

# **Weakly Supervised Learning via Statistical Sufficiency**

**Giorgio Patrini**

A thesis submitted for the degree of  
Doctor of Philosophy  
The Australian National University

December 2016

© Giorgio Patrini 2016

Except where otherwise indicated, this thesis is my own original work.

Giorgio Patrini  
20 December 2016



---

# Abstract

---

The Thesis introduces a novel algorithmic framework for weakly supervised learning, namely, for any any problem in between supervised and unsupervised learning, from the labels standpoint. Weak supervision is the reality in many applications of machine learning where training is performed with partially missing, aggregated-level and/or noisy labels. The approach is grounded on the concept of statistical sufficiency and its transposition to loss functions. Our solution is problem-agnostic yet constructive as it boils down to a simple two-steps procedure. First, estimate a sufficient statistic for the labels from weak supervision. Second, plug the estimate into a (newly defined) linear-odd loss function and learn the model by any gradient-based solver, with a simple adaptation. We apply the same approach to several challenging learning problems: (i) learning from label proportions, (ii) learning with noisy labels for both linear classifiers and deep neural networks, and (iii) learning from feature-wise distributed datasets where the entity matching function is unknown.



---

# Contents

---

<b>Abstract</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Thesis summary . . . . .	1
1.2 Organization and originality . . . . .	4
1.3 First-author publications included in the Thesis . . . . .	5
1.4 Contributed publications . . . . .	5
<b>2 Background</b>	<b>7</b>
2.1 Preliminary notation . . . . .	7
2.2 The supervised learning problem . . . . .	8
2.3 Learning Theory . . . . .	9
2.3.1 Generalization bounds and Rademacher complexity . . . . .	10
2.3.2 Calibrated losses . . . . .	13
2.4 Maximum likelihood, exponential family and sufficient statistics . . . . .	13
2.4.1 Sufficient statistics . . . . .	14
2.5 Weakly supervised learning . . . . .	15
2.5.1 Empirical risk minimization under weak supervision . . . . .	17
2.6 Appendix: proofs . . . . .	19
2.6.1 Proof of Theorem 5 . . . . .	19
2.6.2 Proof of Theorem 7 . . . . .	20
<b>3 Weakly supervised learning and loss factorization</b>	<b>23</b>
3.1 Linear-odd losses and Loss Factorization . . . . .	23
3.1.1 The extent of linear-odd losses . . . . .	25
3.2 Generalization bounds . . . . .	27
3.3 A two-step procedure for weakly supervised algorithms . . . . .	29
3.4 Discussion . . . . .	32
3.5 Appendix: proofs . . . . .	33
3.5.1 Proof of Corollary 20 . . . . .	33
3.5.2 Proof of Lemma 21 . . . . .	33
3.5.3 Proof of Lemma 22 . . . . .	33
3.5.4 Proof of Theorem 23 . . . . .	35
3.5.5 Proof of Lemma 24 . . . . .	39
3.5.6 Proof of Theorem 26 . . . . .	39
3.6 Appendix: additional formal results . . . . .	42
3.6.1 Mean and covariance operators . . . . .	42

---

3.6.2	The generality of factorization . . . . .	43
3.6.3	Factorization of non linear-odd losses . . . . .	44
3.6.4	More graphs on linear and non-linear-odd losses . . . . .	47
3.6.5	The linear-odd losses of du Plessis et al. [2015] . . . . .	47
3.7	References . . . . .	49
3.7.1	The two-step procedure of Raghunathan et al. [2016] . . . . .	49
3.7.2	Learning reductions . . . . .	50
<b>4</b>	<b>Learning from label proportions</b>	<b>51</b>
4.1	Motivation . . . . .	51
4.2	Learning setting . . . . .	52
4.2.1	Symmetric proper losses . . . . .	53
4.3	Estimating the sufficient statistic . . . . .	54
4.4	Mean Map algorithm of Quadrianto et al. [2009] . . . . .	55
4.5	Laplacian Mean Map . . . . .	55
4.6	Estimation: formal guarantees . . . . .	57
4.7	Alternating Mean Map . . . . .	60
4.8	Generalization bounds . . . . .	61
4.9	Experiments . . . . .	63
4.9.1	Algorithms . . . . .	63
4.9.2	Simulated domains . . . . .	64
4.9.3	UCI domains . . . . .	65
4.10	Discussion . . . . .	66
4.11	Appendix: proofs . . . . .	69
4.11.1	Proof of Lemma 35 . . . . .	69
4.11.2	Proof of Theorem 36 . . . . .	69
4.11.3	Proof of Lemma 37 . . . . .	71
4.11.4	Proof of Lemma 38 . . . . .	74
4.11.5	Proof of lemma 39 . . . . .	77
4.11.6	Proof of Theorem 41 . . . . .	78
4.11.7	Proof of Lemma 42 . . . . .	80
4.11.8	Proof of Theorem 43 . . . . .	81
4.11.9	Proof of Theorem 46 . . . . .	83
4.11.9.1	Proof of Equation 4.33 . . . . .	83
4.11.9.2	Proof of Equation 4.34 . . . . .	89
4.12	Appendix: additional experimental results . . . . .	90
4.12.1	Simulated domain for violation of homogeneity assumption . . . . .	90
4.12.2	Additional tests on alter- $\alpha$ SVM Yu et al. [2013] . . . . .	90
4.12.3	Scalability . . . . .	91
4.12.4	Full results on small domains . . . . .	92
4.13	References . . . . .	102
4.13.1	Ecological inference . . . . .	103



---

<b>5</b>	<b>Learning with noisy labels I: theory for linear models</b>	<b>105</b>
5.1	Motivation . . . . .	105
5.2	Learning setting . . . . .	106
5.3	Estimating the sufficient statistic and $\mu$ SGD . . . . .	106
5.4	Generalization bounds . . . . .	107
5.5	Experiments . . . . .	110
5.6	Discussion . . . . .	111
5.7	Appendix: proofs . . . . .	114
5.7.1	Proof of Theorem 52 . . . . .	114
5.7.2	Proof of Theorem 53 . . . . .	114
5.7.3	Proof of Theorem 56 . . . . .	116
5.7.4	Proof of Corollaries 57 and 58 . . . . .	117
5.8	References . . . . .	118
<b>6</b>	<b>Learning with noisy labels II: deep neural networks, multi-class, noise estimation</b>	<b>119</b>
6.1	Motivation . . . . .	119
6.2	Learning setting . . . . .	120
6.3	Loss correction procedures . . . . .	122
6.3.1	The backward correction . . . . .	122
6.3.2	The forward correction . . . . .	122
6.3.3	Estimating the noise rates . . . . .	124
6.4	Noise free Hessians via ReLU . . . . .	125
6.5	Experiments . . . . .	126
6.5.1	Loss corrections with $T$ known or estimated . . . . .	126
6.5.2	Comparing with other loss functions . . . . .	128
6.5.3	Experiments on Clothing1M . . . . .	129
6.6	Discussion . . . . .	130
6.7	Appendix: proofs . . . . .	133
6.7.1	Proof of Theorem 60 . . . . .	133
6.7.2	Proof of Theorem 62 . . . . .	133
6.7.3	Proof of Theorem 63 . . . . .	133
6.7.4	Proof of Theorem 64 . . . . .	134
6.8	References . . . . .	134
<b>7</b>	<b>Learning from vertically distributed data without entity matching</b>	<b>137</b>
7.1	Motivation . . . . .	137
7.1.1	Entity resolution . . . . .	138
7.2	Learning setting . . . . .	139
7.3	Rademacher observations . . . . .	141
7.4	Building and learning from block rados . . . . .	143
7.4.1	Computation and optimality of block rados . . . . .	145
7.4.2	Learning from all block rados . . . . .	145
7.4.3	A more realistic setting . . . . .	146

7.5	Experiments . . . . .	146
7.5.1	Domain generation . . . . .	147
7.5.2	Metric . . . . .	147
7.5.3	Results . . . . .	148
7.6	Discussion and references . . . . .	150
7.7	Appendix: proofs . . . . .	153
7.7.1	Proof of Theorem 68 . . . . .	153
7.7.2	Proof of Lemma 69 . . . . .	153
7.7.3	Proof of Theorem 72 . . . . .	153
7.7.4	Proof of Theorem 73 . . . . .	154
7.8	Appendix: extension to the more general setting . . . . .	154
7.9	Appendix: additional experimental results . . . . .	157
<b>8</b>	<b>Conclusion</b>	<b>165</b>

---

# Introduction

---

## 1.1 Thesis summary

Supervised learning is by far the most effective application of the machine learning paradigm. However, its success in modern real world challenges is undermined by the fundamental assumption that we have perfect knowledge of the target variable, the label, at training time. Despite the unprecedented pace of accumulation of digital datasets, *labeled data* is still rare, due to several reasons: supervision is often generated by costly human annotation; relevant information may be obfuscated by privacy mechanisms; or, merely, labels are only known at a higher level of instance/temporal granularity with the respect to the one required for training. In addition to those issues, data collection is affected by ubiquitous noise corruption of various nature and therefore supervision is never perfectly reliable. As a consequence, learning is often performed with sparse, aggregated-level and/or noisy training labels. We consider all those scenarios all under the name of *weakly supervised learning*.

In this Thesis, we approach this generic learning problem through the lens of Statistics and Learning Theory. Our solution is problem-agnostic yet constructive and it boils down to a simple two-steps procedure. First, estimate a sufficient statistic of the unknown labels from weak supervision; this quantity is called *mean operator*. Second, plug the estimate into a standard loss function and learn the model by any gradient-based solver. The key to the second step is the definition of a family of losses that we call *linear-odd*. Its elements can be computed without the need of any label, given their sufficient statistic, by virtue of what we term *loss factorization*. Several commonly used losses belong to the family, *e.g.* logistic and square. From the theoretical viewpoint, linear-odd losses shed new light on generalization bounds with Rademacher complexity: the contribution of the supervision is isolated into one single term, which accounts for the deviation of the sufficient statistic from its population mean. From the algorithmic viewpoint, we bypass the need of estimating actual labels and therefore circumvent the difficult bi-convex optimization problem arising from naively modeling labels as latent variables. To put this into practice, a well-behaved estimator of the sufficient statistic remains to be defined. This is achieved depending on the particular nature of the weak supervision and relative assumptions.

We study in detail three scenarios of weak supervision. The first is called *learning*

from label proportions, where nothing is known about the target variable but its proportion over subsets of the training set, the *bags*. This setting is inspired by several applications where individual labels are not available, but ratios and percents are easy to estimate for domain experts, or given by other sources in terms of aggregates, e.g. surveys, census data. Despite the poor label knowledge, it is possible to learn linear classifiers with strong theoretical guarantees and good practical performance. We work under a weak distinguish-ability assumption of bags, which relaxes a similar but stricter condition in literature. In light of our two-step framework, the problem is reduced to estimating the mean operator from the label proportions. We do so by least square minimization with a manifold regularizer, which expresses our geometrical assumption on the data. A finite sample guarantee is given: the estimator is all the better as the maximum feature vectors norm increases. We name this algorithm the Laplacian Mean Map (LMM): after estimation, a standard loss function computed with the mean operator is minimized. The model output of LMM enjoys a data dependent approximation bound with respect to an ideal classifier learned with full supervision.

We then propose an iterative algorithm, the Alternating Mean Map (AMM). It takes the solution of LMM as input and optimizes it further over the set of labelings consistent with the proportions, implementing coordinate descent minimization. This simple idea can highly improve the quality of the final model in practice, yet it does not suffer some of the drawbacks of the popular bi-convex iterative optimization. In fact, the LMM initializer performs significantly better than random vectors and is deterministic by definition. Moreover, computing the latent labels that match with the proportions while minimizing the loss can be efficiently done via sorting. Finally, we also formulate a specialized uniform convergence bound, involving a generalization of Rademacher complexity for learning from label proportions. The result includes a bag-wise surrogate risk for which we show that AMM optimizes a tractable bound. We experiment on UCI domains with up to hundreds of thousands of examples, comparing the algorithms to previous work. We simulate bags and their label proportions and retain labels for test set performance. Results display that AMM and LMM outperform the state of the art and sometimes even compete with the fully supervised learner, while requiring few proportions only. Tests on the largest domains display the scalability of both algorithms.

The second weakly supervised scenario is the case of *asymmetric label noise*. In the binary classification setting that we considered, the noise model assumes that one label is flipped into the other at random by class-dependent noise rates. The setting is a well-studied approximation of real-world corruption of categorical labels, by the effect of human mistakes in data annotation or automated data extraction. A known recipe exists for correcting the loss by making it unbiased under asymmetric label noise. Thanks to loss factorization, we can directly apply it to the mean operator. We demonstrate that any algorithm minimizing a linear-odd loss computed with the unbiased mean operator enjoys a generalization bound that tightens results from prior work. We also characterize the whole family of linear-odd losses with an approximate robustness property: the difference in average risk (under the noise) between

---

ideal model and ours cannot be arbitrarily large and the bound is data dependent by the mean operator. On the algorithmic side, we show how to adapt stochastic gradient descent and proximal algorithms to handle weak supervision. Once the mean operator has been estimated, the modification only requires a change of inputs and to sum up the mean operator to the model update. The theory is validated with experiments on UCI domains, on which we inject artificial noise. We assume to know the noise rates for the first step of estimation. Our approach of loss correction is effective: we obtain a significant gain in performance with respect to standard loss minimization with the same algorithm. Interestingly, in the presence of very high noise (one noise rate close to 50%) we still obtain sensible models, while learning with no loss correction often results in random guessing.

We study further the asymmetric label noise setting and consider multi-class classification with deep neural networks, including recurrent neural networks. Once more, the only component we operate on is the loss function, thus our solution is actually independent from any chosen architecture. Although, with neural networks, our two-step learning procedure requires more care. Here the sufficient statistic is computed upon intermediate feature representations that need to be learned and therefore cannot be estimated *a priori*. Still, even with multiple classes, what is sufficient for loss correction is the probability of flipping each class into any other, namely a transition matrix. Given the matrix, we propose two types of loss corrections; one follows the extension of the above idea to multi-class, the other is instead inspired by prior work on robust Deep Learning. Both corrections amounts at most to one inversion and multiplication of the noise matrix. Therefore, we show how to compute the noise rates from (noisy) data as the first step of learning. In particular, we adapt a recently published technique for noise estimation to the multi-class scenario. The whole learning process is summarized as follows: train the chosen neural network with standard loss using corrupted labels; exploit it to estimate the transition matrix; finally re-train the network with the corrected loss. Experiments on MNIST, IMDB, CIFAR-10, CIFAR-100 and a large scale dataset of clothing images show that a diversity of architectures — stacking dense, convolutional, pooling, dropout, batch normalization, word embedding, LSTM and residual layers — demonstrate noise robustness for image recognition and sentiment analysis. Incidentally, we also prove that, when ReLU is the only non-linearity, the loss curvature is immune to label noise.

Finally, we explore an extension of the framework to the challenging problem of *learning from distributed datasets*, where examples are “vertically” (feature-wise) partitioned and the “who-is-who” correspondence is unknown. As a motivating example, we may imagine two peer institutions, *e.g.* a bank and an insurance company, which aim to exploit the joint predictive power of their data assets (expressed on different feature spaces), based on the knowledge that they share many customers. Our goal is to learn a classifier in the cross product space of the two domains, in the hard case when no shared ID is available; companies may not know which partial views are about the same customer in the two datasets or, more likely, matching is not permitted by privacy regulation. We work under two assumptions: all labels are known

and some of the features are shared between the two domains. This can be seen as a peculiar setting of weakly supervised learning, since the mapping between each feature vector and relative label (although known) is effectively missing.

Traditionally, the problem would be approached by first solving (approximate) entity matching and subsequently learning the classifier in a standard manner. Instead, following the underlying philosophy of the Thesis, we bypass the problem of estimating the unknown variables and look for sufficient statistics. In this setting, we require a different sort of loss factorization, expressed by the recently introduced concept of Rademacher observations (rados). Those statistics closely related to the mean operator and indeed can be thought as sufficient for subsets of examples — not just for labels — in the data. We replace the minimization of a loss over examples, which requires entity resolution, by the *equivalent* minimization of a loss over rados. In general, the number of rados is exponential in the sample size. We show that a large subset of these rados does not require to perform entity matching and can be easily obtained linking together partial views of examples with the same value of shared features. With a focus on square loss, we prove that optimization on those rados has time and space complexities smaller than the algorithm minimizing the equivalent square loss on examples. Last, we relax the key assumption that the data is vertically partitioned among peers — in this case, we would not even know the existence of a solution to entity resolution. In this more general setting, experiments support the possibility of beating the optimal peer in hindsight.

## 1.2 Organization and originality

Except Chapter 2 on the background, the content of the Thesis is original or excerpt from novel results either published or submitted to academic conferences from myself and collaborators. In particular, Chapter 3 elaborates the theory of loss functions discussed in Patrini et al. [2016a] including some formal results from Patrini et al. [2014]; it can be thought as a toolbox of abstract results to be specialized for particular instances of weak supervision. Chapter 4 deals with the special case of learning from label proportions, that is the subject of Patrini et al. [2014]. Chapters 5 and 6 discuss the case of learning with asymmetric noisy labels as in Patrini et al. [2016a] (theory, linear and kernel models) and in Patrini et al. [2017] (neural networks, multi-class, noise estimation). Chapter 7 extends the theory from Chapter 3 to the problem of distributed learning. This is a revised presentation from the original work of Patrini et al. [2016b]; it also requires some background material from Nock et al. [2015], that we have co-authored. Results recalled from prior work are accompanied by the relative citations. Literature reviews and related work are topic-specific and presented at the end of each Chapter, along with most of the proofs and additional formal and experimental results.

---

### **1.3 First-author publications included in the Thesis**

PATRINI, G.; NOCK, R.; RIVERA, P.; AND CAETANO, T., (Almost) no label no cry. In *NIPS*, 2014.

PATRINI, G.; NIELSEN, F.; NOCK, R.; AND CARIONI, M., Loss factorization, weakly supervised learning and label noise robustness. In *ICML*, 2016.

PATRINI, G.; NOCK, R.; HARDY, S.; AND CAETANO, T., Fast learning from distributed datasets without entity matching. In *IJCAI*, 2016.

PATRINI, G.; ROZZA, A.; MENON, A.; NOCK, R.; AND QU, L., Making deep neural networks robust to label noise: a loss correction approach. Submitted to *CVPR*, 2017.

### **1.4 Contributed publications**

NOCK, R.; PATRINI, G.; AND FRIEDMAN, A., Rademacher observations, private data and boosting. In *ICML*, 2015.

MUZELLEC, B.; NOCK, R.; PATRINI, G.; AND NIELSEN, F., Tsallis regularized optimal transport and ecological inference. In *AAAI*, 2017.





---

# Background

---

The Chapter recalls basic notions of supervised Machine Learning, Learning Theory and Statistics, and presents a critical review of weakly supervised learning as tackled by prior work. Notations, formal statements and proofs of this introductory material are loosely inspired by the book of Shalev-Shwartz and Ben-David [2014].

## 2.1 Preliminary notation

We begin by defining some notation that will be used throughout the Thesis.  $\mathbb{R}_+$  is the set of non-negative real numbers. For any  $m \in \mathbb{N}$ , the sequence of positive natural number up to  $m$  is  $[m] \doteq \{1, \dots, m\}$ . Boldfaces like  $v$  indicate (column) vectors;  $v_i$  (boldface) is the  $i$ th element of a sequence of vectors instead;  $v_i$  may either be the  $i$ th component of vector  $v$  or the  $i$ th element of a sequence of scalars  $\{v_1, v_2, \dots\}$ .  $\mathbf{1}$  is the vector of all ones, with size determined by context, and  $e^i$  is the indicator vector, that is 1 only at the  $i$ th position.

Capital letters like  $A$  can indicate matrices, or sometimes constants scalar, depending on the context;  $A_i$  and  $A_{.j}$  are respectively row  $i \in [m]$  and column  $j \in [n]$  of matrix  $A \in \mathbb{R}^{m \times n}$ .  $\mathbf{0}$  is either a vector or a matrix filled with zeros, with size determined by context.  $I$  is the identity matrix.  $\text{diag}(v)$  is diagonal matrix with  $v$  on the diagonal.  $\text{tr}(A)$  is the trace of matrix  $A$ .

For any function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , we denote its first derivative at  $x$  by  $f'(x)$ ; the subdifferential set, that is, the set of all subdifferentials at point  $x$  is denoted by  $\partial f(x)$ . For vector functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , we indicate gradient at point  $x$  by  $\nabla f(x)$ .

We denote with  $1\{p\}$  the indicator function of a predicate  $p$ , which is 1 when TRUE.  $[v]_+$  is  $\max(0, v)$ , for any scalar  $v$ . Inner products are written by angular brackets  $\langle \cdot, \cdot \rangle$ . Sets and sequences are italic capital letter like  $\mathcal{S}$ . The probability of an event is denoted by  $\mathbb{P}(\cdot)$ . The expectation over a distribution  $\mathcal{D}$  is denoted as  $\mathbb{E}_{\mathcal{D}}[\cdot]$ ; the same notation is for empirical averages on a sample  $\mathcal{S}$  as  $\mathbb{E}_{\mathcal{S}}[\cdot] \doteq 1/|\mathcal{S}| \sum \cdot$ . We also refer to the set  $\Sigma_m \doteq \{-1, 1\}^m$ , the  $m$ -times Cartesian product of  $\{-1, 1\}$ .

## 2.2 The supervised learning problem

In Supervised Machine Learning we learn a map between two spaces by looking at examples. Learning Theory studies the role of those objects within the process of learning. Let set  $\mathcal{X}$  be the input space set  $\mathcal{Y}$  the output space. Elements  $x \in \mathcal{X} \subseteq \mathbb{R}^d$  are called *feature vectors*, *observations* or *instances* and elements  $y \in \mathcal{Y}$  are called *labels*. The function  $h : \mathcal{X} \rightarrow \mathcal{Y}$ , named *hypothesis* or *model*, belongs to a hypothesis space  $\mathcal{H}$ . Any pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is an example.

It is usual to frame the problem of learning in the language of probability. We assume that examples are drawn *i.i.d* from an unknown but fixed distribution  $\mathcal{D}$ . A (*learning*) *sample*  $\mathcal{S}$  of size  $m$  is a finite sequence of examples  $\{(x_i, y_i), i \in [m]\}$ . Notice that despite the commonly used “set-like” notation,  $\mathcal{S}$  is actually a sequence, which in particular implies that examples can be repeated. At the same time, since examples are drawn *i.i.d.*, their order does not strictly matter, except sometimes for notational convenience in defining algorithms. In regression, the output space is the real line, *i.e.*  $\mathcal{Y} = \mathbb{R}$ . In classification, the output space is instead discrete and labels are also called *classes*. For instance in binary classification, a main focus of this work, it is common to represent the output variable as taking values in  $\{-1, 1\}$ .

The goal of supervised learning is to find a hypothesis that generalizes well on unseen examples. Thus we first ought to define what is the quality measure for hypotheses. In classification, we refer to the *generalization error*, or *risk*, the average number of labels correctly predicted by the model:  $\mathbb{E}_{\mathcal{D}} 1\{h(x) = y\}$ . We could therefore define a *learner* as any algorithm minimizing the generalization error as:

$$h^* \doteq \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{\mathcal{D}} 1\{h(x) \neq y\} \quad (2.1)$$

The minimizer of this problem is called *Bayes optimal* and its risk is the *Bayes risk*. This formulation is problematic for several reasons. First, we have mentioned that the distribution  $\mathcal{D}$  is unknown to the learner. Learning can only make use of the sample  $\mathcal{S}$ , which provides an empirical version of the objective of Equation 2.1, namely the *empirical risk*  $\mathbb{E}_{\mathcal{S}} 1\{h(x) \neq y\}$ . Second, we need to define what the model space  $\mathcal{H}$  is. The choice of the model space is crucial for the success of learning, as discussed below. For most of the Thesis we consider one of the simplest hypothesis space, *i.e.* the set of linear classifiers. To make this explicit, we will denote the model by a vector  $\theta \in \mathbb{R}^d$  as  $h(\cdot) \doteq \langle \theta, \cdot \rangle$ . This often allows to simplify formal arguments while still capturing essential properties. We introduce more complex model spaces in Chapter 6, that deals with deep neural networks, and Chapter 3, which briefly touches on non-parametric kernel methods. Third, in practice it would be difficult to optimize directly the error function, which is discontinuous and non-convex. Instead, it is common to resort to surrogate objectives called *loss functions*.

A loss  $\ell : \mathcal{X} \times \mathcal{Y} \times \mathcal{H} \rightarrow \mathbb{R}$  measures the disagreement between labels and model predictions, although not necessarily by strict inequality, *e.g.* logistic and square losses (in formulae given below). By extension from the empirical risk, we call *empirical  $\ell$ -risk*, or sometimes *surrogate empirical risk*, the average loss  $\mathbb{E}_{\mathcal{S}} \ell(y, h(x))$ .

We denote  $\ell$ -risks and empirical counterparts by  $R_{\mathcal{D},\ell}(h)$  and  $R_{\mathcal{S},\ell}(h)$  respectively. The error itself can be thought as the risk a certain loss function, the 0/1 loss,  $\ell_{01} = 1\{yh(\mathbf{x}) < 0\}$ . As a result, we now formulate the *empirical risk minimization* (ERM) framework, a fundamental paradigm in Machine Learning:

$$\operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{S},\ell}(h) \quad (2.2)$$

One final important issue is hidden in the ERM formulation, that of *overfitting*. By searching for a model *in the learning sample*, we may end up learning all too well particular variations of  $\mathcal{S}$ , but with no guarantee that the model may generalize on unseen examples of  $\mathcal{D}$ . Intuitively, when the size of the learning sample is small relatively to the model complexity, the model could memorize the whole sample itself, and overfitting is a potential threat. When obtaining more data for learning is not an option, overfitting can be combatted in several ways. We may replace the hypothesis space  $\mathcal{H}$  with one containing model with smaller capacity; or equivalently, we may express a “preference” over the elements of  $\mathcal{H}$  by balancing the empirical risk with an additional term acting as a regularizer  $\Omega : \mathcal{H} \rightarrow \mathbb{R}_+$ . We update the learning framework of Equation 2.2 to the *structural risk minimization* to:

$$\operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{S},\ell}(h) + \lambda \cdot \Omega(h) \quad (2.3)$$

with  $\lambda > 0$  balancing the two contributions. As we show below,  $\Omega$  may also be a function of on the learning sample, depending on the kind of prior knowledge one needs to express in the problem.

## 2.3 Learning Theory

Loss functions have been studied extensively in the Machine Learning and Statistics literature. It is common to assume some preconditions for defining losses that are amenable of study: non-negativity and convexity are definitely the most widely required. In our work, we do not strictly work under those conditions but we will make it explicit when particular results require them. Here we limit ourself to recall some definitions and properties that are relevant for what is discussed in the Thesis.

One first important requirement for loss functions is *properness*.

**Definition 1.** A loss  $\ell(y, h(\mathbf{x}))$  is proper when:

$$\operatorname{argmin}_h \mathbb{E}_{y \sim p(y|\mathbf{x})} \ell(y, h(\mathbf{x})) = p(y|\mathbf{x}) . \quad (2.4)$$

This requirement states that the optimal model should be the conditional probability of the class, given the observation. This way, models fitted via proper losses are actual probability estimators for the class.

We may define the domain of a loss function to be  $\mathcal{X} \times \mathcal{Y} \times \mathcal{H}$ . In practice, this parameterization is seldom necessary. It is common to focus on a subset of loss

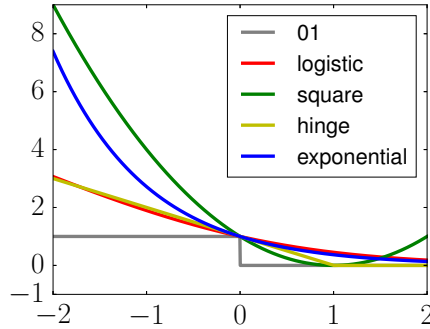


Figure 2.1: Losses as function of  $x = yh(x)$ . Logistic is scaled by  $1/\log(2)$  to be an upper bound of 0/1 loss, which is irrelevant for optimization.

functions called *margin losses*.

**Definition 2.** A margin loss is defined as  $\ell(yh(x)) \doteq \ell(x, y, h)$ .

When clear by the context, we simply use a generic scalar argument  $\ell(x)$ . Margin losses are implicitly *symmetric* since  $\ell(yh(x)) = \ell(-y \cdot (-h(x)))$ . Examples are 0/1 loss  $1\{x < 0\}$ , square loss  $(1 - x)^2$ , logistic loss  $\log(1 + e^{-x})$  (logistic regression), hinge loss  $[1 - x]_+$  (SVM) and exponential loss  $e^{-x}$  (boosting). Figure 2.1 shows how these losses are all designed to be upper bounds of 0/1 loss, a sensible strategy for optimizing the quantity that we ultimately wish to minimize. Among those margin losses, it is well known that hinge loss is not proper and cannot naturally be used as a class probability estimator [Platt, 1999].

Sometimes we will consider losses that are *strongly convex* functions, which is a common assumption to facilitate the derivation of desired upper bounds.

**Definition 3.** Let  $\gamma > 0$ . A differentiable function  $f(x)$  is  $\gamma$ -strongly convex if for any  $x, x' \in \text{Dom}(f)$  we have:

$$f(x) - f(x') \geq \langle \nabla f(x'), x - x' \rangle + \frac{\gamma}{2} \|x - x'\|_2^2. \quad (2.5)$$

Similarly, it is common to refer to  $L$ -Lipschitz functions, that are such that  $|f(x) - f(x')| \leq L\|x - x'\|$  for any  $x, x'$ . A comprehensive discussion on loss functions for binary classification is Reid and Williamson [2010].

### 2.3.1 Generalization bounds and Rademacher complexity

We recall standard results on uniform convergence of learning, taken from Bartlett and Mendelson [2002]; Kakade et al. [2009]; Shalev-Shwartz and Ben-David [2014]. The aim is to provide generalization bounds for ERM. We ground the theory on the concept of Rademacher complexity, a modern data dependent form of complexity of the model space. Several formal results of the Thesis will extend what follows.

**Definition 4.** Let  $\sigma$  be a random variable drawn i.i.d. from  $\{\pm 1\}$  with uniform probability. The empirical Rademacher complexity of a hypothesis space  $\mathcal{H}$  with regard to sample  $\mathcal{S}$  of size  $m$  and loss  $\ell$  is:

$$\mathcal{R}(\ell \circ \mathcal{H} \circ \mathcal{S}) \doteq \mathbb{E}_{\sigma \sim \Sigma_m} \left[ \frac{1}{m} \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \ell(y_i h(x_i)) \right], \quad (2.6)$$

while its population version  $\mathcal{R}(\ell \circ \mathcal{H}) \doteq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \mathcal{R}(\ell \circ \mathcal{H} \circ \mathcal{S})$  is the Rademacher complexity.

We recall two types of bound that we will use in the Thesis. For the sake of completeness, we prove them in the Appendix of this Chapter; proofs rely heavily on the application of McDiarmid's inequality [McDiarmid, 1998].

**Theorem 5.** Assume a bounded loss, i.e.  $|\ell(x)| \leq C, \forall x$ . Then for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  over the choice of  $\mathcal{S}$ , simultaneously for all  $h \in \mathcal{H}$ , we have:

$$R_{\mathcal{D},\ell}(h) - R_{\mathcal{S},\ell}(h) \leq 2\mathcal{R}(\ell \circ \mathcal{H}) + C\sqrt{\frac{2}{m} \log \frac{1}{\delta}}. \quad (2.7)$$

Moreover, let the empirical risk minimizer be  $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{S},\ell}(h)$ . For any  $\delta \in (0, 1)$ , with probability least  $1 - \delta$  over the choice of  $\mathcal{S}$ , we have:

$$R_{\mathcal{D},\ell}(\hat{h}) - \inf_{h \in \mathcal{H}} R_{\mathcal{D},\ell}(h) \leq 4\mathcal{R}(\ell \circ \mathcal{H}) + 2C\sqrt{\frac{2}{m} \log \frac{1}{\delta}}. \quad (2.8)$$

Proof in 2.6.1. This first inequality limits the difference between risk and its empirical approximation; this probabilistic guarantee corroborates the idea of fitting the model on a finite but not too small sample  $\mathcal{S}$ . The second is a bound on the excess  $\ell$ -risk, that is the difference between the risk of ERM model and the smallest risk achievable for models in  $\mathcal{H}$ , the second being unavoidable even for known  $\mathcal{D}$ . Those bounds have a typical shape often encountered in generalization results. Two terms contribute to their significance. The Rademacher complexity quantifies the capacity of the model space  $\mathcal{H}$  of fitting all possible label assignments. While large capacity means better approximation and therefore lower empirical risk, this is not necessarily good for generalization risk on  $\mathcal{D}$ . In fact, low Rademacher complexity tightens the bounds. This is along the line of the classic results on the VC-dimension, which is indeed an upper bound of the Rademacher complexity. The remaining term is a statistical penalty due to requiring uniform convergence. The bound holds true for any  $\mathcal{S}$ , and in particular for samples not being much representative of  $\mathcal{D}$ , that is, drawn from areas with low probability mass on the support of  $\mathcal{D}$ . This is the role of  $\delta$  in the probability inequality, while large  $m$  weights down the possibility of those bad choices of  $\mathcal{S}$ .

Overall the bounds decrease as  $O(1/\sqrt{m})$ . To see this, actually we need to compute the Rademacher complexity for a specific  $\mathcal{H}$  (see Theorem 7 below). Additional conditions can guarantee a *fast rate* of convergence of  $O(1/m)$ , which is also known to

be the fastest rate achievable in learning [Vapnik, 1998]. Those improved guarantees will not be discussed in the Thesis.

Next, we specialize some of the previous results for Lipschitz loss functions and linear classifiers.

**Lemma 6.** *Let  $\ell$  be a  $L$ -Lipschitz function. Then:  $\mathcal{R}(\ell \circ \mathcal{H}) \leq L\mathcal{R}(\mathcal{H})$ , where:*

$$\mathcal{R}(\mathcal{H}) \doteq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \mathbb{E}_{\sigma \sim \Sigma_m} \left[ \frac{1}{m} \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) \right], \quad (2.9)$$

The result is also helpful to clarify what the Rademacher complexity computes: it measures the capacity of hypotheses in  $\mathcal{H}$  by the supremum of correlation with random noise — the random variable  $\sigma$ . We will refer to both  $\mathcal{R}(\mathcal{H})$  and  $\mathcal{R}(\ell \circ \mathcal{H})$  (with dependency on  $\ell$ ) as Rademacher complexity. Second, we can give generalization bounds specialized for linear classifiers.

**Theorem 7.** *Let  $\mathcal{X}$  be a vector space such that  $\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq X < \infty\}$ . Let  $\mathcal{H}$  be the space of bounded linear classifiers  $\mathcal{H} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_2 \leq H < \infty\}$ . Assume a  $L$ -Lipschitz, bounded loss  $|\ell(x)| \leq C, \forall x$ . With probability at least  $1 - \delta$  over the choice of  $\mathcal{S}$ , simultaneously for all  $h \in \mathcal{H}$ , we have:*

$$R_{\mathcal{D}, \ell}(h) - R_{\mathcal{S}, \ell}(h) \leq \frac{2LXH}{\sqrt{m}} + C\sqrt{\frac{2}{m} \log \frac{1}{\delta}}. \quad (2.10)$$

Moreover, let the empirical risk minimizer be  $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{S}, \ell}(h)$ . For any  $\delta \in (0, 1)$ , with probability least  $1 - \delta$  over the choice of  $\mathcal{S}$ , we have:

$$R_{\mathcal{D}, \ell}(\hat{h}) - \inf_{h \in \mathcal{H}} R_{\mathcal{D}, \ell}(h) \leq \frac{4LXH}{\sqrt{m}} + 2C\sqrt{\frac{2}{m} \log \frac{1}{\delta}}. \quad (2.11)$$

Proof in 2.6.2. Finally, we have mentioned that Rademacher complexity is a data dependent measure. Although, previous results account for its population mean, not the empirical complexity. Yet it is simple to show that the empirical Rademacher complexity converges rapidly to its population mean by applying McDiarmid's inequality.

**Lemma 8.** *Assume a bounded loss, i.e.  $|\ell(x)| \leq C, \forall x$ . Let  $\delta > 0$ . We have with probability  $1 - \delta$ :*

$$|\mathcal{R}(\ell \circ \mathcal{H}) - \mathcal{R}(\ell \circ \mathcal{H} \circ \mathcal{S})| \leq C\sqrt{\frac{2}{m} \log \frac{2}{\delta}}. \quad (2.12)$$

It follows that we can express all the above generalization bounds in terms of empirical Rademacher complexity at the price of worsening the statistical penalty.

### 2.3.2 Calibrated losses

The generalization bounds above express uniform convergence guarantees of the ERM with respect to the surrogate loss function chosen for learning. Ultimately, we wish to obtain similar bounds for the error, which is the measure of goodness of the hypothesis. To that objective, we define a subclass of margin losses called *calibrated* [Bartlett et al., 2006], which transforms the previous results in bounds for the generalization error. We introduce them by a sufficient and necessary condition for the case of convex losses.

**Definition 9.** Let  $\ell$  be a non-negative convex loss.  $\ell$  is (classification) calibrated if and only if  $\ell$  is differentiable in 0 and  $\ell'(0) < 0$ .

**Theorem 10.** Let  $\ell$  be calibrated. Then, there exists a convex, non-decreasing, invertible function  $\psi_\ell : [0, 1] \rightarrow \mathbb{R}_+$ , with  $\psi_\ell(0) = 0$ , such that for any  $h \in \mathcal{H}$ :

$$\psi_\ell \left( R_{\mathcal{D},0/1}(h) - \inf_h R_{\mathcal{D},0/1}(h) \right) \leq R_{\mathcal{D},\ell}(h) - \inf_h R_{\mathcal{D},\ell}(h) . \quad (2.13)$$

It follows that any time we can prove a bound for the excess  $\ell$ -risk, we immediately translate it to the error via  $\psi_\ell^{-1}$ . For connections between proper and calibrated losses see Reid and Williamson [2010].

## 2.4 Maximum likelihood, exponential family and sufficient statistics

So far we have seen learning within ERM. An alternative set up for learning problems is the notion of maximum likelihood estimation (MLE), a classic procedure for fitting models in Statistics. We start by assuming that our model belongs to a certain parametric probability distribution. Learning is then accomplished by fitting those parameters to the data by maximizing its likelihood. In particular, we can learn a binary classifier in the conditional (binary) exponential family parameterized by a vector  $\theta \in \mathbb{R}^d$ :

$$p_\theta(y|x) = \exp \left( \langle \theta, yx \rangle - \log \sum_{y \in \mathcal{Y}} \exp \langle \theta, yx \rangle \right) \quad (2.14)$$

Here,  $y$  is the only random variable. The two terms in the exponent are the log-partition function, which normalizes the distribution such that it can sum to one, and the inner product between parameters and the sufficient statistic  $yx$ . Under the *i.i.d.* assumption, we can maximize the likelihood of the data being generated by the

exponential family as:

$$\operatorname{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^m \exp \left( \langle \boldsymbol{\theta}, y_i \mathbf{x}_i \rangle - \log \sum_{y \in \mathcal{Y}} \exp \langle \boldsymbol{\theta}, y \mathbf{x}_i \rangle \right) = \quad (2.15)$$

$$\operatorname{argmax}_{\boldsymbol{\theta}} \exp \left( \sum_{i=1}^m \langle \boldsymbol{\theta}, y_i \mathbf{x}_i \rangle - \sum_{i=1}^m \log \sum_{y \in \mathcal{Y}} \exp \langle \boldsymbol{\theta}, y \mathbf{x}_i \rangle \right) . \quad (2.16)$$

By taking log and a negation, the objective of 2.16 becomes:

$$\sum_{i=1}^m \log \sum_{y \in \mathcal{Y}} \exp \langle \boldsymbol{\theta}, y \mathbf{x}_i \rangle - \sum_{i=1}^m \langle \boldsymbol{\theta}, y_i \mathbf{x}_i \rangle \quad (2.17)$$

$$= \sum_{i=1}^m \log \sum_{y \in \mathcal{Y}} \exp \langle \boldsymbol{\theta}, y \mathbf{x}_i \rangle - \sum_{i=1}^m \log \exp \langle \boldsymbol{\theta}, y_i \mathbf{x}_i \rangle \quad (2.18)$$

$$= \sum_{i=1}^m \log \frac{\exp \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle + \exp \langle \boldsymbol{\theta}, -\mathbf{x}_i \rangle}{\exp \langle \boldsymbol{\theta}, y_i \mathbf{x}_i \rangle} \quad (2.19)$$

$$= \sum_{i=1}^m \log (1 + \exp (-2y_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)) . \quad (2.20)$$

Step 2.20 is true since  $y \in \{\pm 1\}$ . Finally, by re-parameterizing  $\boldsymbol{\theta}$  and normalizing, we obtain logistic loss. These steps prove how the two approaches for learning, ERM and MLE, are equivalent for conditional exponential family and logistic loss. Additionally, Equation 2.17 shows how the loss splits into a linear term aggregating the labels sufficient statistics and another, label free term. We consider this equivalence between conditional exponential family and logistic loss as common knowledge (in spite of being unaware of any published material). One of the key theoretical contribution of the Thesis is to demonstrate how this is not a coincidence: we will introduce a broad new family of losses for which this decomposition is always valid.

### 2.4.1 Sufficient statistics

Why did we call the quantities  $y\mathbf{x}$  and  $\sum_i y_i \mathbf{x}_i$  sufficient statistics? A statistic is any function  $g$  computed from a sample  $\mathcal{S}$ . Intuitively, a statistic  $g(\mathcal{S})$  is sufficient (for the parameter  $\boldsymbol{\theta}$ ) when it aggregates all information of  $\mathcal{S}$ , such that no better model could be learned from  $\mathcal{S}$  in place of  $g$ . We are mostly interested in sufficiency with regard to  $y$ , thus we will only look into this case. In formulas, a statistic  $g$  is sufficient for  $\boldsymbol{\theta}$  with respect to any random variable  $Y$  with outcome  $y$  if:

$$\mathbb{P}(\boldsymbol{\theta}|y) = \mathbb{P}(\boldsymbol{\theta}|g(y)), \forall y \sim Y . \quad (2.21)$$

We provide an equivalent definition that, although somewhat less intuitive, is useful for conceiving some formal results in the Thesis.

**Definition 11** (Sufficient statistic). *Let  $Y$  a random variable.  $g(Y)$  is a sufficient statistic*



for  $\theta$  with regard to  $Y$  when for each pair of outcomes  $y, y'$  of  $Y$  we have:

$$\frac{\mathbb{P}(\theta|g(y))}{\mathbb{P}(\theta|g(y'))} \text{ does not depend on } y \Leftrightarrow g(y) = g(y') . \quad (2.22)$$

An additional equivalent characterization for sufficient statistics is provided by the classic Fisher-Neyman Theorem [Lehmann and Casella, 1998].

**Theorem 12 (Fisher-Neyman).** *Let  $p_{\theta}(y)$  be the probability density function of  $y$ . Then  $g$  is sufficient for  $\theta$  if and only if two non-negative functions  $p^{(i)}$  and  $p^{(ii)}$  can be found such that:*

$$p_{\theta}(y) = p_{\theta}^{(i)}(g(y)) \cdot p^{(ii)}(y) . \quad (2.23)$$

In other words, the probability density factors in two functions, such that  $\theta$  interacts with the  $y$  only through  $g$ . This can be used to show that  $yx$  is indeed sufficient for  $\theta$  with regard to  $y$  in the case of the conditional exponential family (Equation 2.14). It holds that  $p^{(ii)}(y|x) = 1$ ,  $g(y|x) = yx$  and  $p_{\theta}^{(i)}(\cdot|x) = \exp(\langle \theta, \cdot \rangle - \log \sum_{y \in \mathcal{Y}} \exp(\langle \theta, yx \rangle))$ , since the value of  $y$  is not needed for computing  $p_{\theta}^{(i)}$ .

## 2.5 Weakly supervised learning

This Section is halfway between an informal problem statement and a high level literature review for weakly supervised learning problems and relative frameworks for solutions. “Weakly supervised learning” is a non-standard yet widely used terminology to describe scenarios that sit somewhere in between of supervised and unsupervised learning; other literature refers to, for instance, indirect [Raghunathan et al., 2016] or distant supervision [Mintz et al., 2009; Surdeanu et al., 2012]. This class of learning problems relaxes one fundamental assumption of supervised learning: the learner does not have perfect labels, that is, there is no guarantee that each label is fully observable and that each observed label is free from mistakes. This informal definition is deliberately vague as it is meant to encompass a large diversity of learning settings. For example, labels may be missing as with *semi-supervision* [Chapelle et al., 2006] and *positive and unlabeled data* [du Plessis et al., 2015], *noisy* [Natarajan et al., 2013], aggregated as it happens in *multiple instance learning* [Dietterich et al., 1997] and *learning from label proportions* [Kuck and de Freitas, 2005], or given in a candidate set which include the only correct one, as in *partial labels* or *superset label learning* [Cour et al., 2011; Liu and Dietterich, 2014]. We also stress the fact that all these problems share the same objective of supervised learning, which is learning a classifier — or occasionally a regressor —, in contrast with “fully unsupervised” learning.

Those scenarios are rarely considered altogether; formal analysis and algorithmic solutions are usually proposed *ad-hoc*. In contrast, one of the goal of the Thesis is give an unified treatment. We will see how this level of abstraction is beneficial.

Formally, we model a weakly supervised problem by a corruption process that takes the original distribution  $\mathcal{D}$  and produces a corrupted distribution  $\tilde{\mathcal{D}}$ , from which we get a corrupted sample  $\tilde{\mathcal{S}}$ :

$$\mathcal{D} \xrightarrow{\text{corrupt}} \tilde{\mathcal{D}} \xrightarrow{\text{sample}} \tilde{\mathcal{S}} \quad (2.24)$$

A fundamental assumption is that the marginal distribution of the features is untouched by corruption, that is,  $\mathbb{P}_{\mathcal{D}}(x) = \mathbb{P}_{\tilde{\mathcal{D}}}(x)$ . This is all the formalism we need but we give more concrete examples of  $\tilde{\mathcal{S}}$  for some particular cases.

**Example.** In semi-supervised learning, the learner sees two subsamples, one fully labeled and the other unlabeled:  $\tilde{\mathcal{S}} = (\mathcal{S}_L, \mathcal{S}_U)$ , with  $\mathcal{S}_L = \{(x_i, y_i), i \in [m_L]\}$  and  $\mathcal{S}_U = \{x_i, i \in [m_U]\}$ . In positive and unlabeled learning, the scenario is similar but the labeled examples are always positive  $\mathcal{S}_L = \{(x_i, 1), i \in [m_L]\}$ .

**Example.** When labels are noisy, it is usually assumed that the learner sees them all, although without guarantees of their truthfulness. Hence,  $\tilde{\mathcal{S}} = \{(x_i, \tilde{y}_i), i \in [m]\}$ , where  $\tilde{y}$  is drawn from some corrupted label distribution  $p(\tilde{y})$ .

**Example.** It is less obvious how to formalize supervision given at aggregate level; a possible choice is the following. Let us consider the original learning sample given by two ordered sets, one containing observations and the other one the relative labels:  $\mathcal{S} = (\mathcal{S}_X, \mathcal{S}_Y)$ ; observations are mapped to respective labels by indices. The corrupted sample leaves features vectors untouched as usual, yet it partitions them into  $n$  bags as  $\mathcal{S}_X = \bigcup_{j \in [n]} \mathcal{S}_j$  with  $\forall j, j' \mathcal{S}_j \cap \mathcal{S}_{j'} = \emptyset$ . A certain type of supervision is defined at the level of bags by  $\pi \in \mathbb{R}^n$ . The resulting corrupted sample is then  $\tilde{\mathcal{S}} = (\{\mathcal{S}_j\}_{j \in [n]}, \{\pi_j\}_{j \in [n]})$ . In learning from label proportions,  $\pi_j \in [0, 1]$  represents the fraction of positive labels for each bag  $j$ . In multiple instance labels, the supervision is even weaker and  $\pi_j \in \{0, 1\}$  is 1 if at least one label in bag  $j$  is positive<sup>1</sup>.

The difference between those learning settings is not crisp, for two reasons. On one hand, it is plausible that they can be combined by mixing the types of corruption. For instance, we can easily imagine that some labels are missing but those known are noisy; another common scenario is semi-supervised learning with additional prior knowledge on large unlabeled data, *e.g.* as in Bilenko et al. [2004]. On the other hand, some learning settings can be thought to be more general than others, without a total ordering in a hierarchy. For instance, the noisy label distribution  $p(\tilde{y})$  may be modeling the event of “label suppression” so as to encompass missing labels as well [Menon et al., 2015]; multiple instance learning can be thought as a particular, less informative supervision with respecting to label proportions [Kuck and de Freitas, 2005]; at the same time, algorithms for semi-supervised learning have been shown to be apt for multiple instance learning [Zhou and Xu, 2007].

<sup>1</sup>Other encodings of supervision of MIL have been formalized. Also, MIL is sometimes intended as the problem of learning classifiers at bag level; we do not consider this scenario. See Foulds and Frank [2010]; Hernandez-Gonzalez et al. [2016] for more details.

### 2.5.1 Empirical risk minimization under weak supervision

What is the meaning of computing the empirical risk on the corrupted sample,  $R_{\tilde{\mathcal{S}},\ell}(h)$ ? By relaxing the assumption of “full supervision” the framework of ERM loses its sense. Depending on the nature of weak supervision, this quantity may be computable either only partially (semi-supervision) or completely but with little meaning (noisy labels); worse, when supervision is given at multi-instance level, there exists no obvious way to estimate the risk. Mathematically, the problem of searching for risk minimizer becomes ill-posed. We discuss a non-exhaustive overview of families of methods proposed in the past.

Case (i): some observations have individual labels, but some others do not. We adopt the definition of  $\tilde{\mathcal{S}} = \mathcal{S}_L \cup \mathcal{S}_U$  of semi-supervised learning for simplicity. We can formulate an optimization problem such as:

$$\operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{S}_L,\ell}(h) + \lambda \cdot \operatorname{REG}(\tilde{\mathcal{S}}, h) \quad (2.25)$$

The regularizer is intended to exploit the information of the additional unlabeled features vectors to bias the search in the model space. The idea is that the feature geometry, regardless of unknown labels, should be informative about the class distribution. Formally, one of three celebrated assumptions is required for this to be true: the smoothness, the cluster, or the manifold assumption [Chapelle et al., 2006]; the *causal direction* of data generation has also been taken into account for justifying the success of semi-supervised learning [Janzing et al., 2012]. Examples of this framework for learning with missing labels are: transductive SVM [Joachims, 1999], information regularization [Szummer and Jaakkola, 2002], label propagation as a regularizer [Zhu and Ghahramani, 2002; Bengio et al., 2006], entropy regularization [Grandvalet and Bengio, 2004], manifold regularization [Belkin et al., 2006], ladder network [Rasmus et al., 2015] and graph embedding [Weston et al., 2012; Yang et al., 2016].

Case (ii): noisy labels. Design a corrected loss  $\tilde{\ell}$  satisfying a certain property of robustness with respect to the noise and minimize it instead of the original  $\ell$ :

$$\operatorname{argmin}_{h \in \mathcal{H}} R_{\tilde{\mathcal{S}},\tilde{\ell}}(h) \quad (2.26)$$

Often, the robust loss is either non-convex [Masnadi-Shirazi et al., 2010; Ding and Vishwanathan, 2010; Ghosh et al., 2015] or requires certain knowledge of the label corruption  $p(\tilde{y}|y)$  [Stempfel and Ralaivola, 2009; Natarajan et al., 2013], with the noticeable exception of the linear “unhinged” loss [van Rooyen et al., 2015].

Case (iii): aggregate supervision. An alternative to ERM in this case is non trivial, as demonstrated by the many present in literature. An intuitive approach in the event of lack of observation-level labels is to augment the search space to  $\mathcal{Y}$ , namely, by

modeling the unknown labels as latent variables. Although ultimately the learner output is a model, the estimation of the labels become a fundamental intermediate step for learning. This framework closely resembles and it is often implemented as the Expectation Maximization (EM) algorithm, sharing the well known drawbacks. Weak supervision helps the search in the model space by formulating constraints (hard or soft) to be satisfied by the model:

$$\operatorname{argmin}_{h \in \mathcal{H}, y \in \mathcal{Y}} R_{(\mathcal{S}_X, \mathcal{S}_Y), \ell}(h) + \lambda \cdot \text{CONSTR}(\mathcal{S}_X, \mathcal{S}_Y, h) \quad (2.27)$$

Examples of this approach are Yu et al. [2013] for LLP, Andrews et al. [2002] for MIL. This framework easily extends to the case of constraints on (functions of) features vectors, which may be loosely interpreted as weak supervision as well. In this line of work, particularly relevant for applications in NPL, constraints are enforced on particular features, such as the word “amazing” is a strong indicator of a highly rated movie review; a sample of those methods are constraint driven learning Chang et al. [2007], generalized expectation constraints Mann and McCallum [2008], posterior constraints Ganchev et al. [2010], and label regularization [Mann and McCallum, 2010; Ardehaly and Culotta, 2015]. A unifying framework of this kind of EM-based inference with constraints is elaborated by Samdani et al. [2012].

Case (iv): Prior work also attempts an agnostic treatment of weak supervision for ERM, as we envision in the Thesis. A broad, comparative formalization of weakly supervised problems is presented by Garcia-Garcia and Williamson [2011]. Joulin and Bach [2012] attack a version of the generic Problem 2.27 by a convex relaxation and a dedicated optimization algorithm based on semidefinite programming. Li et al. [2013] formulate a tighter and more scalable convex relaxation of a generic weakly supervised max-margin problem. Zantedeschi et al. [2016] propose the  $\beta$ -risk, a surrogate risk augmented with a matrix parameter  $\beta$  which captures the reliability of each instance label; models, labels and  $\beta$  are estimated iteratively. A last important mention is the work of Raghunathan et al. [2016] that, due to the strong similarity with our framework, will be discussed in detail in Section 3.7 of the next Chapter.

Several other ideas are popular in literature. These are for instance self-training and co-training [Blum and Mitchell, 1998; Nigam and Ghani, 2000], multi-view learning [de Sa, 2005] and graph cuts [Blum and Chawla, 2001] for semi-supervised learning, the framework of collective graphical models [Sheldon and Dietterich, 2011; Bernstein and Sheldon, 2016] and, more generally, any Bayesian treatment for incorporating weak supervision into a generative model [Seeger, 2000; Lawrence and Jordan, 2004; Kuck and de Freitas, 2005; Liang et al., 2009; Kingma et al., 2014]. We do not discuss those topics as they do not directly implement a solution for ERM – although some can be shown to link to it. Moreover, readers who are expert in those areas might foresee a conceptual link to multi-armed bandit [Auer et al., 2002] or Reinforcement Learning [Sutton and Barto, 1998]. Admittedly, the need to model the learning process in those settings is the result of the mismatch between observation

(or action) and label (reward). As we do not study in any case neither sequential data nor interactive scenarios, or in other words, we always assume *i.i.d.* samples, those learning setting are as well outside the scope of the Thesis.

## 2.6 Appendix: proofs

### 2.6.1 Proof of Theorem 5

The proof follows Shalev-Shwartz and Ben-David [2014]. Let us first recall McDiarmid's inequality [McDiarmid, 1998].

**Lemma 13.** *Let  $\mathcal{Z}$  be a set and  $f : \mathcal{Z}^m \rightarrow \mathbb{R}$  be a function of  $m$  variables such that for some  $c > 0$ , for all  $i' \in [m]$  and for all  $z_1, \dots, z_m, z_{i'} \in \mathcal{Z}$  we have:*

$$|f(z_1, \dots, z_i, \dots, z_m) - f(z_1, \dots, z_{i'}, \dots, z_m)| \leq c . \quad (2.28)$$

Let now be  $Z_1, \dots, Z_m$   $m$  independent random variables taking values in  $\mathcal{Z}$ . Then, with probability at least  $1 - \delta$  we have:

$$f(Z_1, \dots, Z_m) - \mathbb{E}[f(Z_1, \dots, Z_m)] \leq c \sqrt{\frac{m}{2} \log \frac{1}{\delta}} . \quad (2.29)$$

We also make use of an additional Lemma that connects Rademacher complexity with the maximum deviation of empirical risk from the respective risk.

**Lemma 14.** *Let  $\phi(\mathcal{S}) = \sup_{h \in \mathcal{H}} R_{\mathcal{D}, \ell}(h) - R_{\mathcal{S}, \ell}(h)$ . Then:*

$$\mathbb{E}_{\mathcal{D}} \phi(\mathcal{S}) \leq 2\mathcal{R}(\ell \circ \mathcal{H}) . \quad (2.30)$$

For any bounded loss function  $|\ell(c)| \leq C$ , the function  $\phi$  satisfies McDiarmid's inequality with  $c = 2C/m$ . Therefore:

$$\phi(\mathcal{S}) \leq \mathbb{E}_{\mathcal{D}} \phi(\mathcal{S}) + C \sqrt{\frac{2}{m} \log \frac{1}{\delta}} \leq 2\mathcal{R}(\ell \circ \mathcal{H}) + C \sqrt{\frac{2}{m} \log \frac{1}{\delta}} , \quad (2.31)$$

with probability at least  $1 - \delta$ . Clearly, for all  $h \in \mathcal{H}$ ,  $R_{\mathcal{D}, \ell}(h) - R_{\mathcal{S}, \ell}(h) \leq \phi(\mathcal{S})$ , which proves the first statement. For the second statement, for any  $h^* \in \mathcal{H}$ , we write:

$$\begin{aligned} R_{\mathcal{D}, \ell}(\hat{h}) - R_{\mathcal{D}, \ell}(h^*) &= R_{\mathcal{D}, \ell}(\hat{h}) + \left( -R_{\mathcal{S}, \ell}(\hat{h}) + R_{\mathcal{S}, \ell}(\hat{h}) \right) \\ &\quad + \left( -R_{\mathcal{S}, \ell}(h^*) + R_{\mathcal{S}, \ell}(h^*) \right) - R_{\mathcal{D}, \ell}(h^*) \end{aligned} \quad (2.32)$$

$$\leq \left( R_{\mathcal{D}, \ell}(\hat{h}) - R_{\mathcal{S}, \ell}(\hat{h}) \right) + \left( R_{\mathcal{S}, \ell}(h^*) - R_{\mathcal{D}, \ell}(h^*) \right) \quad (2.33)$$

$$\leq 2\phi(\mathcal{S}) \quad (2.34)$$

Step 2.33 holds because  $\hat{h}$  is the empirical risk minimizer. Finally, we apply again inequality 2.31 and obtain the desired statement.

### 2.6.2 Proof of Theorem 7

We need to compute an upper bound of the Rademacher complexity for linear classifiers.

**Lemma 15.** *Let  $\mathcal{X}$  be a vector space such that  $\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq X < \infty\}$ . Let  $\mathcal{H}$  be the space of bounded linear classifiers  $\mathcal{H} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_2 \leq H < \infty\}$ . Then:*

$$\mathcal{R}(\mathcal{H}) \leq \frac{XH}{\sqrt{m}}. \quad (2.35)$$

**Proof**

$$m \mathcal{R}(\mathcal{H}) = m \mathbb{E}_{\mathcal{D}, \sigma} \sup_{\boldsymbol{\theta} \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle \quad (2.36)$$

$$= \mathbb{E}_{\mathcal{D}, \sigma} \sup_{\boldsymbol{\theta} \in \mathcal{H}} \sum_{i=1}^m \sigma_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle \quad (2.37)$$

$$= \mathbb{E}_{\mathcal{D}, \sigma} \sup_{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_2 \leq B} \sum_{i=1}^m \sigma_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle \quad (2.38)$$

$$= \mathbb{E}_{\mathcal{D}, \sigma} \sup_{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_2 \leq B} \langle \boldsymbol{\theta}, \sum_{i=1}^m \sigma_i \mathbf{x}_i \rangle \quad (2.39)$$

$$\leq H \mathbb{E}_{\mathcal{D}, \sigma} \left[ \left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2 \right]. \quad (2.40)$$

The last Step is due to Cauchy-Schwartz inequality. Next, by Jensen's inequality it holds:

$$\mathbb{E}_{\sigma} \left[ \left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2 \right] = \mathbb{E}_{\sigma} \left[ \left( \left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2^2 \right)^{1/2} \right] \leq \left( \mathbb{E}_{\sigma} \left[ \left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2^2 \right] \right)^{1/2}. \quad (2.41)$$

Finally, since the variables  $\sigma_i$  are independent, we have:

$$\mathbb{E}_{\sigma} \left[ \left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2^2 \right] = \mathbb{E}_{\sigma} \left[ \sum_{i,j} \sigma_i \sigma_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right] \quad (2.42)$$

$$= \sum_{i \neq j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle \mathbb{E}_{\sigma} [\sigma_i \sigma_j] + \sum_{i=1}^m \langle \mathbf{x}_i, \mathbf{x}_i \rangle \mathbb{E}_{\sigma} [\sigma_i^2] \quad (2.43)$$

$$= \sum_{i=1}^m \|\mathbf{x}_i\|_2^2 \quad (2.44)$$

$$\leq mX^2. \quad (2.45)$$

Altogether we have:

$$m \mathcal{R}(\mathcal{H}) \leq \text{HE}_{\mathcal{D}}(mX^2)^{1/2} = \sqrt{m}XH, \quad (2.46)$$

which proves the Lemma. ■

To get Theorem 7, we combine the Theorem 5 and Lemma 6 with the last result.





---

# Weakly supervised learning and loss factorization

---

We introduce a new family of loss functions called linear-odd for the formal study of weak supervision. The central result is the statement of the Loss Factorization Theorem that isolates the effect of supervision into a label sufficient statistic. As a result to this finding, we can formulate a generic two-step approach for solving weakly supervised problem, by calling standard supervised algorithms based on gradient descent. This Chapter is the most abstract and can be thought as a toolbox of formal results and meta-algorithms to be specialized for particular instances of weak supervision. The Chapter, if not the Thesis in a whole, is inspired by the popular quote:

*“One should solve the problem directly and never solve a more general problem as an intermediate step” [Vapnik, 1998].*

## 3.1 Linear-odd losses and Loss Factorization

In Chapter 2 we shown that MLE of the exponential family is equivalent to ERM with logistic loss. In doing so, we proven that logistic loss decomposes into two terms, one of which is linear. We now give a name to that statistical object and show some intriguing properties connecting the idea of sufficiency from Statistics to Learning Theory.

**Definition 16.** *The (empirical) mean operator of a learning sample  $\mathcal{S}$  is  $\mu_{\mathcal{S}} \doteq \mathbb{E}_{\mathcal{S}}[y\mathbf{x}]$ .*

In the following, we drop the dependency of  $\mathcal{S}$  when clear by the context. The name *mean operator*, or mean map, is borrowed from the theory of Hilbert space embedding [Quadrianto et al., 2009]. In this Thesis,  $\mu$  will play the role of sufficient statistic for labels with regard to a set of losses. The next is motivated by Definition 11 by taking log-odd ratio.

**Definition 17.** A function  $T(\mathcal{S})$  is said to be a sufficient statistic for a variable  $y$  with regard to a set of losses  $\mathcal{L}$  and a hypothesis space  $\mathcal{H}$  when for any  $\ell \in \mathcal{L}$ , any  $h \in \mathcal{H}$  and any two samples  $\mathcal{S}$  and  $\mathcal{S}'$  the empirical  $\ell$ -risk is such that<sup>1</sup>:

$$R_{\mathcal{S},\ell}(h) - R_{\mathcal{S}',\ell}(h) \text{ does not depend on } y \Leftrightarrow T(\mathcal{S}) = T(\mathcal{S}') . \quad (3.1)$$

We will often call any such quantity a *label sufficient statistic*. We now state and prove one of the main theoretical contributions of the Thesis, the Loss Factorization Theorem – by reminiscence of Fisher-Neyman Factorization Theorem 12. As a consequence, we establish sufficiency of mean operators for a large set of losses.

**Theorem 18 (Factorization).** Let  $\mathcal{H}$  be the space of linear hypotheses. Assume that a loss  $\ell$  is such that  $\ell_o(x) \doteq (\ell(x) - \ell(-x))/2$  is linear. Then, for any sample  $\mathcal{S}$  and hypothesis  $h \in \mathcal{H}$  the empirical  $\ell$ -risk can be written as:

$$R_{\mathcal{S},\ell}(h) = \frac{1}{2}R_{\mathcal{S}_{2x},\ell}(h) + \ell_o(h(\boldsymbol{\mu})) , \quad (3.2)$$

where  $\mathcal{S}_{2x} \doteq \{(x_i, \sigma), i \in [m], \forall \sigma \in \mathcal{Y}\}$ .

**Proof** We detail the proof as it clearly highlights the relevant elements allowing losses to factor. A key ingredient is an elementary fact from calculus: *any function writes (uniquely) as the sum of an even and an odd function*. Therefore, we can write a loss function  $\ell(x)$  as:

$$\ell(x) = \frac{1}{2} [\ell(x) + \ell(-x) + \ell(x) - \ell(-x)] \quad (3.3)$$

$$= \ell_e(x) + \ell_o(x) , \quad (3.4)$$

where  $\ell_e(x) \doteq \frac{1}{2} [\ell(x) + \ell(-x)]$  and  $\ell_o(x) \doteq \frac{1}{2} [\ell(x) - \ell(-x)]$  are respectively an even and an odd function (see also Figure 3.1). The empirical risk is:

$$R_{\mathcal{S},\ell}(h) = \mathbb{E}_{\mathcal{S}}[\ell(yh(\mathbf{x}))] \quad (3.5)$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{S}} \left[ \ell(yh(\mathbf{x})) + \ell(-yh(\mathbf{x})) + \ell(yh(\mathbf{x})) - \ell(-yh(\mathbf{x})) \right] \quad (3.6)$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{S}} \left[ \sum_{\sigma \in \mathcal{Y}} \ell(\sigma h(\mathbf{x})) \right] + \mathbb{E}_{\mathcal{S}} \left[ \ell_o(yh(\mathbf{x})) \right] \quad (3.7)$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{S}_{2x}} \left[ \ell(\sigma h(\mathbf{x})) \right] + \mathbb{E}_{\mathcal{S}} \left[ \ell_o(h(y\mathbf{x})) \right] . \quad (3.8)$$

Step 3.8 is due to the definition of  $\mathcal{S}_{2x}$  and linearity of  $h$ . Finally, exploiting the properties of the loss assumed by the Theorem and the linearity of expectation, we

<sup>1</sup>Notice that the focus in the Definition is on the variable  $y$  (confront with Definition 11). This is convenient so as to stress that sufficiency is about a certain variable, usually the label, while  $h$  will rarely play a role in the Definition. Yet, the model  $h$  is the object that sufficiency is relevant for, since it is the goal of learning.

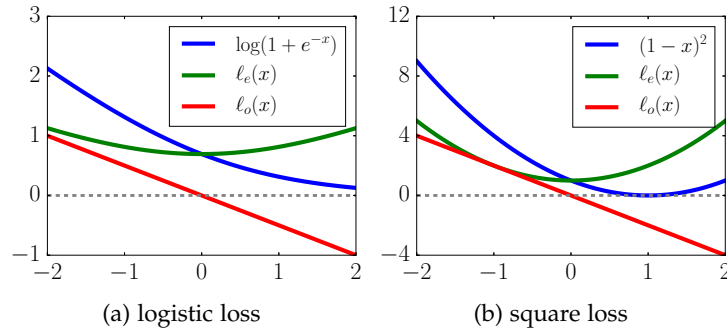


Figure 3.1

have:

$$\mathbb{E}_{\mathcal{S}} \left[ \ell_o(h(yx)) \right] = \ell_o(h(\boldsymbol{\mu})) , \quad (3.9)$$

which concludes the proof. ■

Factorization splits  $\ell$ -risk in two parts. A first term is the  $\ell$ -risk computed *on the same loss* on the “doubled sample”  $\mathcal{S}_{2x}$  that contains each observation twice, labeled with opposite signs, hence it does not require any label knowledge. A second term is a function  $\ell_o$  of  $h$  applied to the mean operator  $\boldsymbol{\mu}$ , which aggregates all sample labels. Also observe that  $\ell_o$  is by construction an odd function, *i.e.* symmetric with respect to the origin. We call the losses satisfying the Theorem *linear-odd*.

**Definition 19** (Linear-odd loss). *A loss  $\ell$  is  $a$ -linear-odd ( $a$ -LOL) when, for any  $a \in \mathbb{R}$ :*

$$\ell_o(x) = (\ell(x) - \ell(-x))/2 = ax . \quad (3.10)$$

Notice how this does not exclude losses that are non-smooth, non-convex, or non-proper. From now on, in this Chapter we consider  $\mathcal{H}$  as the space linear hypotheses  $h(\cdot) = \langle \boldsymbol{\theta}, \cdot \rangle$ . As a consequence of Theorem 18,  $\boldsymbol{\mu}$  is sufficient for all labels.

**Corollary 20.** *The mean operator  $\boldsymbol{\mu}$  is a sufficient statistic for the label  $y$  with regard to LOLs and the space of linear classifiers  $\mathcal{H}$ .*

Proof in 3.5.1. The practical consequence of this Corollary is at the core of the applications in the Thesis: the single vector  $\boldsymbol{\mu} \in \mathbb{R}^d$  summarizes all information concerning the linear relationship between  $y$  and  $x$  for losses that are LOL. But before dealing with the theoretical and practical consequences, we discuss the natural question at this point: how restrictive is the linear-odd condition?

### 3.1.1 The extent of linear-odd losses

Many commonly used losses are linear-odd. We list several examples in Table 3.1 and discuss some in the following.

	loss	even function $\ell_e$	odd function $\ell_o$
generic	$\ell(x)$	$\frac{1}{2}(\ell(x) + \ell(-x))$	$\frac{1}{2}(\ell(x) - \ell(-x))$
LOL	$\ell(x)$	$\frac{1}{2}(\ell(x) + \ell(-x))$	$ax$
$\rho$ -loss	$\rho x  - \rho x + 1$	$\rho x  + 1$	$-\rho x$ ( $\rho \geq 0$ )
unhinged	$1 - x$	1	$-x$
perceptron	$\max(0, -x)$	$x \operatorname{sign}(x)$	$-x$
double-hinge	$\max(-x, 1/2 \max(0, 1 - x))$	$\dagger$	$-x$
SPL	$a_l + l^*(-x)/b_l$	$a_l + \frac{1}{2b_l}(l^*(x) + l^*(-x))$	$-x/(2b_l)$
logistic	$\log(1 + e^{-x})$	$\frac{1}{2} \log(2 + e^x + e^{-x})$	$-x/2$
square	$(1 - x)^2$	$1 + x^2$	$-2x$
Matsushita	$\sqrt{1 + x^2} - x$	$\sqrt{1 + x^2}$	$-x$

Table 3.1: Factorization of losses.  $\dagger$  for reason of space, the even part of double-hinge is written here as:  $\max(-x, 1/2 \max(0, 1 - x)) + \max(x, 1/2 \max(0, 1 + x))$ .

**Example.** For logistic loss it holds that (Figure 3.1a):

$$\ell_o(x) = \frac{1}{2} \log \frac{1 + e^{-x}}{1 + e^x} \quad (3.11)$$

$$= \frac{1}{2} \log \frac{e^{-\frac{x}{2}}(e^{\frac{x}{2}} + e^{-\frac{x}{2}})}{e^{\frac{x}{2}}(e^{-\frac{x}{2}} + e^{\frac{x}{2}})} \quad (3.12)$$

$$= -\frac{x}{2}. \quad (3.13)$$

**Example.** Unhinged loss  $\ell(x) = 1 - x$  of van Rooyen et al. [2015], which is trivially linear-odd. Instead, the standard hinge loss  $\ell(x) = [1 - x]_+$  does not factor in a linear term.

**Example.** Double-hinge and perceptron losses are proven to be linear-odd in du Plessis et al. [2015]. See also Appendix 3.6.5.

