

A supervised learning framework in the context of multiple annotators



Julián Gil González

Universidad Tecnológica de Pereira

This dissertation is submitted for the degree of
Doctor of Engineering

Advisor

Andrés M. Álvarez Meza, PhD.

Committee

Diego Hernán Peluffo-Ordóñez, PhD.

Julian David Arias Londoño, PhD.

Cristian Guarnizo-Lemus, PhD.

Date

August 24, 2021.

To the memory of my grandpa Fabio González.

Acknowledgements

I would firstly like to thank my supervisor Dr. Andrés M. Álvarez for his support, inspiration and guidance during these years. Andrés was always kind and available to discuss every little detail of my research work. Additionally, he trusted me with so many interesting ideas and encourage me to pursuit my own ideas.

Thanks to professor Mauricio Álvarez for allowing me to have the chance to spend highly valuable time at the University of Sheffield. Besides, Thanks to Juan José G. for the academic discussions during my visit to Sheffield.

Thanks to all friends and colleagues at the Automatics research group. It was a privilege the get to know each one of them. Especially, I would like to thank to Cristian T., Hernan G., Cristian G., Hernan V., Carlos Z, Pablo A., Diego A., Jhouben C., Andres L., and Juan L., for their friendship and support during this process.

I would like to acknowledge my source of funding given by Convocatoria 785-2017 from Ministerio de Ciencia, Tecnología e Innovación (MINCIENCIAS).

Finally, I want to thank my family for their love and support. Especially, thanks to Martha González, Bertilda Guerrero, Monica Castro, Alirio Gil, and Sebastián Gil for their never-ending love and patience.

Abstract

The increasing popularity of crowdsourcing platforms, i.e., Amazon Mechanical Turk, is changing how datasets for supervised learning are built. In these cases, instead of having datasets labeled by one source (which is supposed to be an expert who provided the absolute gold standard), we have datasets labeled by multiple annotators with different and unknown expertise. Hence, we face a multi-labeler scenario, which typical supervised learning models cannot tackle. For such a reason, much attention has recently been given to the approaches that capture multiple annotators' wisdom. However, such methods residing on two key assumptions: the labeler's performance does not depend on the input space and independence among the annotators, which are hardly feasible in real-world settings.

This thesis explores several models based on both frequentist and Bayesian perspectives aiming to face multi-labeler scenarios. Our approaches model the annotators' behavior by considering the relationship between the input space and the labelers' performance and coding interdependencies among them.

Regarding the frequentist approaches, we first introduce the Kernel alignment-based annotators relevance analysis-(KAAR), which uses the well-known centered kernel alignment algorithm to measure the similarities between the input features and the labels given by each annotator. Hence, we estimate the labelers' performance by taking into account dependencies among them. KAAR does not consider the relationship between the input features and the labelers' performance; such a drawback is dealt with in the second proposal termed Localized kernel alignment-based annotators relevance analysis-(LKAAR). The main idea is to modify the formulation of KAAR to perform a localized analysis of the similarities between the input space and the labels from multiple annotators; thus, in a multi-labeler scenario, we measure the labelers' performance as a function of the input space while considering dependencies among them. Next, we propose the Regularized Chained Deep Neural Network for Multiple Annotators-(RCDNN), which has as principal purpose to link the parameters of a likelihood from multiple annotators with the outputs of a regularized deep neural network-(RDNN) model. Hence, since the

parameters of a multi-labeler likelihood are related to the labelers' performance, such performance is being modeled as a function of the input features. At the same time, it takes into account annotators' dependencies.

On the other hand, concerning the Bayesian perspective, we first propose the chained Gaussian processes for multiple annotators-(CGPMA), which uses as a basis the Chained Gaussian Processes model aiming to link the parameters of a multi-labeler likelihood to multiple independent latent functions, which follow GP distributions. Thereby, we are modeling each annotator parameter as a function of the input features; however, due to the multiple GPs are assumed to be independent, CGPMA lacks modeling the dependencies among the labelers. To deal with such downside, we propose the correlated chained Gaussian processes for multiple annotators-(CCGPMA), which uses the concepts of multioutput GPs to generate each latent function-(LF) from a semi-parametric latent factor model to code the correlation among the LFs. In consequence, CCGPMA takes into account the annotators' interdependencies. Throughout the study, the proposed frameworks are applied to synthetic and real annotators. Indeed, the introduced frequentist and Bayesian approaches improve, in most cases, supervised learning performances, aiding researchers in their quest to process data from multiple labelers and favor the understanding of their behavior.