**Example.** The class of *symmetric proper losses* (SPLs) [Nock and Nielsen, 2009], e.g. logistic, square and Matsushita losses, satisfies the linear-odd condition. Let  $\phi$  be permissible generator, i.e.  $\phi$  is strictly convex, differentiable and symmetric with respect to  $1/2$  and with  $\operatorname{dom}(\phi) \supseteq [0, 1]$ . SPLs are defined as  $\ell(x) = a_\phi + \phi^*(-x)/b_\phi$ , where  $\phi^*$  is the convex conjugate of  $\phi$ . Then, since  $\phi^*(-x) = \phi^*(x) - x$ , we have:

$$\ell_o(x) = \frac{1}{2} \left( a_\phi + \frac{\phi^*(-x)}{b_\phi} - a_\phi - \frac{\phi^*(x)}{b_\phi} \right) \quad (3.14)$$

$$= \frac{1}{2b_\phi} (\phi^*(x) - x - \phi^*(x)) \quad (3.15)$$

$$= -\frac{x}{2b_\phi}. \quad (3.16)$$

For a broader treatment of SPL, see Subsection 4.2.1 on the next Chapter.

One may question whether SPL and LOL are equivalent. Table 3.1 gives an answer in the negative already: there are examples of non-smooth functions that are LOL. The next result provides a more complete picture by giving an exhaustive characterization of the family of linear-odd losses.

**Lemma 21.** *The exhaustive class of linear-odd losses is in 1-to-1 mapping with a proper subclass of even functions, i.e.  $\ell(x) = \ell_e(x) + ax$ , with  $\ell_e$  any even function.*

Proof in 3.5.2. Interestingly, the proposition also let us engineer losses that always factor: choose any even function  $\ell_e$  with desired properties — it need not be convex nor smooth. The loss is then  $\ell(x) = \ell_e(x) + ax$ , with  $a$  to be chosen. For example, let  $\ell_e(x) = \rho|x| + 1$ , with  $\rho > 0$ .  $\ell(x) = \ell_e(x) - \rho x$  is a LOL; furthermore,  $\ell$  upper bounds 01 loss and intercepts it in  $\ell(0) = 1$ . Non-convex  $\ell$  can be constructed similarly. Yet, not all non-differentiable losses can be crafted this way since they are not LOL.

From the optimization viewpoint, we may want to keep properties of  $\ell$  after factorization. The good news is that we are dealing with the same  $\ell$  plus a linear term. Thus, if the property of interest is closed under summation with linear functions, then it will hold true. An example is convexity: if  $\ell$  is LOL and convex, so is the factored loss. The same is true for differentiability.

In Appendix 3.6, we also elaborate on the relevance of the mean operator relating it to the covariance between  $x$  and  $y$  (Subsection 3.6.1), discuss the generality of Factorization beyond LOLs and linear models (Subsection 3.6.2) and state sufficient and necessary conditions to bound other known losses, e.g. hinge and Huber, by LOLs (Subsection 3.6.3).

## 3.2 Generalization bounds

A consequence of working with LOLs is on generalization bounds. We first derive an improved upper bound to the Rademacher complexity of  $\mathcal{H}$  computed on  $\mathcal{S}_{2x}$ .

**Lemma 22.** *Suppose  $m$  even. Suppose  $\mathcal{X} = \{x : \|x\|_2 \leq X\}$  be the observations space, and  $\mathcal{H} = \{\theta : \|\theta\|_2 \leq H\}$  be the space of linear hypotheses. Let  $\Sigma^{2m} \doteq \times_{j \in [2m]} \mathcal{Y}$ . Then the empirical Rademacher complexity:*

$$\mathcal{R}(\mathcal{H} \circ \mathcal{S}_{2x}) \doteq \mathbb{E}_{\sigma \sim \Sigma^2} \left[ \sup_{\theta \in \mathcal{H}} \frac{1}{2m} \sum_{i \in [2m]} \sigma_i \langle \theta, x_i \rangle \right] \quad (3.17)$$

of  $\mathcal{H}$  on  $\mathcal{S}_{2x}$  satisfies:

$$\mathcal{R}(\mathcal{H} \circ \mathcal{S}_{2x}) \leq v \cdot \frac{XH}{\sqrt{2m}}, \quad (3.18)$$

with  $v \doteq \frac{1}{2} + \frac{1}{2} \sqrt{\frac{1}{2} - \frac{1}{m}}$ .

Proof in 3.5.3. Notice that this holds for the Rademacher complexity as well, since sampling different learning samples from  $D$  does not change the bound. A vanilla

calculation of the complexity of a double sized sample by Lemma 15 would give a coefficient of  $XH/\sqrt{2m}$ . Here we multiply it by  $v < (\sqrt{2} + 1)/(2\sqrt{2}) \approx 0.85$ . This is relevant when combined into the next Theorem. The result sheds new light on excess  $\ell$ -risk bounds on Rademacher complexity with linear hypotheses.

**Theorem 23.** *Assume  $\ell$  is  $a$ -LOL and  $L$ -Lipschitz. Suppose  $\mathbb{R}^d \supseteq \mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq X < \infty\}$  and  $\mathcal{H} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_2 \leq H < \infty\}$ . Let  $c(X, H) \doteq \max_{y \in \mathcal{Y}} \ell(yXH)$  and  $\hat{\boldsymbol{\theta}} \doteq \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} R_{S, \ell}(\boldsymbol{\theta})$ . Then for any  $\delta > 0$ , with probability at least  $1 - \delta$ :*

$$\begin{aligned} R_{\mathcal{D}, \ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}, \ell}(\boldsymbol{\theta}^*) &\leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} + \\ &c(X, H) \cdot \sqrt{\frac{1}{m} \log\left(\frac{1}{\delta}\right)} + 2|a|H \cdot \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2, \end{aligned} \quad (3.19)$$

or more explicitly:

$$\begin{aligned} R_{\mathcal{D}, \ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}, \ell}(\boldsymbol{\theta}^*) &\leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} + \\ &c(X, H) \cdot \sqrt{\frac{1}{m} \log\left(\frac{2}{\delta}\right)} + 2|a|XH \cdot \sqrt{\frac{d}{m} \log\left(\frac{2d}{\delta}\right)}. \end{aligned} \quad (3.20)$$

Proof in 3.5.4. The former expression displays the contribution of the non-linear part of the loss, keeping aside what is missing: a deviation of the empirical mean operator from its population mean. When  $\boldsymbol{\mu}$  is not known because of partial label knowledge, the choice of any estimator would affect the bound only through that norm discrepancy. The second expression highlights the interplay of the two loss components.  $c(X, H)$  is the only non-linear term, which may well be predominant in the bound for fast-growing losses, e.g. strongly convex. Moreover, we confirm that the linear-odd part does not change the complexity and only affects the statistical penalty by a linear factor, with a dependency on  $d$ . We remark that we could as well obtain a bound in the shape of the first statement of Theorem 5; we opt for an excess-risk bound to compare easily with prior work in Chapter 5.

Linear-odd losses are also calibrated under mild conditions.

**Lemma 24.** *Every  $a$ -linear-odd loss that is non-negative, convex, differentiable in 0 and with  $a < 0$  is calibrated.*

Proof in 3.5.5. Consider again Table 3.1. Those losses are all convex and all have  $a < 0$ . They are also differentiable in 0 with the exception of  $\rho$ -loss and “perceptron”. This simple proof can be extended and used for unhinged loss, if we assume bounded models [van Rooyen et al., 2015]. As a consequence of Theorem 10 we know that the previous generalization bound for the  $\ell$ -risk translates into a guarantee for the generalization error.

---

**Meta-Algorithm 1:** Weakly supervised classification via statistical sufficiency — the two-step procedure

---

**Input:**  $\tilde{\mathcal{S}}, \ell$  is  $a$ -LOL,  $\mathcal{H}, \lambda > 0$

(i)  $\hat{\mu} \leftarrow$  estimate from  $(\tilde{\mathcal{S}}, \mathcal{H})$

(ii)  $\hat{\theta} \leftarrow \operatorname{argmin}_{\theta \in \mathcal{H}} \frac{1}{2m} \sum_{i=1}^m \sum_{\sigma \in \mathcal{Y}} \ell(\sigma \langle \theta, x_i \rangle) + a \langle \theta, \hat{\mu} \rangle + \lambda \Omega(\theta)$

**Output:**  $\hat{\theta}$

---

### 3.3 A two-step procedure for weakly supervised algorithms

Recall that in weakly supervised scenarios we learn on  $\tilde{\mathcal{S}}$  with partially observable labels, but aim to generalize to  $\mathcal{D}$ . Let us assume to know an algorithm  $A$  that learns by ERM of a linear-odd loss  $\ell$  over a sample  $\mathcal{S}$  from the non-corrupted distribution  $\mathcal{D}$ . Can we use algorithm  $A$  as is for learning from  $\tilde{\mathcal{S}}$ ? By sufficiency, Corollary 20 leads to a principled two-step approach:

i estimate  $\hat{\mu}$  from  $\tilde{\mathcal{S}}$

ii run algorithm  $A$  with  $\ell$  computed with the estimated  $\hat{\mu}$

Meta Algorithm 1 formalizes this two-step procedure. The framework is as abstract as possible. Step (ii) is the same for every possible weakly supervised scenario, while to implement Step (i) we need to know the particular type of supervision we are dealing with. In other words, any weakly supervised classification problem can be cast into the convex optimization problem of Step (ii), once the label sufficient statistic — or its estimator — is available. This direction has been explored by the work on learning from label proportions of Quadrianto et al. [2009, with logistic loss] and Patrini et al. [2014, symmetric proper losses], and in the setting of noisy labels by Gao et al. [2016, logistic loss]. The approach contrasts with *ad-hoc* optimization methods often aiming to recover the latent labels by coordinate descent and EM (see Chapter 2). Instead, the only difficulty here is to come up with a well-behaved estimator of  $\mu$  — a statistic independent from both model  $h$  and loss  $\ell$ .

Meta Algorithm 1 comes with the generalization guarantees elaborated in this Chapter. Theorem 23 bounds in probability the  $\ell$ -excess risk and, with Lemma 24, the true risk. It is more difficult to derive finite-sample formal results on the output of Meta Algorithm 1, unless we work under additional assumptions. One such result is given in Altun and Smola [2006]. We reformulate it for the whole class of convex

and differentiable linear-odd losses. The Lemma gives an answer to the interesting question of how much model  $\hat{\theta}$  may diverge from the one that would be computed from empirical mean operator  $\mu = \mathbb{E}_{\mathcal{S}}[yx]$ .

**Lemma 25.** *Let  $\ell$  be  $a$ -LOL convex and differentiable and let  $\lambda > 0$ . Call  $\hat{\theta}$  and  $\theta^*$  the minimizer of (ii) in Meta Algorithm 1 when, respectively, the empirical risk is computed with  $\hat{\mu}$  from the estimator of (i) and  $\mu$ . Then*

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{|a|}{\lambda} \|\hat{\mu} - \mu\|_2. \quad (3.21)$$

This Lemma transfers the quality of the sufficient statistic estimation to the quality of the final model. Therefore, by using convex and differentiable linear-odd losses, any well-behaved mean operator estimator would be enough for obtaining reliable loss minimizers. A tighter result than the previous Lemma is the following, a data dependent approximation bound. It requires twice differentiability and strong convexity.

Let us first define some useful quantities. Let  $f_k \in \mathbb{R}^m$  denote the vector encoding the  $k^{\text{th}}$  variable in  $\mathcal{S}$  :  $f_{ki} = x_{ik}$ . For any  $k \in [d]$ , let:

$$\bar{f}_k \doteq \left( \frac{d}{\sum_k \|f_k\|_2^2} \right)^{\frac{d-1}{2d}} f_k \quad (3.22)$$

denote a normalization of vectors  $f_k$  in the sense that:

$$\frac{1}{d} \sum_k \|\bar{f}_k\|_2^2 = \frac{1}{d} \left( \frac{d}{\sum_k \|f_k\|_2^2} \right)^{1-\frac{1}{d}} \sum_k \|f_k\|_2^2 \quad (3.23)$$

$$= \left( \frac{1}{d} \sum_k \|f_k\|_2^2 \right)^{\frac{1}{d}}. \quad (3.24)$$

Let  $\bar{F}$  collect all vectors  $\bar{f}_k$  in column and let  $F$  collect all vectors  $f_k$  in column. Without loss of generality, we assume  $F^\top F \succ 0$ , i.e.  $F^\top F$  positive definite (i.e. no feature is a linear combination of the others), implying that  $\bar{F}^\top \bar{F} \succ 0$  as well, because the columns of  $\bar{F}$  are just positive rescaling of the columns of  $V$ .

**Theorem 26.** *Hold the same conditions of Lemma 25 and additionally assume that  $\ell$  is twice differentiable and  $\gamma$ -strongly convex with  $\gamma > 0$ . Let  $m$  be the learning sample size. Then the following holds:*

$$\|\theta^* - \hat{\theta}\|_2 \leq \frac{|a|}{\lambda + \frac{1}{em} \gamma \text{vol}^2(\bar{F})} \|\mu - \hat{\mu}\|_2, \quad (3.25)$$

where  $\text{vol}(\bar{F}) \doteq \sqrt{\det \bar{F}^\top \bar{F}}$  denote the volume of the (row/column) system of  $\bar{F}$ .

Proof in 3.5.6. To see how large the denominator in 3.25 can be, consider the simple case where all eigenvalues of  $\bar{F}^\top \bar{F}$ ,  $\lambda_k(\bar{F}^\top \bar{F}) \in [\lambda_\circ \pm \delta]$  for small  $\delta$ . In this



**Algorithm 2:**  $\mu$ SGD with  $L_2$  regularization

---

**Input:**  $\mathcal{S}_{2x}, \mu$ ,  $\ell$  is  $a$ -LOL,  $\lambda > 0$ ,  $T > 0$   
 $\theta^0 \leftarrow \mathbf{0}$   
For any  $t = 1, \dots, T$ :  
    Pick  $i \in [|\mathcal{S}_{2x}|]$  uniformly at random  
     $\eta \leftarrow 1/(\lambda t)$   
    Pick any  $v \in \partial \ell(y_i; \langle \theta^t, x_i \rangle)$   
     $\theta^{t+1} \leftarrow (1 - \eta \lambda) \theta^t - \eta (v + a\mu/2)$   
     $\theta^{t+1} \leftarrow \min \left\{ \theta^{t+1}, \theta^{t+1} / (\sqrt{\lambda} \cdot \|\theta^{t+1}\|_2) \right\}$   
**Output:**  $\theta^T$

---

**Algorithm 3:**  $\mu$  proximal algorithm with  $\Theta(\cdot)$  regularization

---

**Input:**  $\mathcal{S}_{2x}, \mu$ ,  $\ell$  is  $a$ -LOL;  $\lambda > 0$ ;  $T > 0$   
 $\theta^0 \leftarrow \mathbf{0}$   
For any  $t = 1, \dots, T$ :  
     $\eta \leftarrow 1/(\lambda t)$  or found by line search  
     $\theta^{t+1} \leftarrow \text{prox}_{\lambda \Theta} \left( \theta^t + \eta \left( \partial R_{\mathcal{S}_{2x}, \ell}(\theta^t) + a\mu/2 \right) \right)$   
**Output:**  $\theta^T$

---

case,  $\text{vol}^2(\bar{F})$  is proportional to the ‘‘average feature norm’’:

$$\frac{\det(\bar{F}^\top \bar{F})}{m} = \frac{\text{tr } V^\top V}{md} + o(\delta) = \frac{\sum_i \|x_i\|_2^2}{md} + o(\delta). \quad (3.26)$$

Theorem 41 in Chapter 4 characterizes further this data dependent results by exploiting additional assumptions.

We now consider two more concrete examples on how to take Step (ii) of the abstract Meta Algorithm 1 for well-known learning algorithms. Algorithm 2,  $\mu$ SGD, adapts SGD to weak supervision. For the sake of presentation, we work on a simple version of SGD based on subgradient descent with  $L_2$  regularization from PEGASO of Shalev-Shwartz et al. [2011]. Given  $\mu$ , changes are trivial: construct  $\mathcal{S}_{2x}$  from  $\tilde{\mathcal{S}}$  and sum  $a\mu/2$  to the subgradients of each example of  $\mathcal{S}_{2x}$ . The only changes are highlighted in grey.

A second noticeable example is given by the family of proximal algorithms [Bach et al., 2012]. The same *modus operandi* leads to Algorithm 3, where the proximal map is  $\text{prox}_{\lambda \Theta}(x) = \text{argmin}_{x'} \lambda \Theta(x') + \frac{1}{2} \|x - x'\|_2^2$  and  $\Theta(\cdot)$  is a regularizer, non necessarily smooth. With some abuse of notation,  $\partial R$  indicates any vector in the subdifferential set of the (decomposable) empirical risk. Once again, the adaptation works by summing  $a\mu/2$  in the gradient step and changing the input to  $\mathcal{S}_{2x}$ .

### **3.4 Discussion**

We have presented novel theoretical tools apt to elaborate a learning theory for weakly supervised problems. The level of abstraction is high enough to let us formulate a Meta Algorithm that is problem agnostic but simply adaptable for exploiting well known gradient-based solvers. Finite sample bounds support the viability of our approach. The rest of the Thesis will put these ideas to use by specialization under different assumptions on the quality of supervision. Applications elsewhere may easily draw from this Chapter as well.

## 3.5 Appendix: proofs

### 3.5.1 Proof of Corollary 20

We need to show the double implication that defines sufficiency for  $y$ .

$\Rightarrow$ ) By Factorization Theorem 18,  $R_{S,\ell}(h) - R_{S',\ell}(h)$  is label independent only if the odd part cancels out.

$\Leftarrow$ ) If  $\mu = \mu'$  then  $R_{S,\ell}(h) - R_{S',\ell}(h)$  is independent from the label, because the label only appears in the mean operator due to Factorization Theorem 18.

### 3.5.2 Proof of Lemma 21

Consider the class of losses satisfying  $\ell(x) - \ell(-x) = 2ax$ . For any element of the class, define  $\ell_e(x) = \ell(x) - ax$ , which is even. In fact we have:

$$\ell_e(-x) = \ell(-x) + ax = \ell(x) - 2ax + ax = \ell(x) - ax = \ell_e(x) . \quad (3.27)$$

### 3.5.3 Proof of Lemma 22

Suppose without loss of generality that  $x_i = x_{m+i}$ . The proof relies on the observation that  $\forall \sigma \in \mathcal{Y}^{2m}$ :

$$\operatorname{argsup}_{\theta \in \mathcal{H}} \mathbb{E}_S[\sigma(\mathbf{x}) \langle \theta, \mathbf{x} \rangle] = \frac{1}{2m} \operatorname{argsup}_{\theta \in \mathcal{H}} \sum_i \sigma_i \langle \theta, \mathbf{x}_i \rangle \quad (3.28)$$

$$= \frac{\sup_{\mathcal{H}} \|\theta\|_2}{\left\| \sum_i \sigma_i \mathbf{x}_i \right\|_2} \sum_i \sigma_i \mathbf{x}_i . \quad (3.29)$$

In fact, this follows from Cauchy-Schwartz inequality – the classifier that maximizes the inner product with the sum is indeed proportional to the normalized sum, times the maximal norm of a classifier. So:

$$\mathcal{R}(\mathcal{H} \circ \mathcal{S}_{2x}) = \mathbb{E}_{\mathcal{Y}^{2m}} \sup_{h \in \mathcal{H}} \{ \mathbb{E}_{\mathcal{S}_{2x}}[\sigma(\mathbf{x}) h(\mathbf{x})] \} \quad (3.30)$$

$$= \frac{\sup_{\mathcal{H}} \|\theta\|_2}{2m} \cdot \mathbb{E}_{\mathcal{Y}^{2m}} \left[ \frac{\left( \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right)^\top \left( \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right)}{\left\| \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right\|_2} \right] \quad (3.31)$$

$$= \sup_{\mathcal{H}} \|\theta\|_2 \cdot \mathbb{E}_{\mathcal{Y}^{2m}} \left[ \frac{1}{2m} \cdot \left\| \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right\|_2 \right] . \quad (3.32)$$

Remark that whenever  $\sigma_i = -\sigma_{m+i}$ ,  $\mathbf{x}_i$  disappears in the sum, and therefore the max norm for the sum may decrease as well. This suggests to split the  $2^{2m}$  assignments into  $2^m$  groups of size  $2^m$ , ranging over the possible number of observations taken into account in the sum. They can be factored by a weighted sum of contributions of

each subset of indices  $\mathcal{I} \subseteq [m]$  ranging over the non-duplicated observations:

$$\mathbb{E}_{\mathcal{Y}^{2m}} \left[ \frac{1}{m} \cdot \left\| \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right\|_2 \right] = \frac{1}{2^{2m}} \sum_{\mathcal{I} \subseteq [m]} \frac{2^{m-|\mathcal{I}|}}{2^m} \cdot \sum_{\sigma \in \mathcal{Y}^{|\mathcal{I}|}} \sqrt{2} \left\| \sum_{i \in \mathcal{I}} \sigma_i \mathbf{x}_i \right\|_2. \quad (3.33)$$

$$= \frac{\sqrt{2}}{2^m} \sum_{\mathcal{I} \subseteq [m]} \frac{1}{2^m} \cdot \underbrace{\frac{1}{2^{|\mathcal{I}|}} \cdot \sum_{\sigma \in \mathcal{Y}^{|\mathcal{I}|}} \left\| \sum_{i \in \mathcal{I}} \sigma_i \mathbf{x}_i \right\|_2}_{u_{|\mathcal{I}|}}. \quad (3.34)$$

The  $\sqrt{2}$  factor appears because of the fact that we now consider only the observations of  $\mathcal{S}$ . For any *fixed*  $\mathcal{I}$ , we renumber its observations in  $[|\mathcal{I}|]$  for simplicity, and observe that, since  $\sqrt{1+x} \leq 1+x/2$ :

$$u_{|\mathcal{I}|} = \frac{1}{2^{|\mathcal{I}|}} \sum_{\sigma \in \mathcal{Y}^{|\mathcal{I}|}} \sqrt{\sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2 + \sum_{i_1 \neq i_2} \sigma_{i_1} \sigma_{i_2} \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}} \quad (3.35)$$

$$= \frac{\sqrt{\sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2}}{2^{|\mathcal{I}|}} \sum_{\sigma \in \mathcal{Y}^{|\mathcal{I}|}} \sqrt{1 + \frac{\sum_{i_1 \neq i_2} \sigma_{i_1} \sigma_{i_2} \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}}{\sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2}} \quad (3.36)$$

$$\leq \frac{\sqrt{\sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2}}{2^{|\mathcal{I}|}} \sum_{\sigma \in \mathcal{Y}^{|\mathcal{I}|}} \left( 1 + \frac{\sum_{i_1 \neq i_2} \sigma_{i_1} \sigma_{i_2} \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2}}{2 \sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2} \right) \quad (3.37)$$

$$= \sqrt{\sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2} + \frac{1}{2^{|\mathcal{I}|} \cdot 2 \sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2} \cdot \sum_{\sigma \in \mathcal{Y}^{|\mathcal{I}|}} \sum_{i_1 \neq i_2} \sigma_{i_1} \sigma_{i_2} \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2} \quad (3.38)$$

$$= \sqrt{\sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2} + \frac{1}{2^{|\mathcal{I}|} \cdot 2 \sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2} \cdot \sum_{i_1 \neq i_2} \mathbf{x}_{i_1}^\top \mathbf{x}_{i_2} \cdot \underbrace{\left( \sum_{\sigma \in \mathcal{Y}^{|\mathcal{I}|}} \sigma_{i_1} \sigma_{i_2} \right)}_{=0} \quad (3.39)$$

$$= \sqrt{\sum_{i \in \mathcal{I}} \|\mathbf{x}_i\|_2^2} \quad (3.40)$$

$$\leq \sqrt{|\mathcal{I}|} \cdot X. \quad (3.41)$$

Plugging this in Equation 3.34 yields:

$$\frac{1}{X} \cdot \mathbb{E}_{\mathcal{Y}^{2m}} \left[ \frac{1}{m} \cdot \left\| \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right\|_2 \right] \leq \frac{\sqrt{2}}{2^m} \sum_{k=0}^m \frac{\sqrt{k}}{2^m} \binom{m}{k}. \quad (3.42)$$

Since  $m$  is even:

$$\mathbb{E}_{\mathcal{Y}^{2m}} \left[ \frac{1}{2m} \cdot \left\| \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right\|_2 \right] \leq \frac{\sqrt{2}}{2^m} \sum_{k=0}^{(m/2)-1} \frac{\sqrt{k}}{2^m} \binom{m}{k} + \frac{\sqrt{2}}{2^m} \sum_{k=m/2}^m \frac{\sqrt{k}}{2^m} \binom{m}{k}. \quad (3.43)$$

Notice that the left one trivially satisfies:

$$\frac{\sqrt{2}}{2^m} \sum_{k=0}^{(m/2)-1} \frac{\sqrt{k}}{2m} \binom{m}{k} \leq \frac{\sqrt{2}}{2^m} \sum_{k=0}^{(m/2)-1} \frac{1}{2m} \cdot \sqrt{\frac{m-2}{2}} \binom{m}{k} \quad (3.44)$$

$$= \frac{1}{2} \cdot \sqrt{\frac{1}{m} - \frac{2}{m^2}} \cdot \frac{1}{2^m} \sum_{k=0}^{(m/2)-1} \binom{m}{k} \quad (3.45)$$

$$\leq \frac{1}{4} \cdot \sqrt{\frac{1}{m} - \frac{2}{m^2}} \quad (3.46)$$

Also, the right one satisfies:

$$\frac{\sqrt{2}}{2^m} \sum_{k=m/2}^m \frac{\sqrt{k}}{2m} \binom{m}{k} \leq \frac{\sqrt{2}}{2^m} \sum_{k=m/2}^m \frac{\sqrt{m}}{2m} \binom{m}{k} \quad (3.47)$$

$$= \frac{1}{\sqrt{2m}} \cdot \frac{1}{2^m} \sum_{k=m/2}^m \binom{m}{k} \quad (3.48)$$

$$= \frac{1}{2} \cdot \frac{1}{\sqrt{2m}}. \quad (3.49)$$

We get:

$$\frac{1}{X} \cdot \mathbb{E}_{y^{2m}} \left[ \frac{1}{m} \cdot \left\| \sum_{i=1}^{2m} \sigma_i \mathbf{x}_i \right\|_2 \right] \leq \frac{1}{4} \cdot \sqrt{\frac{1}{m} - \frac{2}{m^2}} + \frac{1}{2} \cdot \sqrt{\frac{1}{2m}} \quad (3.50)$$

$$= \frac{1}{\sqrt{2m}} \cdot \left( \frac{1}{2} + \frac{1}{2} \sqrt{\frac{1}{2} - \frac{1}{m}} \right). \quad (3.51)$$

And finally:

$$\mathcal{R}(\mathcal{H} \circ \mathcal{S}_{2x}) \leq v \cdot \frac{XH}{\sqrt{2m}}, \quad (3.52)$$

with:

$$v \doteq \frac{1}{2} + \frac{1}{2} \sqrt{\frac{1}{2} - \frac{1}{m}}, \quad (3.53)$$

as claimed.

### 3.5.4 Proof of Theorem 23

We start by proving a helper Lemma, an application of McDiarmid's inequality to evaluate the convergence of the empirical mean operator to its population counterpart.

**Lemma 27.** *Suppose  $\mathbb{R}^d \supseteq \mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq X < \infty\}$  be the observations space. Then for*

any  $\delta > 0$  with probability at least  $1 - \delta$ :

$$\|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2 \leq X \cdot \sqrt{\frac{d}{m} \log \left( \frac{d}{\delta} \right)}. \quad (3.54)$$

**Proof** Let  $\mathcal{S}$  and  $\mathcal{S}'$  be two learning samples that differ for only one example  $(x_i, y_i) \neq (x_{i'}, y_{i'})$ . Let first consider the one-dimensional case. We refer to the  $k$ -dimensional component of  $\boldsymbol{\mu}$  with  $\mu^k$ . For any  $\mathcal{S}, \mathcal{S}'$  and any  $k \in [d]$  it holds:

$$\left| \mu_{\mathcal{S}}^k - \mu_{\mathcal{S}'}^k \right| = \frac{1}{m} \left| x_i^k y_i - x_{i'}^k y_{i'} \right| \quad (3.55)$$

$$\leq \frac{X}{m} |y_i - y_{i'}| \quad (3.56)$$

$$\leq \frac{2X}{m}. \quad (3.57)$$

This satisfies the bounded difference condition of McDiarmid's inequality, which let us write for any  $k \in [d]$  and any  $\epsilon > 0$  that:

$$\mathbb{P} \left( \left| \mu_{\mathcal{D}}^k - \mu_{\mathcal{S}}^k \right| \geq \epsilon \right) \leq \exp \left( -\frac{m\epsilon^2}{2X^2} \right) \quad (3.58)$$

and the multi-dimensional case, by union bound:

$$\mathbb{P} \left( \exists k \in [d] : \left| \mu_{\mathcal{D}}^k - \mu_{\mathcal{S}}^k \right| \geq \epsilon \right) \leq d \exp \left( -\frac{m\epsilon^2}{2X^2} \right). \quad (3.59)$$

Then by negation:

$$\mathbb{P} \left( \forall k \in [d] : \left| \mu_{\mathcal{D}}^k - \mu_{\mathcal{S}}^k \right| \leq \epsilon \right) \geq 1 - d \exp \left( -\frac{m\epsilon^2}{2X^2} \right), \quad (3.60)$$

which implies that for any  $\delta > 0$  with probability  $1 - \delta$ :

$$X \sqrt{\frac{2}{m} \log \left( \frac{d}{\delta} \right)} \geq \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_{\infty} \geq d^{-1/2} \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2. \quad (3.61)$$

This concludes the proof. ■

We now restate and prove Theorem 23.

**Theorem 23** *Assume  $\ell$  is  $\alpha$ -LOL and  $L$ -Lipschitz. Suppose  $\mathbb{R}^d \supseteq \mathcal{X} = \{x : \|x\|_2 \leq X < \infty\}$  be the observations space, and  $\mathcal{H} = \{\theta : \|\theta\|_2 \leq H < \infty\}$  be the space of linear hypotheses. Let  $c(X, H) \doteq \max_{y \in \mathcal{Y}} \ell(yXH)$ . Let  $\hat{\theta} = \operatorname{argmin}_{\theta \in \mathcal{H}} R_{\mathcal{S}, \ell}(\theta)$ . Then for any*

$\delta > 0$ , with probability at least  $1 - \delta$ :

$$\begin{aligned} R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) &\leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} + \\ &c(X, H) \cdot \sqrt{\frac{1}{m} \log\left(\frac{1}{\delta}\right)} + 2|a|H \cdot \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2, \end{aligned} \quad (3.62)$$

or more explicitly:

$$\begin{aligned} R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) &\leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} + \\ &c(X, H) \cdot \sqrt{\frac{1}{m} \log\left(\frac{2}{\delta}\right)} + 2|a|XH \cdot \sqrt{\frac{d}{m} \log\left(\frac{2d}{\delta}\right)} \end{aligned} \quad (3.63)$$

**Proof** Let  $\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{D},\ell}(\boldsymbol{\theta})$ . We have:

$$R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) = \frac{1}{2}R_{\mathcal{D}_{2x},\ell}(\hat{\boldsymbol{\theta}}) + a\langle \hat{\boldsymbol{\theta}}, \boldsymbol{\mu}_{\mathcal{D}} \rangle - \frac{1}{2}R_{\mathcal{D}_{2x},\ell}(\boldsymbol{\theta}^*) - a\langle \boldsymbol{\theta}^*, \boldsymbol{\mu}_{\mathcal{D}} \rangle \quad (3.64)$$

$$= \frac{1}{2} (R_{\mathcal{D}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}_{2x},\ell}(\boldsymbol{\theta}^*)) + a\langle \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*, \boldsymbol{\mu}_{\mathcal{D}} \rangle \quad (3.65)$$

$$= \frac{1}{2} (R_{\mathcal{S}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{S}_{2x},\ell}(\boldsymbol{\theta}^*)) + a\langle \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*, \boldsymbol{\mu}_{\mathcal{D}} \rangle \quad (3.66)$$

$$+ \frac{1}{2} \left( R_{\mathcal{D}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{S}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}_{2x},\ell}(\boldsymbol{\theta}^*) + R_{\mathcal{S}_{2x},\ell}(\boldsymbol{\theta}^*) \right) \} A_1. \quad (3.67)$$

Step 3.64 is obtained by the equality  $R_{\mathcal{D},\ell}(\boldsymbol{\theta}) = \frac{1}{2}R_{\mathcal{D}_{2x},\ell}(\boldsymbol{\theta}) + a\langle \boldsymbol{\theta}, \boldsymbol{\mu}_{\mathcal{D}} \rangle$  for any  $\boldsymbol{\theta}$ . Applying the same equality with regard to  $\mathcal{S}$ , we have:

$$R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) \leq \underbrace{R_{\mathcal{S},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{S},\ell}(\boldsymbol{\theta}^*)}_{A_2} + \underbrace{a\langle \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*, \boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}} \rangle}_{A_3} + A_1. \quad (3.68)$$

$A_2$  is never more than 0 because  $\hat{\boldsymbol{\theta}}$  is the minimizer of  $R_{\mathcal{S},\ell}(\boldsymbol{\theta})$ . From the Cauchy-Schwarz inequality and bounded models it holds true that:

$$A_3 \leq |a| \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 \cdot \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2 \leq 2|a|H \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2. \quad (3.69)$$

We could treat  $A_1$  by calling standard bounds based on Rademacher complexity on a sample with size  $2m$ . Since the complexity  $\mathcal{R}(\mathcal{H} \circ \mathcal{S})$  does not depend on labels, its value would be the same — modulo the change of sample size — for both  $\mathcal{S}$  and  $\mathcal{S}_{2x}$ , as they are computed with same loss and observations. However, the special structure of  $\mathcal{S}_{2x}$  allows us to obtain a tighter structural complexity term, due to

cancellation effects as proven by Lemma 22. In order to exploit it, we first observe:

$$A_1 = \frac{1}{2} \left( R_{\mathcal{D}_{2x,\ell}}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{S}_{2x,\ell}}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}_{2x,\ell}}(\boldsymbol{\theta}^*) + R_{\mathcal{S}_{2x,\ell}}(\boldsymbol{\theta}^*) \right) \quad (3.70)$$

$$\leq \sup_{\boldsymbol{\theta} \in \mathcal{H}} |R_{\mathcal{D}_{2x,\ell}}(\boldsymbol{\theta}) - R_{\mathcal{S}_{2x,\ell}}(\boldsymbol{\theta})| \quad (3.71)$$

$$= \phi(\mathcal{S}_{2x}) . \quad (3.72)$$

Similarly to the proof of Theorem 5, we recall Lemma 6 and notice that for any bounded loss function, the function  $\phi$  satisfies McDiarmid's inequality with  $C = 1/m \cdot c(X, H)$ . Therefore with probability  $\geq 1 - \delta$ :

$$A_1 \leq 2L \mathbb{E}_{\mathcal{D}_{2x}} \mathcal{R}(\mathcal{H} \circ \mathcal{S}_{2x}) + \frac{c(X, H)}{m} \sqrt{m \log \frac{1}{\delta}} \quad (3.73)$$

$$\leq 2v \frac{XHL}{\sqrt{2m}} + c(X, H) \sqrt{\frac{1}{m} \log \frac{1}{\delta}} \quad (3.74)$$

$$\leq \frac{\sqrt{2} + 1}{2} \frac{XHL}{\sqrt{m}} + c(X, H) \sqrt{\frac{1}{m} \log \frac{1}{\delta}} \quad (3.75)$$

where we also applied Lemmas 14 and 22 and the fact that:

$$v = \frac{1}{2} + \frac{1}{2} \sqrt{\frac{1}{2} - \frac{1}{m}} < \frac{\sqrt{2} + 1}{2\sqrt{2}}, \quad \forall m > 1.$$

Finally, we get with probability at least  $1 - \delta$ ,  $\delta > 0$  that:

$$\begin{aligned} R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) &\leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} \\ &+ c(X, H) \cdot \sqrt{\frac{1}{m} \log \left( \frac{1}{\delta} \right)} + 2|a|H \cdot \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2 . \end{aligned} \quad (3.76)$$

This proves the first statement. For the second statement, we apply Lemma 27 that provides the probabilistic bound for the norm discrepancy of the mean operators. We need to combine the two results. For any two events  $E, F$  it holds that:

$$1 - \mathbb{P}(E \wedge F) = \mathbb{P}(\neg(E \wedge F)) \quad (3.77)$$

$$= \mathbb{P}(\neg E \vee \neg F) \quad (3.78)$$

$$\leq \mathbb{P}(\neg E) + \mathbb{P}(\neg F) \quad (3.79)$$

Now consider that events:

$$E \doteq \left\{ 3.76 \right\}, F \doteq \left\{ \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2 \leq X \cdot \sqrt{\frac{d}{m} \log \left( \frac{2d}{\delta} \right)} \right\}$$

are true with probability at least  $1 - \delta/2$ , or equivalently they are false with at most



$\delta/2$  probability. We write:

$$1 - \mathbb{P}(E \wedge F) \leq \delta/2 - \delta/2 = \delta, \quad (3.80)$$

and therefore with  $1 - \delta$  probability:  $\mathbb{P}(E \wedge F) \geq 1 - \delta$ . Finally:

$$\begin{aligned} R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) &\leq \frac{\sqrt{2}+1}{2} \cdot \frac{XHL}{\sqrt{m}} + \\ c(X, H) \cdot \sqrt{\frac{1}{m} \log\left(\frac{2}{\delta}\right)} &+ 2|a|XH \cdot \sqrt{\frac{d}{m} \log\left(\frac{2d}{\delta}\right)}. \end{aligned} \quad (3.81)$$

■

### 3.5.5 Proof of Lemma 24

Consider Definition 9. Since we assume convex and differentiable losses, it suffices to compute the derivative in 0.

$$\ell'(x) = \ell'_e(x) + a \quad (3.82)$$

and because  $\ell'_e$  is an odd function, as derivative of an even function,  $\ell'_e(0) = 0$ . Therefore:

$$\ell'(0) < 0 \iff a < 0 \quad (3.83)$$

### 3.5.6 Proof of Theorem 26

Let us define:

$$R_{\mathcal{S}_X,\ell}(\boldsymbol{\theta}, \boldsymbol{\mu}) = \frac{1}{2m} \left( \sum_{i \in [m]} \sum_{\sigma} \ell(\sigma \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle) \right) + a \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle. \quad (3.84)$$

Define also the regularized loss:

$$R_{\mathcal{S}_X,\ell}(\boldsymbol{\theta}, \boldsymbol{\mu}, \lambda) \doteq R_{\mathcal{S}_X,\ell}(\boldsymbol{\theta}, \boldsymbol{\mu}) + \lambda \|\boldsymbol{\theta}\|_2^2. \quad (3.85)$$

Our proof begins following the same first steps as the proof of Lemma 17 in Altun and Smola [2006] (recalled in Lemma 25) and adding the steps that handle the lower bound on  $\ell''$ . Consider the following auxiliary function  $A_\ell(\boldsymbol{\tau})$ :

$$A_\ell(\boldsymbol{\tau}) \doteq (\nabla R_{\mathcal{S}_X,\ell}(\boldsymbol{\theta}^*, \boldsymbol{\mu}) - \nabla R_{\mathcal{S}_X,\ell}(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\mu}}))^\top (\boldsymbol{\tau} - \hat{\boldsymbol{\theta}}) + \lambda \|\boldsymbol{\tau} - \hat{\boldsymbol{\theta}}\|_2^2, \quad (3.86)$$

where the gradient of  $R_{\mathcal{S}_X,\ell}$  is computed with respect to parameter  $\boldsymbol{\theta}$ . The gradient

of  $A_\ell(\cdot)$  with respect to  $\tau$  is:

$$\nabla A_\ell(\tau) = \nabla R_{\mathcal{S}_X, \ell}(\theta^*, \mu) - \nabla R_{\mathcal{S}_X, \ell}(\hat{\theta}, \hat{\mu}) + 2\lambda(\tau - \hat{\theta}) . \quad (3.87)$$

The gradient is  $\mathbf{0}$  when  $\tau = \theta^*$ :

$$\nabla A_\ell(\theta^*) = \nabla R_{\mathcal{S}_X, \ell}(\theta^*, \mu) - \nabla R_{\mathcal{S}_X, \ell}(\hat{\theta}, \hat{\mu}) + 2\lambda(\theta^* - \hat{\theta}) \quad (3.88)$$

$$= \nabla R_{\mathcal{S}_X, \ell}(\theta^*, \mu, \lambda) - \nabla R_{\mathcal{S}_X, \ell}(\hat{\theta}, \hat{\mu}, \lambda) \quad (3.89)$$

$$= \mathbf{0} , \quad (3.90)$$

since both gradients are  $\mathbf{0}$  because of the optimality of  $\theta^*$  and  $\hat{\theta}$  with respect to  $R_{\mathcal{S}_X, \ell}(\cdot, \mu, \lambda)$  and  $R_{\mathcal{S}_X, \ell}(\cdot, \hat{\mu}, \lambda)$ . Moreover, the Hessian  $H$  of  $A_\ell$  is  $HA_\ell(\tau) = 2\lambda I \succeq 0$ , thus  $A_\ell$  is convex and therefore we can state that it is minimal at  $\tau = \theta^*$ . Additionally,  $A_\ell(\hat{\theta}) = 0$  by definition. It comes thus  $A_\ell(\theta^*) \leq 0$ , which yields equivalently:

$$0 \geq (\nabla R_{\mathcal{S}_X, \ell}(\theta^*, \mu) - \nabla R_{\mathcal{S}_X, \ell}(\hat{\theta}, \hat{\mu}))^\top (\theta^* - \hat{\theta}) + \lambda \|\theta^* - \hat{\theta}\|_2^2 \quad (3.91)$$

$$= \left( \frac{1}{2m} \sum_y \sum_i \nabla \ell(y \langle \theta^*, x_i \rangle) + a\mu - \frac{1}{2m} \sum_y \sum_i \nabla \ell(y \langle \hat{\theta}, x_i \rangle) - a\hat{\mu} \right)^\top (\theta^* - \hat{\theta}) + \lambda \|\theta^* - \hat{\theta}\|_2^2 \quad (3.92)$$

$$= \frac{1}{2m} \underbrace{\left( \sum_y \sum_i \ell(y \langle \theta^*, x_i \rangle) - \sum_y \sum_i \ell(y \langle \hat{\theta}, x_i \rangle) \right)^\top}_{=\beta} (\theta^* - \hat{\theta}) + a(\mu - \hat{\mu})^\top (\theta^* - \hat{\theta}) + \lambda \|\theta^* - \hat{\theta}\|_2^2 . \quad (3.93)$$

Let us lower bound  $\beta$ . We have  $\ell(y \langle \theta^*, x \rangle) = y \ell'(y \langle \theta^*, x \rangle) x$ , and a Taylor expansion brings that for any  $\theta^*, \hat{\theta}$ , there exists some  $\alpha \in [0, 1]$  such that, defining:

$$u_{\alpha, i} \doteq y \langle \alpha \theta^* + (1 - \alpha) \hat{\theta}, x_i \rangle , \quad (3.94)$$

we have:

$$\ell'(y \langle \theta^*, x_i \rangle) = \ell'(y \langle \hat{\theta}, x_i \rangle) + y \langle \theta^* - \hat{\theta}, x_i \rangle \ell''(u_{\alpha, i}) . \quad (3.95)$$

Thus we get:

$$\beta = \left( \sum_y \sum_i \ell(y \langle \boldsymbol{\theta}^*, \mathbf{x}_i \rangle) - \sum_y \sum_i \ell(y \langle \hat{\boldsymbol{\theta}}, \mathbf{x}_i \rangle) \right)^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}) \quad (3.96)$$

$$= \left( \sum_y \sum_i y (\ell'(y \langle \boldsymbol{\theta}^*, \mathbf{x}_i \rangle) - \ell'(y \langle \hat{\boldsymbol{\theta}}, \mathbf{x}_i \rangle)) \mathbf{x}_i \right)^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}) \quad (3.97)$$

$$= \left( \sum_y \sum_i (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}})^\top \mathbf{x}_i \ell''(u_{\alpha,i}) \mathbf{x}_i \right)^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}) \quad (3.98)$$

$$= 2 \sum_i \left( (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}})^\top \mathbf{x}_i \right)^2 \ell''(u_{\alpha,i}) \quad (3.99)$$

$$\geq 2\gamma \sum_i \left( (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}})^\top \mathbf{x}_i \right)^2 \quad (3.100)$$

$$= 2\gamma (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}})^\top S S^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}), \quad (3.101)$$

where matrix  $S \in \mathbb{R}^{d \times m}$  is formed by the observations of  $\mathcal{S}_X$  in columns. Inequality 3.100 comes from the fact that we require  $\gamma$ -strongly convex losses that are also twice differentiable, which implies that the second derivative is lower bounded by  $\gamma$ .

Define  $T \doteq (d / \sum_i \|\mathbf{x}_i\|_2^2) S S^\top$ . Its trace satisfies  $\text{tr } T = d$ . Let  $\lambda_d \geq \lambda_{d-1} \geq \dots \geq \lambda_1 > 0$  denote eigenvalues of  $T$ , with  $\lambda_1$  strictly positive because  $S S^\top = F^\top F \succ 0$ . The Arithmetic Mean-Geometric inequality brings:

$$\prod_2^d \lambda_k \leq \left( \frac{1}{d-1} \sum_{k=2}^d \lambda_k \right)^{d-1} \quad (3.102)$$

$$= \left( \frac{\text{tr } T - \lambda_1}{d-1} \right)^{d-1} \quad (3.103)$$

$$= \left( \frac{d - \lambda_1}{d-1} \right)^{d-1} \quad (3.104)$$

$$\leq \left( \frac{d}{d-1} \right)^{d-1}. \quad (3.105)$$

Multiplying both side by  $\lambda_1$  and rearranging yields:

$$\lambda_1 \geq \left( \frac{d-1}{d} \right)^{d-1} \det T \quad (3.106)$$

Let  $\lambda_\circ > 0$  denote the minimal eigenvalue of  $S S^\top$ . It satisfies:

$$\lambda_\circ = \left( \frac{1}{d} \sum_i \|\mathbf{x}_i\|_2^2 \right) \lambda_1 \quad (3.107)$$

and thus it comes from inequality 3.106:

$$\lambda_{\circ} \geq \left(\frac{d-1}{d}\right)^{d-1} \left(\frac{d}{\sum_i \|x_i\|_2^2}\right)^{d-1} \det SS^{\top} \quad (3.108)$$

$$= \left(\frac{d-1}{d}\right)^{d-1} \det \left[ \left(\frac{d}{\sum_i \|x_i\|_2^2}\right)^{1-\frac{1}{d}} SS^{\top} \right] \quad (3.109)$$

$$= \left(\frac{d-1}{d}\right)^{d-1} \det \bar{F}^{\top} \bar{F} \quad (3.110)$$

$$= \left(\frac{d-1}{d}\right)^{d-1} \text{vol}^2(\bar{F}) \quad (3.111)$$

$$\geq \frac{1}{e} \text{vol}^2(\bar{F}) . \quad (3.112)$$

We have used notation  $\text{vol}(\bar{F}) \doteq \sqrt{\det \bar{F}^{\top} \bar{F}}$ . Since:

$$(\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}})^{\top} SS^{\top} (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}) \geq \lambda_{\circ} \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2^2 , \quad (3.113)$$

combining 3.101 with 3.112 yields the following lower bound on  $\beta$ :

$$\beta \geq \frac{2}{e} \gamma \text{vol}^2(\bar{F}) \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2^2 . \quad (3.114)$$

Going back to 3.93, we get:

$$\lambda \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2^2 + a (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^{\top} (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}) + \frac{1}{em} \gamma \text{vol}^2(\bar{F}) \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2^2 \leq 0 . \quad (3.115)$$

Since  $a (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^{\top} (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}) \leq |a| \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}\|_2 \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2$ , after chaining the inequalities and solving for  $\|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2$  we get:

$$\|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2 \leq \frac{|a|}{\lambda + \frac{1}{em} \gamma \text{vol}^2(\bar{F})} \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}\|_2 , \quad (3.116)$$

as claimed.

## 3.6 Appendix: additional formal results

### 3.6.1 Mean and covariance operators

The intuition behind the relevance of the mean operator becomes clear once we rewrite it as follows.

**Lemma 28.** Let  $\boldsymbol{\mu} \doteq \mathbb{E}[yx]$  and  $\text{Cov}[x, y] \doteq \mathbb{E}_{\mathcal{S}} \left[ (yx - \boldsymbol{\mu})^2 \right]$ , respectively the mean and the

covariance operators. Let  $\pi_+ \doteq \mathbb{E}_{\mathcal{S}} 1\{y > 0\}$  be the positive label proportion of  $\mathcal{S}$ . Then:

$$\boldsymbol{\mu} = \text{Cov}[\mathbf{x}, y] + (2\pi_+ - 1)\mathbb{E}_{\mathcal{S}}[\mathbf{x}] . \quad (3.117)$$

Moreover, if observations are centered ( $\mathbb{E}_{\mathcal{S}}[\mathbf{x}] = 0$ ) or the labels are balanced ( $\pi_+ = 1/2$ ), then  $\boldsymbol{\mu} = \text{Cov}[\mathbf{x}, y]$ .

**Proof**

$$\text{Cov}[\mathbf{x}, y] = \mathbb{E}_{\mathcal{S}}[y\mathbf{x}] - \mathbb{E}_{\mathcal{S}}[y]\mathbb{E}_{\mathcal{S}}[\mathbf{x}] \quad (3.118)$$

$$= \boldsymbol{\mu} - \left( \frac{1}{m} \sum_{i:y_i>0} 1 - \frac{1}{m} \sum_{i:y_i<0} 1 \right) \mathbb{E}_{\mathcal{S}}[\mathbf{x}] \quad (3.119)$$

$$= \boldsymbol{\mu} - (2\pi_+ - 1) \mathbb{E}_{\mathcal{S}}[\mathbf{x}] . \quad (3.120)$$

The second statement follows immediately. ■

We have come to the unsurprising fact that — when observations are centered — the covariance  $\text{Cov}[\mathbf{x}, y]$  is what we need to know about the labels for learning linear models. The rest of the loss (in light of Factorization) may be seen as a data dependent regularizer.

### 3.6.2 The generality of factorization

Factorization is ubiquitous for any (margin) loss, beyond the theory seen so far. The decomposition on even and odd functions is actually all we need to factor  $\ell$ .

**Theorem 29** (Factorization). *For any sample  $\mathcal{S}$  and hypothesis  $h$  the empirical  $\ell$ -risk can be written as:*

$$R_{\mathcal{S}, \ell}(h) = \frac{1}{2} \mathbb{E}_{\mathcal{S}} \left[ \sum_{\sigma \in \mathcal{Y}} \ell(\sigma h(\mathbf{x})) \right] + \mathbb{E}_{\mathcal{S}} \left[ \ell_o(yh(\mathbf{x})) \right] \quad (3.121)$$

where  $\ell_o(\cdot)$  is odd and  $\ell_e(\cdot) \doteq \sum_{\sigma \in \mathcal{Y}} \ell(\sigma h(\cdot))$  is even and both uniquely defined.

Its range of validity is exemplified by 01 loss, a non-convex discontinuous piecewise linear function, which factors as:

$$\ell_e(x) = \begin{cases} \frac{1}{2} & x \neq 0 \\ 1 & \text{otherwise} \end{cases}, \quad \ell_o(x) = -\frac{1}{2} \text{sign}(x) . \quad (3.122)$$

It follows immediately that  $\mathbb{E}_{\mathcal{S}}[\ell_o(\cdot)]$  is sufficient for  $y$ . However,  $\ell_o$  is a function of model  $\theta$ . This defeats the purpose of a sufficient statistic, which we aim to be computable from data only and it is the main reason to define LOLS.

Moreover, Factorization goes beyond simple linear models, and can also be formulated for RKHS. To show that, notice that we satisfy all hypotheses of the Representer Theorem [Schölkopf and Smola, 2002].

**Theorem 30.** Let  $h(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{H}$  be a feature map into a Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{H}$  with symmetric positive definite kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , such that  $h : \mathbf{x} \rightarrow k(\cdot, \mathbf{x})$ . For any learning sample  $\mathcal{S}$ , the empirical  $\ell$ -risk  $R_{\mathcal{S}, \ell}(h)$  with  $\Omega : \|h\|_{\mathcal{H}} \rightarrow \mathbb{R}^+$  regularization can be written as:

$$\frac{1}{2} \mathbb{E}_{\mathcal{S}} \left[ \sum_{\sigma \in \mathcal{Y}} \ell(\sigma h(\mathbf{x})) \right] + \mathbb{E}_{\mathcal{S}} \left[ \ell_o(yh(\mathbf{x})) \right] + \Omega(\|h\|_{\mathcal{H}}) \quad (3.123)$$

and the optimal hypothesis admits a representation of the form  $h(\mathbf{x}) = \sum_{i \in [m]} \alpha_i k(\mathbf{x}, \mathbf{x}_i)$ .

It follows that all results of this Chapter may be read in the context of non-parametric models, with the *kernel* mean operator as sufficient statistic.

Finally, we have formulated the Factorization Theorem for classification problems. We can show a similar statement for regression with square loss  $\ell(\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle, y) = (\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle - y)^2$ :

$$\mathbb{E}_{\mathcal{S}}[(\langle \boldsymbol{\theta}, \mathbf{x} \rangle - y)^2] = \mathbb{E}_{\mathcal{S}}[\langle \boldsymbol{\theta}, \mathbf{x} \rangle^2] + \mathbb{E}_{\mathcal{S}}[y^2] - 2\langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle. \quad (3.124)$$

Taking the minimizers on both sides we obtain:

$$\operatorname{argmin}_{\boldsymbol{\theta}} \mathbb{E}_{\mathcal{S}}[(\langle \boldsymbol{\theta}, \mathbf{x} \rangle - y)^2] = \operatorname{argmin}_{\boldsymbol{\theta}} \mathbb{E}_{\mathcal{S}}[\langle \boldsymbol{\theta}, \mathbf{x} \rangle^2] - 2\langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle \quad (3.125)$$

$$= \operatorname{argmin}_{\boldsymbol{\theta}} \left\| X^{\top} \boldsymbol{\theta} \right\|_2^2 - 2\langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle. \quad (3.126)$$

This last result opens further applications yet, in the Thesis, we keep our focus on classification problems.

### 3.6.3 Factorization of non linear-odd losses

When  $\ell_o$  is not linear, we can find upper bounds in the form of affine functions. It suffices to be continuous and have asymptotes at  $\pm\infty$ .

**Lemma 31.** Let the loss  $\ell$  be continuous. Suppose that it has asymptotes at  $\pm\infty$ , i.e. there exist  $c_1, c_2 \in \mathbb{R}$  and  $d_1, d_2 \in \mathbb{R}$  such that:

$$\lim_{x \rightarrow +\infty} \ell(x) - c_1 x - d_1 = 0, \quad \lim_{x \rightarrow -\infty} \ell(x) - c_2 x - d_2 = 0 \quad (3.127)$$

then there exists  $q \in \mathbb{R}$  such that:  $\ell_o(x) \leq \frac{c_1 + c_2}{2} x + q$ .

**Proof** One can compute the limits at infinity of  $\ell_o$  to get:

$$\lim_{x \rightarrow +\infty} \ell_o(x) - \frac{c_1 + c_2}{2} x = \frac{d_1 - d_2}{2} \quad (3.128)$$

and:

$$\lim_{x \rightarrow -\infty} \ell_o(x) - \frac{c_1 + c_2}{2} x = \frac{d_2 - d_1}{2}. \quad (3.129)$$

Then  $q \doteq \sup\{\ell_o(x) - \frac{c_1+c_2}{2}x\} < +\infty$  as  $\ell_o$  is continuous. Thus  $\ell_o(x) - \frac{c_1+c_2}{2}x \leq q$ . ■

The Lemma covers many cases of practical interest outside the class of LOLs, e.g. hinge, absolute and Huber losses. Exponential loss is the exception since  $\ell_o(x) = -\sinh(x)$  cannot be bounded. Consider for instance hinge loss:  $\ell(x) = [1-x]_+$  is not differentiable in 1, however it is continuous with asymptotes at  $\pm\infty$ . Therefore, for any  $\theta$  its empirical risk is bounded as:

$$R_{S,hinge}(\theta) \leq \frac{1}{2}R_{S_{2x},hinge}(\theta) - \frac{1}{2}\langle\theta, \mu\rangle + q, \quad (3.130)$$

since  $c_1 = 0$  and  $c_2 = 1$ . An alternative proof of this result on hinge is provided next, giving the exact value of  $q = 1/2$ . The odd term for hinge loss is:

$$\ell_o(x) = \frac{1}{2}([1-x]_+ - [1+x]_+) \quad (3.131)$$

$$= \frac{1}{4}(-2x + |1-x| - |1+x|) \quad (3.132)$$

due to an arithmetic trick for the max function:

$$\max(a, b) = (a + b)/2 + |b - a|/2. \quad (3.133)$$

Then for any  $x$ :

$$|1-x| \leq |x| + 1, \quad (3.134)$$

$$|1+x| \geq |x| - 1 \quad (3.135)$$

and therefore:

$$\ell_o(x) \leq \frac{1}{4}(-2x + |x| + 1 - |x| + 1) = \frac{1}{2}(1-x). \quad (3.136)$$

We also provide a "if-and-only-if" version of Lemma 31 fully characterizing which losses can be upper bounded by a LOL.

**Lemma 32.** *Let  $\ell : \mathbb{R} \rightarrow \mathbb{R}$  a continuous function. Then there exists  $c_1, d_1, d_2 \in \mathbb{R}$  such that:*

$$\limsup_{x \rightarrow +\infty} \ell_o(x) - c_1x - d_1 = 0 \quad (3.137)$$

and:

$$\limsup_{x \rightarrow -\infty} \ell_o(x) - c_1x - d_2 = 0, \quad (3.138)$$

if and only if there exists  $q, q' \in \mathbb{R}$  such that  $\ell_o(x) \leq q'x + q$  for every  $x \in \mathbb{R}$ .

**Proof**

⇒) Suppose that such limits exist and they are zero for some  $c_1, d_1, d_2$ . Let prove

that  $\ell_o$  is bounded from above by a line. It holds that:

$$q = \sup_{x \in \mathbb{R}} \{\ell_o(x) - c_1 x\} < \infty , \quad (3.139)$$

because  $\ell_o$  is continuous. So for every  $x \in \mathbb{R}$ :

$$\ell_o(x) \leq c_1 x + q . \quad (3.140)$$

In particular we can take  $c_1$  as the angular coefficient of the line.

$\Leftarrow$ ) Vice versa we proceed by contradiction. Suppose that there exists  $q, q' \in \mathbb{R}$  such that  $\ell_o$  is bounded from above by  $\ell(x) = q'x + q$ . Suppose in addition that the conditions on the asymptotes (3.137) and (3.138) are false. This implies either the existence of a sequence  $x_n \rightarrow +\infty$  such that:

$$\lim_{n \rightarrow \infty} \ell_o(x_n) - q'x_n \rightarrow \pm\infty , \quad (3.141)$$

or the existence of another sequence  $x'_n \rightarrow -\infty$ :

$$\lim_{n \rightarrow \infty} \ell_o(y_n) - q'x'_n \rightarrow \pm\infty . \quad (3.142)$$

On one hand, if at least one of these two limits is  $+\infty$  then we already reach a contradiction, because  $\ell_o(x)$  is supposed to be bounded from above by  $\ell(x) = q'x + q$ . Suppose on the other hand that  $x_n \rightarrow +\infty$  is such that:

$$\lim_{n \rightarrow +\infty} \ell_o(x_n) - q'x_n \rightarrow -\infty . \quad (3.143)$$

Then defining  $x'_n = -x_n$  we have:

$$\lim_{n \rightarrow +\infty} \ell_o(w_n) - mx'_n \rightarrow +\infty , \quad (3.144)$$

and for the same reason as above we reach a contradiction. ■



### 3.6.4 More graphs on linear and non-linear-odd losses

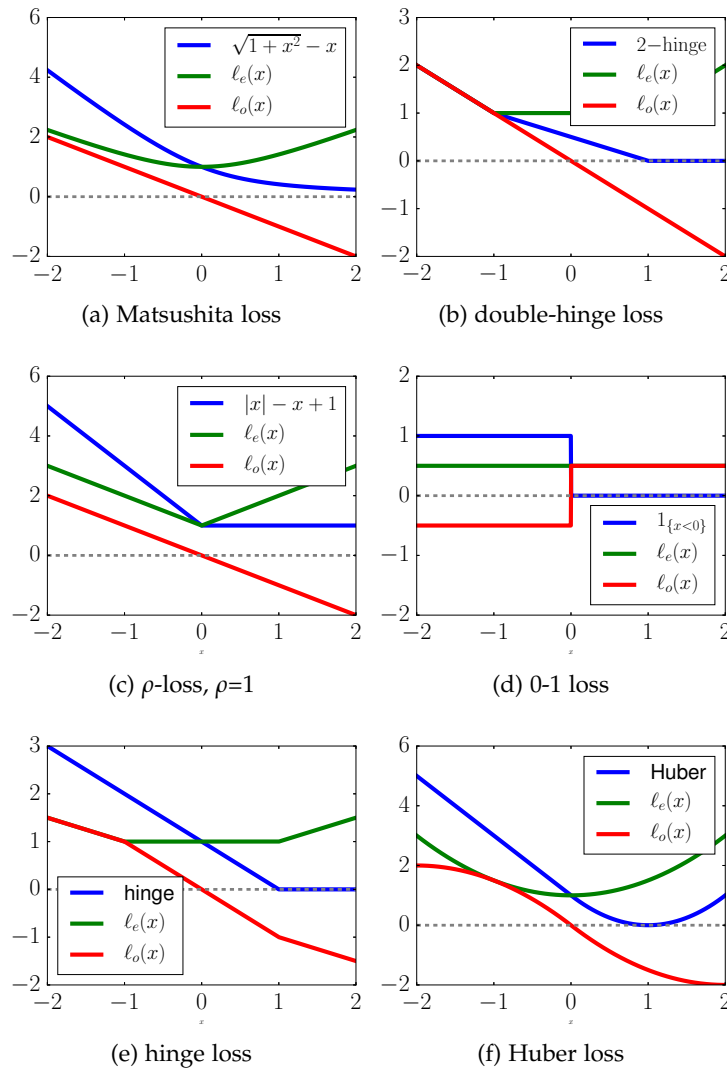


Figure 3.2: Linear-odd: Matsushita, double-hinge,  $\rho$ -loss. Non-linear-odd: 0-1, hinge, Huber.

### 3.6.5 The linear-odd losses of du Plessis et al. [2015]

The work of du Plessis et al. [2015] shows that a linear-odd condition on a *convex*  $\ell$  allows one to derive a tractable, *i.e.* still convex, loss for learning with *positive and unlabeled data*. The approach is similar to ours as it isolates a label-free term in the loss, with the goal of leveraging on the unlabeled examples too. Their manipulation of the loss *is not* equivalent to Theorem 18 though, as we explain here. Beside that, since we reason at the higher level of weakly supervised learning, we may frame

a solution for this setting by calling  $\mu$ SGD on  $\hat{\mu}$  defined above or by building on estimators derived from geometrical ideas discussed in Chapter 4.

For the sake of completeness, we review the use of convex LOLs in the cited work. Let  $\pi_+ \doteq \mathbb{P}(y = 1)$  and let  $\mathcal{D}_+$  and  $\mathcal{D}_-$  respectively the set of positive and negative examples in  $\mathcal{D}$ . Consider first:

$$\mathbb{E}_{(x,\cdot)\sim\mathcal{D}} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.145)$$

$$= \pi_+ \mathbb{E}_{(x,\cdot)\sim\mathcal{D}_+} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] + (1 - \pi_+) \mathbb{E}_{(x,\cdot)\sim\mathcal{D}_-} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.146)$$

Then, it is also true that:

$$\mathbb{E}_{(x,y)\sim\mathcal{D}} [\ell(y\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.147)$$

$$= \pi_+ \mathbb{E}_{(x,y)\sim\mathcal{D}_+} [\ell(y\langle\boldsymbol{\theta},\mathbf{x}\rangle)] + (1 - \pi_+) \mathbb{E}_{(x,y)\sim\mathcal{D}_-} [\ell(y\langle\boldsymbol{\theta},\mathbf{x}\rangle)] . \quad (3.148)$$

Solve Equation 3.145 for:

$$(1 - \pi_+) \mathbb{E}_{(x,y)\sim\mathcal{D}_-} [\ell(y\langle\boldsymbol{\theta},\mathbf{x}\rangle)] = (1 - \pi_+) \mathbb{E}_{(x,y)\sim\mathcal{D}_-} [-\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.149)$$

and substitute it into Equation 3.147 so as to obtain:

$$\mathbb{E}_{(x,y)\sim\mathcal{D}} [\ell(y\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.150)$$

$$= \pi_+ \mathbb{E}_{(x,y)\sim\mathcal{D}_+} [\ell(y\langle\boldsymbol{\theta},\mathbf{x}\rangle)] + \mathbb{E}_{(x,\cdot)\sim\mathcal{D}} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] - \pi_+ \mathbb{E}_{(x,\cdot)\sim\mathcal{D}_+} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.151)$$

$$= \pi_+ \left( \mathbb{E}_{(x,y)\sim\mathcal{D}_+} [\ell(+\langle\boldsymbol{\theta},\mathbf{x}\rangle)] - \mathbb{E}_{(x,\cdot)\sim\mathcal{D}_+} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \right) + \mathbb{E}_{(x,\cdot)\sim\mathcal{D}} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.152)$$

$$= \frac{\pi_+}{2} \mathbb{E}_{(x,y)\sim\mathcal{D}_+} [\ell_o(+\langle\boldsymbol{\theta},\mathbf{x}\rangle)] + \mathbb{E}_{(x,\cdot)\sim\mathcal{D}} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] , \quad (3.153)$$

by our usual definition of  $\ell_o(x) = \frac{1}{2}(\ell(x) - \ell(-x))$ . Recall that one of the goals of the authors is to conserve the convexity of this new crafted loss function. Then, du Plessis et al. [2015, Theorem 1] proceed stating that when  $\ell_o$  is convex, it must also be linear. And therefore they must focus on LOLs. Theorem 1 of du Plessis et al. [2015] is immediate from the point of view of our theory: in fact, an odd function can be convex or concave only if it also linear. The resulting expression based on the fact  $\ell(x) - \ell(-x) = 2ax$  simplifies into:

$$\mathbb{E}_{(x,y)\sim\mathcal{D}} [\ell(y\langle\boldsymbol{\theta},\mathbf{x}\rangle)] = a\pi_+ \mathbb{E}_{(x,y)\sim\mathcal{D}_+} [y\langle\boldsymbol{\theta},\mathbf{x}\rangle] + \mathbb{E}_{(x,\cdot)\sim\mathcal{D}} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] \quad (3.154)$$

$$= a\pi_+ \boldsymbol{\mu}_{\mathcal{D}_+} + \mathbb{E}_{(x,\cdot)\sim\mathcal{D}} [\ell(-\langle\boldsymbol{\theta},\mathbf{x}\rangle)] . \quad (3.155)$$

where  $\boldsymbol{\mu}_{\mathcal{D}_+}$  is a mean operator computed on positive examples only. Notice how the second term is instead label independent, although it is not an even function as in the Factorization Theorem.

## 3.7 References

The name *mean operator* in this context was originally used by Quadrianto et al. [2009]. In the theory of Hilbert space embedding, its importance is due to the injectivity of the map — under conditions on the kernel — which is used in applications such as two-sample and independence tests, feature extraction and covariate shift [Smola et al., 2007].

Factorization of logistic loss has been derived in other contexts in literature [Jaakkola and Jordan, 2000; Gao et al., 2016]. The work of Zantedeschi et al. [2016] implicitly utilizes Loss Factorization for formulating the  $\beta$ -risk, an augmented version of surrogate risks incorporating a confidence parameter for each individual label.

In most of this Chapter, we kept the lighter notation of linear classifiers, but nothing should prevent the extension of our results to non-parametric models, exchanging  $x$  with an implicit feature map  $h(x)$  recalling Theorem 30. A kernelized version of the mean-covariance operators Lemma 28 is given in Song et al. [2009]. A version of Theorem 23 may be derived for RKHS on top of Bartlett and Mendelson [2002]; Kakade et al. [2009].

### 3.7.1 The two-step procedure of Raghunathan et al. [2016]

The publication of our Patrini et al. [2016a] happened concurrently with Raghunathan et al. [2016], which shares many insights with our methodology. The authors design the same two-step procedure that amounts to sufficient statistic estimation followed by convex optimization. In fact, the underlying motivation is to avoid non-convex optimization and EM when labels are only known by indirect supervision. The learning framework of choice is MLE of graphical models, allowing an easy extension of the approach to structured label spaces — this is, in particular, an aspect not touched by our Thesis.

The two examples of weak supervision focus of Raghunathan et al. [2016] are close to what we consider in the next Chapters. The scenario of *learning under local privacy* can be mapped to the one of noisy labels; in particular, label corruption is dictated by a randomization schema that achieves differential privacy. The label sufficient statistic is estimated in a way that resembles Equation 5.2 in Chapter 5, with the same guarantee of unbiasedness. The setting of *learning with lightweight annotation* is a form of aggregated-level supervision, although in the context of structured prediction. The sufficient statistic is obtained by solving a linear system, in a similar fashion of Equation 4.7 and followers in Chapter 4.

Finally, Raghunathan et al. [2016] elaborate on the statistical efficiency of the proposed estimators and respective models, when compared to a model fitted by MLE. The two-step procedure is also studied by a geometrical viewpoint. The problem is mapped to the minimization of a KL divergence on an exponential family defined over a supervision function of  $y$ , not  $y$  itself, which makes the problem of MLE non-convex. It is shown that the two-step approach first optimizes over a relaxed set — the sufficient statistic estimation — and then projects onto the space of the exponen-

tial family, therefore bypassing the non-convex objective by an approximate solution.

### 3.7.2 Learning reductions

We have described the two-step algorithm as a procedure for *casting* a weakly supervised learning problem into an optimization problem, that is equivalent to a fully supervised problem. This is achieved via the equivalence expressed by Factorization. Solving a Machine Learning problem by solutions to other learning problems is a *learning reduction* [Beygelzimer et al., 2015]. Our work can be interpreted as such.

Following Beygelzimer et al. [2005], we define a supervised classification task as a triple  $(\mathcal{K}, \mathcal{Y}, \ell)$ , with supervised advice  $\mathcal{K}$ , predictions space  $\mathcal{Y}$  and loss  $\ell$ . A learning reduction is a procedure that casts a task  $(\mathcal{K}, \mathcal{Y}, \ell)$  into a simpler “reduced” task  $(\mathcal{K}', \mathcal{Y}', \ell')$ . We assume that  $\ell$  is linear-odd and  $\mathcal{K}$  represents any sort of weak supervision. Then, our two-step procedure is a reduction to binary classification  $(\mathcal{Y}, \mathcal{Y}, \ell)$ . The reduction of Meta algorithm 1 is somehow simple, in the sense that  $\mathcal{Y}$  does not change *and neither does*  $\ell$ . Algorithms 2 and 3 are examples of learning reductions from any weakly supervised classification problem to binary classification. Yet we actually modify the internal code of the “reduced learner”, the binary classifier, which contrasts with the concept of reduction. However, take as example Algorithm 2. We could as well write subgradients step as

$$\frac{1}{2} (\partial\ell(\langle \boldsymbol{\theta}^t, \mathbf{x}_i \rangle) + \partial\ell(-\langle \boldsymbol{\theta}^t, \mathbf{x}_i \rangle)) + a\boldsymbol{\mu} , \quad (3.156)$$

which equals  $\partial\ell$ , and thus the fully supervised SGD would be effectively untouched.

---

# Learning from label proportions

---

This Chapter is devoted to the problem of learning from label proportions. Loss Factorization is put to use by showing that we can learn from label proportions only via the estimation of the label sufficient statistic. Several formal guarantees are provided along with two practical algorithms which fill the gap left open in the last Chapter. In particular, we build on the Mean Map of [Quadrianto et al., 2009], relaxing its assumptions and yet proving finite sample guarantees for the estimation of the sufficient statistic. We also introduce a generalization of Rademacher complexity that is more meaningful in this learning setting. We show that our approach is particularly effective by extensive experiments simulating the construction of the instance groups and relative label proportions. The performance success is measured both against prior work and to the closeness to a fully supervised “Oracle”.

## 4.1 Motivation

We are interested in learning a binary classifier with supervision at the level of groups of observations, called *bags*. The type of information we assume available is the *label proportions* per bag, indicating the fraction of positive binary labels of its observations. We refer to this framework as *learning from label proportions* (LLP) [Quadrianto et al., 2009; Yu et al., 2014b]. Several applications fit the LLP abstraction:

(i) Only aggregated labels can be obtained due to the physical limits of measurement tools or constraints of communication and storage, as in particle mass spectrometry and high energy physics [Chen et al., 2006; Musicant et al., 2007], quality control in metallurgy [Stolpe and Morik, 2011] and embryos implantation for human assisted reproduction [Hernández-González et al., 2016]. Similarly, in sentiment analysis positive/negative rates are given to whole documents, but we may be interested in attributing the sentiment to individual sentences [Kotzias et al., 2015].

(ii) Labels existed once but they are now given in an aggregated fashion for privacy-preserving reasons, as in medical databases [Bhowmik et al., 2015], fraud detection [Rüping, 2010], banking [Ma et al., 2016], house price market, election results, census data, etc.

(iii) The problem is semi- or unsupervised but it is simple to obtain additional supervision for the unlabeled observations in the form of expectation from other data source or by domain experts [Quadrianto et al., 2009; Liang et al., 2009]. This is sometimes called *distant supervision*, where individual labels are not only unknown but also virtually never measured or recorded. Supervision is given by pairing two or more data domains: on the one side, the individual feature vectors and, on the other side, aggregated variables from an additional related source of information. For instance, Twitter data can be paired with constraints from county demographics, trends in first names and exemplar Twitter accounts strongly associated with a class label [Mohammady and Culotta, 2014; Ardehaly and Culotta, 2015]. To give another example, advertising experts can estimate the rate of successful ads at a campaign level, while the learning problem concerns the effectiveness of each individual banner and links shown to users [Wager et al., 2015]. Aggregate labels are frequently computed in terms of ratios or percents, *i.e.* proportions.

(iv) In Computer Vision applications, supervision is frequently only available to a higher level than the one interesting for prediction. For example in object detection, image segmentation and region recognition, a predictor should classify image patches or pixels, but those are seldom labeled in annotated image datasets [Kuck and de Freitas, 2005; Chen et al., 2009]. The problem of detecting events in videos by classifying singular frames is characterized by the same kind of weak supervision [Lai et al., 2014]. An example is the case of predicting per-pixel ice concentration on satellite images, when only region-level concentration is available from experts annotation [Li and Taylor, 2015]. Methods for LLP have also been used for learning “human nameable properties” of images [Yu et al., 2014a].

## 4.2 Learning setting

In learning from label proportions, we do not observe directly  $\mathcal{S}$  but  $\mathcal{S}_X$ , which denotes the learning sample without any label. We are given its partition in  $n > 1$  bags, that is  $\mathcal{S}_X = \cup_j \mathcal{S}_j, j \in [n]$  and  $\forall j \neq j', \mathcal{S}_j \cap \mathcal{S}_{j'} = \emptyset$ . This is by no means restrictive: the case when bags constitute a *cover* instead of a partition of  $\mathcal{S}_X$  can be reduced to the defined setting by copying examples among bags and by reweighing the loss accordingly. The “bag assignment function” that partitions  $\mathcal{S}$  is unknown but fixed. We do not assume any knowledge about how observations have been partitioned, in spite of the fact that it is easy to imagine several applications for which something about the process is known. Each bag  $j$  is provided with its respective *label proportion*:

$$\pi_j \doteq \mathbb{E}_{\mathcal{S}_j}[y = 1] \tag{4.1}$$

and bag proportions  $p_j \doteq m_j/m$  with  $m_j = |\mathcal{S}_j|$ . For the sake of clarity, in this Chapter we denote use  $i, j$  and  $k$  to refer respectively to examples, bags and features.

loss name	$\ell(x)$	$-\phi(x)$
logistic loss	$\log(1 + \exp(-x))$	$-x \log x - (1 - x) \log(1 - x)$
square loss	$(1 - x)^2$	$x(1 - x)$
Matsushita loss	$-x + \sqrt{1 + x^2}$	$\sqrt{x(1 - x)}$

Table 4.1: Correspondence between permissible function  $\phi$  and loss  $\ell$ .

$\mathcal{H}$  is be the space of linear classifiers — for a kernelized version of our algorithms we point to the formulation of Quadrianto et al. [2009] that constitutes the basis of our algorithmic contributions.

#### 4.2.1 Symmetric proper losses

We briefly discuss a particular class of losses that are called *symmetric proper losses* (SPL). While not necessary for the statement of the algorithms, some theoretical guarantees derived in this Chapter assume to work with the convenient functional shape of SPLs. Conceptually, we are not moving away from the framework presented in Chapter 3, since the class linear-odd losses has been shown to strictly include the one defined here. Symmetric proper losses are axiomatized in Nock and Nielsen [2009].

**Definition 33** (Symmetric proper losses). *A loss  $\ell(y, h(\mathbf{x}))$  is called symmetric proper (SPL) if it is:*

- (i) lower bounded
- (ii) proper
- (iii) symmetric, that is,  $\ell(y, h(\mathbf{x})) = \ell(-y, -h(\mathbf{x}))$
- (iv) twice differentiable

Notice that since we consider margin losses (Definition 2), condition (iii) is always satisfied. There exists a bijection between the set of symmetric proper losses and a set of *permissible* functions [Kearns and Mansour, 1996]  $\phi$  which are differentiable, strictly convex, symmetric about  $1/2$  and with  $\text{dom}(\phi) \supseteq [0, 1]$ . We also let  $b_\phi \doteq \phi(1/2) - a_\phi > 0$  and  $\phi^*$  be the convex conjugate of  $\phi$ . The bijection is the property that we need, formalized in the next Lemma from Nock and Nielsen [2009].

**Lemma 34.** *A function  $\ell$  is a symmetric proper loss if and only if there exists a permissible function  $\phi$  such that:*

$$\ell(x) = a_\phi + \frac{\phi^*(-x)}{b_\phi} . \quad (4.2)$$

This representation of SPL, which is originally obtained through Bregman divergences in Nock and Nielsen [2009], is sufficient to prove that SPLs factor the same way as linear-odd losses; we have shown that in Chapter 3. The original, more involved proof of Factorization for SPL is contained Patrini et al. [2014], Lemma 1.

### 4.3 Estimating the sufficient statistic

In a direct implementation of our two-step framework (Meta Algorithm 1), the estimation of the mean operator appears to be the learning bottleneck for LLP. The mean operator is in fact unknown since we cannot recover it simply from the label proportions. We start by expanding the mean operator in its bag-wise label-wise components:

$$\boldsymbol{\mu} = \mathbb{E}_{\mathcal{S}} [y\mathbf{x}] \quad (4.3)$$

$$= \sum_{j=1}^n p_j \mathbb{E}_{\mathcal{S}_j} [y\mathbf{x}] \quad (4.4)$$

$$= \sum_{j=1}^n p_j \left( \mathbb{E}_{\mathcal{S}_j} [\mathbf{x}, y = 1] + \mathbb{E}_{\mathcal{S}_j} [\mathbf{x}, y = -1] \right) \quad (4.5)$$

$$= \sum_{j=1}^n p_j \left( \pi_j \mathbb{E}_{\mathcal{S}_j} [\mathbf{x}|y = 1] - (1 - \pi_j) \mathbb{E}_{\mathcal{S}_j} [\mathbf{x}|y = -1] \right) . \quad (4.6)$$

The quantities  $p_j$  and  $\pi_j$  are known. So our problem is turned into the estimation of  $2n$  vectors of unknowns  $\mathbb{E}_{\mathcal{S}_j} [\mathbf{x}|y] \in \mathbb{R}^d$ . Those are the average feature vector — the center — for each bag conditioned to the label sign. How can we use the information available to compute those centers? We come up with a linear system of equations of which they are the only unknowns. By law of total probability:

$$\mathbb{E}_{\mathcal{S}} [\mathbf{x}] = \sum_{y \in \mathcal{Y}} \pi_j \mathbb{E}_{\mathcal{S}_j} [\mathbf{x}|y] , \quad (4.7)$$

where every  $\mathbb{E}_{\mathcal{S}_j} [\mathbf{x}]$  can be computed even without label knowledge as they are simply the bag-wise average feature vectors. For convenience, we rewrite the system in matrix form. Let  $\mathbf{b}_j^y = \mathbb{E}_{\mathcal{S}_j} [\mathbf{x}|y]$  and  $\mathbf{b}_j = \mathbb{E}_{\mathcal{S}_j} [\mathbf{x}]$ . The  $2n$   $\mathbf{b}_j^y$ 's are solutions of:

$$B - \Pi^{\top} B^{\pm} = \mathbf{0} , \quad (4.8)$$

where  $B \doteq [\mathbf{b}_1 | \mathbf{b}_2 | \dots | \mathbf{b}_n]^{\top} \in \mathbb{R}^{n \times d}$ ,  $\Pi \doteq [\text{diag}(\boldsymbol{\pi}) | \text{diag}(\mathbf{1} - \boldsymbol{\pi})]^{\top} \in \mathbb{R}^{2n \times n}$  and  $\boldsymbol{\pi}$  stacks together all label proportions in a vector.  $B^{\pm} \in \mathbb{R}^{2n \times d}$  is the matrix of unknowns:

$$B^{\pm} \doteq \left[ \underbrace{\mathbf{b}_1^{+1} | \mathbf{b}_2^{+1} | \dots | \mathbf{b}_n^{+1}}_{(B^+)^{\top}} \mid \underbrace{\mathbf{b}_1^{-1} | \mathbf{b}_2^{-1} | \dots | \mathbf{b}_n^{-1}}_{(B^-)^{\top}} \right]^{\top} . \quad (4.9)$$

System 4.8 is under-determined, as it is made of  $n \times d$  equations on  $2n \times d$  unknowns. This fact should not be a surprise as it expresses the ill-posed nature of the problem of learning a classifier with label proportions only. To solve the system we need to resort to additional assumptions. We remark that solving System 4.8 accounts for the first estimation step, followed by Equation 4.3 to recover the mean operator. We now focus on solving the linear system.



**Algorithm 4:** Mean Map (MM)

---

**Input:**  $\mathcal{S}_j, \pi_j, p_j, \forall j \in [n], \ell$  is  $a$ -LOL,  $\lambda > 0$   
 $\hat{B}^\pm \leftarrow \Pi^\dagger B$   
 $\hat{\mu} \leftarrow \sum_{y \in \mathcal{Y}} y p(y) \hat{\mathbf{b}}^y$   
 $\hat{\theta} \leftarrow \operatorname{argmin}_{\theta} \frac{1}{2m} \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \ell(y \langle \theta, \mathbf{x}_i \rangle) + a \langle \theta, \hat{\mu} \rangle + \lambda \|\theta\|_2^2$   
**Output:**  $\hat{\theta}$

---

#### 4.4 Mean Map algorithm of Quadrianto et al. [2009]

Quadrianto et al. [2009] originally propose this framework for LLP. It enforces a *homogeneity assumption*, that is a statement of conditional independence for  $j$ :

$$(A4.0) \quad \forall j, \quad \mathbb{E}_{\mathcal{S}_j}[\mathbf{x}|y] = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|y] .$$

Under this hypothesis, the number of unknowns falls to  $2d$ . This is enough to obtain a system that is not under-determined, granted that  $n \geq 2$ . In this case, the linear system is simplified by defining  $B^\pm \doteq [\mathbf{b}^+, \mathbf{b}^-]^\top \in \mathbb{R}^{n \times d}$  and  $\Pi \doteq [\pi | \mathbf{1} - \pi] \in \mathbb{R}^{n \times 2}$ , and the solution is found by pseudo-inversion as  $\hat{B}^\pm = \Pi^\dagger B$ . Let  $\hat{\mathbf{b}}^y$  denote the rows of  $\hat{B}^\pm$  in lieu of the true  $\mathbf{b}^y = \mathbb{E}_{\mathcal{S}}[\mathbf{x}|y]$ . The mean operator is then estimated by:

$$\hat{\mu} = \sum_{y \in \mathcal{Y}} y p(y) \hat{\mathbf{b}}^y , \quad (4.10)$$

where  $p(y)$  is the probability of the label over the whole  $\mathcal{S}$  and can be derived from  $\pi$ . This is effectively the first example of implementation of the Meta Algorithm 1. We state it in Algorithm 4. The last step of convex optimization is expressed with  $L_2$  regularization.

#### 4.5 Laplacian Mean Map

Our proposal, the *Laplacian Mean Map* (LMM) algorithm, aims to solve System 4.8 at a larger extent. In order to do that, we relax the homogeneity assumption as:

$$(A4.1) \quad \forall j, j' \text{ if } j \approx j' \text{ then } \mathbb{E}_{\mathcal{S}_j}[\mathbf{x}|y] \approx \mathbb{E}_{\mathcal{S}_{j'}}[\mathbf{x}|y] .$$

That is, instead of setting every  $\mathbb{E}_{\mathcal{S}_j}[\mathbf{x}|y]$  equal for each bag  $j$  for a given  $y$ , we assume them close when bags are similar, while none is constrained to be equal. The similarity between bags  $j$  and  $j'$  is a domain-specific parameter of the algorithm and we will refer to it by  $v_{j,j'} \geq 0$ ; we discuss the choice of  $v$  below. To incorporate the assumption into the estimation, we solve System 4.8 by least square minimization

with regularization as:

$$\begin{aligned} \operatorname{argmin}_{\mathbf{b}_j^1, \mathbf{b}_j^{-1}} \sum_j \left( \mathbf{b}_j - \pi_j \mathbf{b}_j^1 + (1 - \pi_j) \mathbf{b}_j^{-1} \right)^2 \\ + \gamma \sum_{j,j'} v_{j,j'} \left[ \left( \mathbf{b}_j^1 - \mathbf{b}_{j'}^1 \right)^2 + \left( \mathbf{b}_j^{-1} - \mathbf{b}_{j'}^{-1} \right)^2 \right]. \end{aligned} \quad (4.11)$$

Depending on the regularization strength  $\gamma \geq 0$ , the more the bags are similar — larger  $v_{j,j'}$  — the more their relative estimate of  $\mathbf{b}_j^\pm$  are close to each other. We can restate the problem in matrix form by the Laplacian of the symmetric matrix  $V \in \mathbb{R}^{n \times n}$ , which is the adjacency matrix of the graph induced by the similarity  $v_{j,j'}$ . The Laplacian is defined as  $L_a = D - V$ , and  $D$  is a diagonal matrix such that  $D_j = \sum_{j'} v_{j,j'}$ . For any vector  $\mathbf{u} \in \mathbb{R}^n$ , it holds that:

$$\mathbf{u}^\top L_a \mathbf{u} = \mathbf{u}^\top D \mathbf{u} - \mathbf{u}^\top V \mathbf{u} \quad (4.12)$$

$$= \sum_j D_j u_j^2 - \sum_{j,j'} v_{j,j'} u_j u_{j'} \quad (4.13)$$

$$= \frac{1}{2} \left( \sum_j D_j u_j^2 - 2 \sum_{j,j'} v_{j,j'} u_j u_{j'} + \sum_{j'} D_{j'} u_{j'}^2 \right) \quad (4.14)$$

$$= \sum_{j,j'} v_{j,j'} (u_j - u_{j'})^2 \quad (4.15)$$

Expression 4.15 is obtained by applying the definition of  $D_j$ . This is a standard result in Spectral Graph Theory [Von Luxburg, 2007]. We make use of this equivalence for both  $\mathbf{b}_j^1$  and  $\mathbf{b}_j^{-1}$ , but separately, as the two vectors are subject to distinct constraints. Thanks to the structure of matrix  $B^\pm$ , we can define:

$$L \doteq \varepsilon I + \begin{bmatrix} L_a & | & \mathbf{0} \\ \mathbf{0} & | & L_a \end{bmatrix} \in \mathbb{R}^{2n \times 2n}, \quad (4.16)$$

such that<sup>1</sup>

$$\hat{B}^\pm = \operatorname{argmin}_{X \in \mathbb{R}^{2n \times d}} \operatorname{tr} \left( B - \Pi^\top X \right)^\top D_w \left( B - \Pi^\top X \right) + \gamma \operatorname{tr} X^\top L X. \quad (4.17)$$

$D_w \doteq \operatorname{diag}(w)$  is a user-fixed bias matrix with non-negative element on the diagonal  $w$ . For example, we can re-weight the importance of the linear equations by the size of the respective bags  $m_j$ , by  $w = \mathbf{p}$ . The second term has the form of a manifold regularizer of Belkin et al. [2006], which allows us to reinterpret our assumption from a geometrical standpoint: the bag label-wise feature averages  $\mathbb{E}_{S_j}[\mathbf{x}|y]$  live on a low-dimensional manifold parameterized by the unidimensional similarity function  $v_{j,j'}$ . The Laplacian matrix  $L_a$  is an empirical, discrete approximation of the manifold.

<sup>1</sup>with  $\varepsilon > 0$  to assure non-singularity and numerical stability, see the proof of Theorem 35 in Appendix 4.11.2.

**Algorithm 5:** Laplacian Mean Map (LMM)

**Input:**  $\mathcal{S}_j, \pi_j, p_j, \forall j \in [n], \ell$  is  $a$ -LOL,  $\gamma > 0, D_w, L, \lambda > 0$

$$\hat{B}^\pm \leftarrow (\Pi D_w \Pi^\top + \gamma L)^{-1} \Pi D_w B$$

$$\hat{\mu} \leftarrow \frac{1}{m} \sum_{j=1}^n p_j (\pi_j \hat{\mathbf{b}}_j^+ - (1 - \pi_j) \hat{\mathbf{b}}_j^-)$$

$$\hat{\theta} \leftarrow \operatorname{argmin}_{\theta} \frac{1}{2m} \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \ell(y \langle \theta, \mathbf{x}_i \rangle) + a \langle \theta, \hat{\mu} \rangle + \lambda \|\theta\|_2^2$$

**Output:**  $\hat{\theta}$

The size of the Laplacian is  $O(n^2)$ , which is small compared to  $O(m^2)$  if there are not many bags. This is in contrast with traditional approaches for semi-supervised learning, where the Laplacian matrix is formed on top of similarity between examples, instead of bags [Belkin et al., 2006].

Problem 4.17 admits global optimum in closed form.

**Theorem 35.** *The solution to Problem 4.17*

$$\operatorname{argmin}_{X \in \mathbb{R}^{2n \times d}} \operatorname{tr} \left( B - \Pi^\top X \right)^\top D_w \left( B - \Pi^\top X \right) + \gamma \operatorname{tr} X^\top L X \quad (4.18)$$

is

$$\hat{B}^\pm \doteq \left( \Pi D_w \Pi^\top + \gamma L \right)^{-1} \Pi D_w B \quad (4.19)$$

Proof in 4.11.1. This result explains the role of the penalty  $\varepsilon I$  in 4.16 as  $\Pi D_w \Pi^\top$  and  $L$  have respectively  $n$ - and  $(\geq 1)$ -dim null spaces, so the inversion may not be possible. Even when this does not happen exactly, this may incur numerical instabilities in computing the inverse. For domains where this risk exists, picking a small  $\varepsilon > 0$  solves the problem.

Let  $\hat{\mathbf{b}}_j^y$  denote the rows of  $\hat{B}^\pm$ , the solution of Problem 4.17, in lieu of the true  $\mathbf{b}_j^y = \mathbb{E}_{\mathcal{S}_j}[x|y]$ . Those statistics are used to estimate the mean operator by Equation 4.6. We state the Laplacian Mean Map algorithm (LMM) in Algorithm 5, a second example of our two-step estimation/minimization procedure.

## 4.6 Estimation: formal guarantees

We have shown how to estimate the sufficient statistic. Although, the set up of Problem 4.17 is strongly dependent on our assumptions. It is therefore relevant to study formal guarantees of our estimator. We compare  $\mu_j \doteq \pi_j \mathbf{b}_j^+ - (1 - \pi_j) \mathbf{b}_j^-$  to our estimates  $\hat{\mu}_j \doteq \pi_j \hat{\mathbf{b}}_j^+ - (1 - \pi_j) \hat{\mathbf{b}}_j^-$ ,  $\forall j \in [n]$ , granted that  $\mu = \sum_j p_j \mu_j$  and  $\hat{\mu} = \sum_j p_j \hat{\mu}_j$ .

**Theorem 36.** *Suppose that  $\gamma$  satisfies  $\gamma \sqrt{2} \leq (\varepsilon(2n)^{-1} + \max_{j \neq j'} v_{jj'}) / \min_j w_j$ . Let  $M \doteq$*

$[\boldsymbol{\mu}_1 | \boldsymbol{\mu}_2 | \dots | \boldsymbol{\mu}_n]^\top \in \mathbb{R}^{n \times d}$ ,  $\hat{M} \doteq [\hat{\boldsymbol{\mu}}_1 | \hat{\boldsymbol{\mu}}_2 | \dots | \hat{\boldsymbol{\mu}}_n]^\top \in \mathbb{R}^{n \times d}$  and:

$$\zeta(V, B^\pm) \doteq \left( \frac{\varepsilon}{2n} + \max_{j \neq j'} v_{jj'} \right)^2 \|B^\pm\|_F .$$

The following holds:

$$\|M - \hat{M}\|_F \leq \frac{\sqrt{n/2}}{\min_j w_j^2} \cdot \zeta(V, B^\pm) . \quad (4.20)$$

Proof in 4.11.2. The multiplicative factor to  $\zeta$  in 4.20 is roughly  $O(n^{5/2})$  when there is no large discrepancy in the bias matrix  $D_w$ , so the upper bound is driven by  $\zeta(\cdot, \cdot)$  when there are not many bags. We study its variations when the ‘‘distinguishability’’ between bags increases. This setting is interesting because in this case we may kill two birds with one stone, with the estimation of  $M$  and the subsequent learning problem potentially easier, in particular for linear separators. We consider two examples for  $v_{jj'}$ , the first being the normalized association [Shi and Malik, 2000]:

$$v_{jj'}^{nc} \doteq \frac{1}{2} \left( \frac{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j)}{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j \cup \mathcal{S}_{j'})} + \frac{\text{ASSOC}(\mathcal{S}_{j'}, \mathcal{S}_{j'})}{\text{ASSOC}(\mathcal{S}_{j'}, \mathcal{S}_j \cup \mathcal{S}_{j'})} \right) = \text{NASSOC}(\mathcal{S}_j, \mathcal{S}_{j'}) , \quad (4.21)$$

$$v_{jj'}^{G,s} \doteq \exp(-\|\mathbf{b}_j - \mathbf{b}_{j'}\|_2/s) , s > 0 . \quad (4.22)$$

Here,  $\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_{j'}) \doteq \sum_{\mathbf{x} \in \mathcal{S}_j, \mathbf{x}' \in \mathcal{S}_{j'}} \|\mathbf{x} - \mathbf{x}'\|_2^2$ . To put these two similarity measures in the context of Theorem 36, consider the setting where we can make assumption:

$$(A4.2) \exists \kappa > 0 \text{ such that } \forall j, j' \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2 \geq \kappa \max_{y,j} \|\mathbf{b}_j^y\|_2^2 .$$

This is a weak distinguish-ability property as if no such  $\kappa$  exists, then the centers of distinct bags may just be confounded. Consider also the additional assumption that:

(A4.3)  $\exists \kappa' > 0$  such that  $\max_j d_j^2 \leq \kappa', \forall j \in [n]$ , where  $d_j \doteq \max_{\mathbf{x}_i, \mathbf{x}_{i'} \in \mathcal{S}_j} \|\mathbf{x}_i - \mathbf{x}_{i'}\|_2$  is a bag’s diameter.

In the following Lemma, the *little-o* notation is with respect to the ‘‘largest’’ unknown in Equation 4.11, i.e.  $\max_{y,j} \|\mathbf{b}_j^y\|_2$ .

**Lemma 37.** *There exists  $\varepsilon_* > 0$  such that  $\forall \varepsilon \leq \varepsilon_*$ , the following holds:*

- (i)  $\zeta(V^{nc}, B^\pm) = o(1)$  under assumptions (A4.2 + A4.3);
- (ii)  $\zeta(V^{G,s}, B^\pm) = o(1)$  under assumption (A4.2),  $\forall s > 0$ .
- (iii)  $\zeta(V^{G,s}, B^\pm)$  converges faster than  $\zeta(V^{nc}, B^\pm)$ .

Proof in 4.11.3. Hence, provided a weak (A4.2) or stronger (A4.2+A4.3) distinguish-ability assumption holds, the divergence between  $M$  and  $\hat{M}$  gets smaller with the

increase of the norm of the unknowns  $\mathbf{b}_j^y$ . The following Lemma shows that both similarities also partially encode the hardness of solving the classification problem with linear separators, so that the manifold regularizer “limits” the distortion of the  $\hat{\mathbf{b}}_j^y$ s between two bags that tend not to be linearly separable.

**Lemma 38.** *Take  $v_{jj'} \in \{v_{jj'}^{G,s}, v_{jj'}^{nc}\}$ . There exists  $0 < \kappa_l < \kappa_n < 1$  such that:*

- (i) *if  $v_{jj'} > \kappa_n$  then  $\mathcal{S}_j, \mathcal{S}_{j'}$  are not linearly separable;*
- (ii) *if  $v_{jj'} < \kappa_l$  then  $\mathcal{S}_j, \mathcal{S}_{j'}$  are linearly separable.*

Proof in 4.11.4. This Lemma is an advocacy to fit  $s$  of  $v_{jj'}^{G,s}$  in a data dependent way.

A question may be raised as to whether finite sample approximation results like Theorem 36 can be proven for the Mean Map estimator. The following answers in the negative.

**Lemma 39.** *For any  $\gamma > 0$ , the Mean Map estimator (Equation 4.10) cannot guarantee:*

$$\frac{\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2}{\max_{y,j} \|\mathbf{b}_j^y\|_2} \leq 2 - \gamma, \quad (4.23)$$

*even when (A4.2 + A4.3) hold.*

Proof in 4.11.5. It is not hard to check that a randomized procedure that builds  $\hat{\boldsymbol{\mu}} \doteq yx$  for some random  $(x, y) \in \mathcal{S}$  guarantees  $O(2 + \gamma)$  approximability when some bags are close to the convex hull of  $\mathcal{S}$ , for small  $\gamma > 0$ . In contrast, the Mean Map estimation may be very poor in that respect.

We now move to consider the quality of the model. The following Theorem shows a data dependent approximation bound that refines Theorem 26 under the assumption that we use SPLs that satisfy an additional requirement.

**Definition 40.** *Let  $\phi'(x)$  be the derivative of  $\phi(x)$  with regard to scalar  $x$ . We say that a symmetric proper loss  $\ell$  is proper scoring compliant (PSC) when:*

$$\forall \mathbf{x}_i \in \mathcal{S}, \langle \boldsymbol{\theta}^*, \mathbf{x}_i \rangle \text{ and } \langle \hat{\boldsymbol{\theta}}, \mathbf{x}_i \rangle \in \phi'([0, 1]), \ .$$

This condition always holds for logistic and Matsushita losses for which  $\phi'([0, 1]) = \mathbb{R}$ . For other losses such as square loss for which  $\phi'([0, 1]) = [-1, 1]$ , shrinking the observations in a ball of sufficiently small radius is sufficient to ensure this. Proof in 4.11.6.

**Theorem 41.** *Let  $X \doteq \max_i \|\mathbf{x}_i\|_2$ . Let  $\ell$  be SPL and PSC, and let  $\phi', \phi''$  be first and second derivative of  $\phi$ . Let  $\mathbf{f}_k \in \mathbb{R}^m$  denote the vector encoding the  $k^{\text{th}}$  feature variable in  $\mathcal{S}$ :  $f_{ki} = x_{ik}$ , with  $k \in [d]$ . Let  $\bar{\mathbf{F}}$  denote the feature matrix with column-wise normalized feature vectors:*

$$\bar{\mathbf{f}}_k \doteq \left( \frac{d}{\sum_{k'} \|\mathbf{f}_{k'}\|_2^2} \right)^{(d-1)/(2d)} \cdot \mathbf{f}_k. \quad (4.24)$$

Without loss of generality, we assume  $\bar{F}^\top \bar{F}$  positive definite. Call  $\hat{\boldsymbol{\theta}}$  and  $\boldsymbol{\theta}^*$  the minimizer of problem (ii) in Meta Algorithm 1 when, respectively, the empirical risk is computed with  $\hat{\boldsymbol{\mu}}$  from the estimator of (i) and  $\boldsymbol{\mu} = \mathbb{E}_{\mathcal{S}}[y\mathbf{x}]$ . We have that:

$$\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2^2 \leq \frac{1}{\lambda + q} \|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2^2, \quad (4.25)$$

with:

$$q \doteq \frac{1}{em} \cdot \frac{\det(\bar{F}^\top \bar{F})}{b_\phi \phi''\left(\phi^{-1}\left(\frac{\bar{q}}{b_\phi \lambda}\right)\right)} (> 0), \quad (4.26)$$

for some  $\bar{q}$  in the interval  $[\pm(X + \max\{\|\boldsymbol{\mu}\|_2, \|\hat{\boldsymbol{\mu}}\|_2\})]$ .

## 4.7 Alternating Mean Map

We design a second algorithm for LLP called Alternating Mean Map (AMM). This algorithm does not belong to the family of algorithms defined by our two-step procedure stated in Meta Algorithm 1. Instead, it follows the more conventional strategy of coordinate descent for weakly supervised learning. Estimation of sufficient statistic and model learning are interleaved. Once we learn a new model, it can be applied to provide a better estimation of the mean operator. We do so by estimating the unknown labels under the constraints given by the label proportions. LMM can be used to initialize the alternating procedure. In the experiments of Section 4.9 we also demonstrate that is the best choice for initialization.

Let us denote the set of labelings that are *consistent* with the proportions  $\boldsymbol{\pi}$  as:

$$\Sigma_\pi \doteq \left\{ \boldsymbol{\sigma} \in \Sigma_m : \sum_{i: x_i \in \mathcal{S}_j} \sigma_i = (2\pi_j - 1)m_j, \forall j \in [n] \right\}, \quad (4.27)$$

and the (possibly biased) mean operator computed as  $\boldsymbol{\mu}(\boldsymbol{\sigma}) \doteq (1/m) \sum_i \sigma_i x_i$  from some  $\boldsymbol{\sigma} \in \Sigma_\pi$ . Notice that the true mean operator is  $\boldsymbol{\mu} = \boldsymbol{\mu}(\boldsymbol{\sigma})$  for at least one  $\boldsymbol{\sigma} \in \Sigma_\pi$ . The objective function augmented with latent variables  $\boldsymbol{\sigma}$  is the following:

$$\frac{1}{2m} \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \ell(y \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle) + a \langle \boldsymbol{\theta}, \boldsymbol{\mu}(\boldsymbol{\sigma}) \rangle + \lambda \|\boldsymbol{\theta}\|_2^2. \quad (4.28)$$

This approach fits into the framework of Problem 2.27 for ERM under weak supervision. An important insight is that, by Factorization,  $\boldsymbol{\sigma}$  affects the empirical risk only through the value of  $\langle \boldsymbol{\theta}, \boldsymbol{\mu}(\boldsymbol{\sigma}) \rangle$ . Therefore, for a fixed model  $\boldsymbol{\theta}$ , the optimum is found by optimizing a linear function in the space of consistent labelings  $\Sigma_\pi$ .

The Alternating Mean Map algorithm (Algorithm 6) starts with the output of LMM and then optimizes it further over the set of consistent labelings. We give two possible implementations, corresponding to either a min-max or a min-min learning

**Algorithm 6:** Alternating Mean Map (AMM<sup>OPT</sup>)

---

**Input:** LMM params, optimization strategy  $\text{OPT} \in \{\min, \max\}$ , convergence predicate  $\text{PR}$   
 $\theta_0 \leftarrow \text{LMM}(\text{params})$   
**repeat** for  $t = 0, 1, \dots$   
 (a)  $\sigma_t \leftarrow \arg \text{OPT}_{\sigma \in \Sigma_\pi} \langle \theta, \mu(\sigma) \rangle$   
 (b)  $\theta_t \leftarrow \arg \min_{\theta} \frac{1}{2m} \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \ell(y \langle \theta, x_i \rangle) + a \langle \theta, \mu(\sigma) \rangle + \lambda \|\theta\|_2^2$   
**until** predicate  $\text{PR}$  is true  
**Output:**  $\hat{\theta}$  (see text)

---

strategy. At each iteration, in Step (a) we first pick a consistent labeling in  $\Sigma_\pi$  that is the best ( $\text{OPT} = \min$ ) or the worst ( $\text{OPT} = \max$ ) for the current classifier, and then in Step (b) we fit the classifier  $\hat{\theta}$  on the given set of labels. We iterate until a convergence predicate is met, which tests whether the empirical risk difference with the previous iteration is small enough (AMM<sup>min</sup>), or the number of iterations exceeds a user-specified limit (AMM<sup>max</sup>). The returned classifier  $\hat{\theta}$  is the one with the smallest empirical risk; in the case of AMM<sup>min</sup>, the output is the last learned model since the risk cannot increase.

Step (b) is a convex minimization with no technical difficulty. The label inference Step (a) is combinatorial as we optimize a linear function for integer variables. Yet, it can be solved in time almost linear in  $m$ , exploiting a trick similar to one elaborated in Yu et al. [2013]. The insight is that the search for labels consistent with each label proportion can be done via sorting. Moreover, we are optimizing a linear objective that decomposes by bag. Therefore, bag by bag, we can sort observations by the value of the inner product with the current model and label them as positive until we match the label proportion for the bag. The next Lemma is formally proven in 4.11.7.

**Lemma 42.** *The running time of Step (a) in AMM is  $O(m \log m)$ .*

## 4.8 Generalization bounds

We now study generalization bounds for LLP. We are not aware of much prior work on the topic. The only known generalization theory for LLP is in Yu et al. [2014b] that formulate uniform convergence bounds but with a focus on coarser grained problems — the estimation of bag label proportions – not directly the classifiers.

We obtain a first bound by applying Theorem 23 and utilizing Theorem 36 for evaluating how generalization degrades when we estimate the mean operator with LMM. The next Theorem holds for linear classifiers.

**Theorem 43.** *Let  $\varepsilon = 0$ . Hold the same assumptions and notation from Theorems 23 and 36. Let  $\ell$  be  $a$ -LOL. Let  $\text{Harm}(m) \doteq \frac{n}{\sum_{j \in [m]} 1/m_j}$  be the harmonic mean of the bag sizes. For*

any  $\delta > 0$ , with probability at least  $1 - \delta$ :

$$\begin{aligned} R_{\mathcal{D},\ell}(\hat{\theta}) - R_{\mathcal{D},\ell}(\theta^*) &\leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} + c(X, H) \cdot \sqrt{\frac{1}{m} \log\left(\frac{2}{\delta}\right)} \\ &+ 2|a|H \left( X \sqrt{\frac{d}{m} \log\left(\frac{2d}{\delta}\right)} + X^{1/4} \frac{(\max_{j \neq j'} v_{jj'})^2}{\min_j w_j^2} n^{3/2} \sqrt{\frac{1}{\text{Harm}(m)}} \right). \end{aligned} \quad (4.29)$$

Proof in 4.11.8. The Theorem highlights the effect of the mean operator as estimated by LMM, which trades off number of bags and (the harmonic mean of) their sizes. In particular, if we assume all bags with same size, we have  $\forall j, m_j = m/n$  and therefore:

$$n^{3/2} \sqrt{\frac{1}{n} \sum_{j \in [n]} \frac{n}{m}} = n^{3/2} \sqrt{\frac{n}{m}} = \frac{n^2}{\sqrt{m}}. \quad (4.30)$$

Once again, the rate of convergence is  $O(\sqrt{m})$ . The bound can be improved as we did above for Theorem 36 once we assume (A4.2) or (A4.3). The condition  $\varepsilon = 0$  is not a requirement but it simplifies the bound shape. The result confirms that novel bounds can be derived by simply plugging in any known guarantee on the mean operator estimators. We will show the same convenience in Chapter 5 analyzing the case of noisy labels.

We formulate an additional, more generic bound in the following. This result fits well with AMM and various algorithms that work by minimizing the *bag empirical  $\ell$ -risk*  $\mathbb{E}_{\sigma \sim \Sigma_\pi} \mathbb{E}_{\mathcal{S}}[\ell(\sigma h(x))]$ . This quantity is often part of the objective of LLP algorithms inspired by EM.  $\text{AMM}^{\min}$  and  $\text{AMM}^{\max}$  respectively minimize a lower bound and an upper bound of this risk. Here the model space  $\mathcal{H}$  is generic and not restricted to linear classifiers. The bound relies on a generalization of Rademacher complexity that better suits the LLP setting.

**Definition 44.** *The bag empirical Rademacher complexity of a hypothesis space  $\mathcal{H}$  with regard to sample  $\mathcal{S}$  of size  $m$  and loss  $\ell$  is:*

$$\mathcal{R}^b(\ell \circ \mathcal{H} \circ \mathcal{S}) \doteq \mathbb{E}_{\sigma \sim \Sigma_m} \left[ \sup_{h \in \mathcal{H}} \mathbb{E}_{\sigma' \sim \Sigma_\pi} \left[ \frac{1}{m} \sum_{i=1}^m \sigma' \ell(\sigma' h(x_i)) \right] \right]. \quad (4.31)$$

and the respective bag Rademacher complexity is  $\mathcal{R}^b(\ell \circ \mathcal{H}) \doteq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \mathcal{R}^b(\ell \circ \mathcal{H} \circ \mathcal{S})$ .

Notice that the usual empirical Rademacher complexity (Definition 4) equals  $\mathcal{R}^b$  when  $|\Sigma_\pi| = 1$ . It is also useful to define another type of complexity for the LLP, related with the variation of the label proportions per bag. Let  $\mathcal{I}_1^{l/2}$  and  $\mathcal{I}_2^{l/2}$  be two  $m$ -size *i.i.d.* samples of  $[2m]$  and  $\mathcal{S}(\mathcal{I}_1^{l/2})$  and  $\mathcal{S}(\mathcal{I}_2^{l/2})$  be the size- $m$  subsets of  $\mathcal{S}$  corresponds to those indices. Take  $l = 1, 2$  and any  $x_i \in \mathcal{S}$ . If  $i \notin \mathcal{I}_l^{l/2}$  then  $\pi_{|l}^s(x_i) = \pi_{|l}^l(x_i)$  is  $x_i$ 's bag's label proportion measured on  $\mathcal{S} \setminus \mathcal{S}(\mathcal{I}_l^{l/2})$ . Else,  $\pi_{|2}^s(x_i)$  is its bag's label proportion measured on  $\mathcal{S}(\mathcal{I}_2^{l/2})$  and  $\pi_{|1}^l(x_i)$  is its label, *i.e.* a bag's label proportion that would contain only  $x_i$ . Finally,  $\sigma_1(x) \doteq 2 \cdot 1\{x \in \mathcal{S}(\mathcal{I}_1^{l/2})\} - 1 \in \Sigma_1$ .



**Definition 45.** *The label proportion complexity of  $\mathcal{H}$  is:*

$$L(\mathcal{H}) \doteq \mathbb{E}_{\mathcal{D}_{2m}} \mathbb{E}_{\mathcal{I}_1^2, \mathcal{I}_2^2} \sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \left[ \sigma_1(x_i) \left( \pi_{|2}^s(x_i) - \pi_{|1}^l(x_i) \right) h(x_i) \right] . \quad (4.32)$$

$L(\mathcal{H})$  tends to be smaller as classifiers in  $\mathcal{H}$  have small magnitude on bags whose label proportion is close to  $1/2$ . Despite similar shapes,  $\mathcal{R}^b$  and  $L$  behave differently: when bags are pure,  $\pi_j \in \{0, 1\}, \forall j, L = 0$ . When bags are most impure,  $\pi_j = 1/2, \forall j, \mathcal{R}^b = 0$ . As bags are impure, the bag empirical  $\ell$ -risk also tends to increase. We can formulate a generalization bound for the bag empirical  $\ell$ -risk involving a balance of the two complexity measures.

**Theorem 46.** *Suppose  $\exists H \geq 0, |h(x)| \leq H$ . Then, for any SPL  $\ell$  and any  $0 < \delta \leq 1$ , with probability  $> 1 - \delta$ , the following bound holds over all  $h \in \mathcal{H}$ :*

$$\begin{aligned} \mathcal{R}_{\mathcal{D}, \ell}(h) - \mathbb{E}_{\sigma \sim \Sigma_\pi} \mathbb{E}_{\mathcal{S}}[\ell(\sigma h(x))] &\leq 2\mathcal{R}^b(\ell \circ \mathcal{H} \circ \mathcal{S}) \\ &+ L(\mathcal{H}) + 4 \left( \frac{2H}{b_\phi} + 1 \right) \sqrt{\frac{1}{2m} \log \frac{2}{\delta}} . \end{aligned} \quad (4.33)$$

Furthermore, under PSC, we have for any SPL  $\ell$ :

$$\mathcal{R}^b(\ell \circ \mathcal{H} \circ \mathcal{S}) \leq 2b_\phi \mathbb{E}_{\Sigma_m} \sup_{h \in \mathcal{H}} \{ \mathbb{E}_{\mathcal{S}}[\sigma(x)(\pi(x) - 1/2)h(x)] \} . \quad (4.34)$$

Proof in 4.11.9. This bound is the LLP equivalent to the first statement of Theorem 5, limiting the difference between  $\ell$ -risk and its empirical counterpart, which is here averaged in the set of labels consistent with the given proportions.

## 4.9 Experiments

### 4.9.1 Algorithms

We compare LMM, AMM ( $\ell$  = logistic loss) to the original MM [Quadrianto et al., 2009], InvCal [Rüping, 2010], conv- $\alpha$ SVM and alter- $\alpha$ SVM [Yu et al., 2013] with linear kernels. To obtain strong baselines, we test several additional initializations for AMM: the edge mean map estimator ( $AMM_{EMM}$ ):

$$\hat{\mu} = \frac{1}{m^2} \left( \sum_i y_i \right) \left( \sum_i x_i \right) , \quad (4.35)$$

the constant estimator ( $AMM_1$ ):

$$\hat{\mu} = \mathbf{1} , \quad (4.36)$$

and finally  $AMM_{10\text{ran}}$  which runs 10 random initial models ( $\|\theta_0\|_2 \leq 1$ ) and selects the one with smallest risk. This procedure is akin to alter- $\alpha$ SVM.

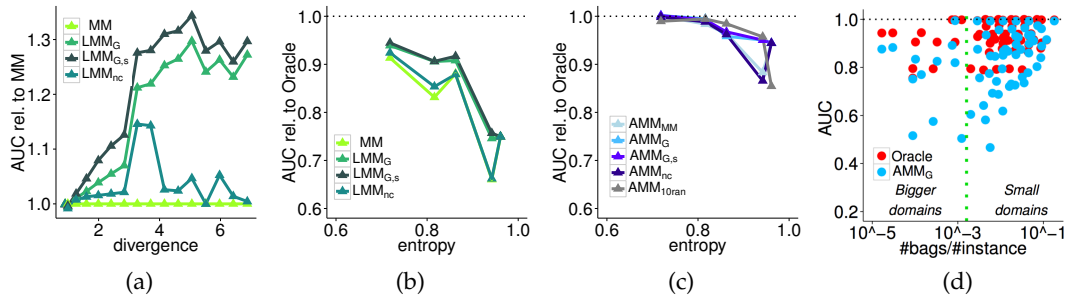


Figure 4.1: Relative AUC (w.r.t. MM) as homogeneity assumption is violated (a). Relative AUC (w.r.t. Oracle) vs entropy on *heart* for LMM(b),  $\text{AMM}^{\min}$ (c). Relative AUC vs  $n/m$  for  $\text{AMM}_{G,s}^{\min}$  (d).

Matrix  $V$  (Equations 4.21, 4.22) used is indicated in subscript:  $\text{LMM}/\text{AMM}_G$ ,  $\text{LMM}/\text{AMM}_{G,s}$ ,  $\text{LMM}/\text{AMM}_{nc}$  respectively denote  $v^{G,s}$  with  $s = 1$ ,  $v^{G,s}$  with  $s$  learned on cross validation and  $v^{nc}$ . We split the algorithms in two groups, *one-shot* and *iterative*. The latter, including AMM and (conv/alter)- $\alpha$ SVM, iteratively optimize a cost over labelings. The former (LMM, InvCal) do not and are thus much faster.

The range of hyper-parameters for cross validation are  $\lambda = \lambda' m$  with  $\lambda' \in \{0\} \cup 10^{\{0,1,2\}}$ ,  $\gamma \in 10^{-\{2,1,0\}}$ ,  $\sigma \in 2^{-\{2,1,0\}}$  for mean operator algorithms. We run all experiments with  $D_w = I$  and  $\varepsilon = 0$ . For  $\alpha$ SVM and InvCal, we used an implementation provided by the authors of [Yu et al., 2013]. Since we test on similar domains (6 are actually the same), hyper-parameters ranges for InvCal and  $\alpha$ SVM are taken from [Yu et al., 2013]. To avoid an additional source of complexity in the analysis, we cross-validate all hyper-parameters using the knowledge of all labels of the validation sets — information that generally would not be accessible in real world applications. For all our experiments, the testing metric is the AUC.

## 4.9.2 Simulated domains

We generate 16 domains that gradually move away the  $b_j^\sigma$  away from each other, thus increasingly violating the homogeneity assumption (A4.0). The degree of violation is measured as  $\|B^\pm - \bar{B}^\pm\|_F$ , where  $\bar{B}^\pm$  is the “homogeneity assumption matrix”, that replaces all  $b_j^\sigma$  by  $b^\sigma$  for  $\sigma \in \{-1, 1\}$ , see Equation 4.8. Figure 4.1 (a) displays the AUC ratios LMM with respect to MM. It shows that LMM is all the better with respect to MM as the homogeneity assumption is violated. Furthermore, learning the width parameter  $s$  in LMM improves the results.

The MM algorithm was shown to learn a model with zero accuracy on a toy domain crafted in Yu et al. [2013]. We reproduce the experiment and test all our methods. In Table 4.2 we report performance of all measured in transductive setting, *i.e.* on the same training set. Although none of the distances used in our experiments in LMM leads reasonable accuracy in the toy dataset (not reported in the Table),  $\text{AMM}^{\max}$  initialized with *any* starting point learns *in one step* a model which perfectly

	AMM <sup>min</sup>	AMM <sup>max</sup>
EMM	100.00	100.00
MM	8.46	100.00
LMM <sub>G</sub>	8.46	100.00
LMM <sub>G,s</sub>	8.46	100.00
LMM <sub>nc</sub>	8.46	100.00
1	8.46	100.00
10ran	100.00	100.00

Table 4.2: AUC on the toy dataset of Yu et al. [2013]

classifies all the instances. We also notice that EMM returns an optimal classifier by itself (not reported in Table 4.2).

### 4.9.3 UCI domains

We convert 10 small domains ( $m \leq 1000$ ) and 4 bigger ones ( $m > 8000$ ) from UCI [Bache and Lichman, 2013] into the LLP setting. We cast to one-against-all classification when the problem is multi-class. On large domains, the bag assignment function is inspired by [Yu et al., 2014b]: we craft bags according to a selected categorical feature, and then we remove that feature from the data. This conforms to the idea that bag assignment is structured and non random in real-world problems. Most of our small domains however do not have many features, so instead of clustering on one feature and then discarding it, we run  $k$ -MEANS on the whole data to make the bags, for  $k = n \in 2^{[5]}$ .

We perform 5-folds nested CV comparisons on the 10 domains = 50 AUC values for each algorithm. Full experimental results are given in Appendix 7.9, including runtime and details by domain. Table 4.3 synthesizes the results, splitting one-shot and iterative algorithms:

- (1) LMM<sub>G,s</sub> outperforms all one-shot algorithms.
- (2) LMM<sub>G</sub> and LMM<sub>G,s</sub> are competitive with many iterative algorithms, but lose against their AMM counterpart, which proves that additional optimization over labels is beneficial.
- (3) AMM<sub>G</sub> and AMM<sub>G,s</sub> are confirmed as the best variant of AMM, the first being the best in this case.
- (4) Surprisingly, all mean map algorithms, even one-shots, are clearly superior to  $\alpha$ SVMs. Further results reveal that  $\alpha$ SVM performances are dampened by learning classifiers with the “inverted polarity” — *i.e.* flipping the sign of the classifier improves its performances. See Appendix 4.12.2.

(5) Figure 4.1 (b, c) presents the AUC relative to the Oracle — which learns the classifier knowing all labels and minimizing the logistic loss —, as a function of the Gini entropy of bag assignment,  $gini(\mathcal{S}) \doteq 4\mathbb{E}_j[\pi_j(1 - \pi_j)]$ . For an entropy close to 1, we were expecting a drop in performances. The unexpected is that on some domains, large entropies ( $\geq .8$ ) do not prevent AMM<sup>min</sup> to compete with the Oracle. No such

pattern clearly emerges for  $\alpha$ SVM and  $\text{AMM}^{\max}$ .

We now consider the 4 bigger datasets. We adopt a 1/5 hold-out method. Scalability results display that every method using  $v^{nc}$  and  $\alpha$ SVM is not scalable to big domains. Table 4.4 presents the results on the big domains, distinguishing the feature used for bag assignment.  $\text{AMM}_{10\text{ran}}$  does not appear because of clearly inferior performance. Big domains confirm the efficacy of LMM+AMM. No approach clearly outperforms the rest, although  $\text{LMM}_{G,s}$  is often the best one-shot and  $\text{AMM}_G^{\min}$  and  $\text{AMM}_{G,s}^{\min}$  outperform all the methods several times.

Figure 4.1 (d) gives the AUC of  $\text{AMM}_G^{\min}$  over the Oracle for *all* domains, as a function of the “degree of supervision”,  $n/m$  ( $=1$  if the problem is fully supervised). Noticeably, on 90% of the runs,  $\text{AMM}_G^{\min}$  gets an AUC representing at least 70% of the Oracle’s. Results on big domains can be remarkable: on the census domain with bag assignment on *race*, 5 proportions are sufficient for an AUC 5 points below the Oracle’s — which learns with 200K labels.

## 4.10 Discussion

We have shown methods that can learn with label proportions successfully. By sufficiency we resort to standard learning procedures for binary classification. This is implemented as the LMM algorithm which estimates the mean operator via a Laplacian-based manifold regularizer relaxing the independence assumption of Quadrianto et al. [2009]. We show that under a weak distinguish-ability assumption between bags, our estimation of the mean operator is all the better as the maximal observation norm increase. This, as we show, cannot hold for the MM algorithm of Quadrianto et al. [2009]. Generalization bound are easily derived specializing the theory of Chapter 3.

We have also provided an iterative algorithm, AMM, that takes as input the solution of LMM and optimizes it further over the set of consistent labelings. We ground the algorithm in a uniform convergence result involving a generalization of Rademacher complexities for the LLP setting. The bound involves a bag surrogate risk for which we show that AMM optimizes tractable bounds.

Experiments display results that are superior to the state of the art, with algorithms that scale to big domains at affordable computational costs. Performances sometimes compete with the Oracle’s — that learns knowing all labels —, even on big domains. On one side, such experimental finding are encouraging for Machine Learning applications where the given supervision is only coarse-grain. On the other side, we uncover severe implications on the reliability of privacy preserving aggregation techniques with simple group statistics like proportions, still common for the public release of summaries on sensitive data. Future work shall study the extent to which LLP methods may threaten anonymity of people and sensitivity of their attributes when polls-like results are released on the Internet and when additional individual information on those individuals is available.



algorithm	mushroom: 8124 × 108			adult: 48842 × 89			marketing: 45211 × 41			census: 299285 × 381		
	I(6)	II(7)	III(10)	IV(5)	V(16)	VI(42)	V(4)	VII(4)	VIII(12)	IV(5)	VIII(9)	VI(42)
EMM	55.61	59.80	76.68	43.91	47.50	66.61	<b>63.49</b>	<b>54.50</b>	44.31	56.05	56.25	57.87
MM	51.99	<b>98.79</b>	5.02	80.93	76.65	74.01	54.64	50.71	49.70	75.21	<b>90.37</b>	75.52
LMMG	73.92	98.57	14.70	81.79	78.40	78.78	54.66	51.00	51.93	75.80	71.75	<b>76.31</b>
LMMG <sub>s</sub>	<b>94.91</b>	98.24	<b>89.43</b>	<b>84.89</b>	<b>78.94</b>	<b>80.12</b>	49.27	51.00	<b>65.81</b>	<b>84.88</b>	60.71	69.74
AMM <sub>EMM</sub>	85.12	99.45	69.43	49.97	56.98	70.19	61.39	55.73	43.10	87.86	87.71	40.80
AMM <sub>MM</sub>	89.81	99.01	15.74	<b>83.73</b>	77.39	80.67	52.85	<b>75.27</b>	58.19	89.68	84.91	68.36
AMM <sub>G</sub>	89.18	99.45	50.44	83.41	<b>82.55</b>	<b>81.96</b>	51.61	75.16	57.52	87.61	88.28	76.99
AMM <sub>G<sub>s</sub></sub>	89.24	<b>99.57</b>	3.28	81.18	78.53	<b>81.96</b>	52.03	75.16	53.98	<b>89.93</b>	83.54	52.13
AMM <sub>I</sub>	<b>95.90</b>	98.49	97.31	81.32	75.80	80.05	65.13	64.96	66.62	89.09	<b>88.94</b>	56.72
AMM <sub>EMM</sub>	93.04	3.32	26.67	54.46	69.63	56.62	51.48	55.63	57.48	71.20	77.14	66.71
AMM <sub>MM</sub>	59.45	55.16	<b>99.70</b>	82.57	71.63	81.39	48.46	51.34	56.90	50.75	66.76	58.67
AMM <sub>G</sub>	95.50	65.32	99.30	82.75	72.16	81.39	50.58	47.27	34.29	48.32	67.54	<b>77.46</b>
AMM <sub>G<sub>s</sub></sub>	95.84	65.32	84.26	82.69	70.95	81.39	<b>66.88</b>	47.27	34.29	80.33	74.45	52.70
AMM <sub>I</sub>	95.01	73.48	1.29	75.22	67.52	77.67	66.70	61.16	<b>71.94</b>	57.97	81.07	53.42
Oracle	99.82	99.81	99.8	90.55	90.55	90.50	79.52	75.55	79.43	94.31	94.37	94.45

Table 4.4: AUC on big domains (*name*: #instances×#features). I=*cap-shape*, II=*habitat*, III=*cap-color*, IV=*race*, V=*education*, VI=*country*, VII=*outcome*, VIII=*job* (number of bags); for each feature, the best result over one-shot, and over iterative algorithms is bold faced.

## 4.11 Appendix: proofs

### 4.11.1 Proof of Lemma 35

Using the fact that  $D_w$  and  $L$  are symmetric, we have:

$$\frac{\partial \ell(L, X)}{\partial X} \quad (4.37)$$

$$= -2 \frac{\partial}{\partial X} \text{tr} B^\top D_w \Pi^\top X + \frac{\partial}{\partial X} \text{tr} X^\top \Pi D_w \Pi^\top X + \gamma \frac{\partial}{\partial X} \text{tr} X^\top L X \quad (4.38)$$

$$= -2 \Pi D_w B + 2 \Pi D_w \Pi^\top X + 2 \gamma L X = 0, \quad (4.39)$$

out of which  $\hat{B}^\pm$  follows.

### 4.11.2 Proof of Theorem 36

We let  $\Pi_o \doteq [\text{diag}(\boldsymbol{\pi}) | \text{diag}(\boldsymbol{\pi} - \mathbf{1})]^\top N$  an orthonormal system ( $n_{jj} = (\pi_j^2 + (1 - \pi_j)^2)^{-1/2}, \forall j \in [n]$  and 0 otherwise). Let  $\mathbb{K}_{\Pi_o}$  be the  $n$ -dim subspace of  $\mathbb{R}^d$  generated by  $\Pi_o$ . The proof of Theorem 36 exploits the following Lemma, which assumes that  $\varepsilon$  is any  $> 0$  real for  $L$  in 4.16 to be  $\succ 0$ . When  $\varepsilon = 0$ , the result of Theorem 36 still holds but follows a different proof.

**Lemma 47.** *Let  $A \doteq \Pi D_w \Pi^\top$  and  $L$  defined as in 4.16. Denote for short:*

$$U \doteq \left( L^{-1} A + \gamma^{-1} I \right)^{-1}. \quad (4.40)$$

Suppose there exists  $\xi > 0$  such that for any  $\mathbf{x} \in \mathbb{R}^{2n}$ , the projection of  $U\mathbf{x}$  in  $\mathbb{K}_{\Pi_o}$ ,  $\mathbf{x}_{U,o}$ , satisfies:

$$\|\mathbf{x}_{U,o}\|_2 \leq \xi \|\mathbf{x}\|_2. \quad (4.41)$$

Then:

$$\|M - \hat{M}\|_F \leq \gamma \xi \|B^\pm\|_F. \quad (4.42)$$

**Proof** Combining Lemma 35 and Equation 4.8, we get

$$B^\pm - \hat{B}^\pm = - \left( (A + \gamma L)^{-1} A - I \right) B^\pm \quad (4.43)$$

$$= \left( (\gamma L)^{-1} A + I \right)^{-1} B^\pm. \quad (4.44)$$

Define the following permutation matrix:

$$C \doteq \left[ \begin{array}{c|c} 0 & I \\ \hline I & 0 \end{array} \right] \in \mathbb{R}^{2n \times 2n}. \quad (4.45)$$

$A \doteq \Pi D_w \Pi^\top$  is not invertible but diagonalizable. Its (orthonormal) eigenvectors can

be partitioned in two matrices  $P_o$  and  $P$  such that:

$$P_o \doteq [\text{diag}(\boldsymbol{\pi} - \mathbf{1}) \mid \text{diag}(\boldsymbol{\pi})]^\top N = C\Pi_o \in \mathbb{R}^{2n \times n} \text{ (eigenvalues 0) } , \quad (4.46)$$

$$P \doteq \Pi N \in \mathbb{R}^{2n \times n} \text{ (eigenvalues } w_j(\pi_j^2 + (1 - \pi_j)^2), \forall j) . \quad (4.47)$$

We have:

$$M - \hat{M} = P_o^\top C B^\pm - P_o^\top C \hat{B}^\pm \quad (4.48)$$

$$= P_o^\top C \left( (\gamma L)^{-1} A + I \right)^{-1} B^\pm \quad (4.49)$$

$$= \Pi_o^\top \left( (\gamma L)^{-1} A + I \right)^{-1} B^\pm \quad (4.50)$$

$$= \gamma \Pi_o^\top \left( L^{-1} A + \gamma^{-1} I \right)^{-1} B^\pm . \quad (4.51)$$

Equation 4.50 follows from the fact that  $C$  is idempotent. Plugging Frobenius norm in 4.51, we obtain:

$$\|M - \hat{M}\|_F^2 = \gamma^2 \left\| \Pi_o^\top \left( L^{-1} A + \gamma^{-1} I \right)^{-1} B^\pm \right\|_F^2 \quad (4.52)$$

$$= \gamma^2 \sum_{k=1}^d \left\| \Pi_o^\top \left( L^{-1} A + \gamma^{-1} I \right)^{-1} \mathbf{b}_k^\pm \right\|_2^2 \quad (4.53)$$

$$\leq \gamma^2 \zeta^2 \sum_{k=1}^d \|\mathbf{b}_k^\pm\|_2^2 \quad (4.54)$$

$$= \gamma^2 \zeta^2 \|B^\pm\|_F^2 , \quad (4.55)$$

which yields 4.42. In 4.54,  $\mathbf{b}_k^\pm$  denotes *column*  $k$  in  $B^\pm$ . Inequality 4.54 makes use of assumption 4.41. ■

To ensure  $\|\mathbf{x}_{U,o}\|_2 \leq \zeta \|\mathbf{x}\|_2$ , it is sufficient that  $\|U\mathbf{x}\|_2 \leq \zeta \|\mathbf{x}\|_2$ , and since  $\|U\mathbf{x}\|_2 \leq \|U\|_F \|\mathbf{x}\|_2$ , it is sufficient to show that:

$$\left\| U_\zeta^{-1} \right\|_F^2 \leq 1 , \quad (4.56)$$

with  $U_\zeta \doteq L_\zeta^{-1} A + \zeta \gamma^{-1} I$ , for relevant choices of  $\zeta$ . We have let  $L_\zeta \doteq (1/\zeta)L$ . Let  $0 \leq \lambda_1(\cdot) \leq \dots \leq \lambda_{2n}(\cdot)$  denote the ordered eigenvalues of a positive-semidefinite matrix in  $\mathbb{R}^{2n \times 2n}$ . It follows that, since  $L$  is symmetric positive definite, we have:

$$\lambda_j(L_\zeta^{-1} A) \geq \frac{\lambda_j(A)}{\lambda_{2n}(L_\zeta)} (\geq 0) , \forall j \in [2n] . \quad (4.57)$$

We have used Equation 4.46. Weyl's Theorem [Horn and Johnson, 2012, Chapter 4]



then brings:

$$\lambda_j(U_{\xi}^{-1}) \leq \frac{\lambda_{2n}(L_{\xi})}{\lambda_j(A) + \xi\gamma^{-1}\lambda_{2n}(L_{\xi})} \leq \begin{cases} \xi^{-1}\gamma & \text{if } j \in [n] \\ \frac{\lambda_{2n}(L_{\xi})}{\lambda_j(A)} & \text{otherwise} \end{cases}. \quad (4.58)$$

Gershgorin's Theorem [Horn and Johnson, 2012, Chapter 6] brings  $\lambda_{2n} \leq (1/\xi)(\varepsilon + \max_j \sum_{j'} |l_{jj'}|)$ , and furthermore the eigenvalues of  $A$  satisfy  $\lambda_j \geq w_j/2, \forall j \geq n+1$ . We thus have:

$$\|U_{\xi}^{-1}\|_F^2 \leq \frac{n\gamma^2}{\xi^2} + \frac{4n \left( \varepsilon + \max_j \sum_{j'} |l_{jj'}| \right)^2}{\xi^2 \min_j w_j^2}. \quad (4.59)$$

In 4.58 and 4.59, we have used the eigenvalues of  $A$  given in Equations 4.46 and 4.47. Assuming:

$$\gamma \leq \frac{\xi}{\sqrt{2n}}, \quad (4.60)$$

a sufficient condition for the right-hand side of 4.59 to be  $\leq 1$  is that:

$$\xi \geq \frac{\varepsilon + \max_j \sum_{j'} |l_{jj'}|}{2\sqrt{n} \min_j w_j}. \quad (4.61)$$

To finish up the proof, recall that  $L = D - V$  with  $d_{jj} \doteq \sum_{j'} v_{jj'}$  and the coordinates  $v_{jj'} \geq 0$ . Hence,

$$\sum_{j'} |l_{jj'}| = 2 \sum_{j \neq j'} v_{jj'} \quad (4.62)$$

$$\leq 2n \max_{j \neq j'} v_{jj'}, \forall j \in [n]. \quad (4.63)$$

The proof is concluded by plugging this upper bound in 4.61 to choose  $\xi$ , then taking the maximal value for  $\gamma$  in 4.60 and finally solving the upper bound in 4.42.

### 4.11.3 Proof of Lemma 37

We first consider the normalized association criterion in 4.21:

$$v_{jj'}^N \doteq \frac{1}{2} \left( \frac{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j)}{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j \cup \mathcal{S}_{j'})} + \frac{\text{ASSOC}(\mathcal{S}_{j'}, \mathcal{S}_{j'})}{\text{ASSOC}(\mathcal{S}_{j'}, \mathcal{S}_j \cup \mathcal{S}_{j'})} \right), \quad (4.64)$$

$$\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_{j'}) \doteq \sum_{x \in \mathcal{S}_j, x' \in \mathcal{S}_{j'}} \|x - x'\|_2^2. \quad (4.65)$$

Remark that:

$$\|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2 = \left\| \frac{1}{m_j} \sum_{x_i \in \mathcal{S}_j} \mathbf{x}_i - \frac{1}{m_{j'}} \sum_{x_{i'} \in \mathcal{S}_{j'}} \mathbf{x}_{i'} \right\|_2^2 \quad (4.66)$$

$$\begin{aligned} &= \frac{1}{m_j^2} \left\| \sum_{x_i \in \mathcal{S}_j} \mathbf{x}_i \right\|_2^2 + \frac{1}{m_{j'}^2} \left\| \sum_{x_{i'} \in \mathcal{S}_{j'}} \mathbf{x}_{i'} \right\|_2^2 - \frac{2}{m_j m_{j'}} \left( \sum_{x_i \in \mathcal{S}_j} \mathbf{x}_i \right)^\top \left( \sum_{x_{i'} \in \mathcal{S}_{j'}} \mathbf{x}_{i'} \right) \\ &= \frac{1}{m_j^2} \left\| \sum_{x_i \in \mathcal{S}_j} \mathbf{x}_i \right\|_2^2 + \frac{1}{m_{j'}^2} \left\| \sum_{x_{i'} \in \mathcal{S}_{j'}} \mathbf{x}_{i'} \right\|_2^2 - \frac{2}{m_j m_{j'}} \sum_{x_i \in \mathcal{S}_j, x_{i'} \in \mathcal{S}_{j'}} \mathbf{x}_i^\top \mathbf{x}_{i'} \quad (4.67) \end{aligned}$$

$$\leq \frac{1}{m_j} \sum_{x_i \in \mathcal{S}_j} \|\mathbf{x}_i\|_2^2 + \frac{1}{m_{j'}} \sum_{x_{i'} \in \mathcal{S}_{j'}} \|\mathbf{x}_{i'}\|_2^2 - \frac{2}{m_j m_{j'}} \sum_{x_i \in \mathcal{S}_j, x_{i'} \in \mathcal{S}_{j'}} \mathbf{x}_i^\top \mathbf{x}_{i'} \quad (4.68)$$

$$= \frac{1}{m_j m_{j'}} \sum_{x_i \in \mathcal{S}_j, x_{i'} \in \mathcal{S}_{j'}} \|\mathbf{x}_i - \mathbf{x}_{i'}\|_2^2 \quad (4.69)$$

$$\begin{aligned} &+ \underbrace{\frac{m_{j'} - 1}{m_j m_{j'}} \sum_{x_i \in \mathcal{S}_j} \|\mathbf{x}_i\|_2^2 + \frac{m_j - 1}{m_j m_{j'}} \sum_{x_{i'} \in \mathcal{S}_{j'}} \|\mathbf{x}_{i'}\|_2^2 - \frac{1}{m_j m_{j'}} \sum_{x_i \in \mathcal{S}_j, x_{i'} \in \mathcal{S}_{j'}} \mathbf{x}_i^\top \mathbf{x}_{i'}}_{\doteq a} \\ &\leq \frac{2}{m_j m_{j'}} \sum_{x_i \in \mathcal{S}_j, x_{i'} \in \mathcal{S}_{j'}} \|\mathbf{x}_i - \mathbf{x}_{i'}\|_2^2 \quad (4.70) \end{aligned}$$

$$= \frac{2}{m_j m_{j'}} \text{ASSOC}(\mathcal{S}_j, \mathcal{S}_{j'}) . \quad (4.71)$$

Equation 4.68 exploits the fact that  $\left(\sum_{j=1}^n a_j\right)^2 \leq n \left(\sum_{j=1}^n a_j^2\right)$  and Equation 4.70 is due to  $a \leq (m_j m_{j'})^{-1} \sum_{x_i \in \mathcal{S}_j, x_{i'} \in \mathcal{S}_{j'}} \|\mathbf{x}_i - \mathbf{x}_{i'}\|_2^2$ . We thus have:

$$\frac{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j)}{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j \cup \mathcal{S}_{j'})} = \frac{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j)}{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j) + \text{ASSOC}(\mathcal{S}_j, \mathcal{S}_{j'})} \quad (4.72)$$

$$\leq \frac{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j)}{\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j) + \frac{m_j m_{j'}}{2} \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} \quad (4.73)$$

$$\leq \frac{\kappa' m_j}{\kappa' m_j + \frac{m_j m_{j'}}{2} \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} \quad (4.74)$$

$$= \frac{1}{1 + \frac{m_{j'}}{2\kappa'} \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} . \quad (4.75)$$

Equation 4.73 uses 4.71 and 4.74 uses assumption **(A4.3)**. Equation 4.74 also holds when permuting  $j$  and  $j'$ , so we get:

$$\begin{aligned} \zeta(V^{\text{NC}}, B^\pm) &\leq \max_{j \neq j'} \left( \frac{\varepsilon}{2n} + \frac{1}{1 + \frac{m_j}{2\kappa'} \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} + \frac{1}{1 + \frac{m_{j'}}{2\kappa'} \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} \right)^2 \|B^\pm\|_F \\ &\leq \left( \frac{\varepsilon}{2n} + \frac{1}{1 + \frac{\min_j m_j}{2\kappa'} \min_{j,j'} \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} \right)^2 \|B^\pm\|_F \end{aligned} \quad (4.76)$$

$$\leq \left( \frac{\varepsilon^2}{2n^2} + 2 \left( \frac{1}{1 + \frac{\min_j m_j}{2\kappa'} \min_{j,j'} \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} \right)^2 \right) \|B^\pm\|_F \quad (4.77)$$

$$\leq \frac{\varepsilon^2}{2n^2} d \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2 + \frac{4\kappa' d \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2}{\min_{j,j'}^2 \|\mathbf{b}_j - \mathbf{b}_{j'}\|_2^2} \quad (4.78)$$

$$\leq \frac{\varepsilon^2}{2n^2} d \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2 + \frac{4\kappa' d}{\kappa^2 \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2} \quad (4.79)$$

$$= f^{\text{NC}} \left( \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2 \right) \quad (4.80)$$

$$= o(1) , \quad (4.81)$$

where the last inequality uses assumption **(A4.2)**, and 4.77 uses the property that  $(a + b)^2 \leq 2a^2 + 2b^2$ . We have let:

$$f^{\text{NC}}(x) \doteq \frac{\varepsilon^2}{2n^2} dx + \frac{4\kappa' d}{\kappa x} , \quad (4.82)$$

which is indeed  $o(1)$  if  $\varepsilon = o(n^2/\sqrt{x})$ . This proves the Lemma for  $\zeta(V^{\text{NC}}, B^\pm)$ . The case of  $\zeta(V^{G,s}, B^\pm)$  is easier, as:

$$\exp \left( -\frac{\|\mathbf{b}_j - \mathbf{b}_{j'}\|_2}{s} \right) \leq \exp \left( -\frac{\min_{j'',j'''} \|\mathbf{b}_{j''} - \mathbf{b}_{j'''}\|_2}{s} \right) \quad (4.83)$$

$$\leq \exp \left( -\frac{\kappa}{s} \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2 \right) , \quad (4.84)$$

from assumption (A4.2) alone, which gives:

$$\zeta(V^{G,s}, B^\pm) \leq \|B^\pm\|_F \left( \frac{\varepsilon}{2n} + \exp\left(-\frac{\kappa}{s} \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2\right) \right)^2 \quad (4.85)$$

$$\leq \|B^\pm\|_F \left( \frac{\varepsilon^2}{2n^2} + 2 \exp\left(-\frac{2\kappa}{s} \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2\right) \right) \quad (4.86)$$

$$\leq d \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2 \left( \frac{\varepsilon^2}{2n^2} + 2 \exp\left(-\frac{2\kappa}{s} \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2\right) \right) \quad (4.87)$$

$$= f^G \left( \max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2 \right) \quad (4.88)$$

$$= o(1) , \quad (4.89)$$

as claimed. We have let  $f^G(x) \doteq \frac{\varepsilon^2}{2n^2} dx + dx \exp(-2\kappa x/s)$ , which is indeed  $o(1)$  if  $\varepsilon = o(n^2/\sqrt{x})$ .