From the developed approaches, we found the following strengths and drawbacks.

Regarding KAAR, we stand out that it is a non-parametric approach, which allows modeling the annotators' performance without assuming any labels' distribution, and it can be applied with different supervised learning algorithms to solve classification or regression problems. Further, KAAR can deal with scenarios where the labelers are not entirely independent. However, KAAR is built under the assumption that the annotator's performance does not depend on the input space, which is hardly feasible in real-world multi-labeler problems, as established in this dissertation.

Aiming to solve the KAAR's drawbacks, we propose the LKAAR, which codes labelers' behavior under two considerations. First, the annotators perform the labeling process based on their knowledge and the features observed from raw data. Second, the decision from one annotator is influenced by other annotators. LKAAR is also flexible, and it can be applied with any supervised learning algorithm; notwithstanding, like KAAR, our LKAAR solves the supervised learning problem as a convex combination of classifiers/regressors (one per labeler), which can be problematic if the number of labelers grows.

Conversely to KAAR and LKAAR, our RCDNN jointly estimates the annotators'

parameters and solves the supervised learning problem. RCDNN allows dealing with multi-labeler data where the annotators are not independent, and their performance depends on the input space. Moreover, we highlight that is flexible, and it can be built from different type of layers (convolutional, sequential), and using any activation functions (RELU, ELU, Softmax, Sigmoid, exponential). However, due to RCDNN is based on a frequentist approach, the estimations for the annotators' parameters and the supervised learning framework are deterministic, and it is not possible to capture the uncertainty related to them.

Contrariwise, KAAR, LKAAR, and RCDNN, CCGPMA is built from a Bayesian framework, which allows estimating the ground truth and the annotators' parameters. Remarkably, the estimation performed by our CCGPMA is probabilistic, and it allows capturing the uncertainty in the data. However, kernel selection is a critical step for this approach; hence, an improper kernel can lead to unsatisfactory outcomes.

As future lines of work, we analyze that it would be interesting to extend our models for multi-task learning in the context of multiple annotators.

Contents

List of Figures	xi
List of Tables	xvii
Notation	xx
1 Preliminaries	1
1.1 Motivation	1
1.2 Problem Statement	2
1.3 Mathematical preliminaries	4
1.3.1 Methods for supervised learning	5
1.3.2 Learning from multiple annotators	11
1.4 Literature review on supervised learning from multiple annotators	11
1.5 Objectives	18
1.5.1 General objective	18
1.5.2 Specific objectives	19
1.6 Outline and contributions	19
1.6.1 Kernel alignment-based annotator relevance analysis (KAAR)	20
1.6.2 Localized kernel alignment-based annotator relevance analysis (LKAAR)	21
1.6.3 Regularized chained deep neural network for multiple annotators (RCDNN)	24
1.6.4 Chained Gaussian processes for multiple annotators (CGPMA) and correlated chained Gaussian processes for multiple annotators (CCGPMA)	25
1.6.5 Thesis structure	26

2	Kernel alignment-based annotator relevance analysis	28
2.1	Centered kernel alignment fundamentals	29
2.2	Kernel alignment-based annotator relevance analysis	30
2.2.1	KAAR for classification and regression	31
2.3	Experimental set-up	32
2.3.1	Classification	33
2.3.2	Regression	36
2.4	Results and discussion	38
2.4.1	Classification	38
2.4.2	Regression	42
2.5	Summary	46
3	Localized kernel alignment-based annotator relevance analysis	48
3.1	Localized Kernel alignment fundamentals	49
3.2	Localized kernel alignment-based annotator relevance analysis	50
3.2.1	LKAAR for classification and regression	52
3.3	Experimental set-up	52
3.3.1	Classification	53
3.3.2	Regression	55
3.4	Results and discussion	56
3.4.1	Classification	56
3.4.2	Regression	66
3.5	Summary	70
4	Regularized chained deep neural network for multiple annotators	72
4.1	Chained Deep Neural Network	73
4.2	Regularized Chained Deep Neural Network for classification with multiple annotators	74
4.3	Experimental set-up	75
4.3.1	Tested datasets	75
4.3.2	Provided and simulated annotations	76
4.3.3	Method comparison and quality assessment	76
4.3.4	RCDNN detailed architecture and training	77
4.4	Results and discussion	78
4.5	Summary	87