For the last statement, remark that we shall have in general  $f^G(x) \leq f^{NC}(x)$  and even  $f^G(x) = o(f^{NC}(x))$  if  $\varepsilon = 0$ , so we may expect better convergence in the case of  $V^{G,s}$  as  $\max_{\sigma,j} \|\mathbf{b}_j^\sigma\|_2$  grows.

#### 4.11.4 Proof of Lemma 38

We first restate the Lemma in a more explicit way, that shall provide explicit values for  $\kappa_l$  and  $\kappa_n$ .

**Lemma 48.** *There exist  $\kappa_{jj'}$  and  $s_{jj'}$  depending on  $d_j, d_{j'}$ , and  $\kappa'_{jj'} > 1$  depending on  $m_j, m_{j'}$ , such that:*

- If  $v_{jj'}^{G,s_{jj'}} > \exp(-1/4)$  then  $\mathcal{S}_j, \mathcal{S}_{j'}$  are not linearly separable;
- If  $v_{jj'}^{G,s_{jj'}} < \exp(-64)$  then  $\mathcal{S}_j, \mathcal{S}_{j'}$  are linearly separable;
- If  $v_{jj'}^{NC} > \kappa_{jj'}$  then  $\mathcal{S}_j, \mathcal{S}_{j'}$  are not linearly separable;
- If  $v_{jj'}^{NC} < \kappa_{jj'} / \kappa'_{jj'}$  then  $\mathcal{S}_j, \mathcal{S}_{j'}$  are linearly separable.

**Proof** We first consider the normalized association criterion in 4.21, and we prove the Lemma for the following expressions of  $\kappa_{jj'}$  and  $\kappa'_{jj'}$ :

$$\kappa_{jj'} \doteq \frac{16}{2 + \frac{d_{j'}^2}{2d_j^2}} + \frac{16}{2 + \frac{d_j^2}{2d_{j'}^2}} , \quad (4.90)$$

$$\kappa'_{jj'} \doteq 512 \max\{m_j, m_{j'}\} , \quad (4.91)$$

with  $d_{jj'} \doteq \max\{d_j, d_{j'}\}$  and  $d_j \doteq \max_{\mathbf{x}, \mathbf{x}' \in \mathcal{S}_j} \|\mathbf{x} - \mathbf{x}'\|_2, \forall j \neq j' \in [n]$ . For any bag  $\mathcal{S}_j$ , we let  $(\mathbf{b}_j^*, r_j) \doteq \text{MEB}(\mathcal{S}_j)$  denote the minimum enclosing ball (MEB) for bag  $\mathcal{S}_j$  and

distance  $L_2$ , that is,  $r_j$  is the smallest unique real such that:

$$\exists! \mathbf{b}_j^* : d(\mathbf{x}, \mathbf{b}_j^*) \doteq \|\mathbf{x} - \mathbf{b}_j^*\|_2 \leq r_j, \forall \mathbf{x} \in \mathcal{S}_j . \quad (4.92)$$

We have let  $d(\mathbf{x}, \mathbf{b}_j^*) \doteq \|\mathbf{x} - \mathbf{b}_j^*\|_2$ . We are going to prove a first result involving the MEBs of  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$ , and then will translate the result to the Lemma's statement. The following properties follows from standard properties of MEBs and the fact that  $d(.,.)$  is a distance (they hold for any  $j \neq j'$ ):

- (a)  $d(\mathbf{x}, \mathbf{x}') \leq 2r_j, \forall \mathbf{x}, \mathbf{x}' \in \mathcal{S}_j$ ;
- (b) If bags  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  are linearly separable, then  $\forall \mathbf{x} \in \text{co}(\mathcal{S}_j), \exists \mathbf{x}' \in \mathcal{S}_{j'}$  such that  $d(\mathbf{x}, \mathbf{x}') \geq \max\{r_j, r_{j'}\}$ ; here, "co" denotes the convex closure;
- (c) If bags  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  are linearly separable, then  $d(\mathbf{b}_j, \mathbf{b}_{j'}) \geq \max\{r_j, r_{j'}\}$ , where  $\mathbf{b}_j$  and  $\mathbf{b}_{j'}$  are the bags average;
- (d)  $\forall \mathbf{x} \in \mathcal{S}_j, \exists \mathbf{x}' \in \mathcal{S}_j$  s.t.  $d(\mathbf{x}, \mathbf{x}') \geq r_j$ ;
- (e)  $d(\mathbf{x}, \mathbf{x}') \leq 2 \max\{r_j, r_{j'}\} + d(\mathbf{b}_j^*, \mathbf{b}_{j'}^*), \forall \mathbf{x} \in \text{co}(\mathcal{S}_j), \forall \mathbf{x}' \in \text{co}(\mathcal{S}_{j'})$ .

Let us define:

$$\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_{j'}) \doteq \sum_{\mathbf{x} \in \mathcal{S}_j, \mathbf{x}' \in \mathcal{S}_{j'}} d^2(\mathbf{x}, \mathbf{x}') . \quad (4.93)$$

We remark that, assuming that each bag contains at least two elements without loss of generality:

$$v_{jj'}^{\text{NC}} = \frac{1}{2} \left( \frac{1}{1 + \frac{\text{ASSOC}(\mathcal{B}_j, \mathcal{B}_{j'})}{\text{ASSOC}(\mathcal{B}_j, \mathcal{B}_j)}} + \frac{1}{1 + \frac{\text{ASSOC}(\mathcal{B}_j, \mathcal{B}_{j'})}{\text{ASSOC}(\mathcal{B}_{j'}, \mathcal{B}_{j'})}} \right) . \quad (4.94)$$

We have  $\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j) \leq 4m_j r_j^2$  and  $\text{ASSOC}(\mathcal{S}_{j'}, \mathcal{S}_{j'}) \leq 4m_{j'} r_{j'}^2$  (because of (a)), and also  $\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_{j'}) \geq \max\{m_j, m_{j'}\} \max\{r_j^2, r_{j'}^2\}$  when  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  are linearly separable (because of (b)), which yields in this case:

$$v_{jj'}^{\text{NC}} \leq \frac{1}{2 + \frac{\max\{m_j, m_{j'}\} \max\{r_j^2, r_{j'}^2\}}{2m_j r_j^2}} + \frac{1}{2 + \frac{\max\{m_j, m_{j'}\} \max\{r_j^2, r_{j'}^2\}}{2m_{j'} r_{j'}^2}} \quad (4.95)$$

$$\leq \frac{1}{2 + \frac{\max\{r_j^2, r_{j'}^2\}}{2r_j^2}} + \frac{1}{2 + \frac{\max\{r_j^2, r_{j'}^2\}}{2r_{j'}^2}} . \quad (4.96)$$

Let us name  $\kappa_{jj'}^\circ$  the right-hand side of 4.96. It follows that when  $v_{jj'}^{\text{NC}} > \kappa_{jj'}^\circ$ ,  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  are not linearly separable.

On the other hand, we have  $\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_j) \geq m_j r_j^2$  and  $\text{ASSOC}(\mathcal{S}_{j'}, \mathcal{S}_{j'}) \geq m_{j'} r_{j'}^2$  (because of (d)), and also:

$$\text{ASSOC}(\mathcal{S}_j, \mathcal{S}_{j'}) \leq m_j m_{j'} (2 \max\{r_j, r_{j'}\} + d(\mathbf{b}_j^*, \mathbf{b}_{j'}^*))^2 \quad (4.97)$$

$$\leq m_j m_{j'} (4 \max\{r_j^2, r_{j'}^2\} + 2d^2(\mathbf{b}_j^*, \mathbf{b}_{j'}^*)) , \quad (4.98)$$

because of (e) and the fact that  $(a + b)^2 \leq 2a^2 + 2b^2$ . It follows that  $\forall j \neq j'$ :

$$v_{jj'}^{\text{NC}} \geq \frac{1}{2 + \frac{2m_{j'}(4 \max\{r_j^2, r_{j'}^2\} + 2d^2(\mathbf{b}_j^*, \mathbf{b}_{j'}^*))}{r_j^2}} + \frac{1}{2 + \frac{2m_j(4 \max\{r_j^2, r_{j'}^2\} + 2d^2(\mathbf{b}_j^*, \mathbf{b}_{j'}^*))}{r_{j'}^2}} . \quad (4.99)$$

For any  $j \neq j'$ , when  $d^2(\mathbf{b}_j^*, \mathbf{b}_{j'}^*) \leq 4 \max\{r_j^2, r_{j'}^2\}$ , then we have from 4.99:

$$v_{jj'}^{\text{NC}} \geq \frac{1}{2 + \frac{16m_{j'} \max\{r_j^2, r_{j'}^2\}}{r_j^2}} + \frac{1}{2 + \frac{16m_j \max\{r_j^2, r_{j'}^2\}}{r_{j'}^2}} \quad (4.100)$$

$$> \kappa_{jj'}^\circ / (32 \max\{m_j, m_{j'}\}) . \quad (4.101)$$

Hence,  $v_{jj'}^{\text{NC}} \leq \kappa_{jj'}^\circ / (32 \max\{m_j, m_{j'}\})$  implies  $d(\mathbf{b}_j^*, \mathbf{b}_{j'}^*) > 2 \max\{r_j, r_{j'}\}$ , implying  $d(\mathbf{b}_j^*, \mathbf{b}_{j'}^*) > r_j + r_{j'}$ , which is a sufficient condition for the linear separability of  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$ .

We can relate the linear separability of  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  to the value of  $v_{jj'}^{\text{NC}}$  with respect to  $\kappa_{jj'}^\circ$  defined in 4.96. To remove the dependence in the MEB parameters and obtain the statement of the Lemma, we just have to remark that  $d_j^2/4 \leq r_j^2 \leq 4d_j^2, \forall j \in [n]$ , which yields  $\kappa_{jj'}/16 \leq \kappa_{jj'}^\circ \leq \kappa_{jj'}$ . Therefore, when  $v_{jj'}^{\text{NC}} > \kappa_{jj'}$ , it follows that  $v_{jj'}^{\text{NC}} > \kappa_{jj'}^\circ$  and  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  are not linearly separable. On the other hand, when  $v_{jj'}^{\text{NC}} \leq \kappa_{jj'}/(16 \times 32 \max\{m_j, m_{j'}\}) = \kappa_{jj'}/\kappa_{jj'}'$ , then  $v_{jj'}^{\text{NC}} \leq \kappa_{jj'}^\circ / (32 \max\{m_j, m_{j'}\})$  and the bags  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  are linearly separable. This achieves the proof of Lemma 38 for the normalized association criterion in 4.21.

The proof for  $v_{jj'}^{\text{G},s}$  is shorter, and we prove it for:

$$s_{j,j'} = \max\{d_j, d_{j'}\} . \quad (4.102)$$

We have  $(1/2) \max\{d_j, d_{j'}\} \leq \max\{r_j, r_{j'}\} \leq 2 \max\{d_j, d_{j'}\}$ . Hence, because of (c) above, if  $\mathcal{S}_j$  and  $\mathcal{S}_{j'}$  are linearly separable, then  $v_{jj'}^{\text{G},s} \leq 1/e^{1/4}$ ; so, when  $v_{jj'}^{\text{G},s} > 1/e^{1/4}$ , the two bags are not linearly separable. On the other hand, if  $d(\mathbf{b}_j^*, \mathbf{b}_{j'}^*) \leq 2 \max\{r_j, r_{j'}\}$ , then because of (e) above  $d(\mathbf{b}_j, \mathbf{b}_{j'}) \leq 4 \max\{r_j, r_{j'}\} \leq 8 \max\{d_j, d_{j'}\}$ , and so  $v_{jj'}^{\text{G},s} \geq 1/e^{64}$ . This implies that if  $v_{jj'}^{\text{G},s} < 1/e^{64}$ , then  $d(\mathbf{b}_j^*, \mathbf{b}_{j'}^*) > 2 \max\{r_j, r_{j'}\} \geq r_j + r_{j'}$ , and thus the two bags are linearly separable, as claimed. This achieves the proof of Lemma 48.  $\blacksquare$

This achieves the proof of Lemma 38.

#### 4.11.5 Proof of lemma 39

Let  $x > 0, \epsilon \in (0, 1), p \in (0, 1), p \neq 1/2$ . We create a dataset from four observations,  $\{(x_1 = 0, 1), (x_2 = 0, -1), (x_3 = x, 1), (x_4 = x, -1)\}$ . There are two bags,  $\mathcal{S}_1$  takes  $1 - \epsilon$  of  $x_2$  and  $\epsilon$  of  $x_1$ .  $\mathcal{S}_2$  takes  $\epsilon$  of  $x_4$  and  $1 - \epsilon$  of  $x_3$ . The label-wise estimators  $\hat{\mathbf{b}}^\sigma$  of Quadrianto et al. [2009] are solution of:

$$\begin{bmatrix} \hat{\mathbf{b}}^+ \\ \hat{\mathbf{b}}^- \end{bmatrix} = \left( \begin{bmatrix} 1 - \epsilon & \epsilon \\ \epsilon & 1 - \epsilon \end{bmatrix}^\top \begin{bmatrix} 1 - \epsilon & \epsilon \\ \epsilon & 1 - \epsilon \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 - \epsilon & \epsilon \\ \epsilon & 1 - \epsilon \end{bmatrix}^\top \begin{bmatrix} x \\ 0 \end{bmatrix} \quad (4.103)$$

$$= \frac{1}{1 - 2\epsilon} \begin{bmatrix} (1 - \epsilon)x \\ \epsilon x \end{bmatrix} \quad (4.104)$$

On the other hand, the true quantities are:

$$\begin{bmatrix} \mathbf{b}^+ \\ \mathbf{b}^- \end{bmatrix} = \begin{bmatrix} (1 - \epsilon)x \\ \epsilon x \end{bmatrix}. \quad (4.105)$$

We now mix classes in  $\mathcal{S}$  and pick bag proportions  $q \doteq \mathbb{P}_{\mathcal{S}}[x \in \mathcal{S}_1]$  and  $1 - q = \mathbb{P}_{\mathcal{S}}[x \in \mathcal{S}_2]$ . We have the class proportions defined by  $\mathbb{P}_{\mathcal{S}}[y = +1] = \epsilon q + (1 - \epsilon)(1 - q) \doteq p$ . Then:

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2 = \left\| p(1 - \epsilon) \left( \frac{1}{1 - 2\epsilon} - 1 \right) x - (1 - p)\epsilon \left( \frac{1}{1 - 2\epsilon} - 1 \right) x \right\|_2 \quad (4.106)$$

$$= \frac{2\epsilon|p - \epsilon|}{1 - 2\epsilon} x \quad (4.107)$$

$$= 2\epsilon(1 - q)x. \quad (4.108)$$

Furthermore,  $\max_i \|\mathbf{b}_i^\sigma\|_2 = x$ . We get:

$$\frac{\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2}{\max_i \|\mathbf{b}_i^\sigma\|_2} = 2\epsilon(1 - q). \quad (4.109)$$

Picking  $\epsilon$  and  $(1 - q)$  both  $> \sqrt{1 - (\gamma/2)}$  is sufficient to have Equation 4.109  $> 2 - \gamma$  for any  $\gamma > 0$ . Remark that both assumptions (A4.2) and (A4.3) hold for any  $\kappa < 1$  and any  $\kappa' > 0$ .

#### 4.11.6 Proof of Theorem 41

Let us start by recalling some notation.

$$R_{S_{X,\ell}}(\boldsymbol{\theta}, \boldsymbol{\mu}) = \frac{1}{2m} \left( \sum_{i \in [m]} \sum_{\sigma} \ell(\sigma \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle) \right) + a \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle . \quad (4.110)$$

Also define the regularized loss:

$$R_{S_{X,\ell}}(\boldsymbol{\theta}, \boldsymbol{\mu}, \lambda) \doteq R_{S_{X,\ell}}(\boldsymbol{\theta}, \boldsymbol{\mu}) + \lambda \|\boldsymbol{\theta}\|_2^2 . \quad (4.111)$$

The proof builds upon Lemma 26, which bound is refined for SPLs. First, we prove a helper Lemma for Lipschitz losses.

**Lemma 49.** *Let  $\ell$  be  $L$ -Lipschitz and  $a$ -LOL. Fix  $\lambda > 0$ , and let  $X \doteq \max_i \|\mathbf{x}_i\|_2$ . Let  $\boldsymbol{\theta}' \doteq \operatorname{argmin}_{\boldsymbol{\theta}} R_{S_{X,\ell}}(\boldsymbol{\theta}, \boldsymbol{\mu}', \lambda)$  where  $\boldsymbol{\mu}'$  is any vector in  $\mathbb{R}^d$ . Then:*

$$\|\boldsymbol{\theta}'\|_2 \leq \frac{LX + |a| \|\boldsymbol{\mu}'\|_2}{\lambda} . \quad (4.112)$$

**Proof** Let us define a shrinking of the optimal solution  $\boldsymbol{\theta}'$ ,  $\boldsymbol{\theta}_\alpha \doteq \alpha \boldsymbol{\theta}'$  for  $\alpha \in (0, 1)$ . We have:

$$R_{S_{X,\ell}}(\boldsymbol{\theta}_\alpha, \boldsymbol{\mu}', \lambda) = \frac{1}{2m} \left( \sum_i \sum_{\sigma} \ell(\sigma \langle \boldsymbol{\theta}_\alpha, \mathbf{x}_i \rangle) \right) + a \langle \boldsymbol{\theta}_\alpha, \boldsymbol{\mu}' \rangle + \lambda \|\boldsymbol{\theta}_\alpha\|_2^2 \quad (4.113)$$

$$= \frac{1}{2m} \left( \sum_i \sum_{\sigma} \ell(\sigma \alpha \langle \boldsymbol{\theta}', \mathbf{x}_i \rangle) \right) + a \alpha \langle \boldsymbol{\theta}', \boldsymbol{\mu}' \rangle + \lambda \alpha^2 \|\boldsymbol{\theta}'\|_2^2 \quad (4.114)$$

$$\leq \frac{1}{2m} \left( \sum_i \sum_{\sigma} \ell(\sigma \langle \boldsymbol{\theta}', \mathbf{x}_i \rangle) + L |\sigma \alpha \langle \boldsymbol{\theta}', \mathbf{x}_i \rangle - \sigma \langle \boldsymbol{\theta}', \mathbf{x}_i \rangle| \right) + a \alpha \langle \boldsymbol{\theta}', \boldsymbol{\mu}' \rangle + \lambda \alpha^2 \|\boldsymbol{\theta}'\|_2^2 \quad (4.115)$$

$$= \frac{1}{2m} \left( \sum_i \sum_{\sigma} \ell(\sigma \langle \boldsymbol{\theta}', \mathbf{x}_i \rangle) \right) + \frac{L(1-\alpha)}{m} \sum_i |\langle \boldsymbol{\theta}', \mathbf{x}_i \rangle| + a \alpha \langle \boldsymbol{\theta}', \boldsymbol{\mu}' \rangle + \lambda \alpha^2 \|\boldsymbol{\theta}'\|_2^2 , \quad (4.116)$$

where 4.115 holds because  $\ell$  is  $L$ -Lipschitz. To have Equation 4.116 smaller than  $R_{S_{X,\ell}}(\boldsymbol{\theta}', \boldsymbol{\mu}', \lambda)$ , we need equivalently:

$$\frac{L(1-\alpha)}{m} \sum_i |\langle \boldsymbol{\theta}', \mathbf{x}_i \rangle| + a \alpha \langle \boldsymbol{\theta}', \boldsymbol{\mu}' \rangle + \lambda \alpha^2 \|\boldsymbol{\theta}'\|_2^2 \leq a \langle \boldsymbol{\theta}', \boldsymbol{\mu}' \rangle + \lambda \|\boldsymbol{\theta}'\|_2^2 , \quad (4.117)$$

that is:

$$\frac{L(1-\alpha)}{m} \sum_i |\langle \boldsymbol{\theta}', \mathbf{x}_i \rangle| - a(1-\alpha) \langle \boldsymbol{\theta}', \boldsymbol{\mu}' \rangle \leq \lambda(1-\alpha^2) \|\boldsymbol{\theta}'\|_2^2 , \quad (4.118)$$



and to find an  $\alpha \in (0, 1)$  such that this holds, because of Cauchy-Schwartz inequality, it is sufficient that:

$$(1 - \alpha)(LX - a\boldsymbol{\mu}') \leq (1 - \alpha)(LX + |a|\|\boldsymbol{\mu}'\|_2) \leq \lambda(1 - \alpha^2)\|\boldsymbol{\theta}'\|_2, \quad (4.119)$$

that is

$$\|\boldsymbol{\theta}'\|_2 \geq \frac{LX + |a|\|\boldsymbol{\mu}'\|_2}{\lambda(1 + \alpha)}. \quad (4.120)$$

Finally, whenever  $\|\boldsymbol{\theta}'\|_2 > (LX + |a|\|\boldsymbol{\mu}'\|_2)/\lambda$ , there is a shrinking of the optimal solution of Equation 4.110 that further decreases the empirical risk, thus contradicting its optimality. This ends the proof of Lemma 49.  $\blacksquare$

We now combine Theorem 26 and Lemma 49. Recall that we assume SPLs. Since they are differentiable by definition, we can verify that they are also Lipschitz by taking the first derivative. We have:

$$\ell'(x) = -\frac{1}{b_\phi}(\phi^*)'(-x) = -\frac{1}{b_\phi}(\phi')^{-1}(-x) \in [-1/b_\phi, 0], \quad (4.121)$$

for any  $x \in \phi'([0, 1])$  and thus  $\ell$  is  $1/b_\phi$ -Lipschitz. Also, we know that  $a = -\frac{1}{2b_\phi}$  with  $b_\phi > 0$  for SPLs (Table 3.1). Therefore, Lemma 49 tells us that:

$$\|\boldsymbol{\theta}'\|_2 \leq \frac{X + \|\boldsymbol{\mu}'\|_2}{b_\phi \lambda}. \quad (4.122)$$

We can obtain an explicit value instead of the generic  $\gamma$  denoting strong convexity in Theorem 26. Let us compute the second derivate for any SPL. By the rule of derivate of an inverse function, we differentiate Equation 4.121:

$$\ell''(x) = -\frac{1}{b_\phi} \left( (\phi')^{-1} \right)'(-x) \quad (4.123)$$

$$= -\frac{1}{b_\phi \phi''((\phi')^{-1}(-x))} \quad (4.124)$$

Recall Equation 3.100 in the proof of Theorem 26. We compute a lower bound of the second derivative  $\ell''$  by bounding its argument for any  $\alpha \in [0, 1]$ :

$$|\pm \langle \alpha \boldsymbol{\theta}^* + (1 - \alpha) \hat{\boldsymbol{\theta}}, \mathbf{x}_i \rangle| \leq (\alpha \|\boldsymbol{\theta}^*\|_2 + (1 - \alpha) \|\hat{\boldsymbol{\theta}}\|_2) X \quad (4.125)$$

$$\leq \frac{X + \alpha \|\boldsymbol{\mu}\|_2 + (1 - \alpha) \|\hat{\boldsymbol{\mu}}\|_2}{b_\phi \lambda} \quad (4.126)$$

$$\leq \frac{X + \max\{\|\boldsymbol{\mu}\|_2, \|\hat{\boldsymbol{\mu}}\|_2\}}{b_\phi \lambda}, \quad (4.127)$$

where Inequality 4.125 follows from Cauchy-Schwartz inequality. 4.126 uses Equ-

tion 4.122.  $\boldsymbol{\mu}$  and  $\hat{\boldsymbol{\mu}}$  are the parameters of  $R_{S_{X,\ell}}(\cdot, \boldsymbol{\mu}, \lambda)$  and  $R_{S_{X,\ell}}(\cdot, \hat{\boldsymbol{\mu}}, \lambda)$  in Lemma 26, and  $\hat{\boldsymbol{\theta}}$  and  $\boldsymbol{\theta}^*$  the respective minimizers. This achieves the proof of Theorem 41 for the interval:

$$[\pm(X + \max\{\|\boldsymbol{\mu}\|_2, \|\hat{\boldsymbol{\mu}}\|_2\})] . \quad (4.128)$$

#### 4.11.7 Proof of Lemma 42

We elaborate the proof for optimization strategy  $\text{OPT} = \min$ . The case  $\text{OPT} = \max$  flips the choice of the label in Step 1. For simplicity, let us define, as in the last proof:

$$R_{S_{X,\ell}}(\boldsymbol{\theta}, \boldsymbol{\mu}) = \frac{1}{2m} \left( \sum_{i \in [m]} \sum_{\sigma} \ell(\sigma \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle) \right) + a \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle . \quad (4.129)$$

To minimize  $R_{S_{X,\ell}}(\boldsymbol{\theta}_t, \boldsymbol{\mu}(\sigma))$  over  $\sigma \in \Sigma_{\pi}$ , we just have to find:

$$\sigma_* \in \operatorname{argmax}_{\sigma \in \Sigma_{\pi}} \left\langle \boldsymbol{\theta}, \sum_i \sigma_i \mathbf{x}_i \right\rangle , \quad (4.130)$$

and we can do that bag by bag. Algorithm 7 presents the labeling (notation  $(m) \doteq \{1, 2, \dots, m-1\}$ ). Remark that the time complexity for one bag is  $O(m_j \log m_j)$  due to the ordering (Step 1), so the overall complexity is indeed  $O(m \max_i \log m_i)$ .

**Lemma 50.** *Let  $\sigma_* \doteq \{\sigma_1^*, \sigma_2^*, \dots, \sigma_m^*\}$  be the set of labels obtained after running  $\text{LA}(\boldsymbol{\theta}, S_j, m_j^+)$  for  $j = 1, 2, \dots, n$ . Then  $\sigma_* \in \operatorname{argmax}_{\sigma \in \Sigma_{\pi}} \boldsymbol{\theta}^{\top} \sum_i \sigma_i \mathbf{x}_i$ .*

**Proof** The total edge,  $\langle \boldsymbol{\theta}, \sum_i \sigma_i \mathbf{x}_i \rangle$  (for any  $\sigma \in \Sigma_{\pi}$ ), can be summable bag-wise with regard to the coordinates of  $\sigma$ . Consider the optimal set:

$$\{\sigma^*\}_{\mathcal{B}} \doteq \operatorname{argmax}_{\sigma \in \{-1,1\}^{m'}: \langle \mathbf{1}, \sigma \rangle = 2m^+ - m'} \left\langle \boldsymbol{\theta}, \sum_{\mathbf{x}_i \in \mathcal{B}} \sigma_i \mathbf{x}_i \right\rangle ,$$

for some bag  $\mathcal{B} = \{\mathbf{x}_i, i = 1, 2, \dots, m'\}$ , with constraint  $m^+ \in [m']$ . This set contains the label assignment  $\sigma_*$  returned by Algorithm 7  $\text{LA}(\boldsymbol{\theta}, \mathcal{B}, m^+)$ , a property that follows from two simple observations:

1. Consider any observation  $\mathbf{x}_i$  of bag  $\mathcal{B}$ ; for any optimal labeling  $\sigma^*$  of  $\mathcal{B}$ , let  $m'^+ \doteq m^+ - \mathbb{I}(\sigma_i^* = 1)$ . Define the set  $\{\sigma'^*\}_i$  of optimal labelings of  $\mathcal{B} \setminus \{\mathbf{x}_i\}$  with constraint  $m'^+ \doteq m^+ - \mathbb{I}(\sigma_i^* = 1)$ . Then this set coincides with the set created by taking the elements of  $\{\sigma^*\}_{\mathcal{B}}$  to which we drop coordinate  $i$ . This follows from the per-observation summability of the total edge with regard to labels.
2. Assume  $m^+ \in (m')$ .  $\forall i^* \in \operatorname{argmax}_i |\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle|$ , there exists an optimal assignment  $\sigma^*$  such that  $\sigma_{i^*}^* = \operatorname{sign} \langle \boldsymbol{\theta}, \mathbf{x}_{i^*} \rangle$ . Otherwise, starting from any optimal assignment  $\sigma^*$ , we can flip the label of  $\mathbf{x}_{i^*}$  and the label of any other  $\mathbf{x}_i$  for which

**Algorithm 7:** Label Assignment (LA)

---

**Input:**  $\theta \in \mathbb{R}^d$ , a bag  $\mathcal{B} = \{x_i \in \mathbb{R}^d, i = 1, 2, \dots, m\}$ , bag size  $m^+ \in [m]$   
**If**  $\mathcal{B} = \emptyset$  **then** stop  
**Else if**  $m^+ \notin (m)$   
 $y_i \leftarrow \mathbb{I}(m^+ = m) - \mathbb{I}(m^+ = 0), \forall i = 1, 2, \dots, m$   
**Else**  
1  $i^* \leftarrow \operatorname{argmax}_i |\langle \theta, x_i \rangle|$   
2  $y_{i^*} \leftarrow \operatorname{sign} \langle \theta, x_{i^*} \rangle$   
3  $\text{LA}(\theta, \mathcal{B} \setminus \{x_{i^*}\}, m^+ - \mathbb{I}(y_{i^*} = 1))$

---

$\sigma_i^* \neq \sigma_{i^*}^*$ , and get a label assignment that satisfies constraint  $m^+$  and cannot be worse than  $\sigma^*$ , and is thus optimal, a contradiction.

Hence,  $\text{LA}(\theta, \mathcal{B}, m^+)$  picks at each iteration a label that matches one in a subset of optimal labelings, and the recursive call preserves the subset of optimal labelings. Since when  $m^+ \notin (m)$  the solution returned by  $\text{LA}(\theta, \mathcal{B}, m^+)$  is obviously optimal, we end up when the current  $\mathcal{B}$  is empty with  $\sigma_* \in \operatorname{argmax}_{\sigma \in \Sigma_\pi} \langle \theta, \sum_i \sigma_i x_i \rangle$ , as claimed. ■

#### 4.11.8 Proof of Theorem 43

We have to combine Theorems 23 (first statement) and 36. The link between the two is the following bound on the norm discrepancy between the population mean operator and its estimator by LMM:

$$\|\mu_{\mathcal{D}} - \hat{\mu}_{\mathcal{S}}\|_2 = \|\mu_{\mathcal{D}} + \mu_{\mathcal{S}} - \mu_{\mathcal{S}} - \hat{\mu}_{\mathcal{S}}\|_2 \quad (4.131)$$

$$\leq \|\mu_{\mathcal{D}} - \mu_{\mathcal{S}}\|_2 + \|\mu_{\mathcal{S}} - \hat{\mu}_{\mathcal{S}}\|_2 . \quad (4.132)$$

Here we are interested in bounding the latter norm in Step 4.132, since we can treat the former one as done in the proof of Theorem 23. Let us denote  $\eta \doteq \mu_{\mathcal{S}} - \hat{\mu}_{\mathcal{S}} \in \mathbb{R}^d$  and similarly  $\eta_j, \eta_j^y \in \mathbb{R}^d$  the difference of bag-wise and bag-wise, label-wise mean operators respectively. By Jensen's inequality and the fact that  $p(j) \leq 1$  and  $\pi_j \leq$

1,  $\forall j$  we have:

$$\|\eta\|_2^2 = \left\| \sum_{j \in [n]} p(j) \eta_j \right\|_2^2 \quad (4.133)$$

$$= \left\| \sum_{j \in [n]} p(j) \left( \pi_j \eta_j^+ - (1 - \pi_j) \eta_j^- \right) \right\|_2^2 \quad (4.134)$$

$$\leq \sum_{j \in [n]} p(j) \left\| \pi_j \eta_j^+ - (1 - \pi_j) \eta_j^- \right\|_2^2 \quad (4.135)$$

$$\leq \sum_{j \in [n]} \left\| \pi_j \eta_j^+ - (1 - \pi_j) \eta_j^- \right\|_2^2 \quad (4.136)$$

$$\leq \sum_{j \in [n]} \sum_{y \in \mathcal{Y}} \left\| y \cdot \eta_j^y \right\|_2^2 \quad (4.137)$$

$$= \sum_{j \in [n]} \sum_{y \in \mathcal{Y}} \sum_{k \in [d]} \left| \left( \eta_j^y \right)^k \right|^2 \quad (4.138)$$

$$= \|M - \hat{M}\|_F^2 \quad (4.139)$$

where  $M, \hat{M}$  are defined in Theorem 36. Therefore, by Theorem 36 and taking  $\varepsilon = 0$ :

$$\|\boldsymbol{\mu}_{\mathcal{D}} - \hat{\boldsymbol{\mu}}_{\mathcal{S}}\|_2 \leq \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2 + \|M - \hat{M}\|_F \quad (4.140)$$

$$\leq \|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2 + \sqrt{n/2} \cdot \frac{(\max_{j \neq j'} v_{jj'})^2}{\min_j w_j^2} \|B^\pm\|_F \quad (4.141)$$

We bound:

$$\|B^\pm\|_F = \sqrt{\sum_{j,y} \left\| \mathbf{b}_j^y \right\|_2^2} \quad (4.142)$$

$$= \sqrt{\sum_{j,y} \left\| \frac{1}{m_j} \sum_{i \in \mathcal{S}_j: y_i=y} \mathbf{x}_i \right\|_2^2} \quad (4.143)$$

$$\leq \sqrt{\sum_{j,y} \frac{1}{m_j^2} \cdot m_j \sqrt{X}} \quad (4.144)$$

$$= \sqrt{2n\sqrt{X} \cdot \sum_j \frac{1}{m_j}} \quad (4.145)$$

and thus Equation 4.141 is upper bounded by:

$$\|\boldsymbol{\mu}_{\mathcal{D}} - \boldsymbol{\mu}_{\mathcal{S}}\|_2 + nX^{1/4} \frac{(\max_{j \neq j'} v_{jj'})^2}{\min_j w_j^2} \cdot \sqrt{\sum_{j \in [n]} \frac{1}{m_j}} . \quad (4.146)$$

Therefore, under the assumptions of Theorem 23 it follows the excess risk is bounded with probability  $1 - \delta$  by:

$$\begin{aligned} R_{\mathcal{D}, \ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}, \ell}(\boldsymbol{\theta}^*) &\leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} + c(X, H) \cdot \sqrt{\frac{1}{m} \log \left( \frac{2}{\delta} \right)} \\ &+ 2|a|H \left( X \sqrt{\frac{d}{m} \log \left( \frac{2d}{\delta} \right)} + nX^{1/4} \frac{(\max_{j \neq j'} v_{jj'})^2}{\min_j w_j^2} \sqrt{\sum_{j \in [n]} \frac{1}{m_j}} \right) . \end{aligned} \quad (4.147)$$

To obtain the Theorem notice:

$$\sum_{j \in [n]} \frac{1}{m_j} = \frac{n}{\text{Harm}(m)} . \quad (4.148)$$

#### 4.11.9 Proof of Theorem 46

We prove separately Equations 4.33 and 4.34.

##### 4.11.9.1 Proof of Equation 4.33

Unless explicitly stated, all samples like  $\mathcal{S}$  and  $\mathcal{S}'$  are of size  $m$ . To make the reading of our expectations clear and simple, we shall write  $\mathbb{E}_{\Sigma_m}$  for  $\mathbb{E}_{\sigma \sim \Sigma_m}$ ,  $\mathbb{E}_{\mathcal{S}}$  for  $\mathbb{E}_{(x,y) \sim \mathcal{S}}$ ,  $\mathbb{E}_{\mathcal{D}'_m}$  for  $\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}}$  and  $\mathbb{E}_{\mathcal{D}_m}$  for  $\mathbb{E}_{\mathcal{S} \sim \mathcal{D}}$ . We now proceed to the proof, that follows the same main steps as that of Theorem 5 in Bartlett and Mendelson [2002]. For any  $q \in [0, 1]$ , let us define the convex combination:

$$\ell(q, h(\mathbf{x})) \doteq q\ell(h(\mathbf{x})) + (1 - q)\ell(-h(\mathbf{x})) . \quad (4.149)$$

It follows that:

$$\mathbb{E}_{\Sigma_\pi} \mathbb{E}_{\mathcal{S}}[\ell(\sigma(\mathbf{x})h(\mathbf{x}))] = \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] , \quad (4.150)$$

with  $\pi(\mathbf{x})$  the label proportion of the bag to which  $\mathbf{x}$  belongs in  $\mathcal{S}$ . We also have  $\forall h$ :

$$\mathbb{E}_{\mathcal{D}}[\ell(yh(\mathbf{x}))] \leq \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] + \phi(\mathcal{S}) , \quad (4.151)$$

with:

$$\phi(\mathcal{S}) \doteq \sup_g \{ \mathbb{E}_{\mathcal{D}}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} . \quad (4.152)$$

Let us bound the deviations of  $\phi(\mathcal{S})$  around its expectation on the sampling of  $\mathcal{S}$ , using McDiarmid's inequality. This departs from the well known result of Lemma 14 because we have extended the definition of  $\phi(\mathcal{S})$  to account for the missing labels by  $\mathbb{E}_{\Sigma_{\pi}}$ . We need to upper bound the maximum difference for the supremum term computed over two samples  $\mathcal{S}$  and  $\mathcal{S}'$  of the same size, such that  $\mathcal{S}'$  is  $\mathcal{S}$  with one example replaced. We have:

$$|\phi(\mathcal{S}) - \phi(\mathcal{S}')| \leq |\mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}'}[\ell(\pi'(\mathbf{x}), g(\mathbf{x}))]| , \quad (4.153)$$

with  $\pi$  and  $\pi'$  denoting the corresponding label proportions in  $\mathcal{S}$  and  $\mathcal{S}'$ . Let  $\{x_1\} = \mathcal{S} \setminus \mathcal{S}'$  and  $\{x_2\} = \mathcal{S}' \setminus \mathcal{S}$ . Let  $x_1 \in \mathcal{S}_j$  and  $x_2 \in \mathcal{S}'_{j'}$  for some bags  $j$  and  $j'$ . The bound 4.153 depends only on bags  $j$  and  $j'$ . For any  $\mathbf{x} \in (\mathcal{S}_j \cup \mathcal{S}'_{j'}) \setminus \{x_1, x_2\}$ , with SPLS, due to Lemma 34, we have:

$$\ell(\pi(\mathbf{x}), g(\mathbf{x})) - \ell(\pi'(\mathbf{x}), g(\mathbf{x})) \leq \frac{|\ell(g(\mathbf{x})) - \ell(-g(\mathbf{x}))|}{m(\mathbf{x})} \quad (4.154)$$

$$= \frac{|g(\mathbf{x})|}{b_{\phi} m(\mathbf{x})} \quad (4.155)$$

$$\leq \frac{H}{b_{\phi} m(\mathbf{x})} , \quad (4.156)$$

where  $m(\mathbf{x})$  is the size of the bag to which it belongs in  $\mathcal{S}$ , plus 1 iff it is bag  $j'$  and  $j' \neq j$ , minus 1 iff it is bag  $j$  and  $j' \neq j$ . Furthermore:

$$\begin{aligned} \ell(\pi(\mathbf{x}), g(\mathbf{x})) &= \ell(|g(\mathbf{x})|) + \frac{1}{b_{\phi}} \cdot ((1 - \pi(\mathbf{x})1\{g(\mathbf{x}) > 0\} \\ &\quad + \pi(\mathbf{x})(1 - 1\{g(\mathbf{x}) > 0\})) \cdot |g(\mathbf{x})| \end{aligned} \quad (4.157)$$

$$\leq \ell(0) + \frac{1}{b_{\phi}} ((1 - \pi(\mathbf{x}))1\{g(\mathbf{x}) > 0\} + \pi(\mathbf{x})(1 - 1\{g(\mathbf{x}) > 0\})) \cdot H \quad (4.158)$$

$$\leq \ell(0) + \frac{H}{b_{\phi}} , \forall \mathbf{x} \in \mathcal{S} . \quad (4.159)$$

Also, it comes from Lemma 34 that:

$$\ell(0) = \frac{1}{b_{\phi}} (0 \cdot \phi'^{-1}(0) - \phi(\phi'^{-1}(0))) \quad (4.160)$$

$$= \frac{-\phi(1/2)}{b_{\phi}} = 1 . \quad (4.161)$$

We obtain that:

$$|\phi(\mathcal{S}) - \phi(\mathcal{S}')| \leq \frac{1}{m} \left(1 + \frac{H}{b_\phi} + 1 + \frac{H}{b_\phi}\right) + \frac{1}{m} \sum_{x \in (\mathcal{S}_j \cup \mathcal{S}_{j'}) \setminus \{x_1, x_2\}} \frac{H}{b_\phi m(x)} \quad (4.162)$$

$$\leq \frac{Q_1}{m}, \quad (4.163)$$

where:

$$Q_1 \doteq 2 \left( \frac{2H}{b_\phi} + 1 \right). \quad (4.164)$$

So McDiarmid's inequality yields that with probability  $\leq \delta/2$  over the sampling of  $\mathcal{S}$ :

$$\phi(\mathcal{S}) \geq \mathbb{E}_{\mathcal{D}_m} \sup_g \{ \mathbb{E}_{\mathcal{D}}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} + Q_1 \sqrt{\frac{1}{2m} \log \frac{2}{\delta}} \quad (4.165)$$

We now upper bound the expectation in 4.165. Using the convexity of the supremum, we have:

$$\mathbb{E}_{\mathcal{D}_m} \sup_g \{ \mathbb{E}_{\mathcal{D}}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} \quad (4.166)$$

$$= \mathbb{E}_{\mathcal{D}_m} \sup_g \{ \mathbb{E}_{\mathcal{D}'_m}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} \quad (4.167)$$

$$\leq \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m} \sup_g \{ \mathbb{E}_{\mathcal{S}'_m}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \}. \quad (4.168)$$

Consider any set  $\mathcal{S} \sim \mathcal{D}_{2m}$ , and let  $\mathcal{I}^{/2} \subseteq [2m]$  be a subset of  $m$  indices, picked uniformly at random among all  $\binom{2m}{m}$  possible choices. For any  $\mathcal{I} \subseteq [2m]$ , let  $\mathcal{S}(\mathcal{I})$  denote the subset of examples whose index matches  $\mathcal{I}$ , and for any  $\mathbf{x} \in \mathcal{S}(\mathcal{I})$ , let  $\pi(\mathbf{x}|\mathcal{S}(\mathcal{I}))$  denote its bag proportion in  $\mathcal{S}(\mathcal{I})$ . For any  $\mathcal{I}_l^{/2}$  indexed by  $l \geq 1$  and any  $\mathbf{x} \in \mathcal{S}$ , let:

$$\pi_{|l}^s(\mathbf{x}) \doteq \begin{cases} \pi(\mathbf{x}|\mathcal{S}(\mathcal{I}_l^{/2})) & \text{if } \mathbf{x} \in \mathcal{S}(\mathcal{I}_l^{/2}) \\ \pi(\mathbf{x}|\mathcal{S} \setminus \mathcal{S}(\mathcal{I}_l^{/2})) & \text{otherwise} \end{cases} \quad (4.169)$$

denote the label proportions induced by the split of  $\mathcal{S}$  in two subsamples  $\mathcal{S}(\mathcal{I}_l^{/2})$  and  $\mathcal{S} \setminus \mathcal{S}(\mathcal{I}_l^{/2})$ . Let:

$$\pi_{|l}^\ell(\mathbf{x}) \doteq \begin{cases} y & \text{if } \mathbf{x} \in \mathcal{S}(\mathcal{I}_l^{/2}) \\ \pi(\mathbf{x}|\mathcal{S} \setminus \mathcal{S}(\mathcal{I}_l^{/2})) & \text{otherwise} \end{cases}, \quad (4.170)$$

where  $y$  is the true label of  $\mathbf{x}$ . Let  $\sigma_l(\mathbf{x}) \doteq 2 \times 1\{\mathbf{x} \in \mathcal{S}(\mathcal{I}_l^{/2})\} - 1$ . The Label Proportion Complexity  $L(\mathcal{H})$  quantifies the discrepancy between these two estimators. When each bag in  $\mathcal{S}$  has label proportion zero or one, each term factoring classifier  $h$  in Equation 4.32 is zero, so  $L(\mathcal{H}) = 0$ .

**Lemma 51.** *The following holds true:*

$$\mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m} \sup_g \{ \mathbb{E}_{\mathcal{S}'}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} \quad (4.171)$$

$$\leq 2\mathbb{E}_{\mathcal{D}_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))] \} + L(\mathcal{H}) . \quad (4.172)$$

**Proof** For any  $\sigma \in \Sigma_m$  and any sets  $\mathcal{S} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$  and  $\mathcal{S}' = \{\mathbf{x}'_1, \mathbf{x}'_2, \dots, \mathbf{x}'_m\}$  of size  $m$ , denote:

$$\mathcal{S}_\sigma \doteq \{\mathbf{x}'_i \text{ iff } \sigma_i = 1, \mathbf{x}_i \text{ otherwise}\} , \quad (4.173)$$

$$\mathcal{S}_{\bar{\sigma}} \doteq \{\mathbf{x}'_i \text{ iff } \sigma_i = -1, \mathbf{x}_i \text{ otherwise}\} = (\mathcal{S} \cup \mathcal{S}') \setminus \mathcal{S}_\sigma . \quad (4.174)$$

and:

$$\pi_*(\mathbf{x}) \doteq \begin{cases} \pi_\sigma(\mathbf{x}) & \text{if } \mathbf{x} \in \mathcal{S}_\sigma , \\ \pi_{\bar{\sigma}}(\mathbf{x}) & \text{otherwise} \end{cases} , \quad (4.175)$$

where:  $\pi_\sigma(\cdot)$  denote the label proportions in  $\mathcal{S}_\sigma$  and  $\pi_{\bar{\sigma}}(\cdot)$  denote the label proportions in  $\mathcal{S}_{\bar{\sigma}}$ . Let  $\pi(\cdot)$  denote the label proportions in  $\mathcal{S}$ ,  $\pi'(\cdot)$  denote the label proportions in  $\mathcal{S}'$  (we know each bag to which each example in  $\mathcal{S}'$  belongs to, so we can compute these estimators), We have:

$$\mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m} \sup_h \{ \mathbb{E}_{\mathcal{S}'}[\ell(yh(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] \} \quad (4.176)$$

$$= \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m} \sup_h \left\{ \mathbb{E}_{\mathcal{S}'}[\ell(\pi'(\mathbf{x}), h(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] - \frac{\Delta_1}{b_\phi} \right\} \quad (4.177)$$

$$= \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m} \sup_h \left\{ \mathbb{E}_{\mathcal{S}_\sigma}[\sigma(\mathbf{x})\ell(\pi^l(\mathbf{x}), h(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}_{\bar{\sigma}}}[\sigma(\mathbf{x})\ell(\pi^r(\mathbf{x}), h(\mathbf{x}))] - \frac{\Delta_1}{b_\phi} \right\} , \quad (4.178)$$

with:

$$\Delta_1 \doteq \mathbb{E}_{\mathcal{S}'}[((1 - \pi'(\mathbf{x}))1\{y = 1\} - \pi'(\mathbf{x})1\{y = -1\})h(\mathbf{x})] , \quad (4.179)$$

$$\pi^l(\mathbf{x}) \doteq \frac{1}{2} ((1 + \sigma(\mathbf{x}))\pi'(\mathbf{x}) + (1 - \sigma(\mathbf{x}))\pi(\mathbf{x})) , \quad (4.180)$$

$$\pi^r(\mathbf{x}) \doteq \frac{1}{2} ((1 + \sigma(\mathbf{x}))\pi(\mathbf{x}) + (1 - \sigma(\mathbf{x}))\pi'(\mathbf{x})) . \quad (4.181)$$

We also have from Lemma 34:

$$\mathbb{E}_{\mathcal{S}_\sigma}[\sigma(\mathbf{x})\ell(\pi^l(\mathbf{x}), h(\mathbf{x}))] = \mathbb{E}_{\mathcal{S}_\sigma}[\sigma(\mathbf{x})\ell(\pi_\sigma(\mathbf{x}), h(\mathbf{x}))] - \frac{\Delta_2}{b_\phi} , \quad (4.182)$$

$$\mathbb{E}_{\mathcal{S}_{\bar{\sigma}}}[\sigma(\mathbf{x})\ell(\pi^r(\mathbf{x}), h(\mathbf{x}))] = \mathbb{E}_{\mathcal{S}_{\bar{\sigma}}}[\sigma(\mathbf{x})\ell(\pi_{\bar{\sigma}}(\mathbf{x}), h(\mathbf{x}))] - \frac{\Delta_3}{b_\phi} , \quad (4.183)$$



with:

$$\Delta_2 \doteq \mathbb{E}_{\mathcal{S}_\sigma}[\sigma(\mathbf{x})(\pi^l(\mathbf{x}) - \pi_\sigma(\mathbf{x}))h(\mathbf{x})] , \quad (4.184)$$

$$\Delta_3 \doteq \mathbb{E}_{\mathcal{S}_{\bar{\sigma}}}[\sigma(\mathbf{x})(\pi^r(\mathbf{x}) - \pi_{\bar{\sigma}}(\mathbf{x}))h(\mathbf{x})] . \quad (4.185)$$

We also have:

$$\Delta_3 - \Delta_2 - \Delta_1 = \mathbb{E}_{\mathcal{S}'}[(\pi_*(\mathbf{x}) - 1\{y = 1\})h(\mathbf{x})] + \mathbb{E}_{\mathcal{S}}[(\pi(\mathbf{x}) - \pi_*(\mathbf{x}))h(\mathbf{x})] \quad (4.186)$$

$$\doteq \Delta_4 . \quad (4.187)$$

Putting Equations 4.178, 4.182, 4.183 and 4.187 altogether, we get, after introducing Rademacher variables:

$$\mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}'}[\ell(yh(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] \} \quad (4.188)$$

$$= \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}_\sigma}[\sigma(\mathbf{x})\ell(\pi_\sigma(\mathbf{x}), h(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}_{\bar{\sigma}}}[\sigma(\mathbf{x})\ell(\pi_{\bar{\sigma}}(\mathbf{x}), h(\mathbf{x}))] + \Delta_4 \} \quad (4.189)$$

$$\begin{aligned} &\leq \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}_\sigma}[\sigma(\mathbf{x})\ell(\pi_\sigma(\mathbf{x}), h(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}_{\bar{\sigma}}}[\sigma(\mathbf{x})\ell(\pi_{\bar{\sigma}}(\mathbf{x}), h(\mathbf{x}))] \} \\ &\quad + \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}'}[(\pi_*(\mathbf{x}) - 1\{y = 1\})h(\mathbf{x})] + \mathbb{E}_{\mathcal{S}}[(\pi(\mathbf{x}) - \pi_*(\mathbf{x}))h(\mathbf{x}))] \} \end{aligned} \quad (4.190)$$

$$\begin{aligned} &= \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}'}[\sigma(\mathbf{x})\ell(\pi^l(\mathbf{x}), h(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))] \} \\ &\quad + \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}'}[(\pi_*(\mathbf{x}) - 1\{y = 1\})h(\mathbf{x})] + \mathbb{E}_{\mathcal{S}}[(\pi(\mathbf{x}) - \pi_*(\mathbf{x}))h(\mathbf{x}))] \} \end{aligned} \quad (4.191)$$

$$\begin{aligned} &\leq 2\mathbb{E}_{\mathcal{D}_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))] \} \\ &\quad + \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}'}[(\pi_*(\mathbf{x}) - 1\{y = 1\})h(\mathbf{x})] + \mathbb{E}_{\mathcal{S}}[(\pi(\mathbf{x}) - \pi_*(\mathbf{x}))h(\mathbf{x}))] \} . \end{aligned} \quad (4.192)$$

Equation 4.191 holds because the distribution of the supremum is the same. We also have:

$$\mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}'}[(\pi_*(\mathbf{x}) - 1\{y = 1\})h(\mathbf{x})] + \mathbb{E}_{\mathcal{S}}[(\pi(\mathbf{x}) - \pi_*(\mathbf{x}))h(\mathbf{x}))] \} \quad (4.193)$$

$$= \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m, \Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}}[(\pi(\mathbf{x}) - \pi_*(\mathbf{x}))h(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}'}[(1\{y = 1\} - \pi_*(\mathbf{x}))h(\mathbf{x}))] \} \quad (4.194)$$

$$= \mathbb{E}_{\mathcal{D}_{2m}} \mathbb{E}_{\mathcal{I}_1^{/2}, \mathcal{I}_2^{/2}} \sup_h \mathbb{E}_{\mathcal{S}}[\sigma_1(\mathbf{x})(\pi_{|2}^s(\mathbf{x}) - \pi_{|1}^l(\mathbf{x}))h(\mathbf{x})] \quad (4.195)$$

$$= L(\mathcal{H}) . \quad (4.196)$$

Equation 4.195 holds because swapping the sample does not make any difference in the outer expectation, as each couple of swapped samples is generated with the same probability without swapping. Putting altogether 4.192 and 4.196 ends the proof of

Lemma 51. ■

We now bound the deviations of  $\mathbb{E}_{\Sigma_m} \sup_h \{\mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\}$  with respect to its expectation over the sampling of  $\mathcal{S}$ ,  $\mathbb{E}_{\mathcal{D}_m, \Sigma_m} \sup_h \{\mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\}$ . To do that, we use a third time McDiarmid's inequality and compute an upper bound for:

$$\left| \mathbb{E}_{\Sigma_m} \sup_h \{\mathbb{E}_{\mathcal{S}_1}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} - \mathbb{E}_{\Sigma_m} \sup_h \{\mathbb{E}_{\mathcal{S}_2}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} \right| \quad (4.197)$$

$$\leq \mathbb{E}_{\Sigma_m} \left[ \left| \sup_h \{\mathbb{E}_{\mathcal{S}_1}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} - \sup_h \{\mathbb{E}_{\mathcal{S}_2}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} \right| \right] \quad (4.198)$$

$$\leq \max_{\Sigma_m} \left[ \left| \sup_h \{\mathbb{E}_{\mathcal{S}_1}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} - \sup_h \{\mathbb{E}_{\mathcal{S}_2}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} \right| \right] \leq \frac{Q_1}{m}, \quad (4.199)$$

where  $Q_1$  is defined in Equation 4.164. Equation 4.198 holds because of the triangular inequality. Inequality 4.199 holds because  $|\sigma(\cdot)| = 1$ . So with probability  $\leq \delta/2$  over the sampling of  $\mathcal{S}$ :

$$\begin{aligned} & \mathbb{E}_{\Sigma_m} \sup_h \{\mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} \\ & \leq \mathbb{E}_{\mathcal{D}_m, \Sigma_m} \sup_h \{\mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))]\} - Q_1 \sqrt{\frac{1}{2m} \log \frac{2}{\delta}}, \end{aligned} \quad (4.200)$$

where  $Q_1$  is defined via Equation 4.163. We obtain that with probability  $> 1 - (\delta/2 +$

$\delta/2) = 1 - \delta$ , the following holds  $\forall h$ :

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}}[\ell(yh(\mathbf{x}))] \\ & \leq \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] + \phi(\mathcal{S}) \quad (\text{see 4.151 and 4.152}) \end{aligned} \quad (4.201)$$

$$\begin{aligned} & \leq \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] + \mathbb{E}_{\mathcal{D}_m} \sup_g \{ \mathbb{E}_{\mathcal{D}}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} \\ & \quad + Q_1 \sqrt{\frac{1}{2m} \log \frac{2}{\delta}} \quad (\text{from 4.165}) \end{aligned} \quad (4.202)$$

$$\begin{aligned} & \leq \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] + \mathbb{E}_{\mathcal{D}_m, \mathcal{D}'_m} \sup_g \{ \mathbb{E}_{\mathcal{S}'}[\ell(yg(\mathbf{x}))] - \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} \\ & \quad + Q_1 \sqrt{\frac{1}{2m} \log \frac{2}{\delta}} \quad (\text{from 4.168}) \end{aligned} \quad (4.203)$$

$$\begin{aligned} & \leq \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] + 2\mathbb{E}_{\mathcal{D}_m, \Sigma_m} \sup_g \{ \mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), g(\mathbf{x}))] \} + L(\mathcal{H}) \\ & \quad + Q_1 \sqrt{\frac{1}{2m} \log \frac{2}{\delta}} \quad (\text{Lemma 51}) \end{aligned} \quad (4.204)$$

$$\begin{aligned} & \leq \mathbb{E}_{\mathcal{S}}[\ell(\pi(\mathbf{x}), h(\mathbf{x}))] + 2\mathbb{E}_{\Sigma_m} \sup_h \{ \mathbb{E}_{\mathcal{S}}[\sigma(\mathbf{x})\ell(\pi(\mathbf{x}), h(\mathbf{x}))] \} + L(\mathcal{H}) \\ & \quad + 2Q_1 \sqrt{\frac{1}{2m} \log \frac{2}{\delta}} \quad (\text{from 4.200}) \end{aligned} \quad (4.205)$$

$$\begin{aligned} & = \mathbb{E}_{\Sigma_\pi} \mathbb{E}_{\mathcal{S}}[\ell(\sigma(\mathbf{x})h(\mathbf{x}))] + 2R^b(\ell \circ \mathcal{H} \circ \mathcal{S}) + L(\mathcal{H}) \\ & \quad + 4 \left( \frac{2H}{b_\phi} + 1 \right) \sqrt{\frac{1}{2m} \log \frac{2}{\delta}}, \end{aligned} \quad (4.206)$$

as claimed.

#### 4.11.9.2 Proof of Equation 4.34

We have  $\ell'(x) = -(1/b_\phi)(\phi^*)'(-x) = -(1/b_\phi)(\phi')^{-1}(-x) \in [-1/b_\phi, 0]$ .  $\ell$  is therefore  $1/b_\phi$ -Lipschitz. So Theorem 4.12 in Ledoux and Talagrand [1991] brings:

$$\mathcal{R}^b(F, \eta) = \mathbb{E}_{\sigma \sim \Sigma_m} \sup_{h \in \mathcal{H}} \left\{ \mathbb{E}_{i \sim [m]} [\sigma_i \mathbb{E}_{\sigma' \sim \Sigma_\pi} [\ell(\sigma'_i h(\mathbf{x}_i) - \eta)]] \right\} \quad (4.207)$$

$$\leq b_\phi \mathbb{E}_{\sigma \sim \Sigma_m} \sup_{h \in \mathcal{H}} \left\{ \mathbb{E}_{i \sim [m]} [\sigma_i \mathbb{E}_{\sigma' \sim \Sigma_\pi} [\sigma'_i h(\mathbf{x}_i) - \eta]] \right\} \quad (4.208)$$

$$= b_\phi \mathbb{E}_{\sigma \sim \Sigma_m} \sup_{h \in \mathcal{H}} \left\{ \mathbb{E}_{i \sim [m]} [\sigma_i \mathbb{E}_{\sigma' \sim \Sigma_\pi} [\sigma'_i h(\mathbf{x}_i)]] \right\} \quad (4.209)$$

$$= b_\phi \mathbb{E}_{\sigma \sim \Sigma_m} \sup_{h \in \mathcal{H}} \left\{ \mathbb{E}_{i \sim [m]} [\sigma_i (2\pi(\mathbf{x}_i) - 1) h(\mathbf{x}_i)] \right\}, \quad (4.210)$$

as claimed.

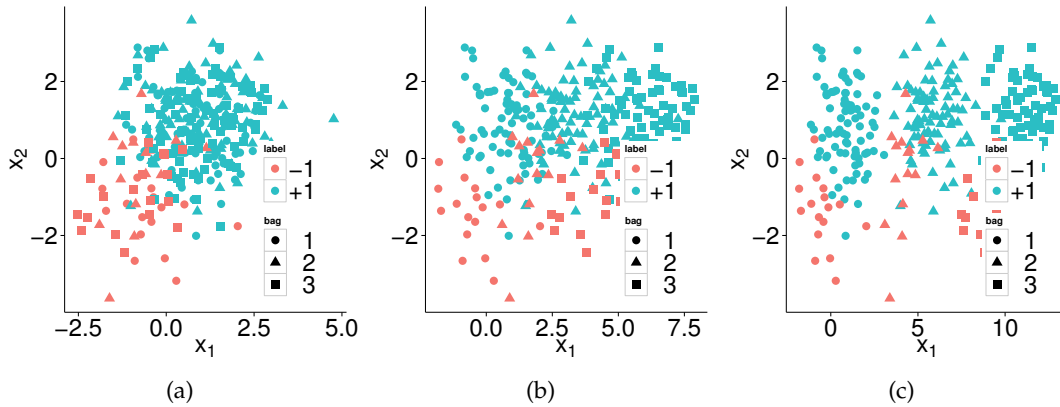
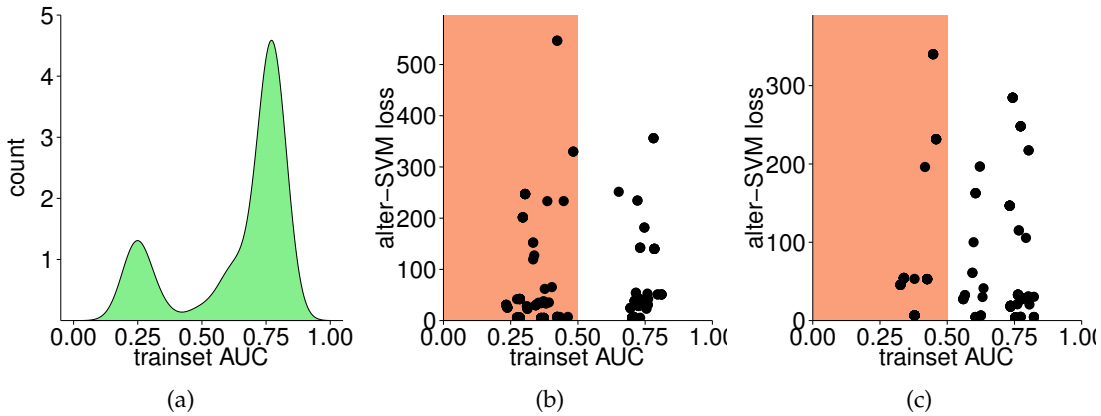


Figure 4.2: Violation of homogeneity assumption

Figure 4.3: alter- $\alpha$ SVM: empirical distribution of AUC (a), and relationship between loss and AUC in two different train split (b)(c)

## 4.12 Appendix: additional experimental results

### 4.12.1 Simulated domain for violation of homogeneity assumption

The synthetic data generated for this test consists on 16 classification problems, each one formed by 16 bags of 100 two-dimensional normal samples. The distribution generating the first dataset satisfies the homogeneity assumption (Figure 4.2 (a)). Then, we gradually change the position of the class-conditional bag-conditional means on one linear direction (to the right on Figure 4.2 (b) and (c)), with different offsets for different bags. In Figure 4.2 we give a graphical explanation of the process on 3 bags.

### 4.12.2 Additional tests on alter- $\alpha$ SVM Yu et al. [2013]

In our experiments, we observe that the AUC achieved by  $\alpha$ SVM can be very high, but it is also often *below* 0.5; in those cases the algorithm outputs models which are

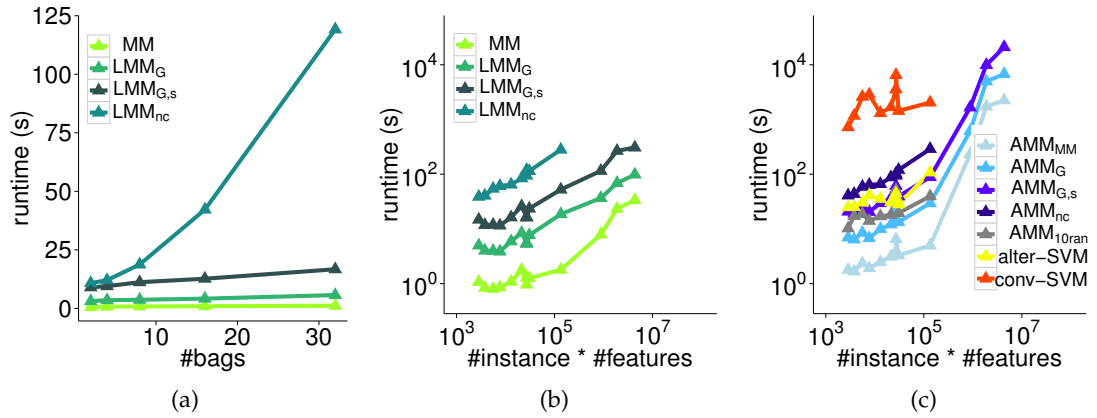


Figure 4.4: Learning runtime of LMM for bags number (a), and for domain size one-shot (b) and iterative methods (c)

worse than random and the average performance over the 5 test splits drops. We are able to reproduce the same behavior on the heart dataset provided by the authors in a demo for alter- $\alpha$ SVM; this also proves our bag assignment for LLP simulation does not introduce the issue. In a first test, we randomly select 3/4 of the dataset, and randomly assign instances to 4 bags of fixed size 64, following Yu et al. [2013]. We repeat the training split 50 times with  $C = C_p = 1$ , as in the demo, and we measure AUC on the same training set. As expected, a consistent number of run (22%) ends up producing AUC smaller than 0.5. We display in Figure 4.3 (a) the AUC density profile, which shows a relevant mass around 0.25; notice also the two distribution modes look symmetric around 0.5.

In a second test, we investigate further measuring pairs of training set AUC and losses obtained by the same execution of the algorithm. In this case, we run over all parameters ranges defined in  $\alpha$ SVM paper, and do not pick the model that minimizes the loss over the 10 random runs, but record losses of all. Figures 4.3 (b) and (c) show scatter plots relative to two chosen training set splits. We observe that loss minimization can lead both to high and low AUC, with few points close to 0.5. A possible explanation might be found in the inverted polarity of the learned linear classifier; inverted polarity in this contest means having a model which would achieve better performance classifying instances labels opposite to the ones predicted. We conclude that optimizing  $\alpha$ SVM loss in some cases might be equivalent to train a max-margin separator of the unlabeled data, exploiting only weakly the information given by the label proportions. This would give a heuristic understanding of the frequent symmetrical behavior of the AUC.

### 4.12.3 Scalability

Figure 4.1 (a) shows runtime of learning (including cross-validation) of MM and LMM with regard to the number of bags – which is the natural parameter of time

dataset	instances	feature
arrhythmia	452	297
australian	690	39
breastw	699	11
colic	368	83
german	1000	27
heart	270	14
ionosphere	351	37
vertebral column	620	9
vote	435	49
wine	178	16

Table 4.5: Small domains size

complexity for our Laplacian-based methods. Although the 3 layers of cross-validation of  $LMM_{G,S}$ ,  $LMM_{nc}$  results the only method clearly not scalable. Figure 4.1 (b) presents how our one-shots algorithms scale on all small domains as a function of problem size. Runtime is averaged over the different bag assignments. The same plot is given in Figure 4.1 (c) for iterative algorithms, in particular  $AMM^{min}$  and (alter/conv)- $\alpha$ SVM. All curves are completed with measurements on bigger domains when available. Runtime of SVMs is not directly comparable with our methods. This is due to both (a) the implementation on different programming languages and (b) to the fact that the code provided implements kernel SVM, even for linear kernels, which is a big overhead in computation and memory access. Nevertheless, the high growth rate of conv- $\alpha$ SVM makes the algorithm not suitable for large datasets. Noticeably, even if alter- $\alpha$ SVM does not show such behavior, we are not able to run it on the bigger domains, since it requires several hours to run on a training set split with fixed parameters.

#### 4.12.4 Full results on small domains

Finally we report details about all experiments run on the 10 small domains (Table 4.5). In the following Tables, columns show the number of bags generated through  $k$ -MEANS. Each cell contains average AUC over 5 test splits and standard deviation; runtime in second is in the separated column. Best performing algorithm and ones not worse than 0.1 AUC are bold faced. Comparisons are made in the respective top/bottom sub-tables, which group one-shot and iterative algorithms. We use  $\uparrow$  to highlight runs which achieve AUC greater or equal than the Oracles.