5	Correlated chained Gaussian processes for multiple annotators	89
5.1	Chained Gaussian processes	90
5.1.1	Correlated chained Gaussian processes	92
5.2	Correlated chained GP for multiple annotators-CCGPMA	94
5.2.1	Classification	94
5.2.2	Regression	96
5.3	Experimental set-up	97
5.3.1	Classification	97
5.3.2	Regression	100
5.4	Results and discussion	101
5.4.1	Classification	101
5.4.2	Regression	108
5.5	Summary	115
6	Final Remarks	116
6.1	Conclusions	116
6.2	Future work	121
6.3	Academic discussion	122
6.3.1	Journal and conference papers	122
6.3.2	Repositories	123
	Appendix A CCGPMA supplementary material	124
A.1	Derivation of CCGPMA lower bounds	124
A.1.1	Gradients w.r.t. the variational parameters	125
A.2	Likelihood functions	127
A.2.1	Multiclass classification with multiple annotators	127
A.2.2	Gaussian distribution for regression with multiple annotators . . .	133
	References	138

List of Figures

1.1	Samples from a GP with SE covariance function over an 1- D problem with $s = 1$ and $l = 0.2$ (Rasmussen and Williams, 2006).	9
1.2	Graphical comparison between datasets for typical supervised learning settings and datasets from multiple annotators. We highlight that the annotators do not have access to the whole dataset; thus, it is common to face missing labels.	12
1.3	Graphical model for the approach proposed in Raykar et al. (2010). \mathbf{x}_n , y_n , and $y_n^{(r)}$ are respectively the input, the ground truth and the label given by the r -th annotator. Shaded nodes represent observed values, while unshaded nodes indicate latent variables.	13
1.4	The graphical model used to code the labelers expertise as a function of the input space. \mathbf{x}_n , y_n , and $y_n^{(r)}$ are respectively the input, the ground truth and the label given by the r -th annotator. Shaded nodes represent observed values, while unshaded nodes indicate latent variables.	16
1.5	Relevant state-of-the-art works for supervised learning with multiple annotators.	19
1.6	Relationship between the aims of this thesis and the developed methodologies. KAAR: Kernel alignment-based annotators relevance analysis, LKAAR: Localized kernel alignment-based annotator relevance analysis, RCDNN: Regularized chained deep neural network for Multiple Annotators, CGPMA: Chained Gaussian Processes for multiple annotators, and CCGPMA: Correlated chained Gaussian Processes for multiple annotators.	20

-
- 1.7 Kernel alignment-based annotator relevance analysis–(KAAR) pipeline. CKA stands for Centered Kernel Alignment. \mathcal{X} indicates the input space, and each \mathcal{Y}_r represents the output space for the r -th annotator. Moreover, $\boldsymbol{\nu} = [\nu_1, \dots, \nu_R]^\top$ is a vector containing the annotators' relevance parameters. Finally, the supervised learning mapping function f is computed as a convex combination of functions g_r that are trained from the r -th annotator's data. 22
- 1.8 Localized kernel alignment-based annotator relevance analysis–(KAAR) pipeline. CKA stands for Centered Kernel Alignment. \mathcal{X} indicates the input space, and each \mathcal{Y}_r represents the output space for the r -th annotator. Moreover, $[q_1(\mathbf{x}_n), \dots, q_R(\mathbf{x}_n)]^\top$ is a vector containing the annotators' relevance parameters for the n -th input sample. Finally, the supervised learning mapping function f is computed as a convex combination of functions g_r that are trained from the r -th annotator's data. 23
- 1.9 Regularized chained Deep Neural Network scheme. \mathbf{x} is an input vector, $\varrho_o(\cdot)$ is the o -th layer of a DNN, $f_j(\mathbf{x})$ is the j -th output neuron for the input \mathbf{x} . Moreover, each $\theta_j(\mathbf{x})$ is the j -th parameter of a given likelihood that is modeled as a function of the input space and J is the number of likelihood's parameters. Finally, $h_j(\cdot)$ is a deterministic function such that $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$ 24
- 1.10 Chained Gaussian processes. Each $f_j(\mathbf{x})$ is a latent function that follows a GP, Moreover, each $\theta_j(\mathbf{x})$ is the j -th parameter of a given likelihood and J is the number of likelihood's parameters. Finally, $h_j(\cdot)$ is a deterministic function such that $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$ 25
- 1.11 Correlated chained Gaussian processes. Each $f_j(\mathbf{x})$ is a latent function computed as a linear combination of functions $\{\mu_q(\mathbf{x})\}_{q=1}^Q$ that follow a GP, Moreover, each $\theta_j(\mathbf{x})$ is the j -th parameter of a given likelihood and J is the number of likelihood's parameters. Finally, $h_j(\cdot)$ is a deterministic function such that $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$ 26

2.1	Iris dataset illustrative results. On the first row from left to right, we show the similarities among samples according to the true labels and the input features. On the second, and third rows, from left to right, we expose the kernels computed from the annotators labels concerning each simulation method (see Table 2.5). Vertical and horizontal axes display the sample index which is sorted regarding the value of the ground truth labels. . . .	39
2.2	Annotators interdependencies. In the left column, from top to bottom, we show the Pearson coefficient computed over the labels from multiple annotators. In the right, we present the annotators' dependencies estimated with our KAAR.	41
2.3	Regression illustrative results. On the first column, from top to bottom, we show the regression results for each of the annotators. On the second column, from top to bottom, we exhibit the result for GPR-GT, the result for our GPR-KAAR; finally, in the last row, we show the Pearson coefficient computed over the labels from multiple annotators and the dependencies estimated with our KAAR.	44
3.1	Fully synthetic dataset results. First column (top to bottom): relevance values $q_r^2(\mathbf{x}_n)$. Second column (top to bottom): decision boundaries produced by a GPC trained over each annotator dataset. Third column (top to bottom): decision boundaries generated by both the GPC-GOLD (gold standard) and the LKAAR. Also, it displays the dependencies among the annotators estimated by LKAAR (from the input samples) vs. the Pearson correlation coefficients (absolute value) from the labelers' annotations.	57
3.2	Annotator dependencies analysis for the Voice dataset (Scale B). First row (left to right): t-sne-based 2D projections holding the annotators labels. red circles: positive class, blue crosses: negative class. Second row (left to right): dependencies among annotators estimated by LKAAR (from the input samples) vs. the Pearson correlation coefficients (absolute value) from the experts' labels.	64
3.3	Number of labeled instances per annotator.	65
3.4	Generalization performance in terms of AUC as a function of the number of annotators.	66