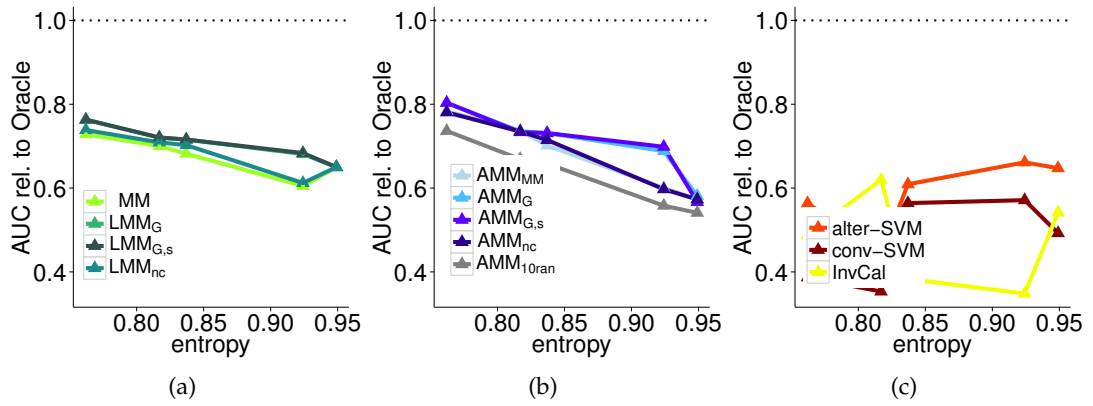


Figure 4.5: Relative AUC (w.r.t. Oracle) vs entropy on arrhythmia

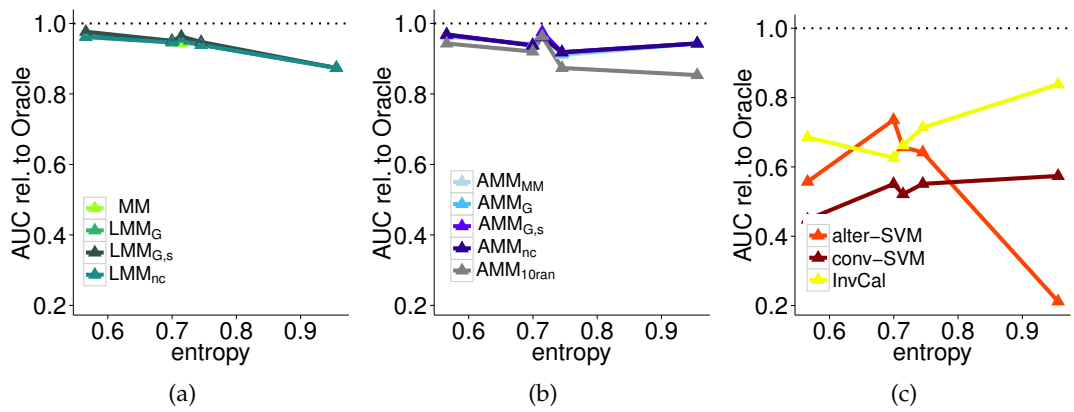


Figure 4.6: Relative AUC (w.r.t. Oracle) vs entropy on australian

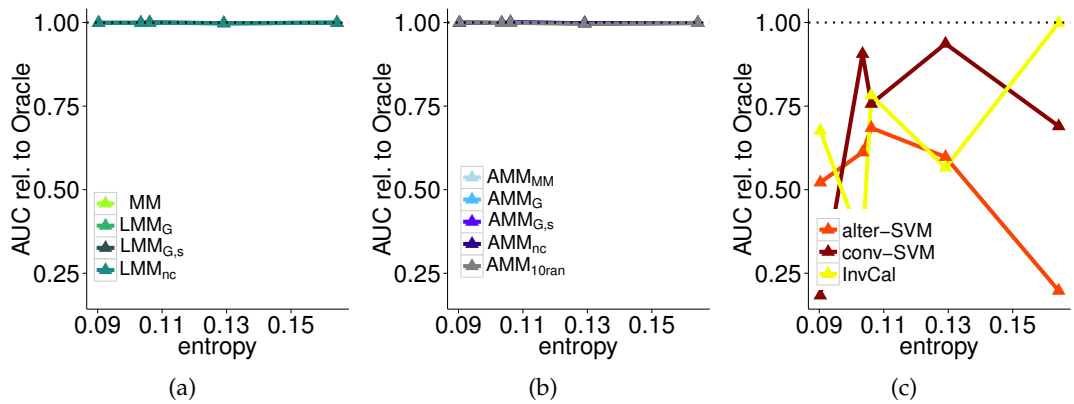


Figure 4.7: Relative AUC (w.r.t. Oracle) vs entropy on breastw

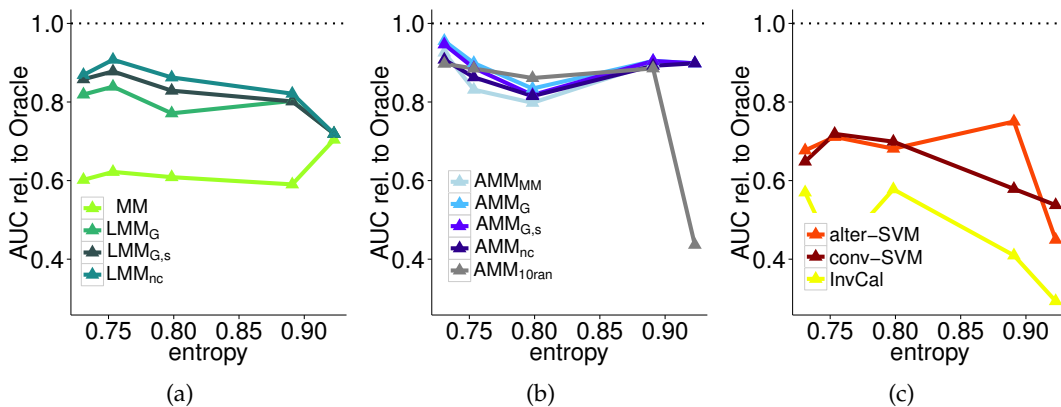


Figure 4.8: Relative AUC (w.r.t. Oracle) vs entropy on colic

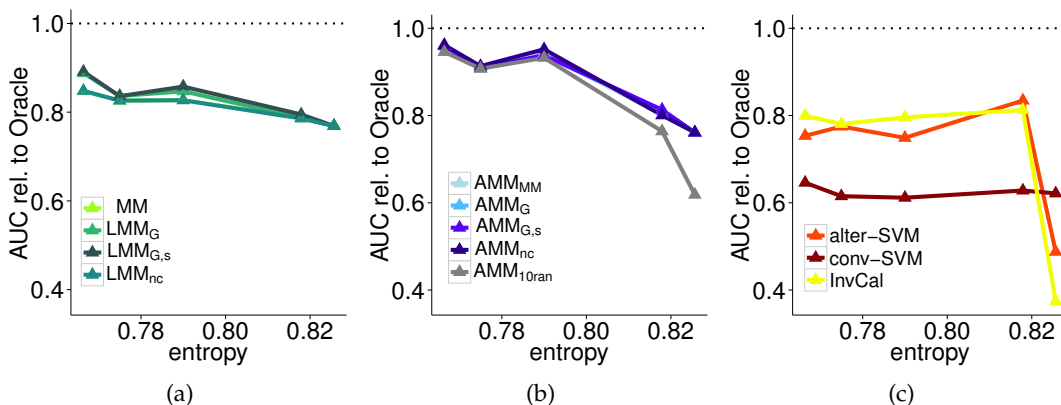


Figure 4.9: Relative AUC (w.r.t. Oracle) vs entropy on german

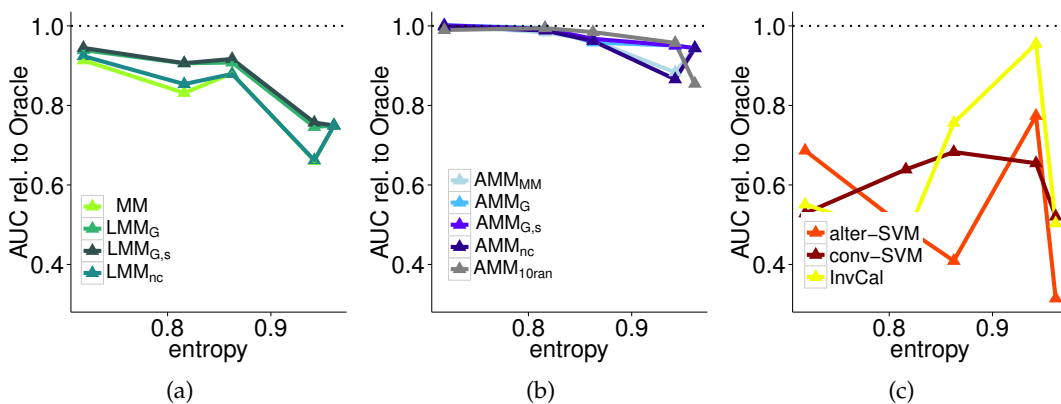


Figure 4.10: Relative AUC (w.r.t. Oracle) vs entropy on heart



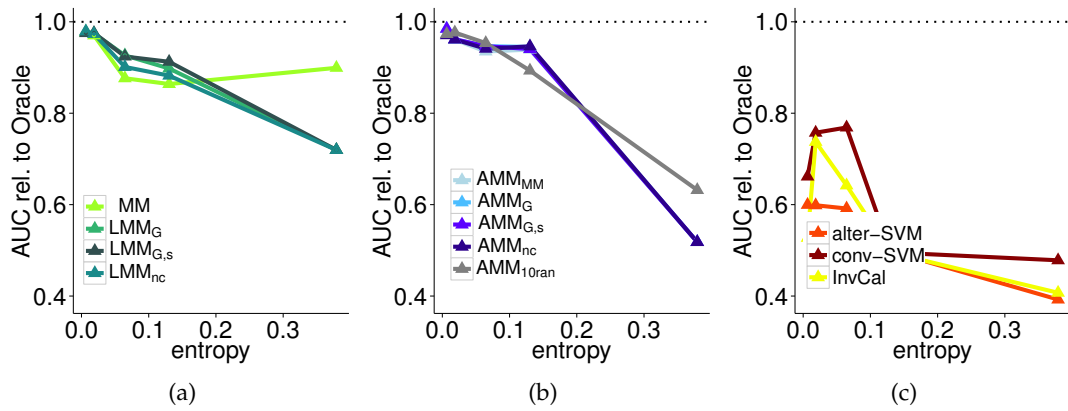


Figure 4.11: Relative AUC (w.r.t. Oracle) vs entropy on ionosphere

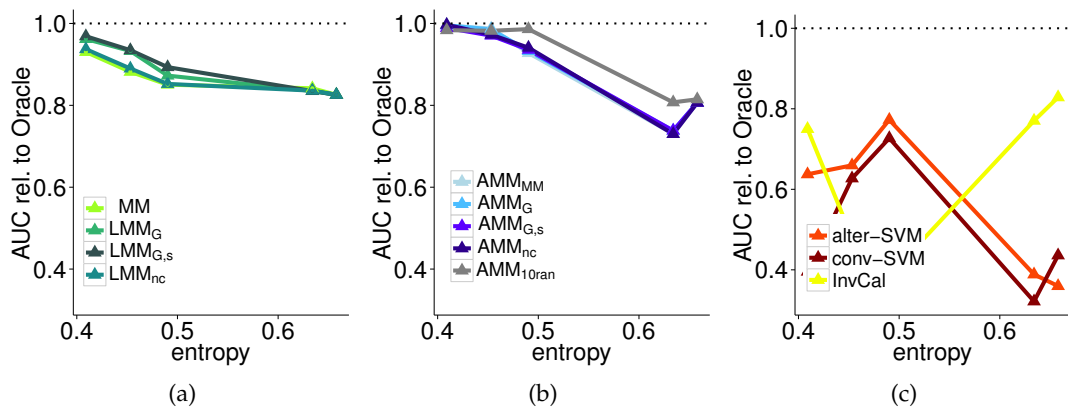


Figure 4.12: Relative AUC (w.r.t. Oracle) vs entropy on vertebral column

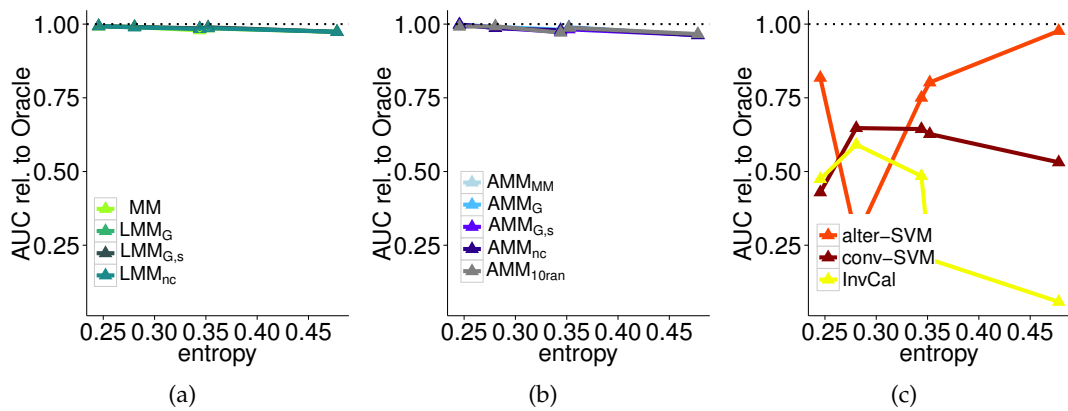


Figure 4.13: Relative AUC (w.r.t. Oracle) vs entropy on vote

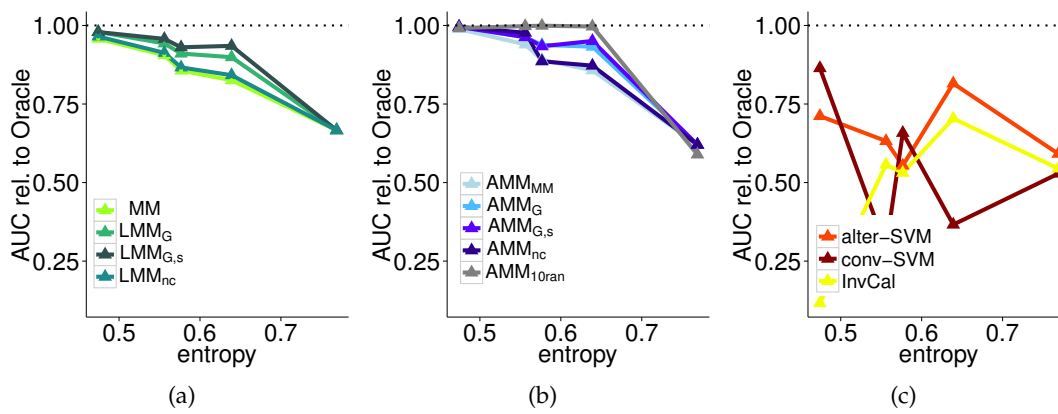


Figure 4.14: Relative AUC (w.r.t. Oracle) vs entropy on wine

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM	$70.91 \pm 6.81$	2	$50.55 \pm 7.54$	2	$50.31 \pm 7.55$	2	$47.03 \pm 6.60$	2	$52.34 \pm 7.25$	2	
MM	$64.99 \pm 2.99$	2	$60.48 \pm 7.28$	1	$68.17 \pm 5.95$	2	$70.01 \pm 9.33$	2	$72.85 \pm 9.49$	2	
LMM <sub>G</sub>	$64.99 \pm 2.99$	18	$68.10 \pm 4.43$	17	<b>71.53 ± 2.36</b>	20	<b>72.06 ± 7.62</b>	18	<b>76.29 ± 7.91</b>	20	
LMM <sub>G,s</sub>	$64.99 \pm 2.99$	49	<b>68.34 ± 3.95</b>	49	<b>71.53 ± 2.36</b>	54	<b>72.06 ± 7.62</b>	52	<b>76.29 ± 7.91</b>	57	
LMM <sub>nc</sub>	$64.99 \pm 2.99$	83	$61.19 \pm 7.53$	83	$70.21 \pm 5.17$	119	$70.89 \pm 9.86$	267	$73.82 \pm 9.29$	854	
InvCal	$64.75 \pm 3.04$	17	$66.12 \pm 2.60$	17	$60.87 \pm 3.54$	17	$44.46 \pm 3.36$	17	$56.36 \pm 5.26$	17	
AMM <sup>min</sup>	AMM <sub>EMM</sub>	$59.54 \pm 7.52$	9	$52.65 \pm 3.10$	8	$63.46 \pm 10.37$	8	$67.85 \pm 9.56$	8	$75.65 \pm 8.81$	8
	AMM <sub>MM</sub>	$57.29 \pm 5.95$	7	$60.00 \pm 7.96$	4	$70.12 \pm 6.46$	4	$73.66 \pm 8.86$	5	$78.36 \pm 8.53$	5
	AMM <sub>G</sub>	$58.15 \pm 6.83$	31	$68.80 \pm 2.15$	28	<b>73.08 ± 2.92</b>	30	$74.54 \pm 7.98$	29	<b>80.32 ± 8.08</b>	30
	AMM <sub>G,s</sub>	$56.67 \pm 4.66$	92	$69.83 \pm 2.69$	84	<b>73.08 ± 2.92</b>	88	$73.34 \pm 7.62$	88	<b>80.32 ± 8.08</b>	91
	AMM <sub>nc</sub>	$57.29 \pm 5.95$	97	$59.71 \pm 8.39$	90	$71.43 \pm 6.21$	126	$73.49 \pm 8.95$	274	$78.04 \pm 8.26$	862
	AMM <sub>1</sub>	<b>65.80 ± 6.92</b>	5	<b>70.00 ± 5.89</b>	4	$68.17 \pm 7.19$	4	$69.93 \pm 4.27$	4	$72.31 \pm 5.02$	5
	AMM <sub>10ran</sub>	$54.09 \pm 12.03$	30	$55.78 \pm 17.36$	32	$66.38 \pm 7.32$	51	$66.89 \pm 6.75$	51	$73.61 \pm 5.15$	57
AMM <sup>max</sup>	AMM <sub>EMM</sub>	$50.59 \pm 5.97$	41	$59.32 \pm 5.82$	41	$60.85 \pm 5.43$	37	$60.38 \pm 4.08$	41	$58.31 \pm 8.40$	40
	AMM <sub>MM</sub>	$62.08 \pm 9.46$	45	$46.86 \pm 3.90$	34	$67.28 \pm 8.92$	33	$74.04 \pm 9.46$	35	$71.00 \pm 7.65$	38
	AMM <sub>G</sub>	$62.08 \pm 9.46$	141	$62.27 \pm 8.14$	128	$65.78 \pm 3.92$	118	$64.64 \pm 10.26$	121	$73.07 \pm 6.72$	124
	AMM <sub>G,s</sub>	$62.08 \pm 9.46$	414	$63.13 \pm 5.17$	380	$63.85 \pm 7.00$	346	$65.49 \pm 10.62$	354	$73.05 \pm 6.70$	374
	AMM <sub>nc</sub>	$62.08 \pm 9.46$	206	$55.57 \pm 6.07$	182	$64.30 \pm 6.24$	207	<b>76.33 ± 3.96</b>	362	$70.82 \pm 4.23$	965
	AMM <sub>1</sub>	$60.53 \pm 9.79$	31	$54.14 \pm 13.28$	34	$67.45 \pm 3.91$	32	$55.85 \pm 8.96$	35	$61.26 \pm 6.95$	38
	AMM <sub>10ran</sub>	$49.79 \pm 8.14$	307	$55.37 \pm 14.62$	370	$53.78 \pm 5.13$	301	$60.62 \pm 8.04$	322	$64.20 \pm 2.84$	338
SVM	alter- $\alpha$	$49.24 \pm 3.92$	96	$57.10 \pm 2.71$	100	$56.38 \pm 2.73$	104	$35.31 \pm 1.30$	114	$38.68 \pm 6.10$	125
	conv- $\alpha$	$54.15 \pm 2.22$	2054	$34.82 \pm 3.20$	2078	$38.31 \pm 8.24$	2168	$61.96 \pm 1.10$	1930	$48.77 \pm 5.73$	2004
Oracle	$99.99 \pm 0.02$	2	$99.98 \pm 0.05$	2	$99.94 \pm 0.13$	2	$100.00 \pm 0.00$	2	$99.97 \pm 0.07$	2	

Table 4.6: arrhythmia

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM	$66.48 \pm 3.16$	<1	$64.67 \pm 4.22$	<1	$63.56 \pm 4.00$	<1	$64.17 \pm 4.80$	<1	$63.14 \pm 5.41$	<1	
MM	<b>81.08 ± 1.66</b>	<1	$87.11 \pm 2.68$	<1	$87.49 \pm 2.86$	1	$87.36 \pm 2.22$	<1	$89.53 \pm 2.13$	2	
LMM <sub>G</sub>	<b>81.08 ± 1.66</b>	4	$87.09 \pm 2.82$	4	<b>87.81 ± 3.16</b>	5	$88.46 \pm 2.50$	6	$89.69 \pm 2.68$	8	
LMM <sub>G,s</sub>	<b>81.08 ± 1.66</b>	14	<b>87.81 ± 3.08</b>	15	<b>87.88 ± 3.21</b>	19	<b>89.18 ± 2.05</b>	20	<b>90.80 ± 2.53</b>	27	
LMM <sub>nc</sub>	<b>81.08 ± 1.66</b>	57	$87.02 \pm 2.72$	49	$87.46 \pm 3.03$	57	$88.06 \pm 2.31$	90	$89.41 \pm 2.41$	217	
Invcal	$19.67 \pm 2.23$	5	$59.50 \pm 5.86$	5	$68.00 \pm 5.27$	5	$60.83 \pm 3.17$	5	$51.81 \pm 4.72$	5	
AMM <sup>min</sup>	AMM <sub>EMM</sub>	$86.65 \pm 2.06$	4	<b>86.59 ± 3.08</b>	4	$86.50 \pm 4.11$	4	$89.51 \pm 2.48$	6	$88.85 \pm 4$	6
	AMM <sub>MM</sub>	<b>87.54 ± 3.84</b>	3	$84.35 \pm 3.63$	4	<b>86.99 ± 3.87</b>	4	$89.43 \pm 1.34$	4	$89.55 \pm 3.18$	5
	AMM <sub>G</sub>	<b>87.54 ± 3.84</b>	10	$84.79 \pm 3.17$	13	$86.78 \pm 4.21$	14	$89.52 \pm 2.18$	14	$89.88 \pm 2.78$	18
	AMM <sub>G,s</sub>	<b>87.54 ± 3.84</b>	30	$85.12 \pm 3.75$	39	$86.75 \pm 4.19$	43	<b>90.37 ± 1.67</b>	43	$89.95 \pm 2.80$	54
	AMM <sub>nc</sub>	<b>87.54 ± 3.84</b>	63	$85.10 \pm 3.55$	57	$86.63 \pm 4.02$	66	$89.00 \pm 1.83$	97	<b>90.11 ± 2.93</b>	227
	AMM <sub>1</sub>	$72.60 \pm 5.70$	2	$85.04 \pm 2.53$	3	<b>86.89 ± 3.73</b>	4	$88.91 \pm 2.32$	4	$88.98 \pm 3.00$	4
	AMM <sub>10ran</sub>	$79.21 \pm 5.07$	27	$80.97 \pm 2.27$	31	$85.08 \pm 3.30$	34	$89.19 \pm 1.81$	46	$87.70 \pm 2.68$	47
AMM <sup>max</sup>	AMM <sub>EMM</sub>	$80.09 \pm 3.99$	17	$71.46 \pm 1.85$	16	$73.41 \pm 6.07$	16	$73.25 \pm 3.33$	18	$81.73 \pm 3.60$	19
	AMM <sub>MM</sub>	$86.83 \pm 4.26$	20	$72.96 \pm 2.30$	15	$70.25 \pm 4.65$	16	$73.89 \pm 5.77$	18	$75.91 \pm 3.50$	21
	AMM <sub>G</sub>	$86.83 \pm 4.26$	61	$73.32 \pm 1.95$	48	$71.16 \pm 4.94$	51	$73.57 \pm 6.86$	55	$75.25 \pm 3.18$	63
	AMM <sub>G,s</sub>	$86.83 \pm 4.26$	181	$73.25 \pm 2.03$	143	$71.19 \pm 4.91$	153	$74.77 \pm 6.85$	163	$75.25 \pm 3.18$	188
	AMM <sub>nc</sub>	$86.83 \pm 4.26$	114	$73.74 \pm 2.48$	92	$70.36 \pm 5.16$	102	$75.16 \pm 5.71$	138	$76.44 \pm 2.74$	272
	AMM <sub>1</sub>	$69.57 \pm 3.99$	15	$73.12 \pm 3.41$	15	$68.25 \pm 2.80$	16	$71.02 \pm 5.46$	17	$81.70 \pm 3.02$	19
	AMM <sub>10ran</sub>	$77.82 \pm 9.12$	192	$68.82 \pm 4.73$	138	$73.58 \pm 4.29$	146	$72.21 \pm 9.35$	164	$74.16 \pm 5.25$	188
SVM	alter- $\alpha$	$53.26 \pm 2.07$	25	$51.08 \pm 2.35$	27	$50.90 \pm 1.63$	31	$48.29 \pm 4.51$	38	$41.66 \pm 5.11$	64
	conv- $\alpha$	$77.80 \pm 6.16$	3924	$66.14 \pm 4.68$	3790	$57.94 \pm 18.54$	3244	$61.37 \pm 21.17$	3327	$63.73 \pm 11.33$	3603
Oracle	$92.81 \pm 2.89$	<1	$92.68 \pm 2.24$	<1	$92.44 \pm 3.01$	1	$92.61 \pm 2.03$	<1	$92.99 \pm 3.58$	<1	

Table 4.7: australian

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM	48.65 ± 7.54	<1	71.45 ± 16.59	<1	61.68 ± 7.47	<1	34.88 ± 12.33	<1	47.50 ± 22.77	<1	
MM	<b>99.42 ± 0.44</b>	2	<b>99.30 ± 0.39</b>	<1	<b>99.28 ± 0.25</b>	<1	<b>99.28 ± 0.37</b>	<1	<b>99.18 ± 0.47</b>	1	
LMM <sub>G</sub>	<b>99.42 ± 0.44</b>	6	<b>99.33 ± 0.38</b>	3	<b>99.28 ± 0.25</b>	3	<b>99.35 ± 0.39</b>	3	<b>99.22 ± 0.46</b>	4	
LMM <sub>G,s</sub>	<b>99.42 ± 0.44</b>	20	<b>99.34 ± 0.39</b>	10	<b>99.37 ± 0.24</b> ↑	11	<b>99.36 ± 0.38</b>	12	<b>99.23 ± 0.44</b>	15	
LMM <sub>nc</sub>	<b>99.42 ± 0.44</b>	41	<b>99.29 ± 0.40</b>	39	<b>99.27 ± 0.25</b>	41	<b>99.30 ± 0.38</b>	59	<b>99.20 ± 0.47</b>	125	
Invcnl	19.67 ± 2.23	5	59.50 ± 5.86	5	68 ± 5.27	5	60.83 ± 3.17	5	51.81 ± 4.72	5	
AMM <sup>min</sup>	AMM <sub>EMM</sub>	<b>99.37 ± 0.42</b>	1	<b>99.33 ± 0.39</b>	1	99.17 ± 0.54	1	<b>99.34 ± 0.40</b>	2	99.29 ± 0.49	2
	AMM <sub>MM</sub>	<b>99.34 ± 0.46</b>	2	<b>99.30 ± 0.37</b>	1	<b>99.36 ± 0.27</b> ↑	2	99.29 ± 0.41	2	99.29 ± 0.48	2
	AMM <sub>G</sub>	<b>99.34 ± 0.46</b>	8	<b>99.30 ± 0.37</b> ↑	5	<b>99.36 ± 0.27</b> ↑	6	99.29 ± 0.41	7	99.30 ± 0.49	8
	AMM <sub>G,s</sub>	<b>99.34 ± 0.46</b>	23	<b>99.30 ± 0.37</b> ↑	16	<b>99.36 ± 0.27</b> ↑	19	99.29 ± 0.41	20	99.30 ± 0.49	25
	AMM <sub>nc</sub>	<b>99.34 ± 0.46</b>	43	<b>99.31 ± 0.35</b>	41	<b>99.36 ± 0.27</b> ↑	44	99.29 ± 0.41	62	99.29 ± 0.48	129
	AMM <sub>1</sub>	<b>99.35 ± 0.45</b>	<1	<b>99.32 ± 0.37</b>	1	99.20 ± 0.45	1	99.30 ± 0.42	1	<b>99.31 ± 0.48</b>	2
	AMM <sub>10ran</sub>	<b>99.36 ± 0.45</b>	8	99.11 ± 0.56	9	99.26 ± 0.35	11	99.28 ± 0.43	11	<b>99.32 ± 0.49</b> ↑	14
	AMM <sub>EMM</sub>	<b>99.42 ± 0.55</b>	6	99.02 ± 0.66	6	<b>99.32 ± 0.25</b> ↑	6	<b>99.43 ± 0.30</b> ↑	7	<b>99.40 ± 0.38</b> ↑	9
AMM <sup>max</sup>	AMM <sub>MM</sub>	99.01 ± 1.12	6	99.00 ± 0.64	6	<b>99.32 ± 0.35</b> ↑	6	<b>99.37 ± 0.38</b>	7	<b>99.39 ± 0.39</b> ↑	9
	AMM <sub>G</sub>	99.01 ± 1.12	20	98.99 ± 0.64	17	<b>99.33 ± 0.35</b> ↑	18	<b>99.37 ± 0.38</b>	21	<b>99.41 ± 0.39</b> ↑	27
	AMM <sub>G,s</sub>	99.01 ± 1.12	60	98.99 ± 0.64	52	99.19 ± 0.45	55	<b>99.37 ± 0.39</b>	63	<b>99.41 ± 0.39</b> ↑	82
	AMM <sub>nc</sub>	99.01 ± 1.12	55	98.99 ± 0.64	53	<b>99.32 ± 0.35</b> ↑	56	<b>99.37 ± 0.39</b>	76	<b>99.40 ± 0.38</b> ↑	148
	AMM <sub>1</sub>	99.09 ± 1.08	5	99.09 ± 0.46	5	<b>99.29 ± 0.26</b>	5	<b>99.37 ± 0.38</b>	6	<b>99.40 ± 0.38</b> ↑	8
	AMM <sub>10ran</sub>	98.97 ± 1.29	47	98.58 ± 0.75	48	<b>99.39 ± 0.27</b> ↑	52	<b>99.37 ± 0.38</b>	61	<b>99.36 ± 0.41</b> ↑	81
	alter- $\alpha$	68.63 ± 17.63	24	93.24 ± 4.43	25	75.17 ± 7.19	33	90.11 ± 2.58	42	18.23 ± 5.67	82
	conv- $\alpha$	<b>99.41 ± 0.48</b>	3346	56.33 ± 4.28	3043	77.71 ± 15.51	2800	32.90 ± 7.24	3036	67.21 ± 8.19	2037
Oracle	99.48 ± 0.41	<1	99.53 ± 0.41	<1	99.31 ± 0.37	<1	99.43 ± 0.39	<1	99.32 ± 0.44	<1	

Table 4.8: breastw

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM	60.69 ± 11.30	<1	51.83 ± 6.36	<1	52.99 ± 5.37	<1	53.83 ± 11.49	<1	52.95 ± 13.28	<1	
MM	<b>62.00 ± 6.44</b>	<1	70.48 ± 7.43	<1	67.13 ± 9.85	2	72.60 ± 9.35	1	72.05 ± 3.38	1	
LMM <sub>G</sub>	<b>62.00 ± 6.44</b>	7	70.37 ± 7.47	6	72.15 ± 8.51	8	75.96 ± 10.38	8	75.47 ± 3.59	9	
LMM <sub>G,s</sub>	<b>62.00 ± 6.44</b>	20	<b>72.10 ± 6.26</b>	20	<b>75.08 ± 7.14</b>	28	<b>78.54 ± 10.20</b>	26	<b>76.43 ± 3.10</b>	27	
LMM <sub>nc</sub>	<b>62.00 ± 6.44</b>	31	70.45 ± 7.46	33	68.38 ± 9.69	52	74.04 ± 10.02	112	72.87 ± 3.20	345	
Invcnl	38.73 ± 5.43	6	65.87 ± 6.70	6	59.30 ± 3.28	6	61.54 ± 4.17	6	59.53 ± 10.00	6	
AMM <sup>min</sup>	AMM <sub>EMM</sub>	59.12 ± 8.86	3	56.23 ± 8.49	3	70.93 ± 10.31	3	<b>78.22 ± 6.00</b>	3	74.22 ± 6.35	4
	AMM <sub>MM</sub>	<b>77.44 ± 3.16</b>	2	78.84 ± 6.95	3	69.46 ± 6.44	4	71.93 ± 7.61	4	81.44 ± 5.18	4
	AMM <sub>G</sub>	<b>77.44 ± 3.16</b>	11	<b>79.41 ± 2.23</b>	12	72.62 ± 5.42	14	77.80 ± 8.11	14	<b>84.05 ± 2.33</b>	16
	AMM <sub>G,s</sub>	<b>77.44 ± 3.16</b>	34	<b>79.41 ± 2.23</b>	36	71.19 ± 5.38	41	76.71 ± 6.70	40	83.27 ± 3.14	47
	AMM <sub>nc</sub>	<b>77.44 ± 3.16</b>	36	78.33 ± 7.35	38	70.95 ± 4.69	57	74.67 ± 9.10	117	79.86 ± 4.87	352
	AMM <sub>1</sub>	38.69 ± 7.18	1	56.07 ± 14.68	2	<b>75.14 ± 4.78</b>	2	75.36 ± 5.64	3	77.51 ± 5.00	3
	AMM <sub>10ran</sub>	37.63 ± 4.19	10	77.75 ± 5.66	12	74.95 ± 5.64	15	76.59 ± 10.81	17	78.94 ± 4.17	23
	AMM <sub>EMM</sub>	50.94 ± 6.54	9	62.44 ± 9.94	9	57.53 ± 13.37	15	53.63 ± 14.71	17	67.63 ± 5.63	19
AMM <sup>max</sup>	AMM <sub>MM</sub>	43.05 ± 14.65	8	75.40 ± 4.64	9	63.72 ± 14.41	16	55.37 ± 10.19	18	69.49 ± 3.17	20
	AMM <sub>G</sub>	43.05 ± 14.65	28	78.19 ± 5.93	31	63.14 ± 7.53	51	61.32 ± 5.69	57	68.21 ± 9.35	62
	AMM <sub>G,s</sub>	43.05 ± 14.65	84	77.91 ± 6.36	91	62.57 ± 6.11	151	64.42 ± 10.77	168	69.47 ± 6.40	184
	AMM <sub>nc</sub>	42.92 ± 14.74	52	73.74 ± 7.21	57	60.39 ± 12.21	94	62.46 ± 15.13	162	68.63 ± 2.37	381
	AMM <sub>1</sub>	51.92 ± 19.91	7	59.89 ± 10.79	8	58.76 ± 12.16	14	62.31 ± 13.32	17	68.25 ± 6.42	18
	AMM <sub>10ran</sub>	56.39 ± 10.26	60	71.28 ± 8.76	68	65.01 ± 13.85	114	69.59 ± 9.96	139	74.40 ± 5.54	159
	alter- $\alpha$	46.33 ± 2.73	18	50.82 ± 1.21	19	60.84 ± 5.51	23	62.20 ± 3.79	32	57.04 ± 10.10	49
	conv- $\alpha$	25.27 ± 3.45	1438	35.96 ± 9.34	1460	50.31 ± 5.57	1439	35.46 ± 9.11	1423	50.13 ± 8.34	1427
Oracle	86.19 ± 4.23	<1	87.80 ± 2.50	<1	87.05 ± 6.05	<1	86.53 ± 7.15	<1	87.97 ± 2.02	<1	

Table 4.9: colic

algorithm AUC	2 bags		4 bags		8 bags		16 bags		32 bags		
	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	
EMM	47.90 ± 4.51	<1	50.11 ± 5.17	<1	46.02 ± 5.88	<1	50.94 ± 1.61	<1	51.02 ± 2.55	<1	
MM	<b>61.07 ± 5.57</b>	<1	62.09 ± 4.00	<1	65.50 ± 6.54	2	65.61 ± 6.05	2	66.96 ± 4.56	2	
LMM <sub>G</sub>	<b>61.07 ± 5.57</b>	4	62.14 ± 4.04	4	67.07 ± 6.36	6	66.43 ± 6.61	6	70.18 ± 4.76	7	
LMM <sub>G,s</sub>	<b>61.07 ± 5.57</b>	11	62.75 ± 3.32	12	<b>67.91 ± 5.80</b>	16	<b>66.40 ± 6.90</b>	19	<b>70.43 ± 5.57</b>	21	
LMM <sub>nc</sub>	<b>61.07 ± 5.57</b>	103	62.04 ± 4.00	87	65.47 ± 6.56	87	65.61 ± 6.06	113	67.01 ± 4.58	209	
Invcal	38.74 ± 5.43	6	<b>65.87 ± 6.70</b>	6	59.30 ± 3.28	6	61.53 ± 4.17	6	59.54 ± 10.00	6	
AMM <sup>min</sup>	AMM <sub>EMM</sub>	53.89 ± 6.82	7	48.63 ± 8.71	7	53.24 ± 8.02	8	57.58 ± 3.44	9	63.64 ± 11.82	11
	AMM <sub>MM</sub>	<b>60.45 ± 5.58</b>	5	63.33 ± 4.99	6	74.58 ± 4.76	6	72.43 ± 1.39	8	75.84 ± 5.24	7
	AMM <sub>G</sub>	<b>60.45 ± 5.58</b>	17	<b>64.16 ± 6.99</b>	18	74.18 ± 4.34	21	72.08 ± 1.24	22	<b>75.94 ± 4.55</b>	24
	AMM <sub>G,s</sub>	<b>60.45 ± 5.58</b>	52	<b>64.20 ± 7.24</b>	57	74.29 ± 4.50	57	72.18 ± 1.37	66	75.77 ± 4.44	74
	AMM <sub>nc</sub>	<b>60.45 ± 5.58</b>	118	63.20 ± 6.09	101	<b>75.37 ± 4.42</b>	100	72.53 ± 1.25	130	<b>75.99 ± 5.26</b>	225
	AMM <sub>1</sub>	37.08 ± 4.42	3	38.53 ± 2.97	3	41.89 ± 2.07	6	41.13 ± 2.58	9	47.09 ± 9.40	10
AMM <sup>max</sup>	AMM <sub>10ran</sub>	49.12 ± 6.50	36	60.31 ± 5.57	38	73.82 ± 4.70	44	72.07 ± 3.22	54	74.73 ± 4.54	72
	AMM <sub>EMM</sub>	46.45 ± 3.30	18	46.31 ± 3.02	19	67.34 ± 13.42	19	72.41 ± 6.17	20	74.58 ± 4.63	22
	AMM <sub>MM</sub>	52.47 ± 8.88	18	58.61 ± 12.19	18	65.14 ± 21.84	19	74.90 ± 4.86	20	74.88 ± 3.75	22
	AMM <sub>G</sub>	52.47 ± 8.88	54	56.12 ± 12.25	53	74.93 ± 8.18	57	73.87 ± 4.55	60	75.43 ± 4.02	67
	AMM <sub>G,s</sub>	52.47 ± 8.88	160	54.79 ± 11.61	158	74.84 ± 8.12	167	73.87 ± 4.55	180	75.40 ± 4.05	197
	AMM <sub>nc</sub>	52.47 ± 8.88	154	49.24 ± 12.68	137	65.11 ± 21.84	137	74.89 ± 4.75	167	74.70 ± 3.71	269
SVM	alter- $\alpha$	58.39 ± 13.20	17	61.04 ± 14.43	17	69.66 ± 16.93	17	<b>76.49 ± 3.29</b>	18	75.44 ± 3.65	20
	conv- $\alpha$	50.47 ± 9.69	168	56.78 ± 10.89	164	60.41 ± 15.48	160	61.62 ± 18.81	170	73.25 ± 6.97	191
	Oracle	49.36 ± 1.68	34	49.59 ± 1.58	37	48.43 ± 2.23	40	48.85 ± 1.55	47	51.05 ± 2.72	64
	29.70 ± 2.03	6031	<b>64.15 ± 5.43</b>	6343	63.01 ± 2.59	6362	62.01 ± 3.61	6765	63.17 ± 3.62	7004	
Oracle	79.43 ± 2.88	<1	78.95 ± 3.99	<1	79.18 ± 1.70	<1	79.42 ± 2.80	<1	79.02 ± 3.62	<1	

Table 4.10: german

algorithm AUC	2 bags		4 bags		8 bags		16 bags		32 bags		
	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	
EMM	51.82 ± 12.39	<1	50.43 ± 23.03	<1	55.09 ± 19.44	<1	49.55 ± 17.47	<1	63.49 ± 18.11	<1	
MM	<b>68.75 ± 6.09</b>	<1	60.24 ± 13.54	<1	80.35 ± 9.42	<1	76.11 ± 6.66	1	83.50 ± 6.22	1	
LMM <sub>G</sub>	<b>68.75 ± 6.09</b>	3	68.04 ± 8.53	3	82.87 ± 6.16	4	<b>82.92 ± 1.28</b>	4	85.85 ± 3.84	6	
LMM <sub>G,s</sub>	<b>68.75 ± 6.09</b>	9	<b>69.04 ± 6.52</b>	12	<b>83.68 ± 5.90</b>	13	<b>82.96 ± 1.79</b>	14	<b>86.36 ± 3.94</b>	17	
LMM <sub>nc</sub>	<b>68.75 ± 6.09</b>	11	60.40 ± 14.18	12	80.24 ± 9.74	189	78.14 ± 4.98	42	84.47 ± 5.06	119	
Invcal	28.84 ± 4.96	4	70.58 ± 6.45	4	37.33 ± 10.31	4	44.96 ± 9.64	4	62.76 ± 15.05	4	
AMM <sup>min</sup>	AMM <sub>EMM</sub>	60.50 ± 30.88	<1	63.36 ± 28.50	1	72.05 ± 19.17	1	80.87 ± 15.51	1	<b>91.63 ± 6.10</b> ↑	2
	AMM <sub>MM</sub>	86.59 ± 6.14	1	80.57 ± 16.72	1	87.96 ± 4.50	2	90.04 ± 5.14	2	91.45 ± 5.70 ↑	2
	AMM <sub>G</sub>	86.59 ± 6.14	5	86.70 ± 5.45	5	87.46 ± 2.67	6	<b>91.06 ± 2.87</b>	7	<b>91.55 ± 5.93</b> ↑	9
	AMM <sub>G,s</sub>	86.59 ± 6.14	15	86.70 ± 5.45	16	88.31 ± 4.00	18	90.86 ± 2.81	21	<b>91.55 ± 5.93</b> ↑	27
	AMM <sub>nc</sub>	86.59 ± 6.14	13	78.97 ± 16.78	14	87.82 ± 4.42	21	90.48 ± 3.53	45	91.25 ± 5.77	125
	AMM <sub>1</sub>	<b>90.62 ± 5.82</b>	<1	89.19 ± 5.90	1	88.64 ± 3.21	1	90.78 ± 2.10	1	91.03 ± 5.82	1
AMM <sup>max</sup>	AMM <sub>10ran</sub>	78.38 ± 30.44	5	87.32 ± 4.71	6	89.85 ± 2.31	7	<b>91.02 ± 2.49</b>	9	90.47 ± 6.39	14
	AMM <sub>EMM</sub>	85.74 ± 13.28	3	84.60 ± 10.87	4	84.60 ± 7.84	3	89.83 ± 2.72	5	71.65 ± 18.52	6
	AMM <sub>MM</sub>	85.35 ± 11.06	4	82.43 ± 9.76	4	<b>90.49 ± 4.75</b>	4	89.92 ± 2.90	4	89.35 ± 6.98	7
	AMM <sub>G</sub>	85.35 ± 11.06	13	87.18 ± 6.56	13	<b>90.49 ± 4.75</b>	13	89.58 ± 2.79	16	88.55 ± 9.71	23
	AMM <sub>G,s</sub>	85.35 ± 11.06	39	<b>90.49 ± 5.05</b>	40	<b>90.58 ± 4.77</b>	40	89.58 ± 2.79	49	89.94 ± 6.63	67
	AMM <sub>nc</sub>	85.35 ± 11.06	20	82.73 ± 9.23	21	89.84 ± 4.24	30	90.06 ± 3.20	54	89.54 ± 6.60	140
SVM	alter- $\alpha$	72.77 ± 37.27	4	89.31 ± 3.99	3	89.68 ± 3.79	3	90.62 ± 3.18	5	87.97 ± 9.42	6
	conv- $\alpha$	89.96 ± 5.62	32	89.93 ± 5.02	31	88.03 ± 3.16	30	90.80 ± 3.61	38	89.61 ± 8.68	54
	Oracle	47.75 ± 17.58	15	59.72 ± 18.21	16	62.32 ± 12.83	20	58.49 ± 10.98	27	48.33 ± 12.77	47
	46.18 ± 43.41	1211	87.13 ± 5.30	1185	69.03 ± 23.18	1197	42.78 ± 23.51	1188	50.34 ± 15.75	1080	
Oracle	91.72 ± 3.95	<1	91.22 ± 4.09	<1	91.27 ± 2.88	<1	91.54 ± 2.76	<1	91.42 ± 5.46	<1	

Table 4.11: heart

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM		44.28 ± 12.13	<1	51.86 ± 8.01	<1	50.69 ± 6.34	<1	44.60 ± 3.91	<1	48.91 ± 11.73	<1
MM		<b>64.81 ± 8.82</b>	<1	77.74 ± 5.23	1	78.95 ± 7.36	1	86.76 ± 2.96	1	88.13 ± 4.16	2
LMM <sub>G</sub>		<b>64.81 ± 8.82</b>	5	80.80 ± 2.32	6	<b>83.46 ± 4.62</b>	5	87.12 ± 2.23	7	<b>88.24 ± 4.41</b>	7
LMM <sub>G,s</sub>		<b>64.81 ± 8.82</b>	14	<b>82.12 ± 2.50</b>	15	83.24 ± 4.84	15	<b>87.23 ± 1.57</b>	17	87.99 ± 4.58	21
LMM <sub>nc</sub>		<b>64.81 ± 8.82</b>	20	79.39 ± 2.12	22	81.18 ± 6.40	32	87.05 ± 2.48	68	<b>88.34 ± 4.32</b>	182
InvCal		35.34 ± 8.76	5	44.78 ± 15.37	5	53.28 ± 9.02	5	53.52 ± 8.51	5	54.08 ± 9.53	5
AMM <sup>min</sup>	AMM <sub>EMM</sub>	56.77 ± 6.42	2	<b>85.07 ± 5.24</b>	2	86.04 ± 5.21	2	86.81 ± 3.81	2	86.71 ± 3.54	3
	AMM <sub>MM</sub>	46.67 ± 8.53	3	84.52 ± 4.60	2	84.23 ± 6.67	2	85.92 ± 4.48	3	87.77 ± 5.56	3
	AMM <sub>G</sub>	46.67 ± 8.53	10	85.05 ± 4.11	9	85.28 ± 6.19	9	85.97 ± 3.19	11	88.85 ± 5.15	12
	AMM <sub>G,s</sub>	46.67 ± 8.53	28	84.63 ± 3.80	26	85.28 ± 6.19	27	86.01 ± 4.37	30	88.85 ± 5.15	36
	AMM <sub>nc</sub>	46.67 ± 8.53	24	<b>85.16 ± 4.39</b>	26	84.77 ± 6.45	36	85.96 ± 4.50	72	87.57 ± 5.23	174
	AMM <sub>1</sub>	51.47 ± 13.46	1	83.65 ± 3.89	2	<b>87.51 ± 4.24</b>	2	86.76 ± 4.07	2	87.83 ± 5.05	2.11
	AMM <sub>10ran</sub>	56.92 ± 22.42	10	80.39 ± 6.36	11	85.89 ± 5.52	12	87.32 ± 3.17	13	87.81 ± 6.52	15
AMM <sup>max</sup>	AMM <sub>EMM</sub>	57.99 ± 8.96	10	76.31 ± 5.29	10	82.07 ± 4.47	11	86.99 ± 7.23	11	87.08 ± 5.86	12
	AMM <sub>MM</sub>	<b>74.57 ± 18.16</b>	10	75.32 ± 4.74	10	78.65 ± 7.93	11	88.84 ± 3.10	12	<b>90.01 ± 5.50</b>	13
	AMM <sub>G</sub>	<b>74.57 ± 18.16</b>	32	78.06 ± 5.11	33	83.24 ± 6.54	35	89.98 ± 3.08 ↑	38	88.41 ± 5.94	41
	AMM <sub>G,s</sub>	<b>74.57 ± 18.16</b>	96	79.21 ± 4.58	98	83.36 ± 6.61	104	<b>90.88 ± 3.11 ↑</b>	112	88.41 ± 5.94	121
	AMM <sub>nc</sub>	<b>74.57 ± 18.16</b>	47	75.80 ± 5.14	50	80.22 ± 6.95	61	88.05 ± 2.47	99	89.19 ± 5.45	198
	AMM <sub>1</sub>	65.53 ± 17.30	10	77.29 ± 6.63	9	82.10 ± 7.95	10	85.45 ± 3.31	11	89.01 ± 7.02	12
	AMM <sub>10ran</sub>	65.05 ± 16.59	85	79.60 ± 6.56	82	78.56 ± 4.77	88	88.44 ± 3.22	94	89.37 ± 6.67	109
SVM	alter- $\alpha$	43.07 ± 6.05	22	44.58 ± 4.95	24	69.24 ± 4.99	27	67.72 ± 12.25	55	59.67 ± 7.01	49
	conv- $\alpha$	36.67 ± 7.44	1316	44.55 ± 9.58	1280	57.84 ± 5.98	1788	65.93 ± 3.90	887	47.58 ± 11.29	1287
Oracle		90.07 ± 5.04	<1	89.99 ± 4.23	<1	90.08 ± 5.50	<1	89.42 ± 6.34	<1	90.22 ± 5.17	<1

Table 4.12: ionosphere

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM		57.91 ± 22.04	<1	59.05 ± 10.46	<1	51.43 ± 17.22	<1	45.39 ± 23.81	<1	61.30 ± 17.86	<1
MM		77.45 ± 6.14	<1	78.97 ± 3.54	<1	79.85 ± 4.14	<1	82.74 ± 2.11	1	87.45 ± 3.57	1
LMM <sub>G</sub>		77.45 ± 6.14	3	78.34 ± 2.82	3	81.93 ± 3.81	3	87.52 ± 2.71	5	90.43 ± 3.20	6
LMM <sub>G,s</sub>		77.45 ± 6.14	9	78.34 ± 2.82	8	<b>83.87 ± 3.63</b>	9	<b>87.71 ± 2.56</b>	13	<b>91.06 ± 3.00</b>	14
LMM <sub>nc</sub>		77.45 ± 6.14	31	78.43 ± 2.74	31	80.02 ± 4.02	35	83.50 ± 2.46	54	88.10 ± 3.57	122
InvCal		33.74 ± 24.95	4	36.46 ± 5.27	4	72.54 ± 5.79	4	61.89 ± 6.25	4	59.91 ± 8.79	4
AMM <sup>min</sup>	AMM <sub>EMM</sub>	<b>81.07 ± 8.12</b>	2	78.56 ± 8.66	2	90.56 ± 3.44	2	92.08 ± 1.78	2	93.14 ± 2.04	3
	AMM <sub>MM</sub>	75.64 ± 5.02	2	68.54 ± 4.90	2	87.10 ± 4.16	2	<b>92.66 ± 1.99</b>	3	93.50 ± 1.93	3
	AMM <sub>G</sub>	75.64 ± 5.02	6	69.27 ± 5.69	7	87.57 ± 4.48	8	92.45 ± 1.89	10	93.59 ± 1.83	11
	AMM <sub>G,s</sub>	75.64 ± 5.02	19	69.27 ± 5.69	22	87.86 ± 4.62	23	91.04 ± 3.82	30	92.97 ± 1.58	32
	AMM <sub>nc</sub>	75.64 ± 5.02	34	68.49 ± 4.86	35	88.33 ± 5.17	39	91.26 ± 3.98	59	<b>93.70 ± 2.09</b>	127
	AMM <sub>1</sub>	74.49 ± 6.08	1	68.66 ± 4.92	1	90.60 ± 3.18	2	92.41 ± 1.58	2	92.95 ± 1.75	2
	AMM <sub>10ran</sub>	76.42 ± 4.80	12	75.75 ± 5.07	16	<b>92.59 ± 0.22</b>	18	92.15 ± 1.44	15	92.46 ± 1.79	19
AMM <sup>max</sup>	AMM <sub>EMM</sub>	76.02 ± 12.70	4	78.42 ± 14.14	5	87.87 ± 1.94	5	87.88 ± 3.29	6	90.71 ± 2.79	8
	AMM <sub>MM</sub>	75.31 ± 13.69	5	<b>87.22 ± 3.13</b>	5	87.43 ± 2.59	6	88.85 ± 2.39	6	90.29 ± 2.47	9
	AMM <sub>G</sub>	75.31 ± 13.69	15	73.91 ± 16.06	17	87.89 ± 1.97	17	87.98 ± 3.27	21	90.29 ± 2.47	28
	AMM <sub>G,s</sub>	75.31 ± 13.69	44	67.48 ± 16.70	50	87.89 ± 1.97	51	87.98 ± 3.27	63	90.18 ± 3.26	82
	AMM <sub>nc</sub>	75.31 ± 13.69	43	82.97 ± 8.05	45	87.85 ± 2.00	49	88.91 ± 2.41	70	90.29 ± 2.47	144
	AMM <sub>1</sub>	77.35 ± 13.61	4	70.14 ± 17.19	5	84.17 ± 2.66	5	89.12 ± 2.31	6	90.94 ± 3.06	8
	AMM <sub>10ran</sub>	72.39 ± 14.33	36	82.49 ± 9.32	47	87.44 ± 1.52	47	85.79 ± 4.54	50	90.87 ± 2.53	69
SVM	alter- $\alpha$	40.88 ± 5.80	21	30.17 ± 7.47	23	68.26 ± 6.40	26	58.84 ± 21.21	33	37.17 ± 17.48	48
	conv- $\alpha$	77.72 ± 6.23	3624	72.28 ± 8.88	2292	36.21 ± 8.38	2328	45.01 ± 14.91	2481	70.49 ± 5.59	2306
Oracle		93.80 ± 1.06	<1	93.83 ± 1.67	<1	93.89 ± 1.89	<1	93.83 ± 1.62	<1	94.00 ± 1.42	<1

Table 4.13: vertebral column

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM		54.32 ± 8.79	<1	45.47 ± 15.63	<1	46.88 ± 6.06	1	55.20 ± 18.03	1	53.93 ± 10.59	1
MM		94.56 ± 2.04	1	95.37 ± 2.62	2	95.65 ± 0.85	2	<b>96.33 ± 1.19</b>	2	96.74 ± 1.50	2
LMM <sub>G</sub>		94.56 ± 2.04	7	95.93 ± 2.47	8	95.87 ± 1.12	8	<b>96.41 ± 1.51</b>	9	<b>96.94 ± 1.67</b>	10
LMM <sub>G,s</sub>		94.56 ± 2.04	20	<b>96.03 ± 2.42</b>	22	<b>96.00 ± 1.18</b>	23	<b>96.38 ± 1.99</b>	25	96.81 ± 2.09	28
LMM <sub>nc</sub>		94.56 ± 2.04	28	95.83 ± 2.34	31	95.71 ± 0.92	43	96.23 ± 1.58	85	96.81 ± 1.50	234
Invcal		<b>94.85 ± 1.71</b>	4	73.10 ± 2.21	4	77.86 ± 4.92	4	26.74 ± 6.82	4	79.77 ± 6.25	4
AMM <sup>min</sup>	AMM <sub>EMM</sub>	93.67 ± 1.84	2	95.04 ± 3.01	2	<b>96.18 ± 0.78</b>	2	96.43 ± 1.31	2	96.94 ± 1.62	3
	AMM <sub>MM</sub>	93.48 ± 2.31	2	95.12 ± 2.89	3	96.10 ± 0.82	3	96.15 ± 1.31	4	<b>97.30 ± 1.58</b>	4
	AMM <sub>G</sub>	93.48 ± 2.31	10	<b>95.61 ± 1.90</b>	12	95.92 ± 1.02	11	96.41 ± 1.12	13	<b>97.36 ± 1.47</b>	15
	AMM <sub>G,s</sub>	93.48 ± 2.31	29	94.87 ± 3.02	33	95.34 ± 0.98	35	96.11 ± 1.30	39	<b>97.36 ± 1.47</b>	46
	AMM <sub>nc</sub>	93.48 ± 2.31	32	95.38 ± 2.38	35	95.81 ± 1.01	46	96.03 ± 1.48	89	<b>97.38 ± 1.45</b>	238
	AMM <sub>1</sub>	93.57 ± 1.99	2	94.32 ± 3.36	2	<b>96.25 ± 0.66</b>	2	96.17 ± 1.20	2	96.83 ± 1.42	2
	AMM <sub>10ran</sub>	<b>93.84 ± 2.23</b>	11	94.59 ± 3.56	11	95.85 ± 0.97	12	<b>96.63 ± 1.32</b>	15	96.66 ± 1.70	18
	AMM <sub>EMM</sub>	91.68 ± 0.81	11	94.97 ± 2.24	12	94.94 ± 1	13	95.83 ± 1.36	14	96.60 ± 1.31	15
	AMM <sub>MM</sub>	92.47 ± 0.38	12	93.43 ± 4.07	13	93.71 ± 1.34	14	95.40 ± 1.10	15	96.77 ± 1.31	17
	AMM <sub>G</sub>	92.47 ± 0.38	40	94.34 ± 2.65	34	94.03 ± 0.81	43	95.65 ± 1.70	48	96.45 ± 1.52	53
AMM <sup>max</sup>	AMM <sub>G,s</sub>	92.47 ± 0.38	124	94.22 ± 2.87	127	94.03 ± 0.81	132	96.01 ± 1.83	142	96.37 ± 1.39	160
	AMM <sub>nc</sub>	92.47 ± 0.38	65	94.96 ± 3.48	66	94.07 ± 0.78	78	95.14 ± 1.18	124	96.74 ± 1.31	275
	AMM <sub>1</sub>	91.60 ± 1.29	11	94.48 ± 2.14	12	94.34 ± 0.82	12	95.36 ± 1.56	13	96.54 ± 1.51	15
	AMM <sub>10ran</sub>	90.49 ± 2.02	101	94.59 ± 2.85	103	94.19 ± 0.73	104	95.73 ± 1.83	112	96.21 ± 1.67	128
	alter- $\alpha$	51.58 ± 3.27	19	62.74 ± 4.27	21	60.88 ± 3.50	25	63.01 ± 9.51	33	41.87 ± 7.12	57
	conv- $\alpha$	5.63 ± 2.03	1848	47.22 ± 4.92	1807	19.62 ± 5.91	1855	57.54 ± 11.22	1598	46.27 ± 9.48	1281
	Oracle	97.11 ± 1.31	<1	97.43 ± 2.25	<1	97.06 ± 0.87	<1	97.33 ± 1.38	<1	97.52 ± 1.49	<1

Table 4.14: vote (feature *physician-fee-freeze* was removed to make the problem more difficult)

algorithm	2 bags		4 bags		8 bags		16 bags		32 bags		
	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	AUC	time(s)	
EMM		<b>70.38 ± 20.39</b>	<1	56.72 ± 29.85	<1	55.42 ± 20.70	<1	65.82 ± 21.45	<1	46.85 ± 16.71	<1
MM		66.45 ± 5.42	1	82.41 ± 6.76	1	85.28 ± 4.80	1	90.35 ± 3.73	1	95.57 ± 2.45	1
LMM <sub>G</sub>		66.45 ± 5.42	4	89.72 ± 3.73	5	90.69 ± 5.30	5	94.09 ± 3.45	5	<b>97.74 ± 0.67</b>	6
LMM <sub>G,s</sub>		66.45 ± 4.412	13	<b>93.32 ± 2.94</b>	13	<b>92.68 ± 6.06</b>	14	<b>95.53 ± 2.40</b>	15	<b>97.69 ± 0.90</b>	19
LMM <sub>nc</sub>		66.45 ± 5.42	9	84.00 ± 5.48	11	86.30 ± 4.18	18	91.10 ± 4.52	40	96.28 ± 2.06	116
Invcal		58.96 ± 5.77	6	81.38 ± 4.59	6	55.18 ± 9.59	6	63.07 ± 12.61	6	71.01 ± 18.19	6
AMM <sup>min</sup>	AMM <sub>EMM</sub>	80.27 ± 18.08	1	90.33 ± 8.87	1	91.46 ± 10.59	1	88.97 ± 6.26	1	88.34 ± 22.79	2
	AMM <sub>MM</sub>	61.84 ± 9.20	2	85.56 ± 7.20	1	88.70 ± 8.31	2	93.78 ± 9.12	2	98.66 ± 1.11	2
	AMM <sub>G</sub>	61.84 ± 9.20	6	93.06 ± 7.88	7	93.42 ± 8.24	7	96.09 ± 8.18	7	<b>99.33 ± 1.01</b>	9
	AMM <sub>G,s</sub>	61.84 ± 9.20	17	94.87 ± 5.68	18	93.00 ± 8.95	20	96.09 ± 8.18	21	<b>99.33 ± 1.01</b>	27
	AMM <sub>nc</sub>	61.84 ± 9.20	10	87.03 ± 3.93	13	88.23 ± 7.90	20	97.49 ± 5.06	43	<b>99.33 ± 1.01</b>	119
	AMM <sub>1</sub>	82.21 ± 11.39	<1	94.12 ± 6.34	1	99.60 ± 0.60	1	96.03 ± 7.57	1	97.03 ± 3.66	1
	AMM <sub>10ran</sub>	58.75 ± 31.30	4	<b>99.47 ± 0.68</b>	5	99.52 ± 0.45	6	<b>99.59 ± 0.54</b>	7	98.95 ± 1.66	10
	AMM <sub>EMM</sub>	74.23 ± 32.62	3	85.52 ± 17.48	4	99.67 ± 0.74	5	98.09 ± 3.09	6	92.00 ± 11.55	7
	AMM <sub>MM</sub>	88.23 ± 18.56	5	97.60 ± 2.40	4	87.42 ± 27.76	6	99.42 ± 0.79	7	98.61 ± 1.69	8
	AMM <sub>G</sub>	88.23 ± 18.56	15	88.41 ± 20	15	<b>100.00 ± 0.00</b> ↑	19	<b>99.63 ± 0.66</b>	20	98.61 ± 1.69	25
AMM <sup>max</sup>	AMM <sub>G,s</sub>	88.23 ± 18.56	44	79.11 ± 23.90	44	<b>100.00 ± 0.00</b> ↑	56	<b>99.63 ± 0.66</b>	59	98.61 ± 1.69	75
	AMM <sub>nc</sub>	88.23 ± 18.56	19	85.44 ± 19.04	21	86.17 ± 27.19	32	99.36 ± 0.74	56	98.61 ± 1.69	135
	AMM <sub>1</sub>	75.24 ± 21.10	3	80.45 ± 10.01	4	91.83 ± 14.63	5	91.79 ± 9.05	5	88.01 ± 9.78	7
	AMM <sub>10ran</sub>	<b>97.54 ± 1.55</b>	30	96.80 ± 3.94	32	99.46 ± 0.82	41	99.21 ± 0.79	47	98.54 ± 1.66	58
	alter- $\alpha$	52.68 ± 2.54	14	36.53 ± 10.97	16	65.54 ± 2.26	19	29.15 ± 9.60	32	86.22 ± 11.93	44
	conv- $\alpha$	54.31 ± 4.63	831	70.23 ± 6.58	794	52.88 ± 13.86	840	55.60 ± 11.29	659	11.58 ± 7.84	495
	Oracle	99.69 ± 0.52	<1	99.80 ± 0.44	<1	99.60 ± 0.43	<1	99.80 ± 0.44	<1	99.78 ± 0.33	<1

Table 4.15: wine

## 4.13 References

To the best of our knowledge, the setting was originally introduced by Kuck and de Freitas [2005], inspired by previous work on multiple instance learning with threshold functions for Computer Vision [Kück et al., 2004]. The seminal work of Kuck and de Freitas [2005] proposes a hierarchical graphical model that generates labels consistent with the proportions. Training is performed in a fully Bayesian fashion by MCMC sampling. Some scaling issues on following this approach are highlighted by Quadrianto et al. [2009]. A similar but simpler graphical model is learned via EM algorithm in Wager et al. [2015].

Somehow independently, Chen et al. [2006] and its follower Musicant et al. [2007] offer a variety of standard machine learning methods such as Bayesian networks, ensemble trees, KNNs, SVMs and neural networks designed to generate self-consistent labels. Along this line, Hernández and Inza [2011]; Hernández-González et al. [2013] rely respectively on naive Bayes and structure learning of Bayesian networks with missing data. Other methods jointly exploit the proportions and geometrical assumptions on the data, by mixing classification with clustering [Chen et al., 2009; Stolpe and Morik, 2011; Cui et al., 2016]. An adaptation of latent discriminant analysis to LLP is given in Pérez-Ortiz et al. [2016].

It seems that the important contribution of Quadrianto et al. [2009] was initially conceived with no connection with previous work until N. De Freitas and A. Smola discussed the topic “while walking along the beach in Australia”<sup>2</sup> at a Machine Learning Summer School in 2008. This piece of work is central to the Thesis, since it is the first published instance of Loss Factorization for weakly supervised learning, with the related estimation of the mean operator. The paper also shows empirically that the Mean Map algorithm outperforms multiple strong baselines as well as the approach of Kuck and de Freitas [2005].

There are several SVM-based algorithms for LLP. Rüping [2010] propose to optimize an inverted calibration function [Platt, 1999] by margin maximization, with the objective of matching bag-wise average prediction with the given proportions. Instead, Yu et al. [2013] follow a more traditional formulation of SVMs, by augmenting the search space with the latent labels. Two alternative approximations are proposed: a convex relaxation, with better performance but hard to scale, and a more practical alternating optimization similar to AMM. Yu et al. [2013] show how their implementation outperform most of other known methods at date. The success of Yu et al. [2013] have inspired several SVM-based solutions [Wang et al., 2015; Ni et al., 2015; Qi et al., 2016; Chen et al., 2016].

Deep neural networks have recently been used for LLP. Fan et al. [2014] solve a density estimation problem with a restricted Boltzmann machines to learn a generative model. Li and Taylor [2015] formulate an alternating algorithm similar to Yu et al. [2013] and our AMM, but utilizing a convolutional neural network as the core model. A rather original insight is contained in the work of Kotzias et al. [2015], with a specific application on sentiment analysis of *sentences* instead of documents. In a

---

<sup>2</sup>Excerpt from a workshop talk by Nando de Freitas at NIPS’15, Montreal.



---

first phase, a convolutional network is trained for predicting document sentiments, which is the level of aggregated supervision that is available; in a second phase, a sentence embedding is obtained from the network by removing the top layer; finally, the architecture is retrained with a loss penalizing mistakes at document level predictions with an additional manifold regularization. The sentence embedding step is essential for constructing a suitable feature representation for the problem and, importantly, it is learned from label proportions only.

Prior work on theory of LLP belongs to two main categories. The first is about estimation guarantee of either the mean operator [Quadrianto et al., 2009] in RKHS or the resulting classifier [Altun and Smola, 2006; Quadrianto et al., 2009]. Those have already been discussed above. The second regards generalization bounds on predicting the bag proportions. Yu et al. [2014b] introduce the framework of *empirical proportion risk estimation*, showing that it is possible to bound generalization error of predicting bag proportions; in turn, the result is used to bound the more interesting error on label prediction by the use of distributional assumptions on the bag assignment process.

LLP has been studied in conjunction with *domain adaptation* in Ardehaly and Cullotta [2016]. A model is learned from label proportions on the source domain and used to generate label proportions for a target domain; a self-training process is repeated to adapt the model from source to target.

Finally, extension of the LLP setting have been recently proposed, where the supervision is weaker but still represented by proportional quantities: *learning from positive-unlabeled proportions* [Hernández-González et al., 2015, 2016], *ballpark learning* [Hope and Shahaf, 2016], *i.e.* learning from bounds on the label proportions, *learning from histograms and order statistics* [Bhowmik et al., 2015]. The last setting is also linked with the broad topic of *data imputation* with the help of auxiliary aggregated information [Park and Ghosh, 2012, 2013] and with the *ecological inference*, presented next.

#### 4.13.1 Ecological inference

We conclude this literature review with a mention to a related yet fundamentally different learning scenario, which originated outside the Machine Learning community. The problem deals with recovering information from aggregate data and thus share the intent of learning from label proportions. Although, in contrast with the topic of this Chapter, the objective is not to classify individual observations, but “disaggregate” variables from an higher level group representation into constituent parts.

The iconic application is inferring electorate behavior: given turnout results for several parties and proportions of some population strata, *e.g.* percentages of ethnic groups, for many geographical regions such as counties, the aim is to recover contingency tables for parties  $\times$  groups for all those counties. In the language of probability the problem is isomorphic to the following: given two random variables and their respective marginal distributions — conditioned to another variable —, reconstruct their conditional joint distribution. Let us simplify the problem in that the feature

vector  $x$  and the label  $y$  are both scalar binary variables. As in LLP, the the problem has another dimension, the one of bags  $j \in [n]$  which often stands for a geographical descriptor. The ecological inference problem is then summarized as: assuming to know the marginals — conditioned to  $j$  —  $p(x|j)$  and  $p(y|j)$ , can we infer the joint distribution  $p(x, y|j)$ ?

The ecological inference arises in a diversity of applied fields such as Econometrics [Cross and Manski, 2002; Cho and Manski, 2008], Sociology and Political Science [King, 1997; King et al., 2004] and Epidemiology [Wakefield and Shaddick, 2006], with a long history [Robinson, 1950]; interestingly, the Empirical Software Engineering community has also explored the idea [Posnett et al., 2011].

As with weak supervision, the problem is fundamentally under-determined and any solution can only provide either loose deterministic bounds [Duncan and Davis, 1953; Cross and Manski, 2002; Cho and Manski, 2008] or needs to enforce additional assumptions and prior knowledge on the data domain [King, 1997]. A decade ago, the problem has witnessed a period of renaissance along with the publication of a diversity of methods from the second family, mostly inspired by several models and related distributional assumptions collected in the book of King et al. [2004]. In contrast, Judge et al. [2004] follow the road of a minimal subset of assumptions and frame the inference as an optimization problem. The approach favors one solution according to some information-theoretic solution, *e.g.* the Cressie-Read power divergence, intended as an entropic measure of the joint distribution. Along this last line of work, we have co-authored a new solution based on *optimal transport* [Villani, 2008] between marginal distributions as presented in Muzellec et al. [2017].

On the one hand, LLP is a strictly more ambitious problem than ecological inference. In fact, in learning from label proportions — the marginals for  $y$  — the prediction happens at the level of each individual. In contrast, the goal of ecological inference is to “disaggregate” data into another, finer-grain level of representation, the one of the bags. From this point of view, a solution for LLP can be used to infer the unknown joint distributions. On the other hand, LLP is assuming a much richer representation of the data as a starting point. Individual features vectors, *e.g.* the individual demographics of voters, are not known and not necessary for solving ecological inference. Yet it has been shown that such information can be beneficially integrated into learning methods for ecological inference [Flaxman et al., 2015].

---

# Learning with noisy labels I: theory for linear models

---

The Chapter specializes our approach to the problem of learning a classifier under asymmetric label noise. The above theory and algorithms are reviewed to adapt to this particular case of weak supervision. A novel estimator for the label sufficient statistic is defined by its property of unbiasedness. We also characterize the whole family of linear-odd losses by an approximate criterion of label noise robustness. Simple experiments with simulated noise corroborate both theoretical and algorithmic statements.

## 5.1 Motivation

Large datasets used in training modern machine learning models are often affected by label noise. The problem is pervasive for a simple reason: manual expert-labeling of each instance at a large scale is not feasible and so researchers and practitioners often resort to cheap but imperfect surrogates. Two such popular surrogates are crowdsourcing using non-expert labelers and (especially for images) the use of a search engine to query instances by a keyword, where it is assumed that the keyword is a valid label for what is collected from the query [Fergus et al., 2010; Schroff et al., 2011; Divvala et al., 2014; Krause et al., 2016]. Both approaches offer the possibility to scale the acquisition of training labels. However, they invariably result in the introduction of label noise, which may adversely affect model training. On top of those issues, label noisy is intrinsic and virtually impossible to eliminate from many data domains, because of inherent ambiguity and task misspecification [Misra et al., 2016].

In the Thesis we adopt the nomenclature of Table 5.1. We focus on label noise that is *asymmetric*, *i.e.* we work in the class-conditional noise setting of Natarajan et al. [2013]. This is clearly a simplification of real-world information corruption, which may also be *instance dependent*. Little is known about learning under such noise, with few exceptions [Xiao et al., 2015; Ghosh et al., 2015; Menon et al., 2016].

name	assumption	example of recent study
symmetric	$p(\tilde{y} y, \mathbf{x}) = p(\tilde{y})$	van Rooyen et al. [2015]
asymmetric	$p(\tilde{y} y, \mathbf{x}) = p(\tilde{y} y)$	Natarajan et al. [2013]
instance dependent	no assumption	Menon et al. [2016]

Table 5.1: Label noise types considered in literature.

## 5.2 Learning setting

In learning with asymmetric noisy labels (ALN),  $\tilde{\mathcal{S}}$  is a set of examples drawn from a distribution  $\tilde{\mathcal{D}}$ , which samples from  $\mathcal{D}$  and flips labels at random. Each example  $(\mathbf{x}_i, \tilde{y}_i)$  is  $(\mathbf{x}_i, -y_i)$  with probability at most  $1/2$  or it is  $(\mathbf{x}_i, y_i)$  otherwise. The *noise rates* are label dependent by  $(p_+, p_-) \in [0, 1/2)^2$  respectively for positive and negative examples. The marginal feature distribution is not affected by noise. In sum, examples are drawn from the distribution:

$$p(\mathbf{x}, \tilde{y}) = p_+ \cdot p(y = 1|\mathbf{x})p(\mathbf{x}) + p_- \cdot p(y = -1|\mathbf{x})p(\mathbf{x}), \quad (5.1)$$

with  $p_+ = p(\tilde{y} = -1|y = 1)$  and  $p_- = p(\tilde{y} = 1|y = -1)$ .

## 5.3 Estimating the sufficient statistic and $\mu$ SGD

The schema of the solution for this learning problem follows closely what we did in the Chapter 4 for LLP and implements our two-step procedure: estimate the sufficient statistic and solve ERM. The estimation step turns out to be much simpler in this scenario. Although, notice that we work under the assumption of knowing the noise rates  $(p_+, p_-)$ . Towards a more realistic treatment, Chapter 6 will deal with the estimation of those quantities as well.

We construct an estimator for the mean operator that is unbiased with respect to the noisy distribution. The result builds on Natarajan et al. [2013, Lemma 1] that provides a recipe for unbiased estimators of *losses*. Instead of estimating the whole  $\ell$ , as consequence of the Factorization Theorem 18 we act on the sufficient statistic.