- 3.5 Regression illustrative results. On the first column, from top to bottom, we show the regression results g_r for each of the annotators. On the second column, from top to bottom, we present the LKAAR-based relevance q_r . In the third column, we exhibit the result for GPR-GOLD, the result for our GPR-LKAAR; finally, in the last row, we show the Pearson coefficient computed over the labels from multiple annotators and the dependencies estimated with our LKAAR. 67
- 4.1 RCDNN architecture details. ϱ_s stands for dense layer. ϱ_1 holds a linear activation, ϱ_2 includes a tanh-based activation, and ϱ_3 and ϱ_4 output the hidden ground truth label and the annotator's reliability fixing a softmax and a sigmoid activation, respectively. 79
- 4.2 RCDNN's decision boundaries for the 2D-PCA Iris dataset (synthetic scenario). AUC= 0.9837. The point's color stands for the Iris dataset classes. PCA₁ and PCA₂ stand for the first and second PCA-based projections. . . 80
- 4.3 RCDNN-based annotators performance (reliability) estimation for the synthetic experiments (2D PCA Iris data). In the first column (from top to bottom), the simulated accuracy for each annotator is presented based on equation (4.10). The second column shows (from top to bottom) the estimated annotators' reliability (λ_r). 81
- 4.4 Target and estimated annotators dependencies for the synthetic 2D PCA Iris dataset. In the left, the Pearson correlation coefficients (absolute value) from simulated accuracies (experts reliability) in matrix \mathbf{P} of equation (4.10) are shown. In the right, the dependencies among the annotators estimated from the RCDNN ϱ_4 layer' weights are displayed. 82
- 4.5 Receiver Operating Characteristic (ROC) plot for the annotators simulated within the spammers and malicious scenario. Blue dots indicate the basis annotators. Red dots show extra annotators with parameters $R_e = 65$, and $p_r = 0.5$. Green dots specify extra labelers with $R_e = 20$, and $p_r = 0.6$. We notice that annotators located in dashed line vicinity are considered Spammers. Similarly, labelers above the dashed line are regarded as good annotators; conversely, labelers located below such a line are malicious annotators. 86
- 4.6 MA-LFC, LKAAR-GPC, and RCDNN performance (AUC) as a function of the number of labelers (spammers and malicious annotators). 87

- 5.1 Fully synthetic dataset results. The PLP is shown, comparing the prediction of our CCGPMA-C($AUC = 1$) and CCGPMA-C($AUC = 0.9999$) against: the theoretical upper bound GPC-GOLD($AUC = 1.0$), the lower bound GPC-MV($AUC = 0.9809$), and the state-of-the-art approaches MA-LFC-C($AUC = 0.9993$), MA-DGRL($AUC = 0.9999$), MA-GPC($AUC = 0.9977$), MA-GPCV($AUC = 0.9515$), MA-DL-MW($AUC = 0.9989$), MA-DL-VW($AUC = 0.9972$), MA-DL-VW+B($AUC = 0.9994$), KAAR(0.9099). Note that the shaded region in GPC-MV, CGPMA-C, and CCGPMA-C indicates the area enclosed by the mean \pm two standard deviations. There is no shaded region for approaches lacking prediction uncertainty. 102
- 5.2 Fully synthetic data reliability results. From top to bottom, the first column exposes the true reliabilities (λ_r). The subsequent columns present the estimation of the reliabilities performed by state-of-the-art models, where the correct values are provided in dashed lines. The shaded region in CGPMA-C and CCGPMA-C indicates the area enclosed by the mean \pm two standard deviations. Also, the accuracy (Acc) is provided. 104
- 5.3 Fully synthetic dataset results. We compare the prediction of our CCGPMA-R($R^2 = 0.9438$), and CGPMA-R($R^2 = 0.9280$) with the theoretical upper bound GPR-GOLD($R^2 = 0.9843$) and lower bound GPR-Av($R^2 = 0.8718$), and state-of-the-art approaches, MA-LFCR($R^2 = -0.0245$), MA-GPR($R^2 = 0.9208$), MA-DL-B($R^2 = 0.7020$), MA-DL-S($R^2 = 0.6559$), MA-DL-B+S($R^2 = 0.5997$). Note that we provided the Gold Standard in dashed lines. The shaded region in GPR-Av, MA-GPR, CGPMA-R, and CCGPMA-R indicates the area enclosed by the mean plus or minus two standard deviations. We remark that there is no shaded region for MA-LFCR, and DLMA since they do not provide information about the prediction uncertainty. 109