**Theorem 52.** *The estimator defined as:*

$$\hat{\boldsymbol{\mu}}_{\mathcal{S}} \doteq \mathbb{E}_{\mathcal{S}} \left[ \frac{y - (p_- - p_+)}{1 - p_- - p_+} \mathbf{x} \right] \quad (5.2)$$

*is unbiased, that is, its expectation over the noise distribution  $\tilde{\mathcal{D}}$  is the population mean operator:*

$$\mathbb{E}_{\tilde{\mathcal{D}}} [\hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}}] = \boldsymbol{\mu}_{\mathcal{D}}. \quad (5.3)$$

*The corresponding risk estimator  $\hat{R}_{\mathcal{S}, \ell}(\boldsymbol{\theta}) \doteq \frac{1}{2} R_{\mathcal{S}_{2x}, \ell}(\boldsymbol{\theta}) + a \langle \boldsymbol{\theta}, \hat{\boldsymbol{\mu}}_{\mathcal{S}} \rangle$  is also unbiased.*

Proof in 5.7.1. We have thus obtained a good candidate as input for any algorithm

**Algorithm 8:**  $\mu$ SGD for asymmetric noisy label

---

**Input:**  $\tilde{\mathcal{S}}, \ell$  is  $a$ -LOL,  $(p_+, p_-)$ ,  $\lambda > 0$ ,  $T > 0$   
 $\mathcal{S}_{2x} \leftarrow \{(\mathbf{x}_i, \sigma), i \in [m], \forall \sigma \in \mathcal{Y}\}$   
 $\hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}} \leftarrow \mathbb{E}_{\tilde{\mathcal{S}}} \left[ \frac{y_-(p_- - p_+)}{1 - p_- - p_+} \mathbf{x} \right]$   
 $\hat{\boldsymbol{\theta}} \leftarrow \mu\text{SGD}(\mathcal{S}_{2x}, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}}, \lambda, T)$   
**Output:**  $\hat{\boldsymbol{\theta}}$

---

**Algorithm 9:**  $\mu$ SGD with  $L_2$  regularization

---

**Input:**  $\mathcal{S}_{2x}, \boldsymbol{\mu}$ ,  $\ell$  is  $a$ -LOL,  $\lambda > 0$ ,  $T > 0$   
 $\boldsymbol{\theta}^0 \leftarrow \mathbf{0}$   
For any  $t = 1, \dots, T$ :  
Pick  $i \in [|\mathcal{S}_{2x}|]$  uniformly at random  
 $\eta \leftarrow 1/(\lambda t)$   
Pick any  $\mathbf{v} \in \partial \ell(y_i \langle \boldsymbol{\theta}^t, \mathbf{x}_i \rangle)$   
 $\boldsymbol{\theta}^{t+1} \leftarrow (1 - \eta \lambda) \boldsymbol{\theta}^t - \eta(\mathbf{v} + a\boldsymbol{\mu}/2)$   
 $\boldsymbol{\theta}^{t+1} \leftarrow \min \left\{ \boldsymbol{\theta}^{t+1}, \boldsymbol{\theta}^{t+1} / (\sqrt{\lambda} \cdot \|\boldsymbol{\theta}^{t+1}\|_2) \right\}$   
**Output:**  $\boldsymbol{\theta}^T$

---

implementing our two-step approach, like  $\mu$ SGD (Algorithm 9). But there is more. On one hand, the estimators of Natarajan et al. [2013] may not be convex even when  $\ell$  is so, but this is never the case with LOLs. In fact,  $\ell(x) - \ell(-x) = 2ax$  may be seen as an alternative sufficient condition to Natarajan et al. [2013, Lemma 4] for convexity, without asking for differentiability.

The estimator of Equation 5.2 is all we need for presenting a complete version of the weakly supervised version of SGD given in Algorithm 8. We also restate Algorithm 9 for the sake of readability. This is another example of implementation of the Meta Algorithm 1, this time for ALN.

## 5.4 Generalization bounds

We now prove that *any* algorithm minimizing LOLs that uses the estimator in Equation 5.2 has a non-trivial generalization bound. We further assume that  $\ell$  is Lipschitz.

**Theorem 53.** Consider the setting of Theorem 23, except that here  $\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} \hat{R}_{\tilde{\mathcal{S}}, \ell}(\boldsymbol{\theta})$ . Then for any  $\delta > 0$ , with probability at least  $1 - \delta$ :

$$R_{\mathcal{D}, \ell}(\hat{\boldsymbol{\theta}}) - \inf_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{D}, \ell}(\boldsymbol{\theta}) \leq \frac{\sqrt{2} + 1}{2} \frac{XHL}{\sqrt{m}} + c(X, H) \sqrt{\frac{1}{m} \log \left( \frac{2}{\delta} \right)} + \frac{2|a|XH}{1 - p_- - p_+} \sqrt{\frac{d}{m} \log \left( \frac{2d}{\delta} \right)}. \quad (5.4)$$

Proof in 5.7.2. The complexity term is tighter than prior work. Natarajan et al. [2013, Theorem 3] prove a factor of  $2L/(1 - p_- - p_+)$  that may even be unbounded due to noise, while our estimate shows a constant of about  $1.2 < 2$  and it is noise free. In fact, LOLS are such that noise affects only the linear component of the bound, as a direct effect of Factorization. Although we are not aware of any other such results, this is intuitive: the Rademacher complexity  $\mathcal{R}(\mathcal{H} \circ \mathcal{S})$  is computed regardless of sample labels and therefore is unchanged by label noise<sup>1</sup>. Furthermore, depending on the loss, the effect of (limited) noise on generalization may also be negligible since  $c(X, H)$  could be very large for losses like strongly convex.

The next Theorem comes in pair with Theorem 53: it holds regardless of algorithm and (linear-odd) loss of choice. In particular, we demonstrate that every learner enjoys a distribution-dependent property of robustness against ALN. No estimate of  $\mu$  is involved and hence the Theorem also applies to any naïve supervised learner *on*  $\tilde{\mathcal{S}}$  via linear-odd losses. We first recall a strong notion of robustness from Manwani and Sastry [2013]; van Rooyen et al. [2015].

**Definition 54.** Let  $(\theta^*, \tilde{\theta}^*)$  respectively be the minimizers of  $(R_{\mathcal{D},\ell}(\theta), R_{\tilde{\mathcal{D}},\ell}(\theta))$  in  $\mathcal{H}$ .  $\ell$  is said ALN robust if for any  $\mathcal{D}, \tilde{\mathcal{D}}$ :

$$R_{\mathcal{D},\ell}(\theta^*) = R_{\mathcal{D},\ell}(\tilde{\theta}^*) . \quad (5.5)$$

The influential work of Long and Servedio [2010], together with van Rooyen et al. [2015], settles the question of whether commonly used loss function can be robust to label noise, in particular in the milder *symmetric* noise scenario. Rather negatively, the literature concludes that no *convex potentials*<sup>2</sup> is immune to such noise. Instead, the already cited unhinged loss  $\ell(x) = 1 - x$ , a linear loss, satisfies Equation 5.5 [van Rooyen et al., 2015].

To study the LOLS under ASL, we take a more pragmatic point of view and extend Definition 5.5 by the weaker formulation of  $\epsilon$ -robustness.

**Definition 55.** In the context of Definition 54,  $\ell$  is said  $\epsilon$ -ALN robust if for any  $\mathcal{D}, \tilde{\mathcal{D}}$ :

$$R_{\tilde{\mathcal{D}},\ell}(\theta^*) - R_{\tilde{\mathcal{D}},\ell}(\tilde{\theta}^*) \leq \epsilon . \quad (5.6)$$

The distance of the two minimizers is measured by  $\ell$ -risk under expected label noise. 0-ALN robust losses are also ALN robust: in fact if  $R_{\tilde{\mathcal{D}},\ell}(\theta^*) = R_{\tilde{\mathcal{D}},\ell}(\tilde{\theta}^*)$  then  $\theta^* \in \operatorname{argmin}_{\theta} R_{\tilde{\mathcal{D}},\ell}(\theta)$ . And hence if  $R_{\tilde{\mathcal{D}},\ell}(\theta)$  has a unique global minimum, that is  $\theta^*$ . More generally we have the following result.

**Theorem 56.** Assume  $\{\theta \in \mathcal{H} : \|\theta\|_2 \leq H\}$ . Then every  $a$ -LOL is  $\epsilon$ -ALN. That is:

$$R_{\tilde{\mathcal{D}},\ell}(\theta^*) - R_{\tilde{\mathcal{D}},\ell}(\tilde{\theta}^*) \leq 4|a|H \max\{p_+, p_-\} \|\mu_{\mathcal{D}}\|_2 \quad (5.7)$$

<sup>1</sup>Strictly speaking, this is true only considering the upper bound  $L \cdot \mathcal{R}(\mathcal{H} \circ \mathcal{S})$  for Lipschitz losses, without dependency on  $\ell$ . In contrast,  $\mathcal{R}(\ell \circ \mathcal{H} \circ \mathcal{S})$  would depend on the sample labels.

<sup>2</sup>Namely, losses  $\ell \in C^1$ , convex, such that  $\ell(0) < 0$  and  $\ell(x) \rightarrow 0$  for  $x \rightarrow \infty$ . Many convex potentials are LOLS but not all. An example is  $e^{-x}$ .

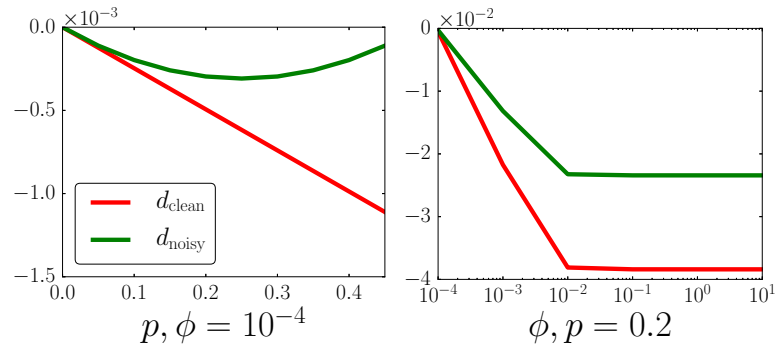


Figure 5.1: Behavior of Theorem 56 on synthetic data. Definition of the axes within the text.

Proof in 5.7.3. The reader might question the bound for the fact that the quantity on the right-hand side can change by rescaling  $\mu_{\mathcal{D}}$  by  $X$ , *i.e.* the max  $L_2$  norm of observations in the space  $\mathcal{X}$ . Although, such transformation would affect  $\ell$ -risks on the left-hand side as well, balancing the effect. With this in mind, we formulate the result without making explicit dependency on  $X$ . Unlike Theorem 53, this bound holds *in expectation* over the noisy risk  $R_{\tilde{\mathcal{D}},\ell}$ . Its shape depends on the population mean operator, a *distribution-dependent* quantity, which empirical approximation is readily available by Equation 5.2. There are two immediate corollaries.

**Corollary 57.** *If  $\|\mu_{\mathcal{D}}\|_2 = 0$  then every LOL is ALN for any  $\tilde{\mathcal{D}}$ .*

**Corollary 58.** *Suppose additionally that  $\ell$  is once differentiable and  $\gamma$ -strongly convex. Then:*

$$\|\theta^* - \tilde{\theta}^*\|_2^2 \leq \frac{2\epsilon}{\gamma}. \quad (5.8)$$

Proofs in 5.7.4. When  $\|\mu_{\mathcal{D}}\|_2 = 0$ , we obtain optimality for all LOLs. Notice that this is not necessarily a condition detrimental to learning. In fact, as remarked in Lemma 28,  $\|\mu_{\mathcal{D}}\|_2 = 0$  does not imply  $\text{Cov}_{\mathcal{D}}[x, y] = 0$ . The second corollary goes further, limiting the minimizers' distance when losses are differentiable and strongly convex. But even more generally, under some compactness assumptions on the domain of  $\ell$ , Theorem 56 tells us more: in the case  $R_{\tilde{\mathcal{D}},\ell}(\theta)$  has a unique global minimizer, the smaller  $\|\mu_{\mathcal{D}}\|_2$ , the closer the minimizer *on noisy data*  $\tilde{\theta}^*$  will be to the minimizer *on clean data*  $\theta^*$ . Therefore, assuming an efficient algorithm that computes a model not far from the global optimum  $\tilde{\theta}^*$ , that will be not far from  $\theta^*$  either. This is true in spite of the presence of local minima and/or saddle points.

Let us compare the significance of Theorem 56 with what is known in literature. Long and Servedio [2010] prove that no convex loss is noise robust, that is, 0-ALN robust. This is not a contradiction. To show the negative statement, the authors craft

a case of  $\mathcal{D}$  breaking all such losses. In fact that choice of  $\mathcal{D}$  does not meet optimality in our bound, because:

$$\|\boldsymbol{\mu}_{\mathcal{D}}\|_2 = \frac{1}{4}(18\gamma^2 + 6\gamma + 1) > 0, \quad (5.9)$$

with  $\gamma \in (0, 1/6)$ . The worst-case result of Long and Servedio [2010], like any extreme-case argument, should be handled with care. It does not give the big picture for all data we may encounter in the real world, but only the most pessimistic. We present such a global view which appears better than expected: learning from noisy data does not necessarily reduce convex losses to a singleton [van Rooyen et al., 2015] but depends on the data at hand — via the mean operator — for a broad set of losses.

We also compare our  $\epsilon$ -ALN robustness with a robustness bound proposed by Ghosh et al. [2015]. In the same setting as Definition 54:

$$R_{\mathcal{D},\ell}(\tilde{\boldsymbol{\theta}}^*) \leq \frac{1}{1 - 2 \max(p_-, p_+)} R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*). \quad (5.10)$$

While relating the (non-noisy)  $\ell$ -risks, this bound is not data dependent and may not be informative for high noise rates.

Finally, Manwani and Sastry [2013] give a sufficient condition for loss function to be ALN robust (Definition 54):

$$\ell(x) + \ell(-x) = \text{const} \Rightarrow \text{0-ALN} \quad (5.11)$$

Quite surprisingly, Loss Factorization also marries two opposite views in one formula:

$$\ell(x) = \frac{1}{2} \left( \underbrace{\ell(x) + \ell(-x)}_{=\text{const} \Rightarrow \text{0-ALN}} + \underbrace{\ell(x) - \ell(-x)}_{=ax \Rightarrow \epsilon\text{-ALN}} \right) \quad (5.12)$$

As a by product, we confirm the peculiar role of unhinged loss, being the only function — modulo linear transformation — that satisfies both conditions.

## 5.5 Experiments

We begin by building a toy planar dataset to probe the behavior of Theorem 56. It is made of four observations:  $(0, 1)$  and  $(\phi/3, 1/3)$  three times, with the first example the only negative, repeated 5 times. We consider this the distribution  $\mathcal{D}$  so as to calculate  $\|\boldsymbol{\mu}_{\mathcal{D}}\|_2 = \phi^2/4$ . We fix  $p_+, p_- = 0.2 = p$  and control  $\phi$  to measure the discrepancy  $d_{\text{noisy}} \doteq R_{\tilde{\mathcal{D}},\ell}(\boldsymbol{\theta}^*) - R_{\tilde{\mathcal{D}},\ell}(\tilde{\boldsymbol{\theta}}^*)$ , its counterpart  $d_{\text{clean}}$  computed on  $\mathcal{D}$ , and how the two minimizers “differ in sign” by  $d_{\text{models}} \doteq \langle \boldsymbol{\theta}^*, \tilde{\boldsymbol{\theta}}^* \rangle / \|\boldsymbol{\theta}^*\|_2 \|\tilde{\boldsymbol{\theta}}^*\|_2$ . The same simulation is run varying the noise rates with constant  $\phi = 10^{-4}$ . We learn with square loss and  $\lambda = 10^{-6}$ . Results are in Figure 5.1. The closer the parameters to 0, the smaller  $d_{\text{clean}} - d_{\text{noisy}}$ , while they are equal when each parameter is individually 0.  $d_{\text{models}}$  is negligible, which is good news for the 01-risk on sight.



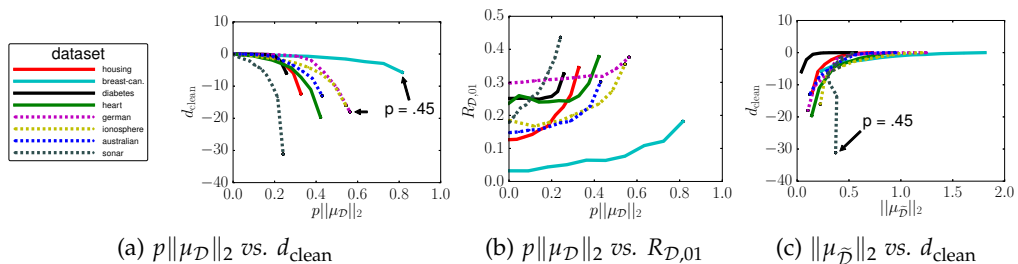


Figure 5.2: How mean operator and noise rate condition risks.  $d_{\text{clean}} \doteq R_{\mathcal{D},\ell}(\theta^*) - R_{\mathcal{D},\ell}(\tilde{\theta}^*)$ .

The next results are obtained by Algorithm 8 on UCI data. We learn with logistic loss, without model intercept and set  $\lambda = 10^{-6}$  and  $T = 4 \cdot 2m$  (4 epochs). We measure  $d_{\text{clean}}$  and  $R_{\mathcal{D},01}$ , injecting symmetric label noise  $p \in [0, 0.45]$  and averaging over 25 runs. Again, we consider *the whole distribution* so as to play with the ingredients of Theorem 56. Figure 5.2a confirms how the combined effect of  $p\|\mu_{\mathcal{D}}\|_2$  can explain most variation of  $d_{\text{clean}}$ . While this is not strictly implied by Theorem 56 that only involves  $d_{\text{noisy}}$ , the observed behavior is expected. A similar picture is given in Figure 5.2b which displays the true risk  $R_{\mathcal{D},01}$  computed on the minimizer  $\tilde{\theta}^*$  of  $\tilde{\mathcal{S}}$ . From 5.2a and 5.2b we also deduce that large  $\|\mu_{\mathcal{D}}\|_2$  is a good proxy for generalization with linear classifiers; see the relative difference between points at the same level of noise. Finally, we have also monitored  $\mu_{\tilde{\mathcal{D}}}$ . Figure 5.2c shows that large  $\|\mu_{\tilde{\mathcal{D}}}\|_2$  indicates small  $d_{\text{clean}}$  as well. This remark can be useful in practice, when the norm can be estimated from  $\tilde{\mathcal{S}}$ , as opposite to  $p$  and  $\mu_{\mathcal{D}}$ , and used to anticipate the effect of noise on the task at hand.

We conclude with a systematic study of hold-out error of  $\mu\text{SGD}$ . The same datasets are now split in 1/5 test and 4/5 training sets once at random. In contrast with the previous experimental setting we perform cross-validation of  $\lambda \in 10^{\{-3, \dots, +3\}}$  on 5-folds in the training set. We compare with vanilla SGD run on corrupted sample  $\tilde{\mathcal{S}}$  and measure the gain from estimating  $\hat{\mu}_{\tilde{\mathcal{S}}}$ . The other parameters  $l, T, \lambda$  are the same for both algorithms; the learning rate  $\eta$  is untouched from Shalev-Shwartz et al. [2011] and not tuned for  $\mu\text{SGD}$ . The only differences are in input and gradient update. Table 5.2 reports test error for SGD and its difference with  $\mu\text{SGD}$ , for a range of values of  $(p_-, p_+)$ .  $\mu\text{SGD}$  beats systematically SGD with large noise rates, and yet performs in pair under low or null noise. Interestingly, in the presence of very intense noise  $p_+ \approx .5$ ,  $\mu\text{SGD}$  still learns sensible models and improves up to 55% relatively to the error of SGD, which is often doomed to random guess.

## 5.6 Discussion

In this Chapter we have shown that both theoretical analysis and algorithms are easily adaptable to the case of noisy labels, granted that we have access of a good estimator of the label sufficient statistic. In this context, we have also derived a

$(p-, p+) \rightarrow$ <i>dataset</i>	(.00, .00)		(.20, .00)		(.20, .10)		(.20, .20)		(.20, .30)		(.20, .40)		(.20, .49)	
	SGD	$\mu$ SGD	SGD	$\mu$ SGD	SGD	$\mu$ SGD	SGD	$\mu$ SGD	SGD	$\mu$ SGD	SGD	$\mu$ SGD	SGD	$\mu$ SGD
australian	0.13	+ .01	0.15	- .01	0.14	$\pm$ .00	0.14	+ .01	0.16	- .01	0.26	- .09	0.45	- .25
breast-can.	0.02	+ .01	0.03	$\pm$ .00	0.03	$\pm$ .00	0.03	$\pm$ .00	0.05	- .01	0.11	- .06	0.17	- .08
diabetes	0.28	- .03	0.29	- .03	0.29	- .03	0.27	- .02	0.28	- .02	0.39	- .13	0.59	- .22
german	0.27	- .02	0.26	$\pm$ .00	0.27	- .02	0.29	- .02	0.31	- .01	0.31	$\pm$ .00	0.31	$\pm$ .00
heart	0.15	+ .01	0.17	- .01	0.16	$\pm$ .00	0.17	$\pm$ .00	0.18	- .01	0.26	- .08	0.35	- .15
housing	0.17	- .03	0.23	- .05	0.22	- .04	0.20	- .02	0.22	- .03	0.34	- .12	0.41	- .13
ionosphere	0.14	+ .05	0.19	- .05	0.20	- .05	0.20	- .03	0.21	- .03	0.35	- .13	0.54	- .29
sonar	0.27	$\pm$ .00	0.29	+ .02	0.29	+ .01	0.34	- .04	0.36	- .03	0.43	- .10	0.45	- .05

Table 5.2: Test error for SGD and  $\mu$ SGD over 25 trials of artificially corrupted datasets.

---

general property of label noisy robustness for linear-odd losses. The next Chapter discusses some practical open issues, in particular the restriction to linear models and the strong hypothesis of known noise rates.

## 5.7 Appendix: proofs

### 5.7.1 Proof of Theorem 52

Natarajan et al. [2013, Lemma 1] provide an unbiased estimator for a loss  $\ell(x)$  computed on  $x$  of the form:

$$\hat{\ell}(y\langle\boldsymbol{\theta}, \mathbf{x}_i\rangle) \doteq \frac{(1 - p_{-y}) \cdot \ell(\langle\boldsymbol{\theta}, \mathbf{x}_i\rangle) + p_y \cdot \ell(-\langle\boldsymbol{\theta}, \mathbf{x}_i\rangle)}{1 - p_- - p_+}. \quad (5.13)$$

We apply the Lemma for estimating the mean operator instead. We are allowed to do so by the very result of the Factorization Theorem, since the noise corruption has only an effect on the linear-odd term of the loss. The estimator of the sufficient statistic of a single example  $y\mathbf{x}$  is:

$$\hat{\mathbf{z}} \doteq \frac{1 - p_{-y} + p_y}{1 - p_- - p_+} y\mathbf{x} \quad (5.14)$$

$$= \frac{1 - (p_- - p_+)y}{1 - p_- - p_+} y\mathbf{x} \quad (5.15)$$

$$= \frac{y - (p_- - p_+)}{1 - p_- - p_+} \mathbf{x}, \quad (5.16)$$

and its average, *i.e.* the mean operator estimator, is:

$$\hat{\boldsymbol{\mu}}_S \doteq \mathbb{E}_S \left[ \frac{y - (p_- + p_+)}{1 - p_- - p_+} \mathbf{x} \right], \quad (5.17)$$

such that in expectation over the noisy distribution it holds  $\mathbb{E}_{\tilde{\mathcal{D}}}[\hat{\mathbf{z}}] = \boldsymbol{\mu}_{\mathcal{D}}$ . Moreover, the corresponding risk estimator  $\hat{R}$  enjoys the same unbiasedness property. In fact:

$$\hat{R}_{\tilde{\mathcal{D}}, \ell}(\boldsymbol{\theta}) = \frac{1}{2} R_{\mathcal{D}_{2x}, \ell}(\boldsymbol{\theta}) + \mathbb{E}_{\tilde{\mathcal{D}}} [a\langle\boldsymbol{\theta}, \hat{\mathbf{z}}\rangle] \quad (5.18)$$

$$= \frac{1}{2} R_{\mathcal{D}_{2x}, \ell}(\boldsymbol{\theta}) + a\langle\boldsymbol{\theta}, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}}\rangle \quad (5.19)$$

$$= \frac{1}{2} R_{\mathcal{D}_{2x}, \ell}(\boldsymbol{\theta}) + a\langle\boldsymbol{\theta}, \boldsymbol{\mu}_{\mathcal{D}}\rangle \quad (5.20)$$

$$= R_{\mathcal{D}, \ell}(\boldsymbol{\theta}), \quad (5.21)$$

by using the independence from labels, hence from noise, of  $R_{\mathcal{D}_{2x}, \ell}$ .

### 5.7.2 Proof of Theorem 53

This is a version of Theorem 23 applied to ALN. The two results differ in three elements. First, we consider the generalization property of a minimizer  $\hat{\boldsymbol{\theta}}$  that is learned on the corrupted sample  $\tilde{\mathcal{S}}$ . Second, the minimizer is computed on the basis of the unbiased estimator of  $\hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}}$  and not barely by the average  $\boldsymbol{\mu}_{\tilde{\mathcal{S}}}$ . Third, as a consequence, Lemma 27 is not valid in this scenario. Therefore, we first prove

a version of the bound for the mean operator norm discrepancy while considering label noise.

**Lemma 59.** *Suppose  $\mathbb{R}^d \supseteq \mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq X < \infty\}$  be the observations space. Let  $\tilde{\mathcal{S}}$  is a learning sample affected by ALN with noise rates  $(p_+, p_-) \in [0, 1/2]^2$ . Then for any  $\delta > 0$  with probability at least  $1 - \delta$ :*

$$\|\hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} - \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}}\|_2 \leq \frac{X}{1 - p_- - p_+} \cdot \sqrt{\frac{d}{m} \log\left(\frac{d}{\delta}\right)}. \quad (5.22)$$

*Proof.* Let  $\tilde{\mathcal{S}}$  and  $\tilde{\mathcal{S}}'$  be two learning samples from the corrupted distribution  $\tilde{\mathcal{D}}$  that differ for only one example  $(\mathbf{x}_i, \tilde{y}_i) \neq (\mathbf{x}_{i'}, \tilde{y}_{i'})$ . Let first consider the one-dimensional case. We refer to the  $k$ -dimensional component of  $\boldsymbol{\mu}$  with  $\mu^k$ . For any  $\tilde{\mathcal{S}}, \tilde{\mathcal{S}}'$  and any  $k \in [d]$  it holds:

$$\left| \hat{\mu}_{\tilde{\mathcal{S}}}^k - \hat{\mu}_{\tilde{\mathcal{S}}'}^k \right| = \frac{1}{m} \left| \left( \frac{\tilde{y}_i - (p_- - p_+)}{1 - p_- - p_+} \right) \mathbf{x}_i^k - \left( \frac{\tilde{y}_{i'} - (p_- - p_+)}{1 - p_- - p_+} \right) \mathbf{x}_{i'}^k \right| \quad (5.23)$$

$$= \frac{1}{m} \left| \frac{\tilde{y}_i \mathbf{x}_i^k}{1 - p_- - p_+} - \frac{\tilde{y}_{i'} \mathbf{x}_{i'}^k}{1 - p_- - p_+} \right| \quad (5.24)$$

$$\leq \frac{X}{m(1 - p_- - p_+)} |\tilde{y}_i - \tilde{y}_{i'}| \quad (5.25)$$

$$\leq \frac{2X}{m(1 - p_- - p_+)}. \quad (5.26)$$

This satisfies the bounded difference condition of McDiarmid's inequality, which let us write for any  $k \in [d]$  and any  $\epsilon > 0$  that:

$$\mathbb{P} \left( \left| \hat{\mu}_{\tilde{\mathcal{D}}}^k - \hat{\mu}_{\tilde{\mathcal{S}}}^k \right| \geq \epsilon \right) \leq \exp \left( -(1 - p_- - p_+)^2 \frac{m\epsilon^2}{2X^2} \right) \quad (5.27)$$

and the multi-dimensional case, by union bound:

$$\mathbb{P} \left( \exists k \in [d] : \left| \hat{\mu}_{\tilde{\mathcal{D}}}^k - \hat{\mu}_{\tilde{\mathcal{S}}}^k \right| \geq \epsilon \right) \leq d \exp \left( -(1 - p_- - p_+)^2 \frac{m\epsilon^2}{2X^2} \right). \quad (5.28)$$

Then by negation:

$$\mathbb{P} \left( \forall k \in [d] : \left| \hat{\mu}_{\tilde{\mathcal{D}}}^k - \hat{\mu}_{\tilde{\mathcal{S}}}^k \right| \leq \epsilon \right) \geq 1 - d \exp \left( -(1 - p_- - p_+)^2 \frac{m\epsilon^2}{2X^2} \right), \quad (5.29)$$

which implies that for any  $\delta > 0$  with probability  $1 - \delta$ :

$$\frac{X}{(1 - p_- - p_+)} \sqrt{\frac{2}{m} \log\left(\frac{d}{\delta}\right)} \geq \|\hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} - \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}}\|_\infty \geq d^{-1/2} \|\boldsymbol{\mu}_{\tilde{\mathcal{D}}} - \boldsymbol{\mu}_{\tilde{\mathcal{S}}}\|_2. \quad (5.30)$$

This concludes the proof.  $\square$

The proof of Theorem 53 follows the structure of Theorem 23's and elements of Theorem 3 of Natarajan et al. [2013]'s. Let  $\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} \hat{R}_{\tilde{\mathcal{D}},\ell}(\boldsymbol{\theta})$  and  $\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathcal{H}} R_{\mathcal{D},\ell}(\boldsymbol{\theta})$ . We have:

$$R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) = \hat{R}_{\tilde{\mathcal{D}},\ell}(\hat{\boldsymbol{\theta}}) - \hat{R}_{\tilde{\mathcal{D}},\ell}(\boldsymbol{\theta}^*) \quad (5.31)$$

$$= \frac{1}{2} R_{\mathcal{D}_{2x},\ell}(\hat{\boldsymbol{\theta}}) + a \langle \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} \rangle - \frac{1}{2} R_{\mathcal{D}_{2x},\ell}(\boldsymbol{\theta}^*) - a \langle \boldsymbol{\theta}^*, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} \rangle \quad (5.32)$$

$$= \frac{1}{2} (R_{\mathcal{D}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}_{2x},\ell}(\boldsymbol{\theta}^*)) + a \langle \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} \rangle \quad (5.33)$$

$$= \frac{1}{2} (R_{\mathcal{S}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{S}_{2x},\ell}(\boldsymbol{\theta}^*)) + a \langle \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} \rangle \quad (5.34)$$

$$+ \frac{1}{2} \left( R_{\mathcal{D}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{S}_{2x},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D}_{2x},\ell}(\boldsymbol{\theta}^*) + R_{\mathcal{S}_{2x},\ell}(\boldsymbol{\theta}^*) \right) \} A_1. \quad (5.35)$$

Step 5.31 is due to unbiasedness. Again, rename Line 5.35 as  $A_1$ , which this time is bounded directly by Theorem 22. Next, we proceed as within the proof of Theorem 22 but now exploiting the fact that  $\frac{1}{2} R_{\mathcal{S}_{2x},\ell}(\boldsymbol{\theta}) = \hat{R}_{\tilde{\mathcal{S}},\ell}(\boldsymbol{\theta}) - a \langle \boldsymbol{\theta}, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}} \rangle$ :

$$R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) \leq \underbrace{\hat{R}_{\tilde{\mathcal{S}},\ell}(\hat{\boldsymbol{\theta}}) - \hat{R}_{\tilde{\mathcal{S}},\ell}(\boldsymbol{\theta}^*)}_{A_2} + \underbrace{a \langle \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*, \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} - \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}} \rangle}_{A_3} + A_1. \quad (5.36)$$

$A_2$  is never more than 0 because  $\hat{\boldsymbol{\theta}}$  is the minimizer of  $\hat{R}_{\tilde{\mathcal{S}},\ell}(\boldsymbol{\theta})$ . From the Cauchy-Schwarz inequality and bounded models it holds true:

$$A_3 \leq |a| \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 \cdot \|\hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} - \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}}\|_2 \leq 2|a|H \|\hat{\boldsymbol{\mu}}_{\tilde{\mathcal{D}}} - \hat{\boldsymbol{\mu}}_{\tilde{\mathcal{S}}}\|_2, \quad (5.37)$$

for which we can call Lemma 59. Finally, by a union bound we get that for any  $\delta > 0$  with probability  $1 - \delta$ :

$$R_{\mathcal{D},\ell}(\hat{\boldsymbol{\theta}}) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) \leq \frac{\sqrt{2} + 1}{2} \cdot \frac{XHL}{\sqrt{m}} + c(X, H) \sqrt{\frac{1}{m} \log\left(\frac{2}{\delta}\right)} + \frac{2|a|XH}{1 - p_+ - p_-} \sqrt{\frac{d}{m} \log\left(\frac{2d}{\delta}\right)}. \quad (5.38)$$

### 5.7.3 Proof of Theorem 56

The proof draws ideas from Manwani and Sastry [2013]. Let us first assume symmetric noise, *i.e.*  $p_+ = p_- = p$ . For any  $\boldsymbol{\theta}$  we have:

$$R_{\tilde{\mathcal{D}},\ell}(\boldsymbol{\theta}^*) - R_{\tilde{\mathcal{D}},\ell}(\boldsymbol{\theta}) = (1 - p) (R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) - R_{\mathcal{D},\ell}(\boldsymbol{\theta})) + p (R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) - R_{\mathcal{D},\ell}(\boldsymbol{\theta}) + 2a \langle \boldsymbol{\theta}^* - \boldsymbol{\theta}, \boldsymbol{\mu}_{\mathcal{D}} \rangle) \quad (5.39)$$

$$\leq (R_{\mathcal{D},\ell}(\boldsymbol{\theta}^*) - R_{\mathcal{D},\ell}(\boldsymbol{\theta})) + 4|a|Hp \|\boldsymbol{\mu}_{\mathcal{D}}\|_2 \quad (5.40)$$

$$\leq 4|a|Hp \|\boldsymbol{\mu}_{\mathcal{D}}\|_2. \quad (5.41)$$

We are working with LOLS, which are such that  $\ell(x) = \ell(-x) + 2ax$  and therefore we can take Step 5.39. Step 5.40 follows from Cauchy-Schwartz inequality and bounded models. Step 5.41 is true because  $\theta^*$  is the minimizer of  $R_{\mathcal{D},\ell}(\theta)$ . We have obtained a bound for any  $\theta$  and so for the supremum with regard to  $\theta$ . Therefore:

$$\sup_{\theta \in \mathcal{H}} \left( R_{\tilde{\mathcal{D}},\ell}(\theta^*) - R_{\tilde{\mathcal{D}},\ell}(\theta) \right) = R_{\tilde{\mathcal{D}},\ell}(\theta^*) - R_{\tilde{\mathcal{D}},\ell}(\tilde{\theta}) . \quad (5.42)$$

To lift the discussion to *asymmetric* label noise, risks have to be split into losses for negative and positive examples. Let  $R_{\mathcal{D}^+,\ell}$  be the risk computed over the distribution of the positive examples  $\mathcal{D}^+$  and  $R_{\mathcal{D}^-,\ell}$  the one of the negatives, and denote the mean operators  $\mu_{\mathcal{D}^+}, \mu_{\mathcal{D}^-}$  accordingly. Also, define the probability of positive and negative labels in  $\mathcal{D}$  as  $\pi_{\pm} = \mathbb{P}(y = \pm 1)$ . The same manipulations for the symmetric case let us write:

$$\begin{aligned} R_{\tilde{\mathcal{D}},\ell}(\theta^*) - R_{\tilde{\mathcal{D}},\ell}(\theta) &= \pi_- (R_{\mathcal{D}^-,\ell}(\theta^*) - R_{\mathcal{D}^-,\ell}(\theta)) + \pi_+ (R_{\mathcal{D}^+,\ell}(\theta^*) - R_{\mathcal{D}^+,\ell}(\theta)) \\ &\quad + 2ap_- \pi_- \langle \theta^* - \theta, \mu_{\mathcal{D}^-} \rangle + 2ap_+ \pi_+ \langle \theta^* - \theta, \mu_{\mathcal{D}^+} \rangle \end{aligned} \quad (5.43)$$

$$\leq (R_{\mathcal{D},\ell}(\theta^*) - R_{\mathcal{D},\ell}(\theta)) + 2a \langle \theta^* - \theta, p_- \mu_{\mathcal{D}^-} + p_+ \mu_{\mathcal{D}^+} \rangle \quad (5.44)$$

$$\leq 4|a|H \cdot \|p_- \mu_{\mathcal{D}^-} + p_+ \mu_{\mathcal{D}^+}\|_2 \quad (5.45)$$

$$\leq 4|a|H \max(p_-, p_+) \cdot \|\pi_- \mu_{\mathcal{D}^-} + \pi_+ \mu_{\mathcal{D}^+}\|_2 \quad (5.46)$$

$$= 4|a|H \max(p_-, p_+) \cdot \|\mu_{\mathcal{D}}\|_2 . \quad (5.47)$$

Then, we conclude the proof by the same argument for the symmetric case.

#### 5.7.4 Proof of Corollaries 57 and 58

The first corollary is immediate by using the additional assumption. For the second, if  $\ell$  is once differentiable and  $\gamma$ -strongly convex in the  $\theta$  argument, so is the risk  $R_{\tilde{\mathcal{D}},\ell}$  by composition with linear functions. Notice also that  $\nabla R_{\tilde{\mathcal{D}},\ell}(\tilde{\theta}^*) = 0$  because  $\tilde{\theta}^*$  is the minimizer. Therefore:

$$\epsilon \geq R_{\tilde{\mathcal{D}},\ell}(\theta^*) - R_{\tilde{\mathcal{D}},\ell}(\tilde{\theta}^*) \quad (5.48)$$

$$\geq \left\langle \nabla R_{\tilde{\mathcal{D}},\ell}(\tilde{\theta}^*), \theta^* - \tilde{\theta}^* \right\rangle + \frac{\gamma}{2} \|\theta^* - \tilde{\theta}^*\|_2^2 \quad (5.49)$$

$$\geq \frac{\gamma}{2} \|\theta^* - \tilde{\theta}^*\|_2^2 , \quad (5.50)$$

$$(5.51)$$

which means that:

$$\|\theta^* - \tilde{\theta}^*\|_2^2 \leq \frac{2\epsilon}{\gamma} . \quad (5.52)$$

## 5.8 References

Learning with noisy labels has been widely investigated in literature; we point any interested reader to Fréney and Verleysen [2014] for a comprehensive survey.

From the theoretical standpoint label noise has been studied in two different regimes, with vastly different conclusions. In the case of low-capacity (typically linear) models, even mild symmetric label noise can produce solutions that are akin to random guessing [Long and Servedio, 2010]. On the other hand, the Bayes-optimal classifier remains unchanged under symmetric [Natarajan et al., 2013; Menon et al., 2015] and even instance dependent label noise [Menon et al., 2016] implying that high-capacity models are robust to essentially any level of such noise, given a sufficient number of samples. A caveat with the latter is that label noise adversely affects the number of samples needed for learning [van Rooyen, 2015, Chapter 3].

In practice, and with more relevance for the mentioned low capacity models, a common strategy for learning with noisy labels is to engineer a loss function that is more suited to handle noise. Suppose one wishes to minimize a loss  $\ell$  on clean data. When the level of noise is known *a priori* — as we have assumed in this Chapter —, the cited Natarajan et al. [2013] provide the general form of a *noise corrected* loss  $\hat{\ell}$  such that minimization of  $\hat{\ell}$  on noisy data is equivalent (by unbiasedness) to minimization of  $\ell$  on clean data. The algorithms proposed in this Chapter can be seen as an application of the method of Natarajan et al. [2013] for correcting linear-odd losses. In the idealized case of symmetric label noise, for certain  $\ell$  one in fact does not need to know the noise rate: Manwani and Sastry [2013]; Ghosh et al. [2015] give a sufficient condition for which  $\ell$  is robust, and several examples of such robust non-convex losses, while van Rooyen et al. [2015] show that the (convex) unhinged loss is its own noise-corrected loss. Other proposed non-convex losses are in Masnadi-Shirazi and Vasconcelos [2009]; Masnadi-Shirazi et al. [2010]; Ding and Vishwanathan [2010].



---

# Learning with noisy labels II: deep neural networks, multi-class, noise estimation

---

In this Chapter, we greatly extend the discussion on label noise from Chapter 5. We propose three important improvements to the idealized learning scenario presented so far. First and foremost, the hypothesis space considered here is the one of deep neural networks. Second, it follows naturally from this choice that we study multi-class classification. Last but not least, we relax the assumption concerning the knowledge of the noisy rates and formulate an estimator for multi-class asymmetric noise; importantly, this last step let us formulate an end-to-end framework. The Chapter is accompanied by a extensive experimental analysis on MNIST, IMDB, CIFAR-10, CIFAR-100, other than on a large scale dataset of clothing images with a variety of convolutional and recurrent architectures.

## 6.1 Motivation

The goal of this Chapter is to effectively train deep neural networks with modern architectures under label noise. In Chapter 5 we discussed one strand of research for combating label noise based on the idea of loss correction. This broad line of work is strongly inspired by a theoretical analysis of this learning setting and is mostly developed within the Machine Learning community. In Chapter 5 we actually interpret this approach within the framework proposed by the Thesis, that is, sufficient statistic estimation followed by learning with linear-odd losses. This interpretation is grounded on the assumption of working with linear (or kernel) models.

Another strand of work is on *ad-hoc deep architectures* tailored to tolerate label noise, primarily developed in Computer Vision [Mnih and Hinton, 2012; Reed et al., 2015; Sukhbaatar et al., 2015; Xiao et al., 2015]. While some such approaches have shown good experimental performance on specific domains, they lack a solid theoretical framework, and often need large a amount of clean labels to obtain acceptable results — in particular, for pre-training or validating hyper-parameters [Xiao et al., 2015; Krause et al., 2016; Reed et al., 2015].

We unify the above research streams by introducing two alternative procedures for loss correction. In doing so, this Chapter slightly departs from the main approach proposed in the Thesis, with the aim of laying a bridge between those two areas.

Both procedures amount to simple linear algebra provided that we know a stochastic matrix  $T$  summarizing the probability of one class being flipped into another under noise. The first procedure, a multi-class extension of Natarajan et al. [2013] applied to neural networks, is called “*backward*” correction as it multiplies the loss by  $T^{-1}$ . The second, inspired by Sukhbaatar et al. [2015], is named “*forward*” correction as it multiplies the network predictions by  $T$ . We prove that both procedures enjoy formal guarantees of robustness with regard to the clean data distribution. Since we only operate on the loss function, the approach is architecture independent and not tied to a particular application domain, other than viable for any loss (even non-linear-odd).

In real applications practitioners may be able to obtain a good estimate of  $T$  by polishing a subset of the available training set [Xiao et al., 2015] — something undoubtedly useful and often necessary for tuning hyper-parameters and testing the model anyway. Nevertheless, we take a further step extending the noise estimator of Menon et al. [2015] to the multi-class setting. Incidentally, we also prove that, when the network only non-linearity is ReLU, the Hessian of the loss is not affected by noise.

A clear motivation of this Chapter is to push for a practical application of our formal work on label noise. We apply our loss corrections to image recognition on MNIST, CIFAR-10, CIFAR-100 and sentiment analysis on IMDB; we simulate corruption by artificially injecting noise on the training labels. In order to show that no architectural choice is the secret ingredient of our robustification recipe, we experiment with a variety of network modules: convolutions and pooling [LeCun et al., 1998], dropout [Srivastava et al., 2014], batch normalization [Ioffe and Szegedy, 2015], word embedding and residual units [He et al., 2016a,b]. Additional tests on LSTM [Hochreiter and Schmidhuber, 1997] confirms that the procedures can be seamlessly applied to recurrent neural networks as well. Comparisons with non-corrected losses and several known methods confirm robustness of our two procedures, with the forward correction dominating the backward. Unsurprisingly, the noise estimator is the bottleneck in obtaining near-perfect robustness, yet in most experiments our approach is often the best compared to prior work. Finally, we experiment with the 1M clothing images dataset of Xiao et al. [2015], establishing a new state of the art.

## 6.2 Learning setting

In supervised  $c$ -class classification, we have a feature space  $\mathcal{X} \subseteq \mathbb{R}^d$  and a label space  $\mathcal{Y} = \{\mathbf{e}^i : i \in [c]\}$ . Note that each  $\mathbf{y}$  only has one non-zero value at the coordinate corresponding to the underlying label.

A  $n$ -layer neural network comprises a transformation  $\mathbf{h} : \mathcal{X} \rightarrow \mathbb{R}^c$ , where  $\mathbf{h} = (\mathbf{h}^{(n)} \circ \mathbf{h}^{(n-1)} \circ \dots \circ \mathbf{h}^{(1)})$  is the composition of a number of intermediate transforma-

tions, the layers, defined by:

$$\begin{aligned} (\forall i \in [n-1]) \mathbf{h}^{(i)}(\mathbf{z}) &= \sigma(W^{(i)}\mathbf{z} + \mathbf{b}^{(i)}) , \\ \mathbf{h}^{(n)}(\mathbf{z}) &= W^{(n)}\mathbf{z} + \mathbf{b}^{(n)} . \end{aligned}$$

where  $W^{(i)} \in \mathbb{R}^{d^{(i)} \times d^{(i-1)}}$  and  $\mathbf{b}^{(i)} \in \mathbb{R}^{d^{(i)}}$  are parameters to be estimated<sup>1</sup>, and  $\sigma$  is any activation function that acts *coordinate-wise*, such as the rectified linear unit (ReLU)  $\sigma(\mathbf{x})_i = \max(0, x_i)$ . Observe that the final layer applies a *linear* projection, unlike all preceding layers. To simplify notation, we write:

$$(\forall i \in [n]) \mathbf{x}^{(i)} \doteq \mathbf{h}^{(i)}(\mathbf{x}^{(i-1)}),$$

with the base case  $\mathbf{x}^{(0)} \doteq \mathbf{x}$ , so that e.g.  $\mathbf{x}^{(1)}$  is exactly the representation in the first layer. Without loss of generality, we assume all layers to be *fully connected*, or dense; for example, convolutions can be represented by dense layers with shared sparse weights. The coordinates of  $\mathbf{h}(\mathbf{x})$  represent the relative weights that the model assigns to each class  $i \in [c]$  to be predicted. The predicted label is thus simply  $\mathbf{y}(\mathbf{x}) = \arg \max_{i \in [c]} h_i(\mathbf{x})$ .

In the training phase, the output of the final layer is contrasted with the true label  $\mathbf{y}$  via two steps. First,  $\mathbf{h}(\cdot)$  passes through the *softmax* function  $e^{h_i(\mathbf{x})} / \sum_{k=1}^c e^{h_k(\mathbf{x})}$ . The softmax output may be interpreted as the vector of class-wise probabilities living in the simplex  $\Delta_{c-1}$ , and hence we may denote it by  $p(\mathbf{y}|\mathbf{x})$ . Next, we measure the discrepancy between label  $\mathbf{y} = \mathbf{e}^i$  and network's output by a loss function  $\ell: \mathcal{Y} \times [0, 1]^c \rightarrow \mathbb{R}$ , given for example by the *cross-entropy*:

$$\ell(\mathbf{e}^i, \mathbf{h}(\mathbf{x})) = -(\mathbf{e}^i)^\top \log p(\mathbf{y}|\mathbf{x}) = -\log p(\mathbf{y}|\mathbf{x})_i . \quad (6.1)$$

With some abuse of notation, we also define the loss in a vector form computed on every possible label:

$$\ell(\mathbf{h}(\mathbf{x})) = \left( \ell(\mathbf{e}^1, \mathbf{h}(\mathbf{x})), \dots, \ell(\mathbf{e}^c, \mathbf{h}(\mathbf{x})) \right)^\top \in \mathbb{R}^c . \quad (6.2)$$

Cross-entropy is essentially a multi-class version of logistic loss when we encode the labels for each class in  $\{0, 1\}$  instead of  $\{\pm 1\}$ . In this Chapter, formal results hold under very mild conditions on generic loss functions. We do not assume to work with linear-odd losses, unless specified as in Theorem 64. Also, notice that here a loss  $\ell$  with a single argument *is not* a margin loss (Definition 2), but it is a notational shortcut to denote losses in vector form in the multi-class setting.

We update the setting of asymmetric label noise of Chapter 5 to multi-class. Each label  $\mathbf{y}$  in the training set is flipped to  $\tilde{\mathbf{y}} \in \mathcal{Y}$  with probability  $p(\tilde{\mathbf{y}}|\mathbf{y})$ . Denote by  $T \in [0, 1]^{c \times c}$  the noise transition matrix specifying the probability of one label being flipped to another, so that  $\forall i, j \quad T_{ij} = p(\tilde{\mathbf{y}} = \mathbf{e}^j | \mathbf{y} = \mathbf{e}^i)$ . The matrix is row-stochastic and not necessarily symmetric across the classes. We aim to modify a loss  $\ell$  so as to

<sup>1</sup>Here,  $d^{(0)} = d$ , the original feature dimensionality, and  $d^{(n)} = c$ , the label dimensionality.

make it robust to label noise; in fact, this is possible if  $T$  is known.

## 6.3 Loss correction procedures

### 6.3.1 The backward correction

We can build an *unbiased estimator* of the loss function in the same sense of Theorem 52. The corrected loss *under expected label noise* equals the original one computed on clean data. This property is stated in the next Theorem, a multi-class generalization of the already cited Lemma 1 Natarajan et al. [2013]. The Theorem is also a particular instance of the more abstract [van Rooyen, 2015, Theorem 3.2].

**Theorem 60.** *Suppose that the noise matrix  $T$  is non-singular. Given a loss  $\ell$ , define the backward corrected loss as:*

$$\ell^{\leftarrow}(\mathbf{h}(\mathbf{x})) = T^{-1}\ell(\mathbf{h}(\mathbf{x})) . \quad (6.3)$$

*Then, the loss correction is unbiased, i.e., its expectation under label noise is exactly the loss:*

$$(\forall \mathbf{y} = \mathbf{e}^i) \quad \mathbb{E}_{p(\tilde{\mathbf{y}}|\mathbf{y})} \ell^{\leftarrow}(\mathbf{h}(\mathbf{x}))_i = \ell(\mathbf{h}(\mathbf{x}))_i , \quad (6.4)$$

*and therefore the minimizers are the same:*

$$\operatorname{argmin}_h \mathbb{E}_D \ell^{\leftarrow}(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \operatorname{argmin}_h \mathbb{E}_D \ell(\mathbf{y}, \mathbf{h}(\mathbf{x})) . \quad (6.5)$$

Proof in 6.7.1. The corrected loss is effectively a linear combination of the loss values for each observable label, which coefficients are due to the probability that  $T^{-1}$  attributes to each possible true label  $\mathbf{y}$ , given the observed one  $\tilde{\mathbf{y}}$ . Intuitively, we are “going one step back” in the noise process described by the Markov chain  $T$ . The corrected loss is differentiable — although not always non-negative — and can be minimized with any off-the-shelf algorithm for back-propagation. Although in practice  $T$  would be invertible almost surely, its condition number may be problematic. A simple solution is to mix  $T$  with the identity matrix before inversion; this may be seen as taking a more conservative noise-free prior.

### 6.3.2 The forward correction

Alternatively, we can correct the model predictions. Following Sukhbaatar et al. [2015], we start by observing that a neural network learned with no loss correction would result in a predictor for noisy labels  $p(\tilde{\mathbf{y}}|\mathbf{x})$ . We can make explicit the depen-

dency on  $T$ . For instance, with cross-entropy we have:

$$\ell(\mathbf{e}^i, \mathbf{h}(\mathbf{x})) = -\log p(\tilde{\mathbf{y}} = \mathbf{e}^i | \mathbf{x}) \quad (6.6)$$

$$= -\log \sum_{j \in [c]} p(\tilde{\mathbf{y}} = \mathbf{e}^i | \mathbf{y} = \mathbf{e}^j) p(\mathbf{y} = \mathbf{e}^j | \mathbf{x}) \quad (6.7)$$

$$= -\log \sum_{j \in [c]} T_{ji} p(\mathbf{y} = \mathbf{e}^j | \mathbf{x}) , \quad (6.8)$$

or in matrix form  $\ell(\mathbf{h}(\mathbf{x})) = -\log T^\top p(\mathbf{y} | \mathbf{x})$ . This loss compares the noisy label  $\tilde{\mathbf{y}}$  to averaged noisy prediction corrupted by  $T$ . We call this procedure “forward” correction.

In order to analyze its behavior, we first need to recall definition and properties of a new family of losses, named *proper composite* [Reid and Williamson, 2010, Section 4]. This is an additional requirement with respect to properness of Definition 1<sup>2</sup>. Many losses are said to be *composite*, in the sense that they can be expressed by the aid of an *link function*.

**Definition 61.** A loss  $\ell_\psi$  is composite with link function  $\psi : \Delta_{c-1} \rightarrow \mathbb{R}^c$ , invertible, if it can be written as:

$$\ell_\psi(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \ell(\mathbf{y}, \psi^{-1}(\mathbf{h}(\mathbf{x}))) . \quad (6.9)$$

Cross-entropy and square are examples of proper composite losses. In the case of cross-entropy, the softmax is the *inverse* link function. When composite losses are also *proper*, their minimizer assumes the particular shape of the link function applied to the class probability:

$$\operatorname{argmin}_h \mathbb{E}_{\mathcal{D}} \ell_\psi(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \psi(p(\mathbf{y} | \mathbf{x})) . \quad (6.10)$$

An intriguing robustness property holds for forward correction of proper composite losses.

**Theorem 62.** Suppose that the noise matrix  $T$  is non-singular. Given a proper composite loss  $\ell_\psi$ , define the forward loss correction as:

$$\ell_{\tilde{\psi}}^{\rightarrow}(\mathbf{h}(\mathbf{x})) = \ell(T^\top \psi^{-1}(\mathbf{h}(\mathbf{x}))) . \quad (6.11)$$

Then, the minimizer of the corrected loss under the noisy distribution is the same as the minimizer of the original loss under the clean distribution:

$$\psi(p(\mathbf{y} | \mathbf{x})) = \operatorname{argmin}_h \mathbb{E}_{\tilde{\mathcal{D}}} \ell_{\tilde{\psi}}^{\rightarrow}(\mathbf{y}, \mathbf{h}(\mathbf{x})) \quad (6.12)$$

$$= \operatorname{argmin}_h \mathbb{E}_{\mathcal{D}} \ell_\psi(\mathbf{y}, \mathbf{h}(\mathbf{x})) . \quad (6.13)$$

<sup>2</sup>Symmetric proper losses of Chapter 4 are also proper composite, with link function equal to the derivative of the generator  $\phi$ . See Nock and Nielsen [2009] and Reid and Williamson [2010] for details.

---

**Algorithm 10: Robust two-stage training**

---

**Input:** the noisy training set  $\tilde{\mathcal{S}}$ , any loss  $\ell$   
 If  $T$  is unknown:  
   Train a network  $h(x)$  on  $\tilde{\mathcal{S}}$  with loss  $\ell$   
   Obtain an unlabeled sample  $\mathcal{S}_X$   
   Estimate  $\hat{T}$  by Equations (6.15)-(6.16) on  $\mathcal{S}_X$   
 Train the network  $h(x)$  on  $\tilde{\mathcal{S}}$  with loss  $\ell^{\leftarrow}$  or  $\ell^{\rightarrow}$   
**Output:**  $h(\cdot)$

---

Proof in 6.7.2. The result expresses a weaker property with respect to unbiasedness of Theorem 60. Robustness applies to the minimizer only: the model learned by forward correction is the minimizer over the *clean* distribution. Yet, Theorem 62 guarantees noise independence without explicitly inverting the noise process, but it does it “behind the scenes” by a “de-noising” link function. This turns out to be an important factor in practice, as shown in Section 6.5 experimentally and discussed in Section 6.6.

### 6.3.3 Estimating the noise rates

A clear limitation of the above procedures is that they require knowing  $T$ . In most applications, the matrix  $T$  would be unknown and needs to be estimated. We present here an extension of the noise estimator of Menon et al. [2015]; Liu and Tao [2016] to the multi-class settings. The estimator is derived under two assumptions.

**Theorem 63.** *Assume  $p(\mathbf{x}, \mathbf{y})$  is such that:*

(i) *there exist “perfect examples” of each of class  $j \in [c]$ , in the sense that:*

$$(\exists \bar{\mathbf{x}}^j \in \mathcal{X}) : p(\bar{\mathbf{x}}^j) > 0 \wedge p(\mathbf{y} = \mathbf{e}^j | \bar{\mathbf{x}}^j) = 1.$$

(ii) *given sufficiently many corrupted samples,  $h$  is rich enough to model  $p(\tilde{\mathbf{y}}|\mathbf{x})$  accurately.*

*It follows that:*

$$\forall i, j \in [c], T_{ij} = p(\tilde{\mathbf{y}} = \mathbf{e}^j | \bar{\mathbf{x}}^i). \quad (6.14)$$

Proof in 6.7.3. Rather surprisingly, Theorem 63 tells us that we can estimate each component of matrix  $T$  *just based on noisy class probability estimates*, that is, the output of the softmax of a network trained with noisy labels. In particular, let  $\mathcal{S}_X$  be any set of features vectors. This can be taken from  $\mathcal{S}$  itself, but not necessarily: we do not require this sample to have *any* label at all and therefore whenever more unlabeled samples are easy to obtain from the same distributions; they could be used in place

loss	correction	$\mathbb{E}_{\hat{\mathcal{D}}}$	Hessian of $\mathbb{E}_{\hat{\mathcal{D}}}$
$\ell$	-	no guarantee	unchanged
$\ell^{\leftarrow}$	$T^{-1}$ .	unbiased estimator of $\ell$	unchanged
$\ell^{\rightarrow}$	$T$ .	same minimizer of $\ell$	no guarantee

Table 6.1: Qualitative comparison of loss corrections.

of  $\mathcal{S}$ . We can approximate  $T$  with two steps:

$$\bar{\mathbf{x}}^i = \operatorname{argmax}_{\mathbf{x} \in \mathcal{S}_X} p(\tilde{\mathbf{y}} = \mathbf{e}^i | \mathbf{x}) \quad (6.15)$$

$$\hat{T}_{ij} = p(\tilde{\mathbf{y}} = \mathbf{e}^j | \bar{\mathbf{x}}^i) . \quad (6.16)$$

In practice, assumption (1) of Theorem 63 might hold true when  $\mathcal{S}_X$  is large enough. Assumption (2) is more difficult to justify: we require the network to perfectly model the probability of *the noisy labels*. Although, in the experiments we can often recover  $T$  close to the ground truth and find that small estimation errors have a mild, not catastrophic effect on the quality of the correction.

Algorithm 10 summarizes the end-to-end approach. If we know  $T$ , for example by cleaning manually a subset of training data, we can train with  $\ell^{\leftarrow}$  or  $\ell^{\rightarrow}$ . Otherwise, we first have to train the network with  $\ell$  on noisy data, and obtain from it estimates of  $p(\tilde{\mathbf{y}} | \mathbf{x})$  for each class via the output of the softmax. After training  $\hat{T}$  is computable in  $O(c^2 \cdot |\mathcal{S}_X|)$ . Finally, we re-train with the corrected loss, while potentially utilizing the first network to help initializing the second one.

## 6.4 Noise free Hessians via ReLU

We now present a result of independent interest in the context of label noise. The ReLU activation function appears to be a good fit for an architecture in our noise model, since it brings the particular convenience that the Hessian of the loss *does not depend on noise*, and hence the local curvature is left unchanged. At the same time, we are assured that backward correction by  $T$  — or any arbitrarily bad estimator of the matrix — has no impact on those second order properties of the loss — something that does not hold for the forward correction though. We stress the fact that other activation functions like the sigmoid do not share this guarantee. The proof in 6.7.4 makes use of the Factorization Theorem 18.

**Theorem 64.** *Assume  $\ell$  is LOL and all network activation functions are ReLUs. Then, the Hessian of  $\ell$  does not change under noise. Moreover, the Hessians of  $\ell^{\leftarrow}$  and  $\ell$  are the same for any  $T$ .*

Theorem 64 does not provide any assurance on minima: indeed, stationary points may change location due to label noise. What it *does* guarantee is that the convergence rate of first-order methods is the same: the loss curvature cannot blow up or flat out and instead it is the same point by point in the model space. The Theorem advocates

for use of ReLU networks, in line with the recent theoretical breakthrough allowing for deep learning with no local minima [Kawaguchi, 2016]. Table 6.1 summarizes the properties of loss correction.

## 6.5 Experiments

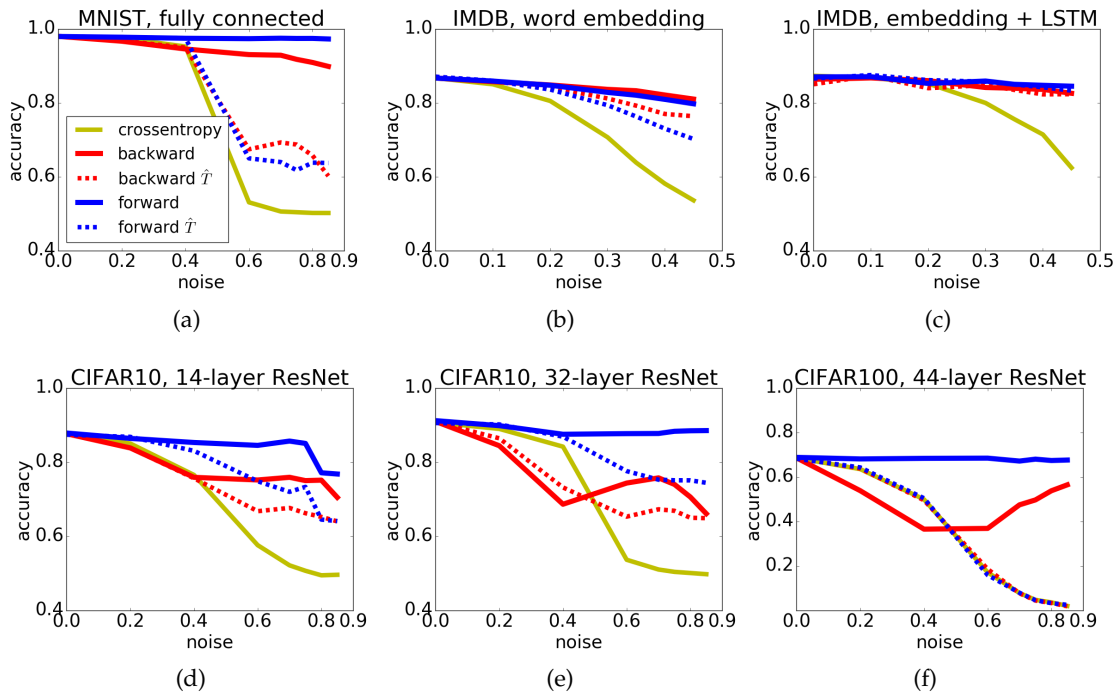


Figure 6.1: Comparison of cross-entropy with its corrections, with known or estimated  $T$ .

We now test the theory on various deep neural networks trained on MNIST [LeCun et al., 1998], IMDB [Maas et al., 2011], CIFAR-10, CIFAR-100 [Krizhevsky and Hinton, 2009] and Clothing1M [Xiao et al., 2015] so as to stress that our approach is independent on both architecture and data domain.

### 6.5.1 Loss corrections with $T$ known or estimated

We artificially corrupt labels by a parametric matrix  $T$ . The rationale is to mimic some of the structure of real mistakes for similar classes, *e.g.* CAT  $\rightarrow$  DOG. Transitions are parameterized by  $N \in [0, 1]$  such that ground truth and wrong class have probability respectively of  $1 - N, N$ . An example of  $T$  used for MNIST with  $N = 0.7$  is on the



left:

$$\begin{array}{c}
 \left[ \begin{array}{cccccccc}
 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & .3 & 0 & 0 & 0 & .7 & 0 \\
 0 & 0 & 0 & .3 & 0 & 0 & 0 & .7 \\
 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & .3 & .7 & 0 \\
 0 & 0 & 0 & 0 & 0 & .7 & .3 & 0 \\
 0 & .7 & 0 & 0 & 0 & 0 & 0 & .3 \\
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
 \end{array} \right], & \left[ \begin{array}{cccccccc}
 1 & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon \\
 \epsilon & 1 & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon \\
 \epsilon & \epsilon & .33 & \epsilon & \epsilon & \epsilon & \epsilon & .67 \\
 \epsilon & \epsilon & \epsilon & .35 & \epsilon & \epsilon & \epsilon & .65 \\
 \epsilon & \epsilon & \epsilon & \epsilon & 1 & \epsilon & \epsilon & \epsilon \\
 \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & .29 & .71 & \epsilon \\
 \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & .73 & .26 & \epsilon \\
 \epsilon & .75 & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & .25 \\
 \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & 1 \\
 \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon \\
 \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & \epsilon & 1
 \end{array} \right]
 \end{array} \tag{6.17}$$

Common to all experiments is what follows. The loss  $\ell$  chosen for comparison is cross-entropy. 10% of training data is hold out for validation. *The loss* is evaluated on it during training: with the corrected losses we can validated *on noisy data*, which is advantageous over other approaches that measure noisy validation accuracy instead. The available standard test sets are used for testing. We use ReLUs for all networks and initialize weights prior to ReLUs as in He et al. [2015], otherwise by uniform sampling in  $[-0.05, 0.05]$ . The mini-batch size is 128. The estimator of  $T$  from noisy labels is applied to  $\mathcal{S}_X$  being training and validation sets together; in fact, preliminary experiments highlighted that a  $\mathcal{S}_X$  of large size sensibly improves the approximation of  $T$ ; after estimation, we row-normalize the matrix. Following Menon et al. [2015], we take a  $\alpha$ -percentile in place of the argmax of Equation 6.15, and we find  $\alpha = 97\%$  to work well for most experiments. Although, such estimator performs very poorly with CIFAR100, possibly due the small number of images per class, and we find it is better to run the argmax instead.

*Fully connected network on MNIST.* In the first set of experiments we consider MNIST. Pixels are normalized in  $[0, 1]$ . Noise flips some of the similar digits:  $1 \rightarrow 7, 2 \rightarrow 7, 3 \rightarrow 8, 5 \leftrightarrow 6$ ; see Equation (6.17, left). We train an architecture with two dense hidden layers of size 128, with probability 0.5 of dropout. AdaGrad [Duchi et al., 2011] is run for 40 epochs with initial learning rate 0.01 and  $\delta = 10^{-6}$ . We repeat each experiment 5 times to account for noise and weight initialization. It is clear from Figure 6.1a that, although the model is somewhat robust to mild noise, high level of corruption has a disrupting effect on  $\ell$ . Instead, our losses do not witness a drastic drop. With  $\hat{T}$  estimated performance lays in between, yet it is significantly better than with no correction. An example of  $\hat{T}$  is in Equation (6.17, right), with  $\epsilon < 10^{-6}$ .

*Word embedding and LSTM on IMDB.* We keep only the top 5000 most frequent words in the corpus. Each review is either truncated or padded to 400-word long. To simulate asymmetric noise in this binary problem, we keep constant noise for the transition  $0 \rightarrow 1$  at 5%, while  $1 \rightarrow 0$  is parameterized as above; 0/1 are the two review’s sentiments. We trained two models inspired by the baselines of Dai and Le [2015]. The first maps words into 50-dimensional embedding, before passing through ReLUs; dropout with probability 0.8 is applied to the embedding output as in Yarin and Ghahramani [2016]. In the second model the embedding has dimension 256 and it is followed by an LSTM with 512 units and by a last 512-dimensional hidden layer with 0.5 dropout. AdaGrad is run for 50 epochs with the same setup as above; results are averages over 5 runs. Figures 6.1b-6.1c display an outcome similar to what previously observed on MNIST, in spite of difference in dataset, number of classes,

architecture and structure of  $T$ . Noticeably, our approach is effective on recurrent networks as well. Correcting with  $\hat{T}$  is in line with the true  $T$  here; we believe this is because estimation is easier on this binary problem.

*Residual networks on CIFAR-10 and CIFAR-100.* For both datasets we perform per-pixel mean subtraction and data augmentation as in He et al. [2016a], by horizontal random flips and  $32 \times 32$  random crops after padding with 4 pixels on each side.  $T$  for CIFAR-10 is described by: TRUCK  $\rightarrow$  AUTOMOBILE, BIRD  $\rightarrow$  AIRPLANE, DEER  $\rightarrow$  HORSE, CAT  $\leftrightarrow$  DOG. In CIFAR-100, the 100 classes are grouped into 20 5-size super-classes, e.g. AQUATIC mammals contain BEAVER, DOLPHIN, OTTER, SEAL and WHALE. Within super-classes, the noise flips each class into the next, circularly.

For this last set of experiments we use deep residual networks (ResNet), the CIFAR10/100 architectures from He et al. [2016a]. In short, residual blocks implements a non-linear operation  $F(x)$  in parallel with an identity shortcut, so as to sum the input with the output of the same block:  $x \rightarrow x + F(x)$ . Conceptually, this allows gradients to propagate more freely and helps convergence in deep architectures.  $F$  is implemented as cascade of twice batch normalization  $\rightarrow$  ReLU  $\rightarrow 3 \times 3$  convolution, following the “pre-activation” recommendation of He et al. [2016b]. Striding is used instead of max pooling; average pooling over channels is placed just before the output layer. We point to He et al. [2016a] for a detailed description. Here we experiment with ResNets of depth 14 and 32 (CIFAR10) and 44 (CIFAR100). By common practice [Huang et al., 2016], we run SGD with 0.9 momentum and learning rate 0.01, and divide it by 10 after 40 and 80 epoch (120 in total) for CIFAR-10 and after 80 and 120 (150 in total) for CIFAR-100; weight decay is set to  $10^{-4}$ . Training deep ResNets is more time consuming and therefore experiments are run only once. Since we use shallower networks than the ones in He et al. [2016a], performance is not necessarily comparable with the original work. Figures 6.1d-6.1e-6.1f forward correction does not show any significant drop at all. Except with the shallowest ResNet, backward correction does not work so well in the low noise regime and it is possibly affected by high variance. Finally, notice how the noise estimation seems particularly difficult on CIFAR100.

### 6.5.2 Comparing with other loss functions

We now compare with other methods. Data, architectures and artificial noise are the same as above. Additionally, we test the case of symmetric noise:  $N$  is the probability of label flip that is spread uniformly among all the other classes. We select methods prescribing changes in the loss function, similarly to ours: unhinged [van Rooyen et al., 2015], sigmoid [Ghosh et al., 2015], Savage [Masnadi-Shirazi and Vasconcelos, 2009] and soft and hard bootstrapping [Reed et al., 2015]; hyper-parameters of the last two methods are set in accordance with their paper.

Unhinged loss is unbounded and cannot be used alone. In the original work  $L_2$  regularization is applied to address the problem, when training non-parametric kernel models. We tried to regularize every layer with little success: learning either does not converge (too little regularization) or converge to very poor solutions (too

much). On preliminary experiments sigmoid loss ran into the opposite issue, namely premature saturation: the loss reaches a plateau too quickly, a well-known problem with sigmoidal activation functions [Glorot and Bengio, 2010]. To make those losses usable for comparison, we stack a layer of batch normalization right before the loss function. Essentially, the network output is normalized around 0 and thus operates in a bounded and non-saturated area of the loss; note that this is never required for linear or kernel models.

Table 6.3 presents the empirical analysis. We list the key findings.

(a) In the absence of artificial noise (first column for each dataset), all losses reach similar accuracy with a spread of 2 points; exceptions are some instances of unhinged and sigmoid, and Savage on CIFAR100 that makes learning impossible. Additionally, with IMDB there are cases († in Table 6.3) of loss correction with noise estimation that perform slightly better than assuming no noise: clearly, the estimator is able to recover the *natural* noise in the sentiment reviews.

(b) With low asymmetric noise (second column) results differ between simple architecture/tasks (datasets on the left) and deep networks/more difficult problems (right): in the former case, the two corrections behave similarly and are not statistically far from the competitors; in the latter case, forward correction with known  $T$  is unbeaten, with no clear winner among the remaining ones.

(c) With asymmetric noise (last two columns) the two loss corrections with known  $T$  are overall the best performing, confirming the practical implications of their formal guarantees; forward is usually the best.

(d) If we exclude CIFAR100, the noise estimation accounts for average accuracy drops between 0 (IMBD with LSTM model) and 27 points (MNIST); nevertheless, our performance is better than any other methods in many occasions.

(e) In the experiment on CIFAR100 we obtain essentially perfect noise robustness with the ideal forward correction. The noise estimation works well except in the very last column, yet it guarantees again better accuracy over competing methods. We discuss this issue in Section 6.6.

### 6.5.3 Experiments on Clothing1M

Finally, we test on Clothing1M [Xiao et al., 2015], consisting of 1M images with noisy labels, with additional 50k, 14k, 10k of clean data respectively for training, validation and testing; we refer to those sets by their size. We aim to classify images within 14 classes, *e.g.* t-shirt, suit, vest. In the original work two AlexNets [Krizhevsky et al., 2012] are trained together via EM; the networks are pre-trained with ImageNet. Two practical tricks are fundamental: a first learning phase with the clean 50k to help EM (#1 in Table 6.2) and a second phase with the mix of 50k bootstrapped to 500k and 1M (#3). Data augmentation is also applied, same as in Section 6.5.1 for CIFAR10.