- 5.4 Estimated values of error-variance for the five annotators in the *fully synthetic* experiment. In the first column, from top to bottom, we expose the error-variances used to simulate the labels from each annotator. Furthermore, the subsequent columns from top to bottom present the estimation of such error-variances performed by state-of-the-art models that include these kinds of parameters in their formulation; moreover, the true error-variances are provided in dashed lines. The shaded region in CGPMA-R and CCGPMA-R indicates the area enclosed by the mean plus or minus two standard deviations. We remark that there is no shaded region for MA-LFCR, and MA-GPR since these approaches perform a fixed-point estimation for the annotators' parameters. Finally, we remark that the R^2 score between the true and estimated error variances are provided. . . 111
- 6.1 Malicious annotators. On the left, we show the Receiver Operating Characteristic (ROC) plot for the simulated annotators. Red dots indicate the basis annotators. Green dots show malicious annotators. On the right, we present the performance (AUC) of DGRL, MA-GPC, MA-DL-VW+B, KAAR, LKAAR, RCDNN, CGPMA, and CCGPMA as a function of the number of malicious annotators. 120

List of Tables

2.1	Datasets description.	33
2.2	A brief overview of the state-of-the-art methods tested. GPC: Gaussian Processes classifier, LRC: logistic regression classifier, MV: majority voting, MA: multiple annotators, MAE: Modelling annotators expertise, LFC: Learning from crowds, DGRL: Distinguishing good from random labelers.	35
2.3	Datasets for regression.	37
2.4	A brief overview of state-of-the-art methods tested for regression tasks. GPR: Gaussian Processes Regression, LR: logistic regression, Av: average, MA: multiple annotators, DL: Deep learning, LFCR: Learning from crowds for regression.	37
2.5	KAAR-based annotator coding results for the iris dataset. Targets: the true labels (ν_{TL}) and the input features (ν_{IF}).	40
2.6	UCI repository classification results. Bold: the highest performance excluding the upper bound (target) classifier GPC-GOLD	41
2.7	Real annotators datasets results. Bold: the method with the highest performance excluding the upper bound (target) classifier GPC-GOLD.	42
2.8	KAAR-based annotator coding results the <i>fully synthetic</i> dataset.	43
2.9	Regression results in terms of R^2 score over <i>semi synthetic datasets</i> . Bold: the highest R^2 excluding the upper bound GPR-GOLD.	43
2.10	Regression results in terms of R^2 score over <i>fully real dataset</i> . Bold: the highest R^2 excluding the upper bound GPR-GOLD.	45

3.1	A brief overview of the state-of-the-art methods tested. GPC: Gaussian Processes classifier, LRC: logistic regression classifier, MV: majority voting, MA: multiple annotators, MAE: Modelling annotators expertise, LFC: Learning from crowds, DGRL: Distinguishing good from random labelers, KAAR: kernel alignment-based annotator relevance analysis.	55
3.2	A brief overview of state-of-the-art methods tested for regression tasks. GPR: Gaussian Processes Regression, LR: logistic regression, Av: average, MA: multiple annotators, DL: Deep learning, LFCR: Learning from crowds for regression.	56
3.3	UCI repository classification results. Bold: the highest AUC excluding the upper bound (target) classifier GPC-GOLD. Marked with *: the highest accuracy (Acc) except the upper bound.	60
3.4	Fully real datasets results. Bold: the method with the highest performance excluding the upper bound (target) classifier GPC-GOLD.	62
3.5	Regression results in terms of R^2 score over <i>semi synthetic datasets</i> . Bold: the highest R^2 excluding the upper bound GPR-GOLD.	69
3.6	Regression results in terms of R^2 score over <i>fully real dataset</i> . Bold: the highest R^2 excluding the upper bound GPR-GOLD.	69
4.1	A short overview of the tested state-of-the-art approaches. GPC: Gaussian Processes classifier, LRC: logistic regression classifier, MV: majority voting, MA: multiple annotators, MAE: Modelling annotators expertise, LFC: Learning from crowds, DGRL: Distinguishing good from random labelers, KAAR: kernel alignment-based annotator relevance analysis, LKAAR: localized kernel alignment-based annotator relevance analysis.	77
4.2	Semi-synthetic datasets results. Bold: the highest AUC excluding the upper bound (target) classifier GPC-GOLD. Marked with *: the highest accuracy (Acc) except the upper bound. The last column presents the average ranking for both the AUC score and the overall accuracy (GPC-GOLD is not considered), the best average ranking for AUC is highlighted in bold, and the accuracy is marked with *.	84
4.3	Fully real-world datasets results. Bold: the method with the highest performance excluding the upper bound (target) classifier GPC-GOLD.	85
5.1	A brief overview of the state-of-the-art methods tested.	100