We learn a 50-layer ResNet pre-trained on ImageNet — the bottleneck architecture of He et al. [2016a] — with SGD with learning rate  $10^{-3}$  and  $10^{-4}$  for 5 epochs each. and 0.9 momentum. Batch size is 32. When we train with 50k we use weight decay  $5 \cdot 10^{-2}$  and data augmentation, while with 1M we use only weight decay of  $10^{-3}$ .

Clothing1M					
#	model	loss	init	training	accuracy
1	AlexNet	cross-.	ImageNet	50k	72.63
2	AlexNet [Sukhbaatar et al., 2015]	cross-.	#1	1M, 50k	76.22
3	AlexNet [Xiao et al., 2015]	cross-.	#1	1M, 50k	78.24
4	50-ResNet	cross-	ImageNet	1M	68.94
5	50-ResNet	backward	ImageNet	1M	69.13
6	50-ResNet	forward	ImageNet	1M	69.84
7	50-ResNet	cross-.	ImageNet	50k	75.19
8	50-ResNet	cross-.	#6	50k	<b>80.38</b>

Table 6.2: Results on the top section are from Xiao et al. [2015]. In #2, #3 the clean 50k are bootstrapped to 500k. Best result #8 is obtained by fine tuning a net trained with forward correction.

The ResNet gives an uplift of about 2.5 points by training with 50k only (#7 vs. #1). However, the large amount of noisy images is essential to compete with #3. Instead of estimating the matrix  $T$  by (6.15)-(6.16), we exploit the curated labels of 50k and their noisy versions in 1M. Forward and backward corrections confirm to work better than cross-entropy (#6, #5 vs. #4), yet cannot reach the state of the art without the additional clean data. Thus, we fine-tune the networks with 50k, with the same learning parameters as in #7; because of time constraints we only fine-tune #6. The new state of the art is #8 that outperforms Xiao et al. [2015] of more than 2 points, which is achieved without time consuming bootstrapping of the 50k.

## 6.6 Discussion

We have proposed two methods for training deep neural networks with noisy labels that boils down to two loss corrections based on modeling the noise by a row-stochastic matrix  $T$ . Test accuracy is consistently only a few percents away from training cross-entropy *on clean data*, while corruption often worsen performance of cross-entropy by 40 points or more. Forward correction often outperforms the backward one. The explicit inversion of  $T$  — sometimes with high condition number — may be the root of the problem; indeed, backward correction is a linear combination of loss values for each possible label, and their coefficients may be far by orders of magnitude, which intuitively makes optimization hard. Instead, forward correction projects model predictions into another distribution in the probability simplex.

The quality of noise estimation is evidently the key factor for successfully obtaining robustness. Estimation works fairly well in most experiments with a median drop of only 10 points of accuracy with respect to using the true  $T$ . The only exception is the very last column of tests on CIFAR100, where estimation destroys most of the gain from correcting the loss. We believe that the combination of high noise and limited number of images per class (500) is detrimental to the proposed estimator. This is confirmed by the sensitivity of the hyper-parameter  $\alpha$ . In fact, the same value

---

used in other experiments led to  $\hat{T}$  with no resemblance to the ground truth, leading the learning process to very poor solutions.

We attempt two explanations. First, the sample size may be too small for estimating  $T$  with 100 classes — it is 10 times smaller per class than on CIFAR10. Second, we assumption (2) of Theorem 63 may be the bottleneck: cross-entropy achieves less than 70 accuracy on clean data versus  $\sim 90$  with CIFAR10; this is an indication that the model may well be misspecified for CIFAR100 and hence the estimator of  $p(\tilde{y}|x)$  may be not good enough.

Future work shall improve the estimation phase by incorporating knowledge of the noise structure, for example assuming low rank  $T$ . Improvements on this direction may also enlarge the applicability of our approach to multi-class with thousands of classes. It remains an open question whether more realistic *instance*-dependent noise may be included in our approach [Xiao et al., 2015; Menon et al., 2016]. Finally, we anticipate the application of our approach as a seamless tool for pre-training models with noisy data from the Web, in the spirit of Krause et al. [2016].



## 6.7 Appendix: proofs

### 6.7.1 Proof of Theorem 60

Simply:

$$\mathbb{E}_{p(\tilde{\mathbf{y}}|\cdot)} \ell^{\leftarrow}(\mathbf{h}(\mathbf{x})) = T \ell^{\leftarrow}(\mathbf{h}(\mathbf{x})) = T T^{-1} \ell(\mathbf{h}(\mathbf{x})) = \ell(\mathbf{h}(\mathbf{x})) . \quad (6.18)$$

The second statement follows from  $\ell(\mathbf{e}^i, \mathbf{h}(\mathbf{x})) = (\mathbf{e}^i)^\top \ell(\mathbf{h}(\mathbf{x}))$ . Therefore, the minimizers are the same.

### 6.7.2 Proof of Theorem 62

First notice that:

$$\ell_{\tilde{\psi}}^{\rightarrow}(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \ell(\mathbf{y}, T^\top \psi^{-1}(\mathbf{h}(\mathbf{x}))) = \ell_{\phi}(\mathbf{y}, \mathbf{h}(\mathbf{x})) \quad (6.19)$$

where we denote  $\phi^{-1} = \psi^{-1} \circ T^\top$ , or equivalently  $\phi = (T^{-1})^\top \circ \psi$  by rule of inverse of composition.  $\phi$  is invertible by composition of invertible functions, its domain is  $[0, 1]$  as of  $\psi$  and its codomain is  $\mathbb{R}$  because of composition of  $T^{-1}$  with  $\psi$ . The last loss in Equation 6.19 is therefore proper composite with link  $\phi$ . Finally, from Equation 6.10, the loss minimizer over the noisy distribution is:

$$\operatorname{argmin}_h \mathbb{E}_{\tilde{\mathcal{D}}} \ell_{\phi}(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \phi(p(\tilde{\mathbf{y}}|\mathbf{x})) \quad (6.20)$$

$$= \psi((T^{-1})^\top p(\tilde{\mathbf{y}}|\mathbf{x})) = \psi(p(\mathbf{y}|\mathbf{x})) , \quad (6.21)$$

that proves the Theorem.

### 6.7.3 Proof of Theorem 63

For any  $j \in [c]$  and any  $\mathbf{x} \in \mathcal{X}$ , we have:

$$p(\tilde{\mathbf{y}} = \mathbf{e}^j | \mathbf{x}) = \sum_{k \in [c]} p(\tilde{\mathbf{y}} = \mathbf{e}^j | \mathbf{y} = \mathbf{e}^k) p(\mathbf{y} = \mathbf{e}^k | \mathbf{x}) \quad (6.22)$$

$$= \sum_{k \in [c]} T_{kj} p(\mathbf{y} = \mathbf{e}^k | \mathbf{x}) . \quad (6.23)$$

When  $\mathbf{x} = \bar{\mathbf{x}}^i$ , we have  $p(\mathbf{y} = \mathbf{e}^k | \bar{\mathbf{x}}^i) = 0$  for  $k \neq i$ .

### 6.7.4 Proof of Theorem 64

We give the proof for cross-entropy for simplicity. One can use the Factorization Theorem 18 to generalize the result to all LOLs. When  $\mathbf{y} = \mathbf{e}^i$  the loss is:

$$\ell(\mathbf{e}^i, \mathbf{h}(\mathbf{x})) = -\log p(\mathbf{y}|\mathbf{x})_i \quad (6.24)$$

$$= -\log \frac{e^{W_i^{(n)} \mathbf{x}^{(n-1)} + \mathbf{b}_i^{(n)}}}{\sum_{k=1}^c e^{W_k^{(n)} \mathbf{x}^{(n-1)} + \mathbf{b}_k^{(n)}}} \quad (6.25)$$

$$= -W_i^{(n)} \mathbf{x}^{(n-1)} + \mathbf{b}_i^{(n)} + \log \sum_{k=1}^c e^{W_k^{(n)} \mathbf{x}^{(n-1)} + \mathbf{b}_k^{(n)}}. \quad (6.26)$$

The only dependence on the true class  $\mathbf{e}^i$  above are the first two terms. The log-partition is *independent* of the precise class  $i$ . Evidently, the noise affects the loss only through  $W_i^{(n)}$  and  $\mathbf{b}_i^{(n)}$ : those are the *only* terms in which  $\ell(\mathbf{y}, \mathbf{h}(\mathbf{x}))$  and  $\ell(\tilde{\mathbf{y}}, \mathbf{h}(\mathbf{x}))$  may differ. Therefore we can rewrite the backward corrected loss as:

$$\ell^{\leftarrow}(\mathbf{e}^j, \mathbf{h}(\mathbf{x})) = \left( T^{-1} \ell(\mathbf{h}(\mathbf{x})) \right)_j \quad (6.27)$$

$$= - \left( T^{-1} W^{(n)} \right)_j \cdot \mathbf{x}^{(n-1)} - \left( T^{-1} \mathbf{b}^{(n)} \right)_j \\ + \log \sum_{k=1}^c e^{W_k^{(n)} \mathbf{x}^{(n-1)} + \mathbf{b}_k^{(n)}}. \quad (6.28)$$

In fact, note that  $T^{-1}$  does not affect the log-partition function. To see this, let  $A(\mathbf{x}) = \log \left( \sum_{k=1}^c e^{W_k^{(n)} \mathbf{x}^{(n-1)} + \mathbf{b}_k^{(n)}} \right)$ , with the (vector) log-partition being  $A(\mathbf{x})\mathbf{1}$ . It follows that its correction is  $T^{-1}A(\mathbf{x})\mathbf{1} = A(\mathbf{x})\mathbf{1}$ , by left-multiplication of  $T$  and because  $T\mathbf{1} = \mathbf{1}$  since  $T$  is row-stochastic. Thus  $\ell^{\leftarrow}(\mathbf{e}^j, \mathbf{h}(\mathbf{x})) = B(\mathbf{x}) + A(\mathbf{x})$ , where  $B(\mathbf{x}) = - \left( T^{-1} W^{(n)} \right)_j \cdot \mathbf{x}^{(n-1)} - \left( T^{-1} \mathbf{b}^{(n)} \right)_j$  is a *piece-wise linear* function of the *model parameters*, and the log-partition  $A(\mathbf{x})$  is non-linear because of the loss and the architecture but *does not depend on noise*. Since the composition of piece-wise linear function is piece-wise linear, the Hessian of  $B(\mathbf{x})$  vanishes, and therefore the Hessian of  $\ell^{\leftarrow}$  is noise independent for any  $T$ . The same holds for  $\ell$  (no correction) by taking  $T = I$  and hence the Hessians are the same.

## 6.8 References

We review relevant literature on Deep Learning. Several works have attempted to deal with noisy labels of late, especially in Computer Vision. This is often achieved by formulating noise-aware models. Mnih and Hinton [2012] build a noise model for binary classification of aerial image patches, which can handle omission and wrong location of training labels. Xiao et al. [2015] construct a more sophisticated mix of symmetric, asymmetric and instance-dependent noise; two networks are learned by EM as model for classifier and noise type. Reed et al. [2015] augment the objective



---

similarly to entropy regularization. In practice, it is often the case that a small set of clean labels is needed in order either to pre-train or fine-tune the model [Xiao et al., 2015; Krause et al., 2016; Reed et al., 2015].

The work of Sukhbaatar et al. [2015] deserves a particular mention. The method augments the architecture by adding a linear layer on top of the network. Once learned, this layer plays the role of our matrix  $T$ . The insight is that, at training time, the effect of the noise is captured by the linear layer, emulating the corruption of the model predictions; at test time, the linear layer must be removed to obtain clean predictions. However, learning this architecture appears problematic; heuristics are necessary, such as trace regularization and a fixed updating schedule for the linear layer. We sidestep those issues by decoupling the two phases: we first estimate  $T$  and then learn with loss correction.

Recent work has provided methods to estimate label flip probabilities directly from noisy samples. (Note that the matrix is fully specified by 2 noise rates in binary classification.) If one has access to clean samples, estimation is reliable Kearns [1998]. With only noisy samples, one can obtain *bounds* on these rates [Laird, 1988, Algorithm 5.8], one can obtain  $T$  under some assumptions on the generating distribution [Scott et al., 2013; Sanderson and C. Scott, 2014; Liu and Tao, 2016; Menon et al., 2015; Ramaswamy et al., 2016]. Typically, it is required that the generating distribution is “weakly separable”: that is, such that for each class, there exists some “perfect” instance, *i.e.* one that is classified with probability equal to one. Proposed estimators involve either the use of kernel mean embedding [Ramaswamy et al., 2016], or post-processing the output of a standard class-probability estimator such as logistic regression using order statistics on the range of scores [Liu and Tao, 2016; Menon et al., 2015] or the slope of the induced ROC curve [Sanderson and C. Scott, 2014]. A common limitation is the focus on the case of binary labels, with the exception of Sanderson and C. Scott [2014].

We are not aware of any other attempt at either applying the noise-corrected loss approach of Natarajan et al. [2013] to neural networks, nor on combining those losses with the above noise rate estimators. Our work sits precisely in this intersection. Note that, even though in principle loss correction should not be necessary for high-capacity models like deep neural networks, owing to aforementioned theoretical results, in practice such correction may offset the sub-optimality of these models arising from training on finite samples.



---

# Learning from vertically distributed data without entity matching

---

This last Chapter deals with a rather particular learning problem. Our goal is to learn classifiers combining features from two different spaces. To a first approximation, we can consider the data as vertically (feature-wise) partitioned, hence distributed. But the scenario is actually more challenging. We do not even know how to match features vectors from one dataset to the other; their identity has been lost, for instance due to anonymization. Therefore, to some extent, the weak supervision derives from the fact that we do not exactly know how to match “parts” of features vectors, while we do assume to know all the relative labels. Once again, we design a solution within the two-step framework by first recovering the missing information via estimation of sufficient statistics. This approach allow us to bypass a more traditional solution, combinatorial in nature, that resorts to entity matching. While conceptually similar to the rest of the Thesis, the peculiar constraints for this learning problem push for a different workaround manipulating loss functions. We require the notion of Rademacher observations, aggregated statistics from which we learn instead of examples. This way, entity resolution can be bypassed to carry out supervised learning almost as accurate as if its solution were known.

## 7.1 Motivation

Consider the following data fusion scenario: two datasets/peers contain the same real-world entities described using partially shared features, for example banking and insurance company records of the same customer base. Our goal is to learn a classifier in the cross product space of the two domains, in the hard case in which no shared ID is available – *e.g.* due to anonymization.

A main motivation towards this objective comes from the reported experience that combining features from different sources leads to better predictive power. For instance, insurance and banking data together can improve fraud detection; shopping

	Peer 1			Peer 2			
		shared			shared		
	$x_1$	$x_3$	$c$		$x_2$	$x_3$	$c$
$e_1$	1	1	1	$e'_1$	-1	1	1
$e_2$	-1	1	1	$e'_2$	1	1	1

Table 7.1: A simple case of vertical partitioning with  $p = 2$  peers, two shared variables  $x_3$  and  $c$  (the class to predict). This toy example has binary description features and a binary shared feature, but this restriction does not need to hold in the general case. For example, each shared feature can be any categorical/ordinal feature, like “post-code”, “age-band”, etc.

records well complement medical history for estimating risk of disease [Tsui et al., 2003]; joining heterogeneous data helps prediction in genomics [Lanckriet et al., 2004; Yamanishi et al., 2004]; security agencies integrate various sources for terrorism intelligence [Sweeney, 2005; Christen, 2006; Sproull et al., 2015]

Typical data fusion methods however rely on a known map between entities [Bleiholder and Naumann, 2008], *i.e.* peers have partially different views of the *same* examples. Instead, we assume the datasets do not share a common ID, as shown in Figure 7.2; that is, for example, the case when data collection is performed independently by each peer, or when sources are deliberately anonymized. Thus, we can think the data as vertically partitioned (VP).

Learning from massively distributed data collections and multiple information sources has become a pivotal problem, yet it faces critical challenges, among which is the fact that it relies on reconstructing consistent examples from diverse features distributed between different data handling *peers*. Exhaustive search to solve this problem is simply not scalable, nor communication efficient, and sometimes not even accurate [Estrada et al., 2010; Zhang et al., 2015].

### 7.1.1 Entity resolution

*Entity resolution* (ER), or entity matching [Christen, 2012], would be the traditional approach for reconciling entities with no shared ID. It approximates a JOIN operation, assuming that some of the attributes are shared, such as *age-band*, *gender*, *post-code* (etc.), and hence can be used as “weak IDs”. Most techniques for ER are based on similarity functions and thresholding: candidate entities are selected as matches when their similarity is above a threshold. Both components can be tuned on some ground truth matches and effectively enhanced with learning techniques [Bilenko and Mooney, 2003; Christen, 2012]. Performance metrics of ER encompass lots of different parameters, including generality, accuracy, soundness, scalability, parallelizability [Rastogi et al., 2011].

The standard pipeline for learning with ER is depicted in Figure 7.1 (left): (1) entities are matched based on similarity and heuristics, (2) they are merged in one unique database and (3) a model is learned on the joint data. Here, we deliberately

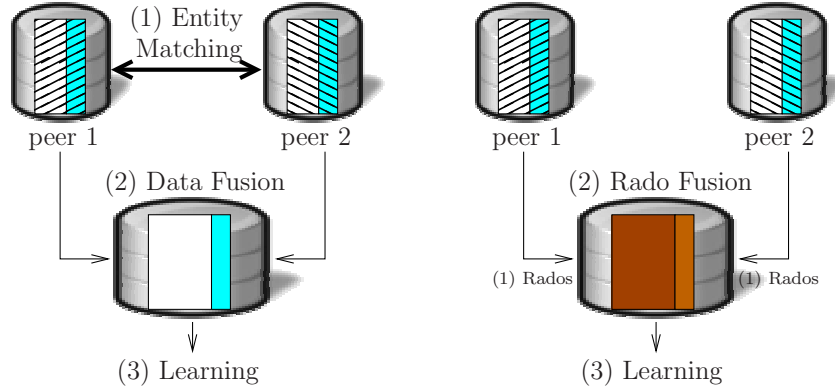


Figure 7.1: Learning on top of ER (left) or with rados (right).

do not consider common issues, such as *conflicts* and *heterogeneity* [Bleiholder and Naumann, 2008].

From a high level view, ER integrates data as a pre-process for other tasks. When it comes to learning from ER'ed data, small changes in ER can have large impact on evaluating classifiers, even for simple classifiers as linear models. To see this, suppose we are in the toy example of Table 7.1. Here, all shared variables have the same values, so entity matching has two potential solutions; notice that one of the shared variable is class  $c$ . One, say ER1, is matching  $e_1$  with  $e'_1$  and  $e_2$  with  $e'_2$ . We denote the examples obtained by

$$e_{11} = ((1, -1, 1), 1), \quad e_{22} = ((-1, 1, 1), 1) .$$

The other solution, say ER2, is matching  $e_1$  with  $e'_2$  and  $e_2$  with  $e'_1$ . We denote the examples obtained by

$$e_{12} = ((1, 1, 1), 1), \quad e_{21} = ((-1, -1, 1), 1) .$$

Consider linear classifier  $\theta = (1, 1, 1) \in \mathbb{R}^3$ ; the predicted class is given by the sign of its inner product with an observation,  $\text{sign}\langle \theta, z \rangle$ . While  $\theta$  classifies perfectly on  $\{e_{11}, e_{22}\}$  (zero error), it classifies no better than random on  $\{e_{12}, e_{21}\}$  (error 50%). The simple demonstration shows how the approach of entity matching is not only computationally expensive but potentially detrimental to the final learning performance. We advocate for an entirely different solution addressing both issues.

## 7.2 Learning setting

In this Chapter we work again on the setting of binary classification with linear models  $\theta$ . We refer to  $\mathcal{S} \doteq \{(x_i, y_i), i \in [m]\}$  as the *total* learning sample. We have  $p$  (sub)samples,  $\mathcal{S}^j$  of size  $m_j$ ,  $j \in [p]$  for some  $p > 1$ . Each one is defined in its own

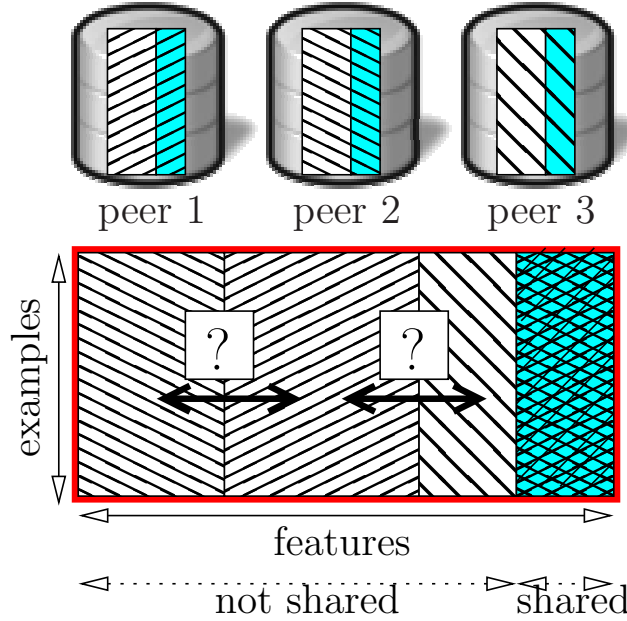


Figure 7.2: Schematic view of our setting, with  $p = 3$  peers (best viewed in color). Some features (cyan) are described in each peer and one these shared features is a class. Non-shared features are split among peers. A so-called *total* sample  $\mathcal{S}$  is figured by the red rectangle. In the case of vertical partition (**VP**) all peers see different views of the same examples, but do not know who is who among their datasets (“?”). Hence, each bit of the total sample is seen by one peer.

feature space

$$\mathcal{X}^j \doteq \times_{k=1}^{d_j} \mathcal{X}_{j,k}, j_k \in [d], \forall k. \quad (7.1)$$

To get a simple case of this framework, shown in Figure 7.2, one may see each  $\mathcal{S}^j \doteq \{(x_i^j, y_i^j), i \in [m_j]\}$  handled by a peer  $P^j$ . We denote this setting by the work “distributed” not necessarily meaning that the data is split on different computing machines or locations — although this is also possible —, but to stress the fact that the partial views of the dataset cannot be simply joined together. To avoid confusion, we reserve the word *entity* for a generic record in a dataset, the object of matching, and *attributes* or *features* to its fields. In the following, subscripts  $i$  will refer to an example or entity, while superscripts  $j$  to a peer.

We rely on the following assumption:

- (A7) The class  $\mathcal{Y}$  and a subset of features  $\mathcal{J} \subseteq \{\mathcal{X}_k\}_{k=1}^d$  are shared by all peers. Each other feature is exclusive to one peer.

Hence, there exists  $\dim(\mathcal{J}) + 1$  columns that represent the same set of attributes among peers, and one of them is the class. Each of the dimensions of  $\mathcal{J}$  is in all  $\mathcal{X}^j$ s. This is a realistic assumption for the features in  $\mathcal{J}$ : in the (**VP**) setting the domain is

vertically partitioned for the non-shared features, implying  $m_j = m_{j'} = m, \forall j, j' \in [p]$ . In this case, there exists an (unknown) one-to-one mapping between the peers' rows. The shared labels may be harder to justify, since they are the attribute we aim to predict. However, we argue in Section 7.6 that if at least one peer has classes than all peers can get their labels as well *without entity resolution*, by methods for LLP from Chapter 4.

### 7.3 Rademacher observations

In order to elaborate our solution to the learning problem of this Chapter, we need to recall elements of the theory on *Rademacher observations* (rados). The notion is introduced in Nock et al. [2015], which we have co-authored. In the standard binary classification setting, a Rademacher observation is a simple transformation of the examples in  $\mathcal{S}$ .

**Definition 65.** Let  $\sigma \in \Sigma_m = \{-1, 1\}^m$ . A Rademacher observation, or rado, is a vector in  $\mathbb{R}^d$  defined as:

$$\pi_\sigma \doteq \sum_{i=1}^m 1\{y_i = \sigma_i\} y_i x_i . \quad (7.2)$$

We also term any element  $yx$  as an *edge* vector. Rados are effectively sums of edge vectors restricted to subsample of  $\mathcal{S}$ . From  $\mathcal{S}$ , we can obtain  $2^m$  rados. For any  $\Sigma'_m \subseteq \Sigma_m \doteq \{-1, 1\}^m$ , we let  $\mathcal{R}_{\mathcal{S}, \Sigma'_m} \doteq \{\pi_\sigma : \pi_\sigma \in \Sigma'_m\}$  which denotes the set of rados that can be crafted from  $\Sigma'_m$  using  $\mathcal{S}$ . One rado in  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$  is a non-normalized version of the mean operator  $\pi = \sum_{i=1}^m y_i \cdot x_i$ , with  $\sigma = \mathbf{y}$ . A comment on notation in this Chapter: for maintaining coherence with published work, we use the greek letter  $\pi$  for rados and the letter  $\mathcal{R}$  for sets of rados. Since there is no mention of neither label proportions (Chapter 4) nor Rademacher complexity, this choice should not be a source of confusion.

One reason behind Rademacher observations is the existence of multiple equivalences between loss functions computed on examples and loss functions computed on rados. The idea was originally introduced by mapping the standard logistic loss on examples to a particular exponential loss on rados [Nock et al., 2015, Lemma 2].

**Theorem 66.** *The following holds true for any  $\theta$  and  $\mathcal{S}$ :*

$$\mathbb{E}_{(x,y) \sim \mathcal{S}} \log \left( 1 + e^{-y\langle \theta, x \rangle} \right) = \log(2) + \frac{1}{m} \log \left( \mathbb{E}_{\pi \sim \mathcal{R}_{\mathcal{S}, \Sigma_m}} e^{-\langle \theta, \pi \rangle} \right) \quad (7.3)$$

By virtue of the Theorem, in order to compute and hence optimize the right-hand side of 7.3, we don't need to know any individual examples. Instead, we can train a model on rados. Importantly, the minimizers of the two sides of the Equation are the same. Therefore, no post-processing of the model is required. This should remind of the Factorization Theorem 18.

There is an caveat though: the equivalence strictly holds true only for the whole, exponentially-sized set  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$ , undermining any direct application. Nevertheless, we can get interesting generalization bounds based on Equation 7.3 when we compute the right-hand side by uniform sampling of rados belonging in  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$  and by limiting ourselves to a set of linear size [Nock et al., 2015, Theorem 3].

Instead of concentrating on the losses of Theorem 66, we formulate an equivalence based on square loss. Such another relationship between the two worlds, examples and rados, suggests the existence of a wide theory underpinning those aggregated statistics. For a more complete view on the topic, the interested reader can consult Nock [2015].

We consider the empirical risk of an  $L_2$  regularized square loss  $\ell(x) = (1 - x)^2$ . Here  $\Gamma$  is a symmetric positive definite matrix.

$$\mathbb{E}_{\mathcal{S}} (1 - y\langle \boldsymbol{\theta}, \mathbf{x} \rangle)^2 + \boldsymbol{\theta}^\top \Gamma \boldsymbol{\theta} . \quad (7.4)$$

The loss admits a simple closed form solution:

$$\boldsymbol{\theta}^* \doteq \arg \min_{\boldsymbol{\theta}} \mathbb{E}_{\mathcal{S}} (1 - y\langle \boldsymbol{\theta}, \mathbf{x} \rangle)^2 + \boldsymbol{\theta}^\top \Gamma \boldsymbol{\theta} = \left( XX^\top + m \cdot \Gamma \right)^{-1} \boldsymbol{\pi} , \quad (7.5)$$

where  $X \doteq [x_1 | x_2 | \dots | x_m]$ , and so,  $XX^\top = \sum_i x_i x_i^\top$ . Remark that the computation of  $\boldsymbol{\theta}^*$  involves  $\boldsymbol{\pi}$ , the rado corresponding with the mean operator; labels do not appear anywhere else in the formula in light of label sufficiency. Let us define the following  $M$ -loss over rados via its empirical risk.

**Definition 67.** *The empirical risk associated with the  $M$ -loss over  $\mathcal{R}_{\mathcal{S}, \Sigma'_m}$  of classifier  $\boldsymbol{\theta}$  is:*

$$\frac{1}{2} \mathbb{V}_{\Sigma'_m} \langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle - \mathbb{E}_{\Sigma'_m} \langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle , \quad (7.6)$$

where expectation and variance are computed with respect to uniform sampling of  $\sigma$  in  $\Sigma'_m$ .

Equation (7.6) resembles a Markowitz mean-variance criterion [Markowitz, 1952] — hence the loss name —, with no coefficient for the risk aversion. What this means is that a good classifier trained on rados should have large “return” and small “risk”, where the risk is the variance of its predictions and the return is its inner product with the expected rado.  $M$ -loss and squared loss are linked via the next Theorem.

**Theorem 68.** *The following holds true for any  $\boldsymbol{\theta}$  and  $\mathcal{S}$ :*

$$\mathbb{E}_{(x,y) \sim \mathcal{S}} (1 - y\langle \boldsymbol{\theta}, \mathbf{x} \rangle)^2 = 1 + \frac{2}{m} \mathbb{V}_{\sigma \sim \Sigma_m} \langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle - \frac{4}{m} \mathbb{E}_{\sigma \sim \Sigma_m} \langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle . \quad (7.7)$$

The Proof in 7.7.1 simplifies the long derivation in Nock [2015]. The above equivalence still holds if both risks are regularized by  $\boldsymbol{\theta}^\top \Gamma \boldsymbol{\theta}$ , for any symmetric positive semi-definite  $\Gamma$ . Hence, minimizing the  $L_2$  regularized square loss over examples is equivalent to minimizing a regularized version of the  $M$ -loss, over the complete set of all rados.



Rados can also be thought as sufficient statistics, similarly to the mean operator.

**Lemma 69.** *The function  $\mathcal{S} \mapsto \mathcal{R}_{\mathcal{S}, \Sigma_m}$  that produces the whole set of rados  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$  is a sufficient statistic for variables  $(\mathbf{x}, y)$  with regard to the hypothesis space of linear classifiers  $\mathcal{H}$  and square loss, in the sense that, for any two sample  $\mathcal{S}$  and  $\mathcal{S}'$  and for any  $\theta \in \mathcal{H}$  we have:*

$$\begin{aligned} \mathbb{E}_{\mathcal{S}}(1 - y\langle\theta, \mathbf{x}\rangle)^2 - \mathbb{E}_{\mathcal{S}'}(1 - y\langle\theta, \mathbf{x}\rangle)^2 \text{ does not depend on } (\mathbf{x}, y) \\ \iff \mathcal{R}_{\mathcal{S}, \Sigma_m} = \mathcal{R}_{\mathcal{S}', \Sigma_m} . \end{aligned} \quad (7.8)$$

Proof in 7.7.2. As we have noticed already,  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$  contains the mean operator and hence sufficiency for the label variable is not a surprise. Although, the claim is about all examples in  $\mathcal{S}$ , not just the labels. Yet, the Lemma is in some sense trivial. In fact, the whole set of rados  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$  contains also all “the sums” of single edge vectors  $y_i x_i$  — when there is only one  $i$  such that  $\sigma_i$  in the definition of rados is equal to  $y_i$ . The knowledge of all edge vectors would be sufficient for all  $\mathcal{S}$  even without the use of rados, in case we use any margin loss with linear classifiers. Therefore, expressing the empirical risk with rados requires many more elements to achieve statistical sufficiency, *i.e.* the whole  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$ .

Therefore, why using rados instead of examples? The main reason is that rados are invariant to the selection of different solutions for entity resolution. For example, consider again Table 7.1. Since all classes are positive, computing a rado is just summing observations. Let  $\pi_{ij,kl}$  be the rado that sums those of examples  $e_{ij}$  and  $e_{kl}$ . Then, surprisingly, regardless of the solution to ER, this rado is the *same*:

$$(E1) \quad \pi_{11,22} = (1, -1, 1) + (-1, 1, 1) \quad (7.9)$$

$$= (0, 0, 2) \quad (7.10)$$

$$= (1, 1, 1) + (-1, -1, 1) = \quad (7.11)$$

$$= \pi_{12,21} (E2) . \quad (7.12)$$

This always holds in the (VP) setting: there a set of rado, potentially of exponential size, that matches the set that could be built *knowing* the true entity resolution. Next, we give the algorithm that builds these rados and it is communication efficient and easily parallelizable.

## 7.4 Building and learning from block rados

In our distributed setting, we extend the definition of rados in the following way. We let  $s \in \mathcal{J}$  denote a *signature*, a vector of shared attributes.

**Definition 70.** *The rado of signature  $s$  and peer  $P^j$  is:*

$$\pi_{(s,y)}^j \doteq \sum_{i=1}^m 1_{\{\text{proj}_{\mathcal{J}}(\mathbf{x}_i^j) = s \wedge y_i^j = y\}} y_i^j \mathbf{x}_i^j , \quad (7.13)$$

where  $\text{proj}_{\mathcal{I}}(\mathbf{z})$  denotes the restriction of a vector  $\mathbf{z}$  to  $\mathcal{I}$ .

In short,  $\pi_{(s,y)}^j$  sums edge vectors local to  $P^j$  whose examples match signature  $\mathbf{s}$  and class  $y$ . Intuitively, we can conceptualize those rados and expressing statistics *locally* sufficient for the examples sharing the same signature  $\mathbf{s}$  in the data of  $P^j$ . Let  $\mathcal{F}(\mathbf{z}) \subseteq \mathcal{X}$  be the set of features of  $\mathbf{z}$ . We also define, for any  $\mathcal{F}' \supseteq \mathcal{F}(\mathbf{z})$ ,  $\text{lift}_{\mathcal{F}'}(\mathbf{z})$  to be the vector  $\mathbf{z}'$  described using  $\mathcal{F}'$  such that  $\text{proj}_{\mathcal{F}(\mathbf{z})}(\mathbf{z}') = \mathbf{z}$  and  $\text{proj}_{\mathcal{F}' \setminus \mathcal{F}(\mathbf{z})}(\mathbf{z}') = \mathbf{0}$ . While  $\text{proj}_{\mathcal{F}}(\mathbf{z})$  removes coordinates of  $\mathbf{z}$ ,  $\text{lift}_{\mathcal{F}'}(\mathbf{z})$  “completes” the coordinates of  $\mathbf{z}$  with zeroes.

By analogy with entity resolution [Whang et al., 2009], we define *block rados* as rados, lifted to  $\mathcal{X}$ , that are the (weighted) sums of examples matching a particular signature and class in all peers.

**Definition 71.** For any  $\mathbf{s} \in \mathcal{J}$ ,  $y \in \{-1, 1\}$ , let  $m_{(s,y)}$  be the number of examples matching signature  $(\mathbf{s}, y)$ . Then a basic block (BB) rado for  $(\mathbf{s}, y)$  is:

$$\pi_{(s,y)} \doteq \sum_{j=1}^p \text{lift}_{\mathcal{X}}(\pi_{(s,y)}^j) - m_{(s,y)}(p-1) \cdot \text{lift}_{\mathcal{X}}(y \cdot \mathbf{s}) . \quad (7.14)$$

We need to subtract the second term to take into account that  $\mathbf{s}$  has already been summed  $m_{(s,y)}$  times by each peer. Let:

$$\mathcal{J}_* \doteq \{(\mathbf{s}, y) \in \mathcal{J} \times \{-1, 1\} : \exists j \in [p], \pi_{(s,y)}^j \neq \mathbf{0}\} . \quad (7.15)$$

This latter set, which can easily be computed from all peers, has cardinal  $m_* \doteq |\mathcal{J}_*| \leq m$ , and even  $m_* \ll m$  when few features are shared. We let:

$$\mathcal{R}_B \doteq \{\pi_{v_i}, \forall i \in [m_*]\} \quad (7.16)$$

denote the ordered set of each BB rado, each coordinate of  $v = (\mathbf{s}, y)$  being in one-one correspondence with an element of  $\mathcal{J}_*$ . A superset of  $\mathcal{R}_B$  is interesting, that considers all sums of vectors from  $\mathcal{R}_B$ :

$$\mathcal{R}_* \doteq \left\{ \sum_{i \in \mathcal{U}} \pi_{v_i}, \forall \mathcal{U} \subseteq [m_*] \right\} . \quad (7.17)$$

We call  $\mathcal{R}_*$  the set of *block rados*. Notice that we may have  $|\mathcal{R}_*| = \Omega(2^{\sum_j |\mathcal{S}^j|})$ . It is therefore intractable in general to *explicitly* compute  $\mathcal{R}_*$ . However,  $|\mathcal{R}_B| = O(\sum_j |\mathcal{S}^j|)$  and to compute it, we just need the set of  $\pi_{(s,y)}^j$ , hence a communication complexity that can be much smaller than  $\sum_j |\mathcal{S}^j|$ .

This set has exponential size. A possibility is to randomly subsample the set, along with proving good uniform convergence bounds for the  $M$ -loss — this can be done in the same way as in [Nock et al., 2015]. However, in the case of the square loss, greed pays twice: learning from all rados in  $\mathcal{R}_*$  may be both computationally cheap and accurate.

**Algorithm 11:** RadoCraft( $P^1, P^2, \dots, P^p$ )

---

**Input:** Peers  $P^1, P^2, \dots, P^p$   
 $\mathcal{R}_B \leftarrow \emptyset$   
**for**  $s \in \mathcal{J}, y \in \{\pm 1\}$   
    Let  $\boldsymbol{\pi}_{(s,y)} \leftarrow \mathbf{0} \in \mathbb{R}^d$   
    **for**  $j \in [p]$   
         $\boldsymbol{\pi}_{(s,y)} \leftarrow \boldsymbol{\pi}_{(s,y)} + \text{lift}_{\mathcal{X}}(\text{CRAFT}(s, y) \rightsquigarrow P^j)$   
     $\mathcal{R}_B \leftarrow \mathcal{R}_B \cup \{\boldsymbol{\pi}_{(s,y)}\}$   
**Output:**  $\mathcal{R}_B$

---

**7.4.1 Computation and optimality of block rados**

We do not have access to all rados because we do not assume to know the entity matching function. Yet, we are going to show a first result which is, in a sense, *stronger*:  $\mathcal{R}_*$  always belongs to  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$ . Therefore  $\mathcal{R}_*$  —potentially exponential-sized— gives us a set of rados that would have been built *from*  $\mathcal{S}$ , *had we known the perfect solution to entity matching*. So, even without carrying out entity matching, we have access to a potentially huge set of “ideal” rados which we can use to learn  $\boldsymbol{\theta}$  via the minimization of the empirical risk on  $M$ -loss.

**Theorem 72.** *In setting (VP), for any  $p \geq 2$ , any  $\mathcal{S}$ , any  $\mathcal{J}$ . Let  $\mathcal{R}_B$  be the output of Algorithm 11 and let  $\mathcal{R}_*$  its superset by Equation (7.17). Then,  $\mathcal{R}_* \subseteq \mathcal{R}_{\mathcal{S}, \Sigma_m}$ .*

Proof in 7.7.3. Furthermore, there exists a simple algorithm to build  $\mathcal{R}_B$ . Algorithm 11 summarizes the protocol. Each peer  $P^j$  crafts rados upon request of a particular signature and label; “ $\text{CRAFT}(s, y) \rightsquigarrow$ ” symbolizes a message sent, expecting  $\boldsymbol{\pi}_{(s,y)}^j$  in return. The computation of each rado for each peer can easily be performed in parallel. Algorithm 11 always provides the basis for the set  $\mathcal{R}_*$  of the “ideal” rados.

**7.4.2 Learning from all block rados**

The questions that remain are how we minimize the regularized  $M$ -loss and, more importantly, what subset of rados from  $\mathcal{R}_*$  we shall use. As already discussed, we choose “greediness” against randomization: instead of picking a (small) random subset of  $\mathcal{R}_*$ , we want to use them *all* because we know that all of them are “ideal” or close to being so via Theorems 72. Recall that  $|\mathcal{R}_*|$  may be of exponential size (in  $m, d, |\mathcal{J}_*|$ , etc.). We now show that if we consider all of  $\mathcal{R}_*$ , the optimal model learned from rados has an analytic expression which depends *only* on the rados of  $\mathcal{R}_B$ . In short, it is even *faster* to compute than  $\hat{\boldsymbol{\theta}}$  from  $\mathcal{S}$  in Equation (7.5), and can be directly computed from the output of Algorithm 11.

**Theorem 73.** *Let  $\hat{\boldsymbol{\theta}}$  be the minimizer of Equation (7.6). Then:*

$$\hat{\boldsymbol{\theta}} = \left( BB^\top + \dim_c(B) \cdot \Gamma \right)^{-1} B\mathbf{1} , \quad (7.18)$$

**Algorithm 12:** DRL( $P^1, P^2, \dots, P^p; \Gamma$ )

---

**Input:** Peers  $P^1, P^2, \dots, P^p, \Gamma, \gamma > 0$   
 $B \leftarrow \text{Column}(\text{RadoCraft}(P^1, P^2, \dots, P^p))$   
 $\theta \leftarrow (BB^\top + \gamma \cdot \Gamma)^{-1} B\mathbf{1}$   
**Output:**  $\theta$

---

where  $B$  stacks in columns the rados of  $\mathcal{R}_B$ , and  $\dim_c(B)$  is the number of columns of  $B$ .

Proof in 7.7.4. When  $m_* = m$ , each element of  $\mathcal{R}_B$  is in fact an example, and we retrieve Equation (7.5). One consequence of Theorem 73 is the following convergence property which we state informally: in the (VP) setting, for any  $\varepsilon \geq 0$ , there exists a minimal size for  $\mathcal{J}_*$  such that  $\hat{\theta}$  will be  $\varepsilon$ -close to  $\theta^*$ , where the closeness can be measured by  $\|\hat{\theta} - \theta^*\|_2$  or  $|\cos(\hat{\theta}, \theta^*)|$ . The statement of Distributed Rado-Learn (DRL) is in Algorithm 12. “column(.)” takes a set of vectors and put them in column in a matrix.

### 7.4.3 A more realistic setting

What happens if we drop the assumption of data being vertically partitioned (VP)? Or equivalently, what if examples are not shared by *all* peers? This is a much more realistic scenario. Since there is no shared ID — and the data may have been anonymized — we are not even in a situation where we can guarantee that a specific client of the bank *is*, or *is not*, a client of the insurance company. Thus, there may be significant unknown data “to reconstruct” the total sample  $\mathcal{S}$ , but we do not know which specific examples have missing features.

In this most general setting (G), it is possible to show that a very simple transformation of the rados, involving only the shared features, has the same properties so far described and for which Theorem 72 holds *in expectation*. We leave the details to Appendix 7.8 to avoid the heavier notation. However, in the next Section, we provide an experimental validation of our approach in both the settings.

## 7.5 Experiments

We evaluate the leverage that DRL provides compared to the peers, that would learn using only their local dataset. Each peer  $P^j$  estimates learns through a 10-folds stratified cross-validation minimization of regularized squared loss (Equation 7.5), where  $\gamma$  is also locally optimized through a ten-folds CV in set  $\{.01, 1.0, 100.0\}$ . DRL solves Equation 7.18 where  $\mathcal{R}_B$  is built using RadoCraft, with the set of all peers as input.

We have carried out a very simple optimization of the regularization matrix of DRL as a diagonal matrix which weights differently the shared features:

$$\Gamma = \text{diag}(\text{lift}_{\mathcal{X}}(\text{proj}_{\mathcal{J}}(\mathbf{1}))) + \gamma \cdot \text{diag}(\text{lift}_{\mathcal{X}}(\text{proj}_{\mathcal{X} \setminus \mathcal{J}}(\mathbf{1}))), \text{ for } \gamma \in \mathcal{G}. \quad (7.19)$$

$\gamma$  is optimized by a 10-folds CV on  $\mathcal{J}_*$ , defined in Equation 7.15. CV is performed on *rados* as follows: first,  $\mathcal{R}_B$  is split in 10 folds,  $\mathcal{R}_{B,\ell}$ , for  $\ell = 1, 2, \dots, 10$ . Then, we repeat for  $\ell = 1, 2, \dots, 10$  (and then average) the following CV routine:

1. DRL is trained using  $\mathcal{R}_B \setminus \mathcal{R}_{B,\ell}$
2. DRL's solution is evaluated on "test rados" by computing the average  $M$ -loss

The expression of  $\Gamma$  for rados exploits the idea that the estimations related to a shared feature can be much more accurate than for another, non shared feature.

### 7.5.1 Domain generation

We ran experiments on a dozen UCI domains [Bache and Lichman, 2013]. Only two are fully detailed here, the rest are presented in Appendix 7.9. They are *mice* ( $m \times d \approx (1K \times 70)$ ) and *musk* ( $\approx (6.5K \times 166)$ ). For each domain, we have varied (i) the number of peers  $p$ , (ii) the number of shared features  $\dim(\mathcal{J})$ , and (iii) the number  $b$  of numeric modalities ("bins") each shared feature was reduced to (it controls the size of  $\mathcal{J}_*$ ). The training sample is split among peers, each keeping record of  $\mathcal{I}$  and its own features (non shared features are evenly partitioned among peers). Finally, for some  $p_s \in [0, 1]$ , each peer  $P^j$  selects a proportion  $p_s$  of its examples index and for each of them, another peer  $P^{j'}$ , chosen at random, gets the example as well (on its own set of features  $\mathcal{X}^{j'}$ ). This policy simulate scenario (G). When  $p_s = 0$ , this is setting (VP). We then run *all* algorithms for *each* value  $p, \dim(\mathcal{J}), b, p_s$ . As we shall see,  $b$  appears to have a relatively small influence compared to the other factors, so we mainly report results combining various values for  $p, \dim(\mathcal{J})$  and  $p_s$ , for the range of values of  $p, \dim(\mathcal{J})$  specified in Table 7.6, and for  $p_s \in \{0.0, 0.2\}$ . We have chosen  $b = 4$  for all domains, except when it is not possible (if for example all features are boolean), in which case we pick  $b = 2$ .

### 7.5.2 Metric

We used two metrics. The first:

$$\Delta \doteq p_{\text{err}}(\text{DRL}) - \min_j p_{\text{err}}(P^j) \in [-1, 1], \quad (7.20)$$

is the test error for DRL minus that of the *optimal* peer *in hindsight* (since we consider the peer's test error). when  $\Delta < 0$ , DRL beats *all* peers. For example, Table 7.3 (left) provides the results obtained on UCI domain mice. We see that for almost all combinations of  $p$  and  $\dim(\mathcal{J})$ , DRL beats all peers.

The smallest test error obtained for a peer among all runs for each domain: this is an indication of the room of improvement for DRL, and it also shows that in general, at least some (and in fact most) peers were always very significantly better than random guessing, a safe-check that DRL is not just beating unbiased coins.

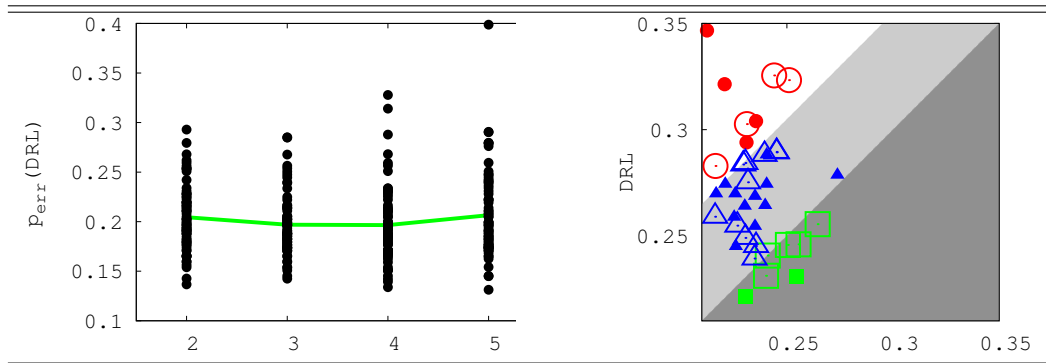


Table 7.2: *Left*: test error of DRL on domain ionosphere, as a function of the number of bins, aggregating all values of the number of peers  $p$  and number of shared features  $\dim(\mathcal{J})$  used; the green line denotes the average values. *Right*: scatter plot of the test error of DRL ( $y$ ) vs that of the Oracle (learning using the complete entity-resolved domain). Points in the dark grey area (green) denote better performances of DRL; points in the light grey area (blue) denote better performances of the Oracle (but not statistically better). Points in the white area (red) denote *statistically* better performances of the Oracle (filled points:  $p_s = 0.2$ ; empty points:  $p_s = 0.8$ ).

To evaluate the statistical significance, we compute:

$$q \doteq \text{proportion of peers } \textit{statistically} \textit{ beaten by DRL} . \quad (7.21)$$

To compute the test, we use the powerful Benjamini-Hochberg procedure [Benjamini and Hochberg, 1995] on top of paired  $t$ -tests with  $q^* = p\text{-val} = 0.05$ .  $q = 0.8$  surface helps see when DRL *statistically beats all peers*. For example, Table 7.3 (right) displays that DRL does not always *statistically* beat all peers when  $\Delta < 0$ , yet it manages to statistically beat all of them in approximately one third to one half of the total tests, which implies that, on this domain, there is a significant chance that DRL improves on the peers, regardless of their number and the number of shared features.

### 7.5.3 Results

All domains display that there exists regimes  $(p, \dim(\mathcal{J}))$  for which DRL improves on all peers, in some cases significantly. Sometimes, the improvement is sparse (phishing, creditcard), but sometimes it is quite spectacular and in fact (almost) systematic (page, ionosphere, steelplates). Domain steelplate’s case is interesting, since the so-called *Oracle*, *i.e.* the learner that learns from the complete training fold *before* it is split among peers — and therefore knows the solution to entity matching —, has for this domain almost optimal error, but local peers are in fact very far from this optimum. This indicates that many features, properly combined, are necessary to attain the best performances. DRL’s performances are close to the Oracle, which accounts for the huge gap in classification compared to peers — sometimes, DRL’s test error is smaller than that of the *best* peer by more than 20% —, and so it seems

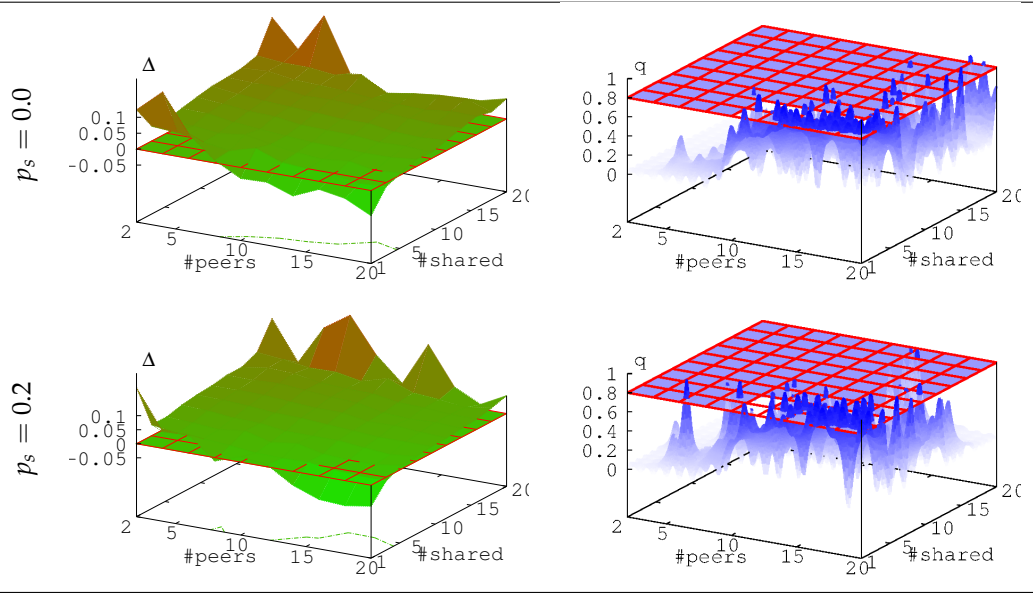


Table 7.3: Results on domain mice: fig/rado of  $\Delta \doteq p_{\text{err}}(\text{DRL}) - \min_j p_{\text{err}}(P^j)$  (left) and  $q = \text{prop. peers simultaneously beaten by DRL}$  (right) as a function of the number of peers  $p$  and the number of shared features  $\dim(\mathcal{J})$ . On mice,  $\min_j p_{\text{err}}(P^j) = 0.30$ . Top: proportion of shared examples  $p_s = 0.0$  (VP); bottom: proportion of shared examples  $p_s = 0.2$  (G). The isoline on the left fig/rado is  $\Delta = 0$ .

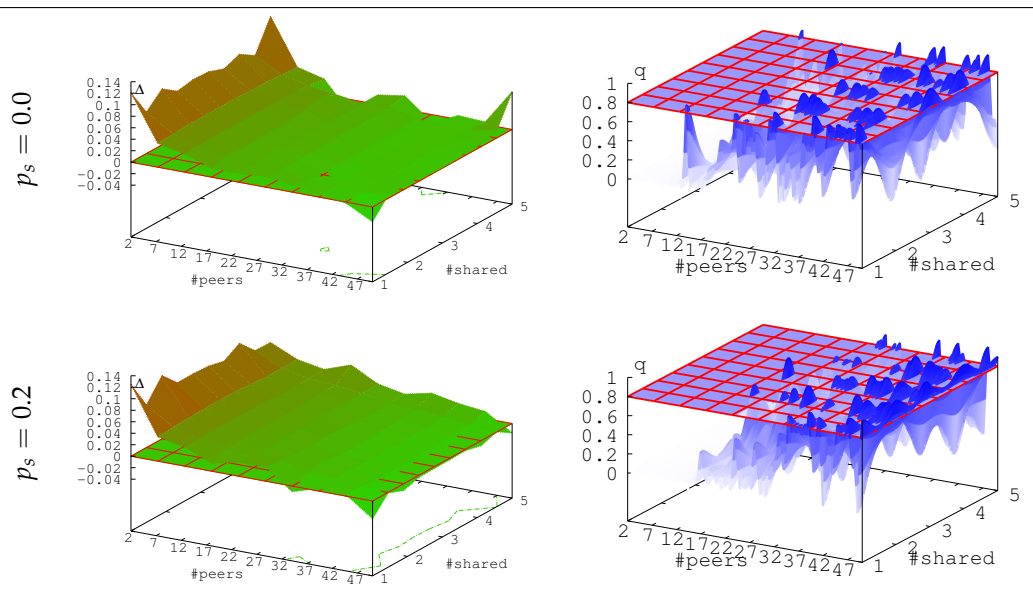


Table 7.4: Results on domain musk, using the same convention as Table 7.3. On musk,  $\min_j p_{\text{err}}(P^j) = 0.25$

that DRL indeed successfully bypasses entity matching to learn a classifier that almost matches the Oracle’s performances, and therefore represents a very significant leverage of each peer’s data.

For more detailed results, Table 7.2 (left) displays that binning indeed does not affect significantly DRL on average, which is also good news, since it means that there is no restriction on the shared features for DRL to perform well: shared features can be binary, or categorical with any number of modalities. Table 7.7 displays that while the CV tuning of  $\Gamma$  offers leverage to DRL (*vs*  $\Gamma = \text{Id}_d$ ) in general (firmteacher), there are some (rare) domains (mice) on which relying on the simplest  $\Gamma = \text{Id}_d$  improves upon the results of CV. This, we believe, comes from the fact that CV as we have carried out is certainly not optimal because one rado can aggregate any number of examples. Finally, Table 7.2 (right) drills down a bit more into the performances of DRL with respect to those of the Oracle on a domain for which DRL obtains somehow “median” performances among all domains, sonar. The Oracle (10-folds CV from the *total* ER’ed  $\mathcal{S}$ ) is *idealistic* since in general we do not know the solution to ER, yet it gives clues on how close DRL may be from the “grail”. Interestingly, DRL comes frequently under the statistical significance radar ( $\alpha = 0.05$ ). In notable cases (more frequent as  $p_s$  increases), DRL beats Oracle — but not significantly. Aside from theory, these are good news as DRL does not assume ER’ed data, and uses an amount of data which can be  $\sim p^2$  times *smaller* than Oracle.

## 7.6 Discussion and references

We remark that our framework is not formally comparable with ER, since the two address different problems. On the one hand, ER has a much broader applicability than the problem we are interested in here; learning on distributed datasets is less general than ER: in fact, we show a solution that bypasses ER. On the other hand, *learning-based* ER [Bilenko and Mooney, 2003] as well as manifold alignment techniques [Lafon et al., 2006] are viable only knowing some ground truth matches — which are not required for working with rados. From another perspective, in concert with the *open issues* in Getoor and Machanavajjhala [2012], we study ER as component of a pipeline for classification, and highlight how matching is not necessary for the purpose of learning.

In spite of those considerations, we can still draw comparisons with methods that learn on top of data merged through ER (Table 7.5). In both settings, no ID is shared between datasets but some attributes must be so, in order to allow entities comparison for matching or for building rados. Obviously, entity matching does not require the labels to be one of those shared attributes, while this is a fundamental hypothesis of our approach. Although, it is not as restrictive as it might seem at first: if just one peer has labels, then *all* can obtain labels on their own data, via learning from label proportions (Chapter 4): the label handling peer computes the label proportions per each block; the “bags” are defined by examples matching a particular signature. Proportions are then shared among all other peers, which can



metric	ER + learning	Algorithms 11+12
Hp: shared IDs	no	no
Hp: shared variables	necessary	necessary
Hp: shared labels	no	may be relaxed
Fusion / Rados crafting	$O(m^2/m^*T_{sim})$	$O(m)$
Communication	$m \times d$	$m^* \times d$
Learning problem	$m \times d$	$m^* \times d$
Privacy	complex	many guarantees

Table 7.5: Multiple metrics of comparison between learning on top of ER and our approach. Time complexity are estimated for 2 peers in the (VP) scenario, assuming all blocks of equal size. “Hp” is short for hypothesis. See Section 7.6 for details.

train a classifier with them so as to estimate a label for each observation.

To discuss time complexity, let us consider a simplified problem with only 2 peers with  $m$  examples each in the (VP) scenario. In terms of complexity of fusion, if we assume that examples are uniformly distributed in the blocks, each block has size  $m/m^*$  ( $m^*$  is the number of blocks). DRL builds each block rado in time  $O(m/m^*)$ , with total cost linear in  $m$ . ER takes  $O((m/m^*)^2 \cdot T_{sim})$  to match entities in each of the  $m^*$  blocks, where  $T_{sim}$  is the cost of evaluation any similarity function; learning-based methods spend additional time for training; advanced blocking strategies can reduce the average complexity [Bilenko et al., 2006; Whang et al., 2009; Whang and Garcia-Molina, 2012].

Most literature on distributed learning is concerned with limiting communication and designing optimal strategies for merging models [Balcan et al., 2012; Liu and Ihler, 2014]; beside that, previous works focus on horizontal split by observations, with few exceptions [Liu and Ihler, 2012]. In contrast, we exploit what is sufficient to merge *about the data*. The communication protocol is extremely simple. Once rados are crafted locally, they are sent to a central learner in one shot. By Theorem 73, only  $d$ -dimensional  $m^*$  block rados are needed. *Data is not accessed anymore* and learning takes place centrally. Moreover, rados help with data compression, being  $m^* \times d$ ,  $m^* \ll m$  the problem size. ER needs to transfer and learn from all entities, for a total size of  $m \times d$ .

Learning on data described by different feature sets is the topic of multiple view learning and co-training [Blum and Mitchell, 1998; Sindhvani et al., 2005]. To the best of our knowledge, co-training with unknown matches has not been addressed before. Brefeld et al. [2006] present a multi-view distributed algorithm with co-regularization; although it requires matches for all unlabeled examples.

In settings with multiple data providers, privacy can be crucial [Balcan et al., 2012]. The peers have to trade off model enhancements and information leaks. A learner receives rados to train the model; this can be done by one of the peers, or by a third party — paralleling multi-party ER scenarios [Christen, 2006]. The only information sent through the channel consists of rados, while examples, with their individual sensitive features, are never shared. Computational complexity results

on reconstruct-ability of examples have been proven in Nock et al. [2015], along with NP-HARD characterizations, and protection in the sense of differential privacy [Dwork, 2011; Dwork and Roth, 2013]. Regarding ER, since matching has the potential of de-anonymizing the entities, privacy is usually a very relevant issue to address [Christen, 2006]. However, solutions are not straightforward, as proven by the vast amount of research on the topic [Vatsalan et al., 2013]. Techniques based on partial share of attributes, anonymization or hashing can severely impair the process.

The key message of the Chapter is that entity matching addresses a very general and difficult problem but, in the comparatively restricted context of supervised learning from distributed datasets, accurate learning evading the pitfalls of entity matching is possible with Rademacher observations. Rados have another advantage: they offer a cheap, easily parallelizable material which somehow “compresses” examples while allowing accurate learning. They also offer a readily available solution for guaranteed private exchange of data in a distributed setting. Finally, some domains display that there is significant room space for improvement of how cross-validation of optimized parameters is performed. This interesting problem comes in part from the fact that statistical properties of cross-validation on rados are *not* the same as when carried out on examples; see the recent study of Nock [2016].

## 7.7 Appendix: proofs

### 7.7.1 Proof of Theorem 68

First, we remark that  $\mathbb{E}_{\Sigma_m}[\langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle] = \langle \boldsymbol{\theta}, \mathbb{E}_{\Sigma_m}[\boldsymbol{\pi}_\sigma] \rangle = (1/2) \cdot \langle \boldsymbol{\theta}, \boldsymbol{\pi} \rangle$ , since each example participates to half of the  $2^m$  rados. Letting  $v \doteq 2^{m+2} \cdot \mathbb{V}_{\Sigma_m}[\langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle]$ , we also have:

$$v = 4 \cdot \sum_{\sigma \in \Sigma_m} \left( \langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle - \frac{1}{2} \cdot \langle \boldsymbol{\theta}, \boldsymbol{\pi} \rangle \right)^2 \quad (7.22)$$

$$= \sum_{\sigma \in \Sigma_m} \left( \sum_i \sigma_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle \right)^2 \quad (7.23)$$

$$= \sum_{\sigma \in \Sigma_m} \left[ \sum_{i=1}^m \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle^2 + \sum_{i=1}^m \sum_{i' \neq i} \sigma_i \sigma_{i'} \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle \langle \boldsymbol{\theta}, \mathbf{x}_{i'} \rangle \right] \quad (7.24)$$

$$= 2^m \cdot \sum_{i=1}^m \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle^2 + \sum_{i=1}^m \sum_{i' \neq i} v_{ii'} \cdot \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle \boldsymbol{\theta}^\top \mathbf{x}_{i'} , \quad (7.25)$$

with  $v_{ii'} \doteq \sum_{\sigma \in \Sigma_m} \sigma_i \sigma_{i'}$ . For any  $i \neq i'$ ,  $\sigma_i \sigma_{i'}$  takes exactly the same number of times value  $+1$  and value  $-1$ , and so  $v_{ii'} = 0, \forall i \neq i'$ . We get from Equation 7.25:

$$\mathbb{V}_{\Sigma_m}[\langle \boldsymbol{\theta}, \boldsymbol{\pi}_\sigma \rangle] = (1/4) \cdot \sum_{i=1}^m \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle^2 = (1/4) \cdot \sum_{i=1}^m (y_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 . \quad (7.26)$$

Finally,

$$1 - \frac{2}{m} \cdot \sum_{i=1}^m y_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle + \frac{1}{m} \cdot \sum_{i=1}^m (y_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 = \frac{1}{m} \cdot \sum_i (1 - y_i \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)_2^2 . \quad (7.27)$$

### 7.7.2 Proof of Lemma 69

As with the sufficiency of mean operator, we need to show the double implication that defines sufficiency for  $(x, y)$ . This is trivial by applying Theorem 68.

### 7.7.3 Proof of Theorem 72

We give a proof sketch. Once three simple facts are established in the (VP) setting, the Theorem follows. (a) The ground truth entity matching exists. (b) Any BB rado for pair  $(s, y)$  would be obtained as a rado summing the contributions of all examples in  $\mathcal{S}$  matching the corresponding signature  $s$  and class  $y$ . (c) We obtain  $\mathcal{R}_B \subseteq \mathcal{R}_{\mathcal{S}, \Sigma_m}$ , from which follows the Theorem's statement with Equation (7.18) and the fact that any sum of a subset of rados in  $\mathcal{R}_B$  would also be in  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$  since an example cannot match two distinct couples (signature, class).

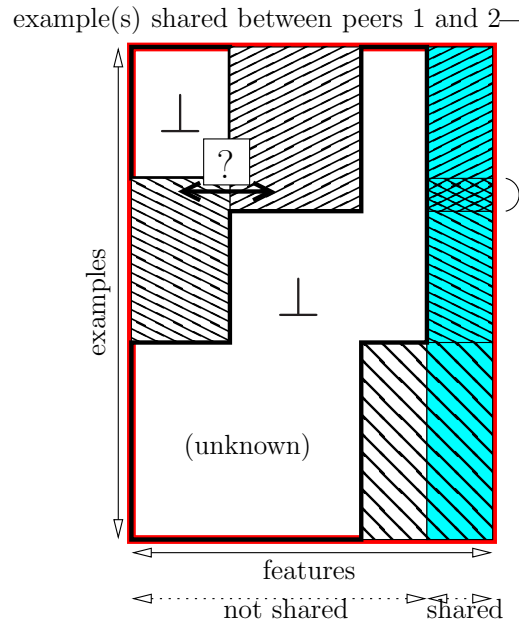


Figure 7.3: Schematic view of our setting, with  $p = 3$  peers. Some features (cyan) are described in each peer and one of these shared features is a class. Non-shared features are split among peers. The *total* sample  $\mathcal{S}$  is figured by the red rectangle. In the more general setting (G), it is not known whether one example, viewed by a peer, also exists in other peers' datasets. In this case, there may be a lot of missing data ( $\perp$ ), but it is not known of which examples.

#### 7.7.4 Proof of Theorem 73

The proof uses the following trick: consider any sample  $\mathcal{S}'$  such that its edge vectors match the basic block rados. Remark that  $XX^\top = \sum_i (y_i x_i)(y_i x_i)^\top$  in Equation (7.5) depends only on edge vectors. Thus, since  $\pi = B\mathbf{1}$ , the optimal square loss classifier on  $\mathcal{S}'$  is  $\hat{\theta}$  in Equation (7.18), which, through Theorem 68, is also the optimal classifier of the empirical risk associated with  $M$ -loss.

### 7.8 Appendix: extension to the more general setting

We extend here the algorithms from Section 7.4 to setting (G). Figure (7.3) depicts this more challenging learning setting. We do not assume that the peers handle the same examples, and therefore  $m \geq m_j, \forall j \in [p]$ . However, hypothesis (7A) still holds, that is, the peers share the same set of features.

The definition of rados itself needs to be upgraded first. In the (VP) setting we could avoid an obvious complication that we face here: in the computation of the basic block rados, we need to rescale the signature in order to take into account differences in the number of examples per block for each peer. First, we redefine the rados, by a projection onto the set of features that are not shared among peers.

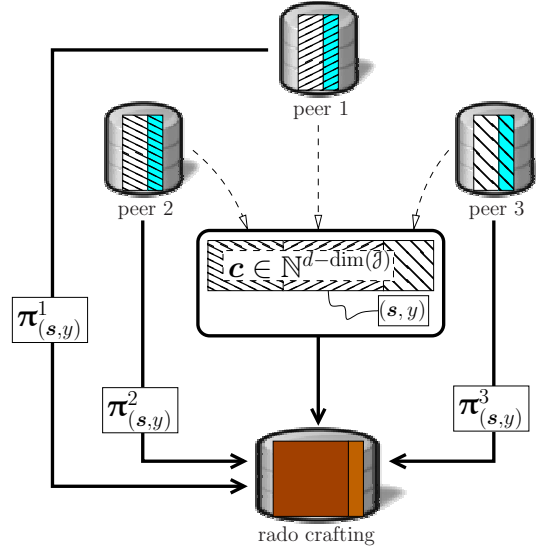


Figure 7.4: Communication for one BB rado, with  $(s, y) \in \mathcal{J}_*$ . Counter  $c$  is defined in Algorithm RadoCraft (see text).

**Definition 74.** For any  $s \in \mathcal{J}$ ,  $y \in \{-1, 1\}$  and  $P^j$ , a rado is:

$$\pi_{(s,y)}^j \doteq \text{proj}_{\mathcal{X}^j \setminus \mathcal{J}} \left( \sum_{i=1}^{m_j} 1_{\text{proj}_{\mathcal{J}}(x_i^j) = s \wedge y_i^j = y} y_i^j \cdot x_i^j \right). \quad (7.28)$$

The definition of  $u$ -basic block rados follows.

**Definition 75.** For any  $s \in \mathcal{J}$ ,  $y \in \{-1, 1\}$ ,  $u \in \mathbb{R}$ , the  $u$ -basic block (BB) rado for pair  $(s, y)$  is:

$$\pi_{(s,y,u)} \doteq u \cdot \text{lift}_{\mathcal{X}}(y \cdot s) + \sum_{j=1}^p \text{lift}_{\mathcal{X}}(\pi_{(s,y)}^j). \quad (7.29)$$

The set of block rados is upgraded accordingly. Recall that  $\mathcal{J}_* \doteq \{(s, y) \in \mathcal{J} \times \{-1, 1\} : \exists j \in [p], \pi_{(s,y)}^j \neq \mathbf{0}\}$  and that  $m_* \doteq |\mathcal{J}_*|$ . Then, for any  $\mathbf{u} \in \mathbb{R}^{m_*}$ , we let:

$$\mathcal{R}_B^{\mathbf{u}} \doteq \{\pi_{v_i}^{u_i}, \forall i \in [m_*]\} \quad (7.30)$$

denote the set of each  $u_i$ -BB rado, each coordinate of  $\mathbf{u}$  being in one-one correspondence with an element of  $\mathcal{J}_*$  (represented by  $v_i$ ). Finally, the set of  $\mathbf{u}$ -block rados is:

$$\mathcal{R}_*^{\mathbf{u}} \doteq \left\{ \sum_{i \in \mathcal{U}} \pi_{v_i}^{u_i}, \forall \mathcal{U} \subseteq [m_*] \right\}. \quad (7.31)$$

We state (without proof) that in the general settings (G) there exists  $\mathbf{u} \in \mathbb{R}^{m_*}$  such that  $\mathcal{R}_*^{\mathbf{u}}$  belongs to  $\mathcal{R}_{\mathcal{S}, \Sigma_m}$  in expectation. This is now obviously more difficult

**Algorithm 13:** RadoCraft( $P^1, P^2, \dots, P^p$ )

---

**Input** Peers  $P^1, P^2, \dots, P^p$   
 $\mathcal{R}_B^u \leftarrow \emptyset$   
**for**  $s \in \mathcal{J}, y \in \{\pm 1\}$   
 $\pi_{(s,y)} \leftarrow \mathbf{0} \in \mathbb{R}^d, \mathbf{c} \leftarrow \mathbf{0} \in \mathbb{N}^{d-\dim(\mathcal{J})};$   
**for**  $j \in [p]$   
 $\pi_{(s,y)} \leftarrow \pi_{(s,y)} + \text{lift}_{\mathcal{X}}(\text{CRAFT}(s, y) \rightsquigarrow P^j)$   
 $u \leftarrow (\mathbf{1}^\top \mathbf{c})(d - \dim(\mathcal{J}))^{-1}$   
 $\mathcal{R}_B^u \leftarrow \mathcal{R}_B^u \cup (u \cdot \text{lift}_{\mathcal{X}}(y \cdot s) + \pi_{(s,y)})$   
**Return**  $\mathcal{R}_B^u$

---

to tackle than in Theorem 72, since there may be a large amount of missing data ( $\perp$  in Figure 7.2) and there would be no one-one correspondence between the peers' examples in general. Yet, there is an interesting property which can be shown in the following (R)andomized model: each peer's features remain fixed but there exists a fixed  $\eta \in [0, 1]^m$  such that example  $i$  has probability  $\eta_i$  to be seen by a peer. Let  $\bar{\mathcal{S}}$  denote the "expected" sample, where each example is weighted by its probability. For any signature  $s$  and class  $y$ ,  $\mathbb{E}[\pi_{(s,y)}]$  denotes the expected rado put in  $\mathcal{R}_B^u$  in Algorithm 13.

**Theorem 76.** Under (R),  $\forall (s, y) \in \mathcal{J}_*, \mathbb{E}[\pi_{(s,y)}] \in \mathcal{R}_{\bar{\mathcal{S}}, \Sigma_m}$ .

Therefore, under setting (G), if examples are "seen" independently at random by peers, the expected output of Algorithm 13 still meets the guarantees of Theorem 72 with respect to the expected sample. The fact that  $\mathcal{R}_B^u \subseteq \mathcal{R}_{\mathcal{S}, \Sigma_m}$  from Theorem 72 is also a consequence of Theorem 76 for  $\eta = \mathbf{1}$ .

Algorithm 13 is a variation of the original Algorithm 11 which takes care of the computation of  $u$ . Specifically,  $P^j$  does the following:

- it computes and return  $\pi_{(s,y)}^j$ ; let  $C_j$  be the number of examples that are counted in the sum in Equation (7.28);
- it updates counter vector  $\mathbf{c}$ : for each feature  $k \notin \mathcal{J}$  it possesses in its local dataset, it does  $c_k \leftarrow c_k + C_j$ ;

Letting  $v_i \doteq (s, y) \in \mathcal{J}_*$ , the corresponding value of  $u_i$  is given by:

$$u_i \doteq (\mathbf{1}^\top \mathbf{c})(d - \dim(\mathcal{J}))^{-1}, \quad (7.32)$$

which is guaranteed to be non-zero since  $v_i \in \mathcal{J}_*$ .

## 7.9 Appendix: additional experimental results

domain	$m$	$d$	$\min_j p_{\text{err}}(P^j)$	$p$	$\dim(\mathcal{J})$	results
Wine	178	12	0.07	$\{2, 3, \dots, 8\}$	$\{1, 2, 3, 4\}$	Table 7.15
Sonar	208	60	0.29	$\{2, 3, \dots, 16\}$	$\{1, 2, \dots, 20\}$	Table 7.9
Ionosphere	351	33	0.20	$\{2, 3, \dots, 9\}$	$\{1, 2, \dots, 9\}$	Table 7.11
Mice	1 080	77	0.30	$\{2, 3, \dots, 20\}$	$\{1, 2, \dots, 20\}$	Table 7.3
Winered	1 599	11	0.26	$\{2, 3, \dots, 7\}$	$\{1, 2, 3, 4\}$	Table 7.12
Steelplates	1 941	33	0.16	$\{2, 3, \dots, 14\}$	$\{1, 2, \dots, 5\}$	Table 7.18
Statlog	4 435	36	0.05	$\{2, 3, \dots, 30\}$	$\{1, 2, \dots, 5\}$	Table 7.17
Winewhite	4 898	11	0.32	$\{2, 3, \dots, 7\}$	$\{1, 2, 3, 4\}$	Table 7.13
Page	5 473	10	0.21	$\{2, 3, \dots, 6\}$	$\{1, 2, 3, 4\}$	Table 7.8
Musk	6 598	166	0.25	$\{2, 3, \dots, 50\}$	$\{1, 2, 3, 5\}$	Table 7.4
Firmteacher	10 800	16	0.26	$\{2, 3, \dots, 7\}$	$\{1, 2, \dots, 7\}$	Table 7.10
Phishing	11 055	30	0.11	$\{2, 3, 4, 5\}$	$\{1, 2, 3, 4\}$	Table 7.14
Credit card	14 599	23	0.32	$\{2, 3, \dots, 18\}$	$\{1, 2, \dots, 5\}$	Table 7.16

Table 7.6: UCI domains used in our experiments Bache and Lichman [2013], with for each the indication of the total number of features ( $d$ ), examples ( $m$ ) and the error of the optimal peer in hindsight obtained in our experiments,  $\min_j p_{\text{err}}(P^j)$ . Two of the right columns present, for each domain, the range of values for the number of peers ( $p$ ) and the number of shared features ( $\dim(\mathcal{J})$ ) considered. Experiments are performed considering *all* possible combinations of values of  $p$  and  $\dim(\mathcal{J})$  within the allocated sets. The rightmost column points to the Table collecting specific results for each domain.

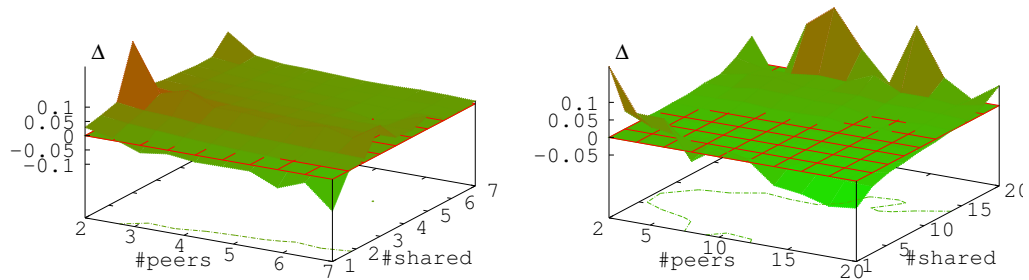


Table 7.7: Results of the dummy regularized DRL ( $\Gamma = \text{Id}_d$ ) on domains firmteacher (left) and mice (right), following the convention of Table 7.11 ( $p_s = 0.2$ ).

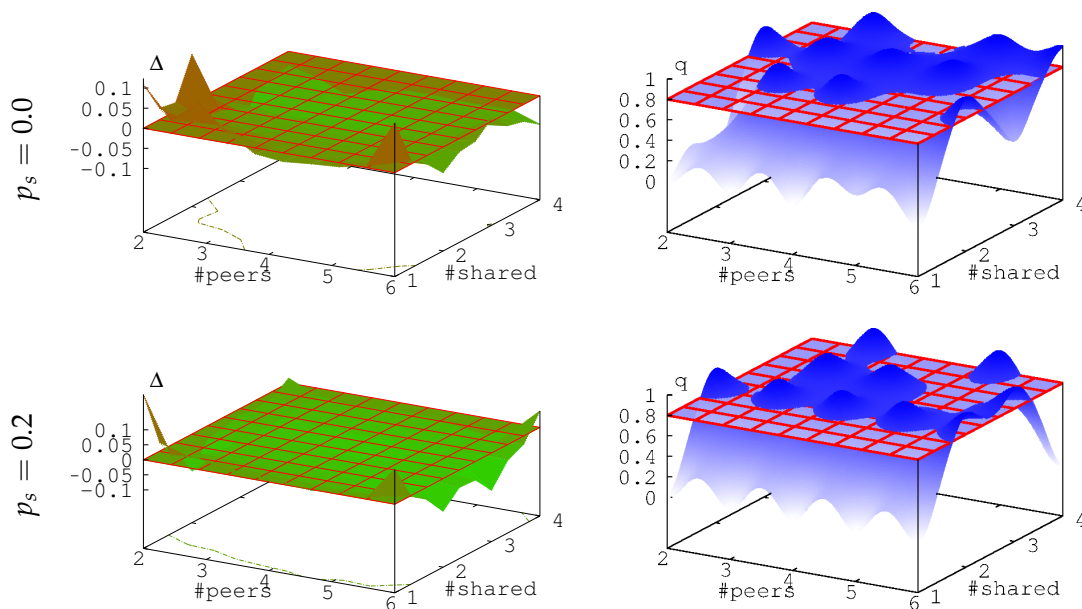


Table 7.8: Results on domain sonar, using the same convention as Table 7.3.



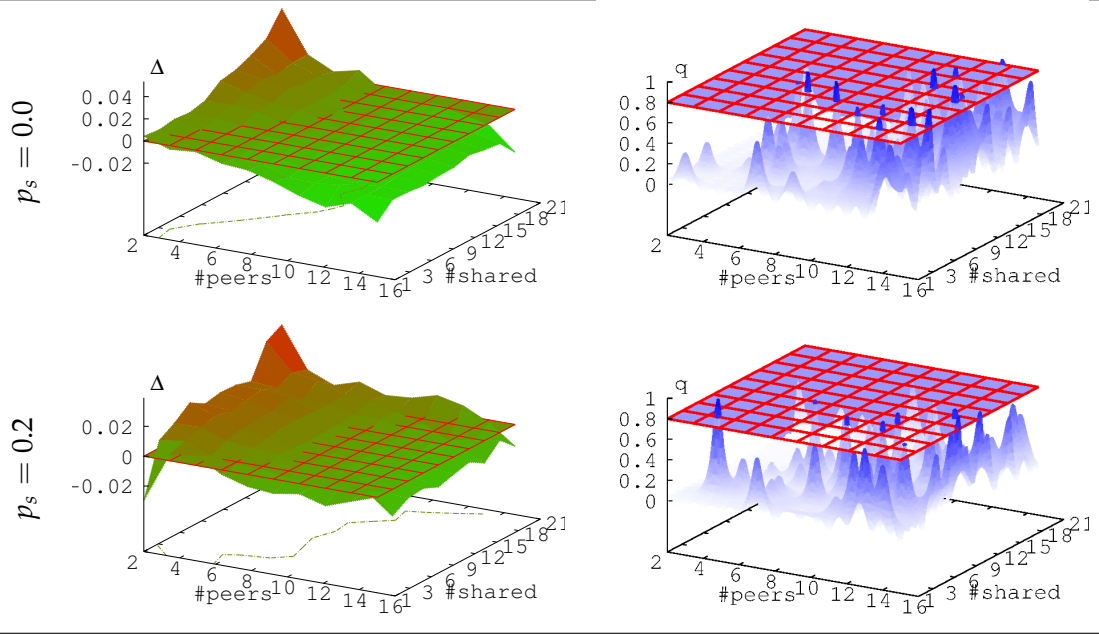


Table 7.9: Results on domain sonar, using the same convention as Table 7.3.

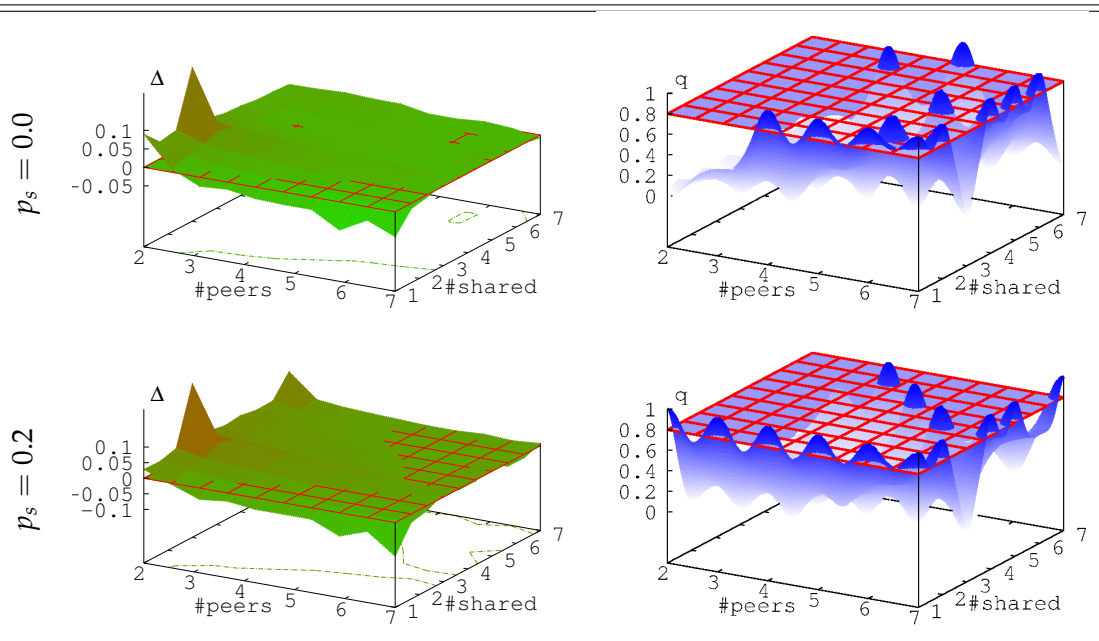


Table 7.10: Results on domain firmteacher, using the same convention as Table 7.3.

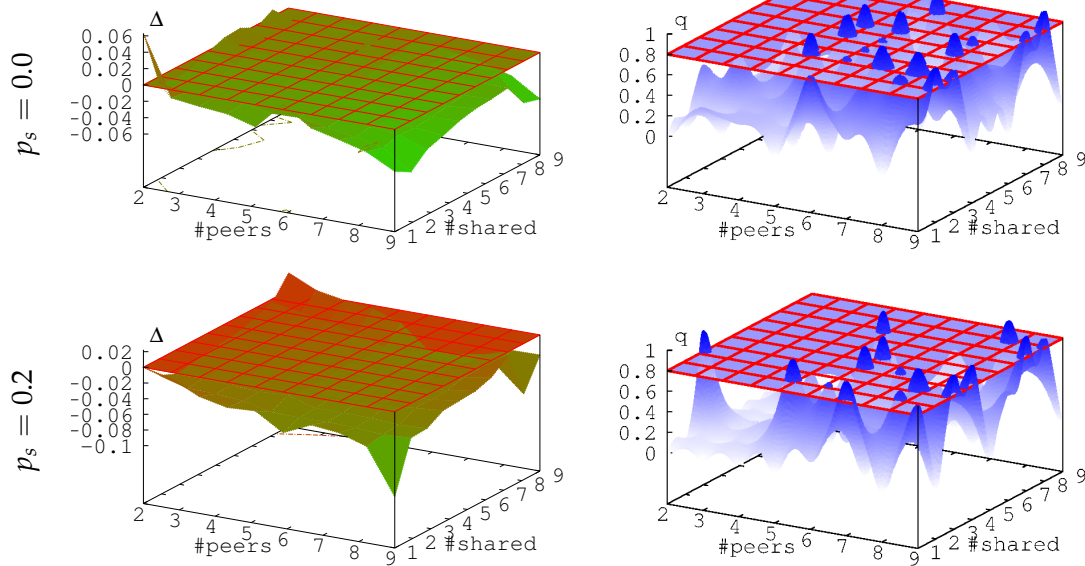


Table 7.11: Results on domain ionosphere, using the same convention as Table 7.3.

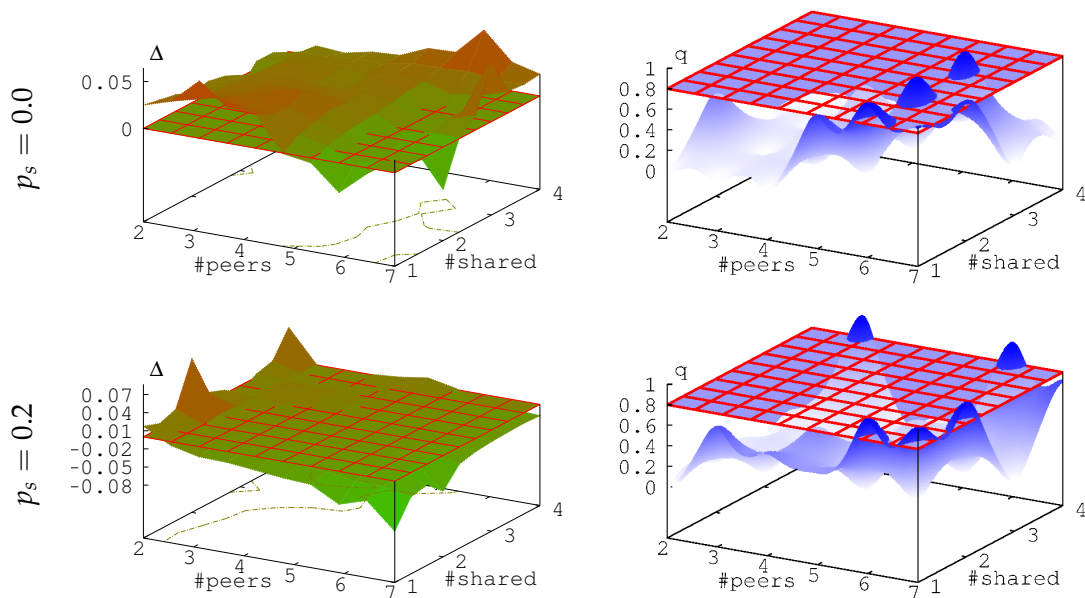


Table 7.12: Results on domain winered, using the same convention as Table 7.3.

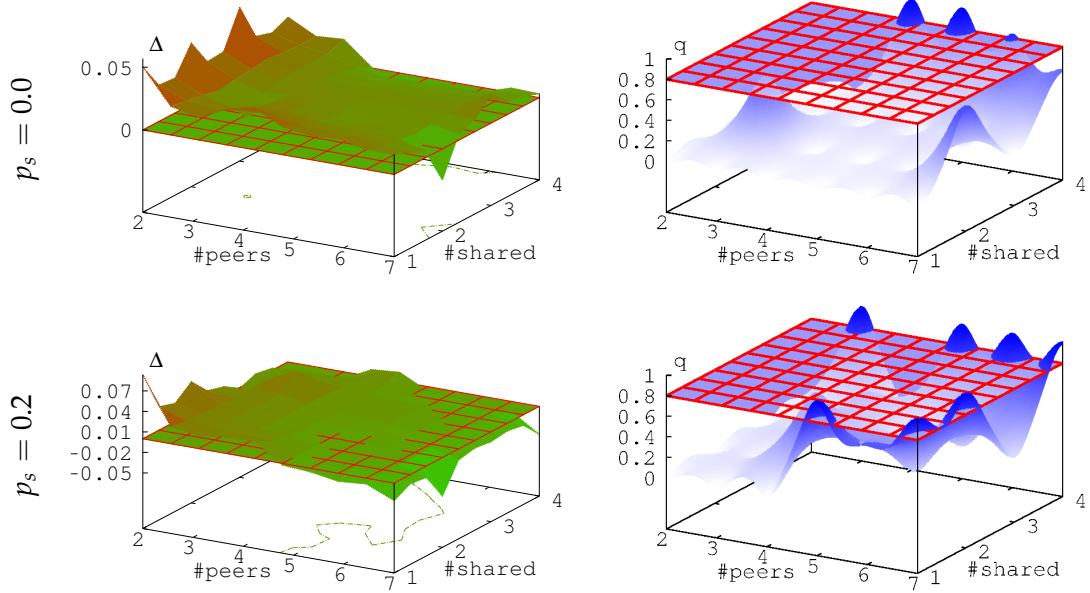


Table 7.13: Results on domain winewhite, using the same convention as Table 7.3.

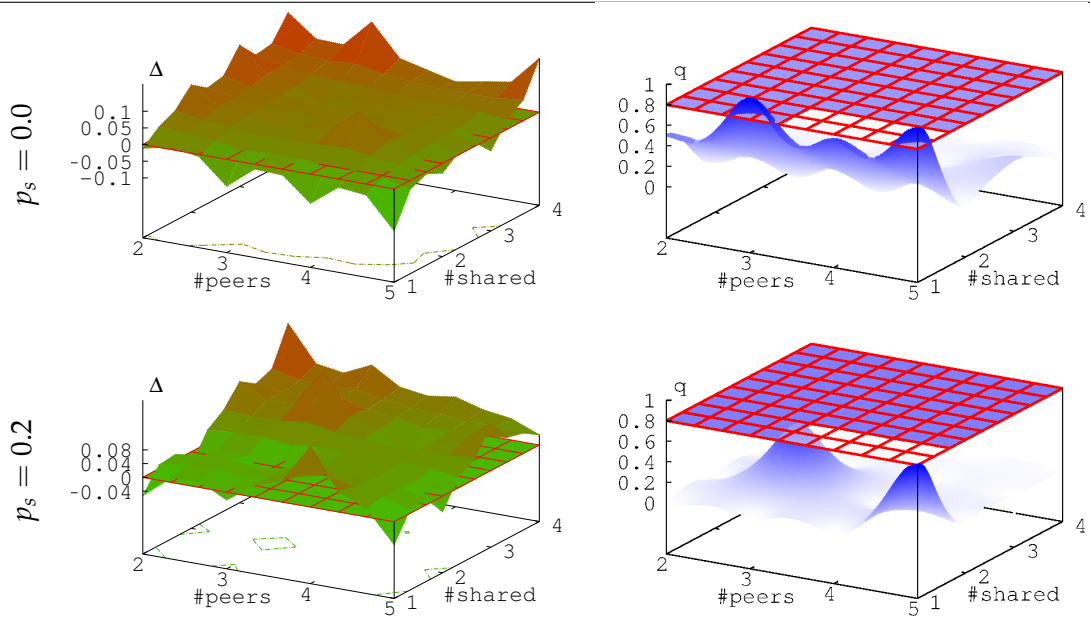


Table 7.14: Results on domain phishing, using the same convention as Table 7.3.

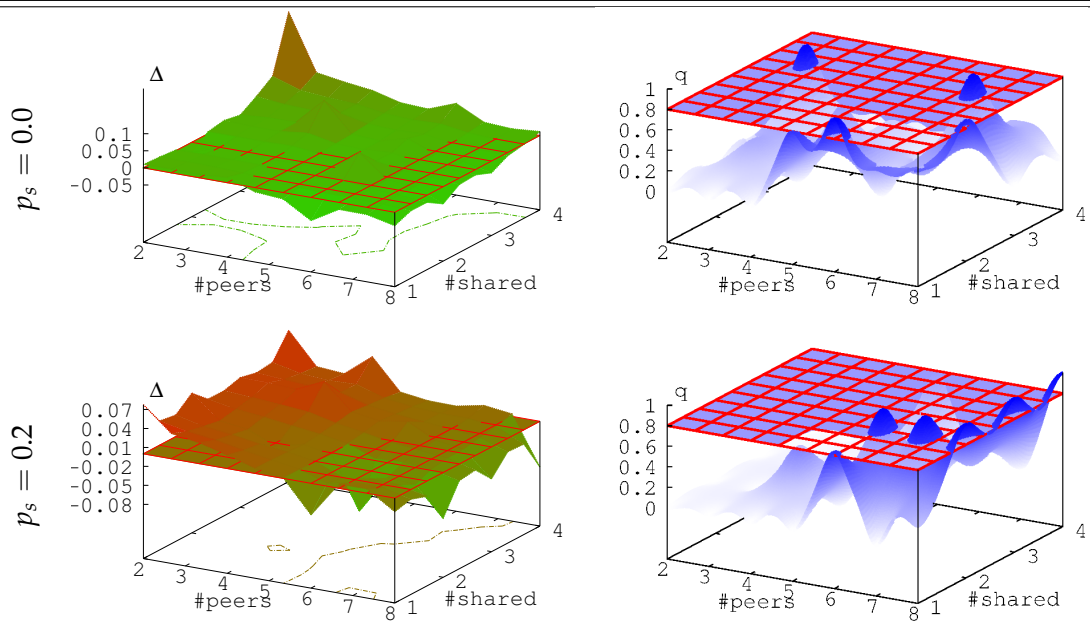


Table 7.15: Results on domain wine, using the same convention as Table 7.3.

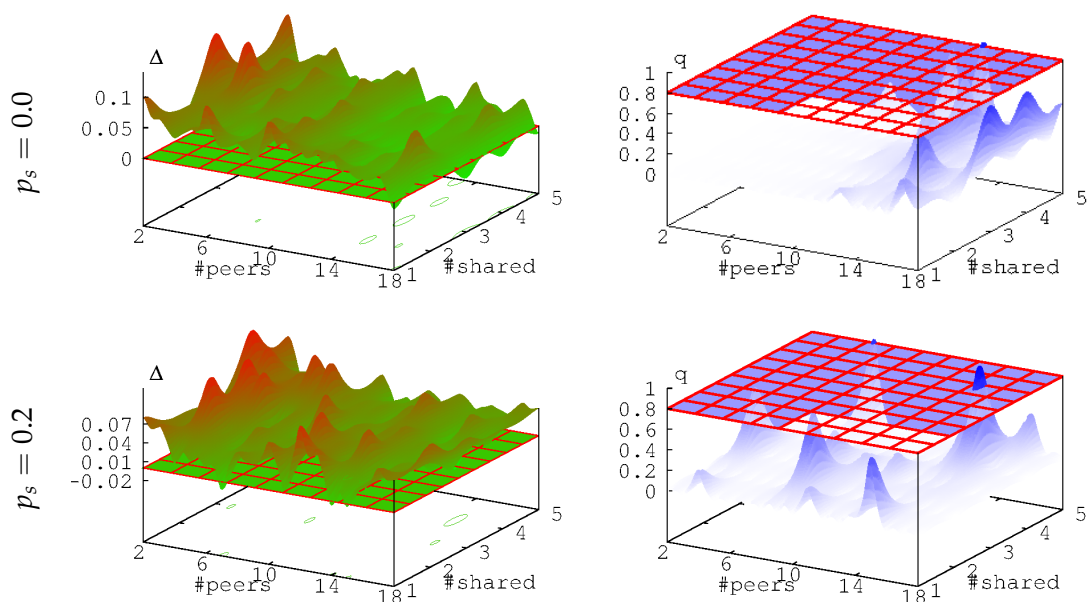


Table 7.16: Results on domain creditcard, using the same convention as Table 7.3.

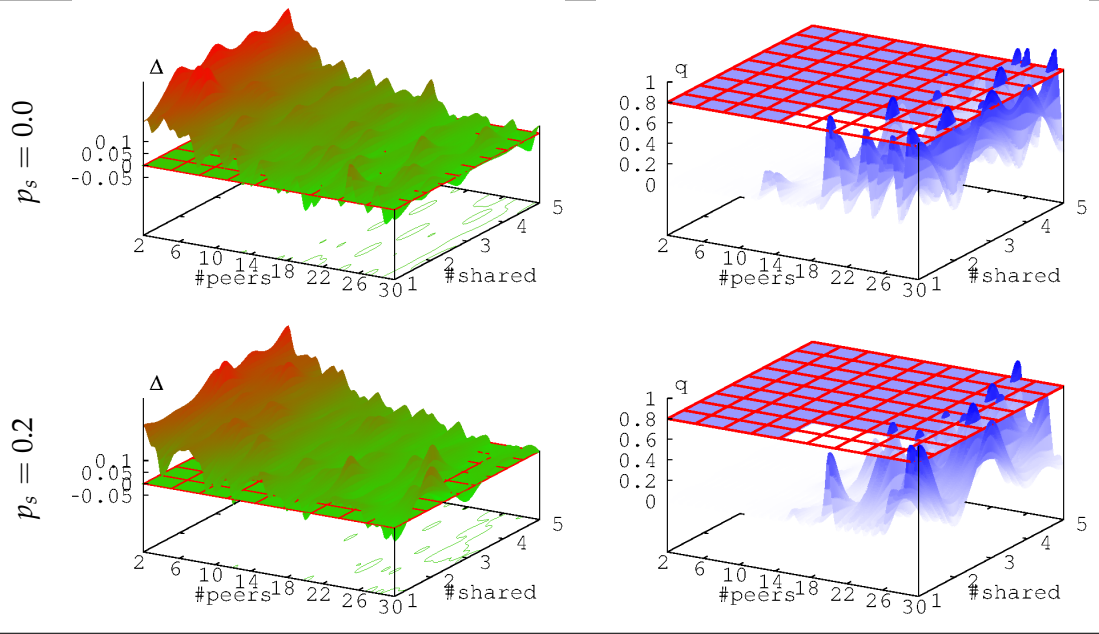


Table 7.17: Results on domain statlog, using the same convention as Table 7.3.

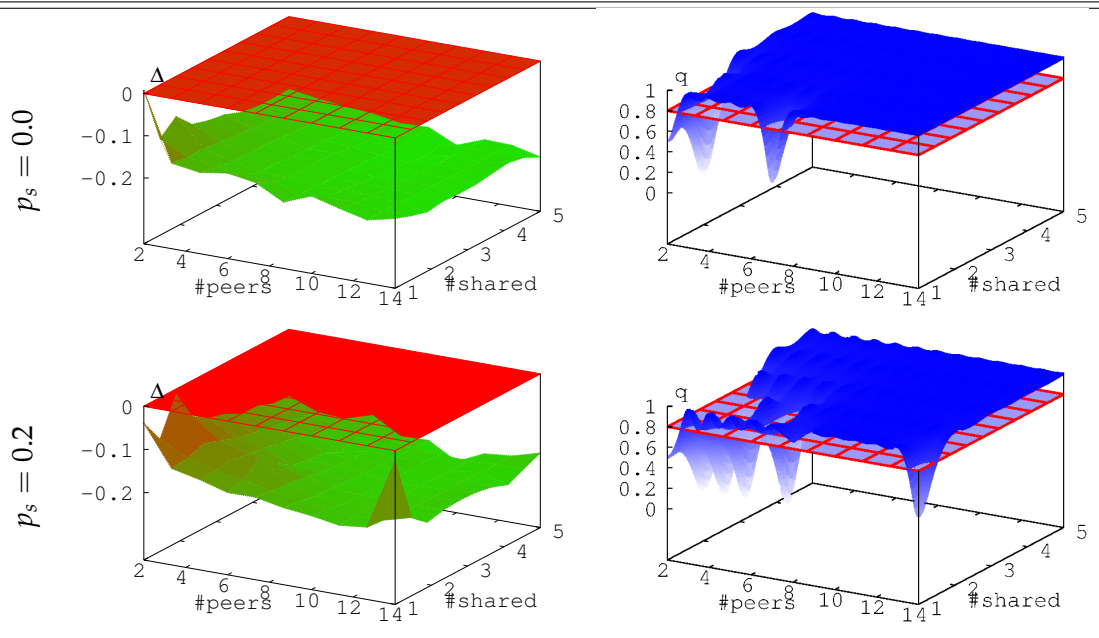


Table 7.18: Results on domain steelplates, using the same convention as Table 7.3.



---

# Conclusion

---

Research is seldom a linear path that proceeds from an open problem to its resolution and this Thesis is no exception. The presentation has not strictly followed the chronological order — by publication date of each respective piece of work.

This research journey originated from reasoning on how a classifier may be trained efficiently from label proportions and to do so without resorting to restrictive previous assumptions. The empirical success of Patrini et al. [2014] brought to question whether the proposed framework was viable at the larger extent of weakly supervised problems. A more abstract and mature view was later published in Patrini et al. [2016a] from which, in retrospective, we could have subsumed the case of LLP itself. The content of Chapter 3 and 4 is taken from those publications.

The scenario of learning with noisy labels was analyzed in Patrini et al. [2016a] to reinforce that our approach is also suited for solving the more common and well studied learning problem (Chapter 5). Patrini et al. [2017] deepens the experimental analysis bringing us closer to mainstream deep learning applications (Chapter 6) and also constitutes the major empirical effort of the Thesis in showing the practical outcome of our work on real world problems – for examples in Vision and Language.

At first sight, the problem of learning from vertically distributed datasets elaborated in Patrini et al. [2016b] has little relation to our framework. The peculiar setting was inspired discussing business needs of some of our collaborators involved in building start-up products. Yet we could carve a solution to this demanding application out of our conceptual framework (Chapter 7).

It is therefore evident that the core contribution of this Thesis is methodological. Most of the algorithmic content has revolved around the idea of casting problems with weak supervision into the two-step approach of sufficient statistic estimation followed by standard “fully” supervised learning. Different assumptions and procedures are required for the estimation of the sufficient statistic and we have detailed how to operate in several settings, as recalled above. The effectiveness of this *modus operandi*, by virtue of abstraction and adaptivity to several learning scenarios, suggests that our insight should open new ways for solving other challenging non-standard Machine Learning problems. In fact, the framework already unifies a growing body of literature [Quadrianto et al., 2009; Gao et al., 2016; Raghunathan et al., 2016].

There are similarities with many principles of traditional Software Engineering.

Decoupling is essential for analyzing a complex problem and attacking it via a *divide and conquer* strategy. We have treated the “label issue” in a modular fashion, separated from the learning process. We do not reinvent the wheel and instead we resort to well known optimization methods: gradient descent/L-BFGS (Algorithm 5), SGD and proximal algorithms (Algorithms 9 and 3), back-propagation (Algorithm 10) and the closed form solution for square loss by matrix inversion (Algorithm 12).

Ultimately, we should never be solving a more general problem than the one we are interested in the first place – *i.e.* estimating the latent variable is not a necessary step [Vapnik, 1998; Joulin and Bach, 2012]. We have given another important example in the introduction of Chapter 7, where we have motivated the need of bypassing the expensive and error prone step of entity resolution. This is precisely the glorified role of sufficient statistics, borrowed from statistical modeling and here fitted into Learning Theory; in the Thesis, we have entitled the mean operator and Rademacher observations of such power. We may also think of those statistics as compressing information that allows us to save computational time for learning.

Those insights challenge the literature which is dominated by problem specific solutions, where either loss functions or optimization algorithms are re-designed to handle the lack of supervision. We believe that our proposal takes the right direction of rethinking Machine Learning, that is historically “more akin to a craft than an engineering discipline” [Williamson, 2009]. We have also interpreted our framework as a family of learning reductions [Beygelzimer et al., 2005, 2015].

Theoretical arguments constitute a large section of the Thesis. The Factorization Theorem 18 underpins many of them. We have seen how problem agnostic generalization bounds (Theorem 23) can be easily tailored to LLP and ALN (Theorems 43 and 53). The shape of linear-odd losses, and in particular of symmetric proper losses in Chapter 4, has allowed further manipulation. We have formulated data dependent finite sample bounds on learning linear models (Theorems 26) and we have characterized a form of distribution dependent noise robustness (Theorem 56).

Aside from those results, Theorem 46 expresses generalization bounds that take into account both a more suited definition of Rademacher complexity and a novel complexity measure of the problem, related to the variance of the label proportions; these ideas are new. Surprisingly, Loss Factorization was also the key ingredient for proving that the loss curvature of ReLU networks is immune to label noise (Theorem 64).

More empirical investigation is required to confirm the practical implications of the algorithms for LLP. The success we have obtained on UCI domains — with synthesized bags — is promising for real world applications, where we believe that the relaxed assumption **(A1)**, Chapter 4, may be well justified. In some cases we have obtained surprisingly good performance with AMM (Algorithm 6). The Algorithm — which admittedly does not fit into the two-step framework, but it takes advantage of it via LMM (Algorithm 5) — should be effective in practice, at least in scenarios where linear models can bring good predictive power, in the vein of Mohammady and Culotta [2014]; Ardehaly and Culotta [2015]. In fact, we believe that LLP will have the largest impact in the future on this class of prediction problems, where



---

features and labels represent people sensitive attributes and behavior.

It is well known that personal attributes such as the one recorded in social networks are highly predictive of sensitive traits [Kosinski et al., 2013]. In turn, such private variables, *e.g.* electoral behavior, sexual preference, likelihood of contracting a disease or committing a crime, are often publicly available given by aggregate as recorded by governments and polling institutions, *e.g.* the census. This can be thought as a scenario for LLP. High predictive performance can be particularly problematic in this context. Owners of large assets of personal data would be able to infer ever more sensitive attributes of their own users or customers, information that has potentially never been shared or recorded elsewhere. Researchers working on LLP have discussed such potentially severe implications [Quadrianto et al., 2009; Yu et al., 2014b]. Yet we are not aware of any experimental evaluation on real world datasets; more research is needed on this front.

Most of the results presented, both from the algorithmics and the theory, is supported by the Factorization Theorem 18 — or the intimately related Theorem 68. It tells us that losses decompose in a way that we can isolate the *contribution of supervision* into a sufficient statistic. An intriguing open question is whether Factorization could help to identify what really matters in learning that is instead *completely unsupervised* [Sutskever et al., 2015].



---

# Bibliography

---

- ALTUN, Y. AND SMOLA, A. J., Unifying divergence minimization and statistical inference via convex duality. In *COLT*, 2006. (cited on pages 29, 39, and 103)
- ANDREWS, S.; TSOCHANTARIDIS, I.; AND HOFMANN, T., Support vector machines for multiple-instance learning. In *NIPS*, 2002. (cited on page 18)
- ARDEHALY, E. M. AND CULOTTA, A., Inferring latent attributes of twitter users with label regularization. In *HLT*, 2015. (cited on pages 18, 52, and 166)
- ARDEHALY, E. M. AND CULOTTA, A., Domain adaptation for learning from label proportions using self-training. In *IJCAI*, 2016. (cited on page 103)
- AUER, P.; CESA-BIANCHI, N.; AND FISCHER, P., Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47, 2-3 (2002), 235–256. (cited on page 18)
- BACH, F.; JENATTON, R.; MAIRAL, J.; AND OBOZINSKI, G., Optimization with sparsity-inducing penalties. *Foundations and Trends® in Machine Learning*, 4, 1 (2012), 1–106. (cited on page 31)
- BACHE, K. AND LICHMAN, M., 2013. UCI machine learning repository. <http://archive.ics.uci.edu/ml>. (cited on pages 65, 147, and 157)
- BALCAN, M. F.; BLUM, A.; FINE, S.; AND MANSOUR, Y., Distributed learning, communication complexity and privacy. *arXiv:1204.3514*, (2012). (cited on page 151)
- BARTLETT, P. L.; I, M. I. J.; AND MCAULIFFE, J. D., Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101, 473 (2006), 138–156. (cited on page 13)
- BARTLETT, P. L. AND MENDELSON, S., Rademacher and gaussian complexities: Risk bounds and structural results. *JMLR*, 3 (2002). (cited on pages 10, 49, and 83)
- BELKIN, M.; NIYOGI, P.; AND SINDHWANI, V., Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *JMLR*, 7 (2006), 2399–2434. (cited on pages 17, 56, and 57)
- BENGIO, Y.; DELALLEAU, O.; AND ROUX, N. L., Label propagation and quadratic criterion. *Semi-supervised learning*, 10 (2006). (cited on page 17)
- BENJAMINI, Y. AND HOCHBERG, Y., Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Stat. Society. Series B*, 57, 1 (1995), 289–300. (cited on page 148)

- BERNSTEIN, G. AND SHELDON, D., Consistently estimating markov chains with noisy aggregate data. In *AISTATS*, 2016. (cited on page 18)
- BEYGEZIMER, A.; DANI, V.; HAYES, T.; LANGFORD, J.; AND ZADROZNY, B., Error limiting reductions between classification tasks. In *ICML*, 2005. (cited on pages 50 and 166)
- BEYGEZIMER, A.; III, H. D.; LANGFORD, J.; AND MINEIRO, P., Learning reductions that really work. *arXiv:1502.02704*, (2015). (cited on pages 50 and 166)
- BHOWMIK, A.; GHOSH, J.; AND KOYEJO, O., Generalized linear models for aggregated data. In *AISTATS*, 2015. (cited on pages 51 and 103)
- BILENKO, M.; BASU, S.; AND MOONEY, R. J., Integrating constraints and metric learning in semi-supervised clustering. In *ICML*, 2004. (cited on page 16)
- BILENKO, M.; KAMATH, B.; AND MOONEY, R. J., Adaptive blocking: Learning to scale up record linkage. In *ICDM*, 2006. (cited on page 151)
- BILENKO, M. AND MOONEY, R. J., Adaptive duplicate detection using learnable string similarity measures. In *KDD*, 2003. (cited on pages 138 and 150)
- BLEIHOLDER, J. AND NAUMANN, F., Data fusion. *ACM Computing Surveys (CSUR)*, 41, 1 (2008), 1. (cited on pages 138 and 139)
- BLUM, A. AND CHAWLA, S., Learning from labeled and unlabeled data using graph mincuts. In *ICML*, 2001. (cited on page 18)
- BLUM, A. AND MITCHELL, T., Combining labeled and unlabeled data with co-training. In *COLT*, 1998. (cited on pages 18 and 151)
- BREFELD, U.; GÄRTNER, T.; SCHEFFER, T.; AND WROBEL, S., Efficient co-regularised least squares regression. In *ICML*, 2006. (cited on page 151)
- CHANG, M. W.; RATINOV, L.; AND ROTH, D., Guiding semi-supervision with constraint-driven learning. In *ACL*, 2007. (cited on page 18)
- CHAPELLE, O.; SCHÖLKOPF, B.; AND ZIEN, A., 2006. *Semi-supervised learning*. MIT press Cambridge. (cited on pages 15 and 17)
- CHEN, B. C.; CHEN, L.; RAMAKRISHNAN, R.; AND MUSICANT, D. R., Learning from aggregate views. In *ICDE*, 2006. (cited on pages 51 and 102)
- CHEN, S.; LIU, B.; QIAN, M.; AND ZHANG, C., Kernel k-means based framework for aggregate outputs classification. In *ICDMW*, 2009. (cited on pages 52 and 102)
- CHEN, Z.; QI, Z.; WANG, B.; CUI, L.; MENG, F.; AND SHI, Y., Learning with label proportions based on nonparallel support vector machines. *Knowledge-Based Systems*, (2016). (cited on page 102)
- CHO, W.-K.-T. AND MANSKI, C.-F., Cross level/ecological inference. *Oxford Handbook of Political Methodology*, (2008), 547–569. (cited on page 104)

- 
- CHRISTEN, P., Privacy-preserving data linkage and geocoding: Current approaches and research directions. In *ICDMW*, 2006. (cited on pages 138, 151, and 152)
- CHRISTEN, P., 2012. *Data Matching Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*. Springer Data-Centric Systems and Applications. (cited on page 138)
- COUR, T.; SAPP, B.; AND TASKAR, B., Learning from partial labels. *JMLR*, (2011), 1501–1536. (cited on page 15)
- CROSS, P.-J. AND MANSKI, C.-F., Regressions, short and long. *Econometrica*, 70, 1 (2002), 357–368. (cited on page 104)
- CUI, L.; QI, Z.; AND MENG, F., A proportion learning algorithms with density peaks. *Procedia Computer Science*, 91 (2016), 841–846. (cited on page 102)
- DAI, A. M. AND LE, Q. V., Semi-supervised sequence learning. In *NIPS*, 2015. (cited on page 127)
- DE SA, V. R., Spectral clustering with two views. In *ICML Workshop on Learning with Multiple Views*, 2005. (cited on page 18)
- DIETTERICH, T. G.; LATHROP, R. H.; AND LOZANO-PÉREZ, T., Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89 (1997), 31–71. (cited on page 15)
- DING, N. AND VISHWANATHAN, S. V. N., t-logistic regression. In *NIPS*, 2010. (cited on pages 17 and 118)
- DIVVALA, S.; FARHADI, A.; AND GUESTRIN, C., Learning everything about anything: Webly-supervised visual concept learning. In *CVPR*, 2014. (cited on page 105)
- DU PLESSIS, M. C.; NIU, G.; AND SUGIYAMA, M., Convex formulation for learning from positive and unlabeled data. In *ICML*, 2015. (cited on pages viii, 15, 26, 47, and 48)
- DUCHI, J.; HAZAN, E.; AND SINGER, Y., Adaptive subgradient methods for online learning and stochastic optimization. *JMLR*, 12 (2011), 2121–2159. (cited on page 127)
- DUNCAN, O. D. AND DAVIS, B., An alternative to ecological correlation. *American sociological review*, (1953), 665–666. (cited on page 104)
- DWORK, C., 2011. Differential privacy. In *Encyclopedia of Cryptography and Security*, 338–340. Springer. (cited on page 152)
- DWORK, C. AND ROTH, A., The algorithmic foundations of differential privacy. *Theoretical Computer Science*, 9, 3-4 (2013), 211–407. (cited on page 152)

- ESTRADA, T.; ARMEN, R.; AND TAUFER, M., Automatic selection of near-native protein-ligand conformations using a hierarchical clustering and volunteer computing. In *International Conference on Bioinformatics and Computational Biology*, 2010. (cited on page 138)
- FAN, K.; ZHANG, H.; YAN, S.; WANG, L.; ZHANG, W.; AND FENG, J., Learning a generative classifier from label proportions. *NeuroComputing*, 139 (2014), 47–55. (cited on page 102)
- FERGUS, R.; FEI-FEI, L.; PERONA, P.; AND ZISSERMAN, A., Learning object categories from internet image searches. *Proceedings of the IEEE*, 98, 8 (2010), 1453–1466. (cited on page 105)
- FLAXMAN, S.-R.; WANG, Y.-X.; AND SMOLA, A.-J., Who supported obama in 2012?: Ecological inference through distribution regression. In *KDD*, 2015. (cited on page 104)
- FOULDS, J. AND FRANK, E., A review of multi-instance learning assumptions. *The Knowledge Engineering Review*, 25, 01 (2010), 1–25. (cited on page 16)
- FRÉNEY, B. AND VERLEYSSEN, M., Classification in the Presence of Label Noise: A Survey. *IEEE Transactions on Neural Networks and Learning Systems*, 25, 5 (May 2014), 845–869. (cited on page 118)
- GANCHEV, K.; GRAÇA, J.; GILLENWATER, J.; AND TASKAR, B., Posterior regularization for structured latent variable models. *JMLR*, 11 (2010), 2001–2049. (cited on page 18)
- GAO, W.; WANG, L.; LI, Y. F.; AND ZHOU, Z. H., Risk minimization in the presence of label noise. In *AAAI*, 2016. (cited on pages 29, 49, and 165)
- GARCIA-GARCIA, D. AND WILLIAMSON, R. C., Degrees of supervision. In *NIPS Workshops*, 2011. (cited on page 18)
- GETOOR, L. AND MACHANAVAJJHALA, A., Entity resolution: theory, practice & open challenges. *Proceedings of the VLDB Endowment*, 5, 12 (2012), 2018–2019. (cited on page 150)
- GHOSH, A.; MANWANI, N.; AND SASTRY, P. S., Making risk minimization tolerant to label noise. *Neurocomputing*, 160 (2015), 93–107. (cited on pages 17, 105, 110, 118, and 128)
- GLOROT, X. AND BENGIO, Y., Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010. (cited on page 129)
- GRANDVALET, Y. AND BENGIO, Y., Semi-supervised learning by entropy minimization. In *NIPS*, 2004. (cited on page 17)

- 
- HE, K.; ZHANG, X.; REN, S.; AND SUN, J., Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015. (cited on page 127)
- HE, K.; ZHANG, X.; REN, S.; AND SUN, J., Deep residual learning for image recognition. In *CVPR*, 2016a. (cited on pages 120, 128, and 129)
- HE, K.; ZHANG, X.; REN, S.; AND SUN, J., Identity mappings in deep residual networks. In *ECCV*, 2016b. (cited on pages 120 and 128)
- HERNÁNDEZ, J. AND INZA, I., Learning naive bayes models for multiple-instance learning with label proportions. In *Conference of the Spanish Association for Artificial Intelligence*, Springer, 2011. (cited on page 102)
- HERNÁNDEZ-GONZÁLEZ, J.; INZA, I.; CRISOL-ORTÍZ, L.; GUEMBE, M. A.; IÑARRA, M. J.; AND LOZANO, J. A., Fitting the data from embryo implantation prediction: Learning from label proportions. *Statistical methods in medical research*, (2016). (cited on page 51)
- HERNANDEZ-GONZALEZ, J.; INZA, I.; AND LOZANO, J., 2016. Weak supervision and other non-standard classification problems: a taxonomy. In *PRL*. Elsevier. (cited on page 16)
- HERNÁNDEZ-GONZÁLEZ, J.; INZA, I.; AND LOZANO, J. A., Learning bayesian network classifiers from label proportions. *Pattern Recognition*, 46, 12 (2013), 3425–3440. (cited on page 102)
- HERNÁNDEZ-GONZÁLEZ, J.; INZA, I.; AND LOZANO, J. A., A novel weakly supervised problem: Learning from positive-unlabeled proportions. In *Conference of the Spanish Association for Artificial Intelligence*, Springer, 2015. (cited on page 103)
- HERNÁNDEZ-GONZÁLEZ, J.; INZA, I.; AND LOZANO, J. A., Learning from proportions of positive and unlabeled examples. *International Journal of Intelligent Systems*, (2016). (cited on page 103)
- HOCHREITER, S. AND SCHMIDHUBER, J., Long short-term memory. *Neural computation*, 9, 8 (1997), 1735–1780. (cited on page 120)
- HOPE, T. AND SHAHAF, D., Ballpark learning: Estimating labels from rough group comparisons. In *ECML-PKDD16*, 2016. (cited on page 103)
- HORN, R. A. AND JOHNSON, C. R., 2012. *Matrix analysis*. Cambridge university press. (cited on pages 70 and 71)
- HUANG, G.; SUN, Y.; LIU, Z.; SEDRA, D.; AND WEINBERGER, K., Deep networks with stochastic depth. In *ECCV*, 2016. (cited on page 128)
- IOFFE, S. AND SZEGEDY, C., Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. (cited on page 120)

- JAANKOLA, T. S. AND JORDAN, M. I., Bayesian parameter estimation via variational methods. *Statistics and Computing*, 10, 1 (2000), 25–37. (cited on page 49)
- JANZING, D.; PETERS, J.; SGOURITSA, E.; ZHANG, K.; MOOIJ, J. M.; AND SCHÖLKOPF, B., On causal and anticausal learning. In *ICML*, 2012. (cited on page 17)
- JOACHIMS, T., Transductive inference for text classification using support vector machines. In *ICML*, 1999. (cited on page 17)
- JOULIN, A. AND BACH, F. R., A convex relaxation for weakly supervised classifiers. In *ICML*, 2012. (cited on pages 18 and 166)
- JUDGE, G. G.; MILLER, D. J.; AND CHO, W. K. T., 2004. An information theoretic approach to ecological estimation and inference. In *Ecological inference: New methodological strategies*, 162–187. Cambridge University Press. (cited on page 104)
- KAKADE, S. M.; SRIDHARAN, K.; AND TEWARI, A., On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *NIPS*, 2009. (cited on pages 10 and 49)
- KAWAGUCHI, K., Deep learning without poor local minima. In *NIPS*, 2016. (cited on page 126)
- KEARNS, M., Efficient noise-tolerant learning from statistical queries. *Journal of the ACM*, 45 (1998). (cited on page 135)
- KEARNS, M. J. AND MANSOUR, Y., On the boosting ability of top-down decision tree learning algorithms. In *STOC*, 1996. (cited on page 53)
- KING, G., 1997. *A solution to the ecological inference problem: reconstructing individual behavior from aggregate data*. Princeton University Press. (cited on page 104)
- KING, G.; TANNER, M.-A.; AND ROSEN, O., 2004. *Ecological inference: New methodological strategies*. Cambridge University Press. (cited on page 104)
- KINGMA, D. P.; MOHAMED, S.; REZENDE, D. J.; AND WELLING, M., Semi-supervised learning with deep generative models. In *NIPS*, 2014. (cited on page 18)
- KOSINSKI, M.; STILLWELL, D.; AND GRAEPEL, T., Private traits and attributes are predictable from digital records of human behavior. *PNAS*, 110, 15 (2013), 5802–5805. (cited on page 167)
- KOTZIAS, D.; DENIL, M.; FREITAS, N. D.; AND SMYTH, P., From group to individual labels using deep features. In *KDD*, 2015. (cited on pages 51 and 102)
- KRAUSE, J.; SAPP, B.; HOWARD, A.; ZHOU, H.; TOSHEV, A.; DUERIG, T.; PHILBIN, J.; AND FEI-FEI, L., The unreasonable effectiveness of noisy data for fine-grained recognition. In *ECCV*, 2016. (cited on pages 105, 119, 131, and 135)



- 
- KRIZHEVSKY, A. AND HINTON, G., Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. (cited on page 126)
- KRIZHEVSKY, A.; SUTSKEVER, I.; AND HINTON, G. E., Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. (cited on page 129)
- KÜCK, H.; CARBONETTO, P.; AND DE FREITAS, N., A constrained semi-supervised learning approach to data association. In *ECCV*, 2004. (cited on page 102)
- KUCK, H. AND DE FREITAS, N., Learning about individuals from group statistics. In *UAI*, 2005. (cited on pages 15, 16, 18, 52, and 102)
- LAFON, S.; KELLER, Y.; AND COIFMAN, R. R., Data fusion and multicue data matching by diffusion maps. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28, 11 (2006), 1784–1797. (cited on page 150)
- LAI, K. T.; YU, F.; CHEN, M. S.; AND CHANG, S. F., Video event detection by inferring temporal instance labels. In *CVPR*, 2014. (cited on page 52)
- LAIRD, P. D., 1988. *Learning from Good and Bad Data*. Kluwer Academic Publishers, Norwell, MA, USA. (cited on page 135)
- LANCKRIET, G. R. G.; BIE, T. D.; CRISTIANINI, N.; JORDAN, M. I.; AND NOBLE, W. S., A statistical framework for genomic data fusion. *Bioinformatics*, 20, 16 (2004), 2626–2635. (cited on page 138)
- LAWRENCE, N. D. AND JORDAN, M. I., Semi-supervised learning via gaussian processes. In *NIPS*, 2004. (cited on page 18)
- LECUN, Y.; BOTTOU, L.; BENGIO, Y.; AND HAFFNER, P., Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86, 11 (1998), 2278–2324. (cited on pages 120 and 126)
- LEDOUX, M. AND TALAGRAND, M., 1991. *Probability in Banach Spaces*. Springer Verlag. (cited on page 89)
- LEHMANN, E. L. AND CASELLA, G., 1998. *Theory of point estimation*, vol. 31. Springer Science & Business Media. (cited on page 15)
- LI, F. AND TAYLOR, G., Alter-cnn: An approach to learning from label proportions with application to ice-water classification. In *NIPS Workshop on Learning and privacy with Incomplete Data and Weak Supervision*, 2015. (cited on pages 52 and 102)
- LI, Y.-F.; TSANG, I. W.; KWOK, J. T.; AND ZHOU, Z.-H., Convex and scalable weakly labeled svms. *JMLR*, (2013), 2151–2188. (cited on page 18)
- LIANG, P.; JORDAN, M. I.; AND KLEIN, D., Learning from measurements in exponential families. In *ICML*, 2009. (cited on pages 18 and 52)

- LIU, L.-P. AND DIETTERICH, T. G., Learnability of the superset label learning problem. In *ICML*, 2014. (cited on page 15)
- LIU, Q. AND IHLER, A., Distributed parameter estimation via pseudo-likelihood. In *ICML*, 2012. (cited on page 151)
- LIU, Q. AND IHLER, A. T., Distributed estimation, information loss and exponential families. In *NIPS*, 2014. (cited on page 151)
- LIU, T. AND TAO, D., Classification with noisy labels by importance reweighting. *IEEE Transactions on pattern analysis and machine intelligence*, 38, 3 (2016), 447–461. (cited on pages 124 and 135)
- LONG, P. M. AND SERVEDIO, R. A., Random classification noise defeats all convex potential boosters. *Machine learning*, 78, 3 (2010), 287–304. (cited on pages 108, 109, 110, and 118)
- MA, F.; SHI, Y.; WANG, B.; AND CHEN, Z., Research on the classification of commercial banks? fund clients based on learning with label proportions. *Procedia Computer Science*, 91 (2016), 988–994. (cited on page 51)
- MAAS, A. L.; DALY, R. E.; PHAM, P. T.; HUANG, D.; NG, A. Y.; AND POTTS, C., Learning word vectors for sentiment analysis. In *ACL*, 2011. (cited on page 126)
- MANN, G. S. AND MCCALLUM, A., Generalized expectation criteria for semi-supervised learning of conditional random fields. In *ACL*, 2008. (cited on page 18)
- MANN, G. S. AND MCCALLUM, A., Generalized expectation criteria for semi-supervised learning with weakly labeled data. *JMLR*, (2010), 955–984. (cited on page 18)
- MANWANI, N. AND SASTRY, P. S., Noise tolerance under risk minimization. *IEEE Transactions on Cybernetics*, 43, 3 (2013), 1146–1151. (cited on pages 108, 110, 116, and 118)
- MARKOWITZ, H., Portfolio selection. *Journal of Finance*, 6 (1952), 77–91. (cited on page 142)
- MASNADI-SHIRAZI, H.; MAHADEVAN, V.; AND VASCONCELOS, N., On the design of robust classifiers for computer vision. In *CVPR*, 2010. (cited on pages 17 and 118)
- MASNADI-SHIRAZI, H. AND VASCONCELOS, N., On the design of loss functions for classification: theory, robustness to outliers, and savageboost. In *NIPS*, 2009. (cited on pages 118 and 128)
- MCDIARMID, C., Concentration. In *Probabilistic Methods for Algorithmic Discrete Mathematics*, 1–54, Springer Verlag, 1998. (cited on pages 11 and 19)

- 
- MENON, A.; ROOYEN, B. V.; ONG, C. S.; AND WILLIAMSON, B., Learning from corrupted binary labels via class-probability estimation. In *ICML*, 2015. (cited on pages 16, 118, 120, 124, 127, and 135)
- MENON, A.; VAN ROOYEN, B.; AND NATARAJAN, N., Learning from binary labels with instance-dependent corruption. *arXiv preprint arXiv:1605.00751*, (2016). (cited on pages 105, 106, 118, and 131)
- MINTZ, M.; BILLS, S.; SNOW, R.; AND JURAFSKY, D., Distant supervision for relation extraction without labeled data. In *ACL*, 2009. (cited on page 15)
- MISRA, I.; ZITNICK, C. L.; MITCHELL, M.; AND GIRSHICK, R., Seeing through the human reporting bias: Visual classifiers from noisy human-centric labels. In *CVPR*, 2016. (cited on page 105)
- MNIH, V. AND HINTON, G. E., Learning to label aerial images from noisy data. In *ICML*, 2012. (cited on pages 119 and 134)
- MOHAMMADY, E. AND CULOTTA, A., Using county demographics to infer attributes of twitter users. In *ACL*, 2014. (cited on pages 52 and 166)
- MUSICANT, D. J.; CHRISTENSEN, J. M.; AND OLSON, J. F., Supervised learning by training on aggregate outputs. In *ICDM*, 2007. (cited on pages 51 and 102)
- MUZELLEC, B.; NOCK, R.; PATRINI, G.; AND NIELSEN, F., Tsallis regularized optimal transport and ecological inference. In *AAAI*, 2017. (cited on page 104)
- NATARAJAN, N.; DHILLON, I. S.; RAVIKUMAR, P. K.; AND TEWARI, A., Learning with noisy labels. In *NIPS*, 2013. (cited on pages 15, 17, 105, 106, 107, 108, 114, 116, 118, 120, 122, and 135)
- NI, T.; CHUNG, F.-L.; AND WANG, S., Support vector machine with manifold regularization and partially labeling privacy protection. *Information Sciences*, 294 (2015), 390–407. (cited on page 102)
- NIGAM, K. AND GHANI, R., Analyzing the effectiveness and applicability of co-training. In *CIKM*, 2000. (cited on page 18)
- NOCK, R., Learning games and Rademacher observations losses. *CoRR*, abs/1512.05244 (2015). (cited on page 142)
- NOCK, R., On regularizing rademacher observation losses. In *NIPS*, 2016. (cited on page 152)
- NOCK, R. AND NIELSEN, F., Bregman divergences and surrogates for learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31 (2009), 2048–2059. (cited on pages 26, 53, and 123)
- NOCK, R.; PATRINI, G.; AND FRIEDMAN, A., Rademacher Observations, private data and boosting. In *ICML*, 2015. (cited on pages 4, 141, 142, 144, and 152)

- PARK, Y. AND GHOSH, J., A probabilistic imputation framework for predictive analysis using variably aggregated, multi-source healthcare data. In *SIGHIT International Health Informatics Symposium*, 2012. (cited on page 103)
- PARK, Y. AND GHOSH, J., Cudia: Probabilistic cross-level imputation using individual auxiliary information. *ACM Transactions on Intelligent Systems and Technology*, 4, 4 (2013), 66. (cited on page 103)
- PATRINI, G.; NIELSEN, F.; NOCK, R.; AND CARIONI, M., Loss factorization, weakly supervised learning and label noise robustness. In *ICML*, 2016a. (cited on pages 4, 49, and 165)
- PATRINI, G.; NOCK, R.; HARDY, S.; AND CAETANO, T., Fast learning from distributed datasets without entity matching. In *IJCAI*, 2016b. (cited on pages 4 and 165)
- PATRINI, G.; NOCK, R.; RIVERA, P.; AND CAETANO, T., (Almost) no label no cry. In *NIPS*, 2014. (cited on pages 4, 29, 53, and 165)
- PATRINI, G.; ROZZA, A.; MENON, A.; NOCK, R.; AND QU, L., Making neural networks robust to label noise: a loss correction approach. In *Submitted to CVPR*, 2017. (cited on pages 4 and 165)
- PÉREZ-ORTIZ, M.; GUTIÉRREZ, P. A.; CARBONERO-RUZ, M.; AND HERVÁS-MARTÍNEZ, C., Learning from label proportions via an iterative weighting scheme and discriminant analysis. In *Conference of the Spanish Association for Artificial Intelligence*, 79–88, Springer, 2016. (cited on page 102)
- PLATT, J., Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10, 3 (1999), 61–74. (cited on pages 10 and 102)
- POSNETT, D.; FILKOV, V.; AND DEVANBU, P., Ecological inference in empirical software engineering. In *International Conference on Automated Software Engineering*, 2011. (cited on page 104)
- QI, Z.; WANG, B.; MENG, F.; AND NIU, L., Learning with label proportions via npsvm. *IEEE Transactions on Cybernetics*, (2016). (cited on page 102)
- QUADRIANTO, N.; SMOLA, A. J.; CAETANO, T. S.; AND LE, Q. V., Estimating labels from label proportions. *JMLR*, 10 (2009), 2349–2374. (cited on pages viii, 23, 29, 49, 51, 52, 53, 55, 63, 66, 77, 102, 103, 165, and 167)
- RAGHUNATHAN, A.; FROSTIG, R.; DUCHI, J.; AND LIANG, P., Estimation from indirect supervision with linear moments. In *ICML*, 2016. (cited on pages viii, 15, 18, 49, and 165)
- RAMASWAMY, H. G.; SCOTT, C.; AND TEWARI, A., Mixture proportion estimation via kernel embedding of distributions. In *ICML*, 2016. (cited on page 135)

- 
- RASMUS, A.; BERGLUND, M.; HONKALA, M.; VALPOLA, H.; AND RAIKO, T., Semi-supervised learning with ladder networks. In *NIPS*, 2015. (cited on page 17)
- RASTOGI, V.; DALVI, N.-N.; AND GAROFALAKIS, M.-N., Large-scale collective entity matching. *VLDB*, 4, 4 (2011), 208–218. (cited on page 138)
- REED, S.; LEE, H.; ANGUELOV, D.; SZEGEDY, C.; ERHAN, D.; AND RABINOVICH, A., Training deep neural networks on noisy labels with bootstrapping. *ICLR Workshops*, (2015). (cited on pages 119, 128, 134, and 135)
- REID, M. D. AND WILLIAMSON, R. C., Composite binary losses. *JMLR*, 11 (2010), 2387–2422. (cited on pages 10, 13, and 123)
- ROBINSON, W.-S., Ecological correlations and the behavior of individuals. *American Sociological Review*, 15, 3 (1950), 351–357. (cited on page 104)
- RÜPING, S., Svm classifier estimation from group probabilities. In *ICML*, 2010. (cited on pages 51, 63, and 102)
- SAMDANI, R.; CHANG, M.-W.; AND ROTH, D., Unified expectation maximization. In *HLT*, 2012. (cited on page 18)
- SANDERSON, T. AND C. SCOTT, C., Class proportion estimation with application to multiclass anomaly rejection. In *AISTATS*, 2014. (cited on page 135)
- SCHÖLKOPF, B. AND SMOLA, A. J., 2002. *Learning with kernels: Support vector machines, regularization, optimization, and beyond*. MIT press. (cited on page 43)
- SCHROFF, F.; CRIMINISI, A.; AND ZISSERMAN, A., Harvesting image databases from the web. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33, 4 (2011), 754–766. (cited on page 105)
- SCOTT, C.; BLANCHARD, G.; AND HANDY, G., Classification with asymmetric label noise : Consistency and maximal denoising. In *COLT*, 2013. (cited on page 135)
- SEEGER, M., Learning with labeled and unlabeled data. Technical report, Institute for ANC, Edinburgh, UK, 2000. (cited on page 18)
- SHALEV-SHWARTZ, S. AND BEN-DAVID, S., 2014. *Understanding machine learning: From theory to algorithms*. Cambridge University Press. (cited on pages 7, 10, and 19)
- SHALEV-SHWARTZ, S.; SINGER, Y.; SREBRO, N.; AND COTTER, A., Pegasos: Primal estimated sub-gradient solver for svm. *Mathematical programming*, 127, 1 (2011), 3–30. (cited on pages 31 and 111)
- SHELDON, D. R. AND DIETTERICH, T. G., Collective graphical models. In *NIPS*, 2011. (cited on page 18)
- SHI, J. AND MALIK, J., Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22 (2000), 888–905. (cited on page 58)

- SINDHWANI, V.; NIYOGI, P.; AND BELKIN, M., A co-regularized approach to semi-supervised learning with multiple views. In *ICML Workshop on Learning with Multiple Views*, 2005. (cited on page 151)
- SMOLA, A.; GRETTON, A.; SONG, L.; AND SCHÖLKOPF, B., A Hilbert space embedding for distributions. In *Algorithmic Learning Theory*, 13–31, Springer, 2007. (cited on page 49)
- SONG, L.; HUANG, J.; SMOLA, A. J.; AND FUKUMIZU, K., Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *ICML*, 2009. (cited on page 49)
- SPROULL, R. F.; DUMOUCHEL, W. H.; KEARNS, M.; LAMPSON, B. W.; LANDAU, S.; LEITER, M. E.; PARKER, E. R.; AND WEINBERGER, P. J., 2015. Bulk collection of signal intelligence: technical options. In *Committee on Responding to Section 5(d) of Presidential Policy Directive 28: The Feasibility of Software to Provide Alternatives to Bulk Signals Intelligence Collection*. National Academy Press. (cited on page 138)
- SRIVASTAVA, N.; HINTON, G. E.; KRIZHEVSKY, A.; SUTSKEVER, I.; AND SALAKHUTDINOV, R., Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 15, 1 (2014), 1929–1958. (cited on page 120)
- STEMPFEL, G. AND RALAIVOLA, L., 2009. Learning SVMs from sloppily labeled data. In *ICANN*, 884–893. Springer. (cited on page 17)
- STOLPE, M. AND MORIK, K., Learning from label proportions by optimizing cluster model selection. In *ECML-PKDD*, 2011. (cited on pages 51 and 102)
- SUKHBAATAR, S.; BRUNA, J.; PALURI, M.; BOURDEV, L.; AND FERGUS, R., Training convolutional networks with noisy labels. In *ICLR Workshops*, 2015. (cited on pages 119, 120, 122, 130, and 135)
- SURDEANU, M.; TIBSHIRANI, J.; NALLAPATI, R.; AND MANNING, C. D., Multi-instance multi-label learning for relation extraction. In *EMNLP-CoNLL*, 2012. (cited on page 15)
- SUTSKEVER, I.; JOZEFOWICZ, R.; GREGOR, K.; REZENDE, D.; LILICRAP, T.; AND VINYALS, O., Towards principled unsupervised learning. *arXiv preprint arXiv:1511.06440*, (2015). (cited on page 167)
- SUTTON, R. S. AND BARTO, A. G., 1998. *Reinforcement learning: An introduction*, vol. 1. MIT press Cambridge. (cited on page 18)
- SWEENEY, L., Privacy-enhanced linking. *ACM SIGKDD Explorations Newsletter*, 7, 2 (2005), 72–75. (cited on page 138)
- SZUMMER, M. AND JAAKKOLA, T. S., Information regularization with partially labeled data. In *NIPS*, 2002. (cited on page 17)

- 
- TSUI, F.; ESPINO, J. U.; DATO, V. M.; GESTELAND, P. H.; HUTMAN, J.; AND WAGNER, M. M., Technical description of rods: a real-time public health surveillance system. *Journal of the American Medical Informatics Association*, 10, 5 (2003), 399–408. (cited on page 138)
- VAN ROOYEN, B., 2015. *Machine Learning via Transitions*. Ph.D. thesis, The Australian National University. (cited on pages 118 and 122)
- VAN ROOYEN, B.; MENON, A. K.; AND WILLIAMSON, R. C., Learning with symmetric label noise: The importance of being unhinged. In *NIPS*, 2015. (cited on pages 17, 26, 28, 106, 108, 110, 118, and 128)
- VAPNIK, V., 1998. *Statistical Learning Theory*. John Wiley. (cited on pages 12, 23, and 166)
- VATSALAN, D.; CHRISTEN, P.; AND VERYKIOS, V. S., A taxonomy of privacy-preserving record linkage techniques. *Information Systems*, 38, 6 (2013), 946–969. (cited on page 152)
- VILLANI, C., 2008. *Optimal transport: old and new*, vol. 338. Springer Science & Business Media. (cited on page 104)
- VON LUXBURG, U., A tutorial on spectral clustering. *Statistics and computing*, 17, 4 (2007), 395–416. (cited on page 56)
- WAGER, S.; BLOCKER, A.; AND CARDIN, N., Weakly supervised clustering: Learning fine-grained signals from coarse labels. *The Annals of Applied Statistics*, 9, 2 (2015), 801–820. (cited on pages 52 and 102)
- WAKEFIELD, J. AND SHADDICK, G., Health-exposure modeling and the ecological fallacy. *Biostatistics*, 7, 3 (2006), 438–455. (cited on page 104)
- WANG, B.; CHEN, Z.; AND QI, Z., Linear twin svm for learning from label proportions. In *International Conference on Web Intelligence and Intelligent Agent Technology*, 2015. (cited on page 102)
- WESTON, J.; RATLE, R.; MOBAHI, H.; AND COLLOBERT, R., 2012. Deep learning via semi-supervised embedding. In *Neural Networks: Tricks of the Trade*, 639–655. Springer. (cited on page 17)
- WHANG, S. E. AND GARCIA-MOLINA, H., Joint entity resolution. In *ICDE*, 2012. (cited on page 151)
- WHANG, S.-E.; MENESTRINA, D.; KOUTRIKA, G.; THEOBALD, M.; AND GARCIA-MOLINA, H., Entity resolution with iterative blocking. In *SIGMOD*, 219–232, 2009. (cited on pages 144 and 151)
- WILLIAMSON, R., Reconceiving machine learning. Technical report, NICTA, 2009. (cited on page 166)

- XIAO, T.; XIA, T.; YANG, T.; HUANG, C.; AND WANG, X., Learning from massive noisy labeled data for image classification. In *CVPR*, 2015. (cited on pages 105, 119, 120, 126, 129, 130, 131, 134, and 135)
- YAMANISHI, Y.; VERT, J. P.; AND KANEHISA, K., Protein network inference from multiple genomic data: a supervised approach. *Bioinformatics*, 20, suppl 1 (2004), i363–i370. (cited on page 138)
- YANG, Z.; COHEN, W.; AND SALAKHUTDINOV, R., Revisiting semi-supervised learning with graph embeddings. In *ICML*, 2016. (cited on page 17)
- YARIN, G. AND GHAHRAMANI, Z., A theoretically grounded application of dropout in recurrent neural networks. In *NIPS*, 2016. (cited on page 127)
- YU, F. X.; CAO, L.; MERLER, M.; CODELLA, N.; CHEN, T.; SMITH, J. R.; AND CHANG, S.-F., Modeling attributes from category-attribute proportions. In *International Conference on Multimedia*, 2014a. (cited on page 52)
- YU, F. X.; KUMAR, S.; JEBARA, T.; AND CHANG, S. F., On learning with label proportions. *CoRR*, abs/1402.5902 (2014). (cited on pages 51, 61, 65, 103, and 167)
- YU, F. X.; LIU, D.; KUMAR, S.; JEBARA, T.; AND CHANG, S. F.,  $\alpha$ SVM for Learning with Label Proportions. In *ICML*, 2013. (cited on pages viii, 18, 61, 63, 64, 65, 90, 91, and 102)
- ZANTEDESCHI, V.; EMONET, R.; AND SEBBAN, M., beta-risk: a new surrogate risk for learning from weakly labeled data. In *NIPS*, 2016. (cited on pages 18 and 49)
- ZHANG, B.; ESTRADA, T.; CICOTTI, P.; BALAJI, P.; AND TAUFER, M., Accurate scoring of drug conformations at the extreme scale. In *International Symposium on Cluster, Cloud and Grid Computing*, 2015. (cited on page 138)
- ZHOU, Z.-H. AND XU, J.-M., On the relation between multi-instance learning and semi-supervised learning. In *ICML*, 2007. (cited on page 16)
- ZHU, X. AND GHAHRAMANI, Z., Learning from labeled and unlabeled data with label propagation. Technical report, CMU-CALD-02-107, 2002. (cited on page 17)