5.2	Semi-synthetic classification results. Bold: the highest AUC excluding the upper bound (target) classifier GPC-GOLD. Marked with *: the highest accuracy (Acc) except the upper bound.	106
5.3	AUC classification results for the fully real datasets. Bold: the highest performance excluding the GPC-GOLD bound.	108
5.4	Semi-synthetic regression results. Bold: the highest R^2 excluding the upper bound (target) classifier GPR-GOLD. Marked with *: the highest Pearson coefficient (Pear) except the upper bound.	112
5.5	Regression results in terms of R^2 score over <i>fully real dataset</i> . Bold: the highest R^2 excluding the upper bound GPR-GOLD.	114
6.1	AUC classification results for the semi-synthetic and fully real datasets. Bold: the highest performance excluding the GPC-GOLD bound. The last column presents the average ranking for AUC score, in bold the best average ranking. The Friedman test returns a Chi-square value of 26.84 ($p - value = 0.0081$).	119

Notation

Mathematical notation

Generalities

P	input space dimensionality
N	total number of data points
R	number of annotators
K	number of classes
Q	number of latent functions
J	number of likelihood parameters
M	number of inducing points
S	number of points in a Gaussian Hermite quadrature
N_r	a set containing the indexes of samples labeled by the n -th annotator $\{1, \dots, N\}$
$ N_r $	cardinality of set N_r
R_n	a set containing the indexes of annotators that labeled the n -th sample $\{1, \dots, R\}$
$ R_n $	cardinality of set R_n
\mathbb{R}	set of the real numbers
ψ	regularization parameter
\mathcal{X}	input space
\mathcal{Y}	output space
\mathcal{H}_κ	reproducing kernel Hilbert space
l	length-scale for squared exponential kernel
s	output-scale amplitude for squared exponential kernel
σ^2	variance parameter in a Gaussian distribution
ι	variance parameter in a Gaussian distribution

ψ	regularization parameter
ϵ	random number sampled from a Gaussian distribution
τ	random number sampled from a Bernoulli distribution
\mathcal{M}_j	The space regarding the parameter $\theta(\mathbf{x})$
w	combination parameter in a semi-parametric latent factor model

Operators

$\ \cdot\ _{\mathcal{H}_\kappa}$	the norm in a reproducing kernel Hilbert space
$\langle \cdot, \cdot \rangle_{\mathcal{H}_\kappa}$	the inner product in a reproducing kernel Hilbert space
$\ \cdot\ _{\mathbf{F}}$	Frobenius norm
$\langle \cdot, \cdot \rangle_{\mathbf{F}}$	Frobenius product
$\mathbb{E}[\cdot]$	expected value
$\text{tr}(\cdot)$	trace of a matrix
$ \cdot $	determinant of a matrix

Functions

$f(\cdot)$	function estimated from a supervised learning algorithm
$\kappa(\cdot, \cdot)$	kernel function
$u_q(\mathbf{x})$	q -th input modeled as a Gaussian processes
$\xi_d(\cdot)$	d -th Basis function
$\varsigma(\cdot)$	Sigomid function
$\Xi(\cdot)$	Softmax function
$\varrho_o(\cdot)$	o -th layer's representation if a Deep neural network.
$\theta_j(\mathbf{x})$	represents the j -th parameter in a given likelihood
$\mu(\cdot)$	latent function in a semi-parametric latent factor model

Vectors and matrices

\mathbf{x}	a P -dimensional input vector
\mathbf{k}_*	a kernel matrix with entries $[\kappa(\mathbf{x}_1, \mathbf{x}_*), \dots, \kappa(\mathbf{x}_N, \mathbf{x}_*)]^\top \in \mathbb{R}^{N \times 1}$
\mathbf{K}	a kernel matrix formed by the evaluation of $\kappa(\cdot, \cdot)$ over \mathbf{X}
\mathbf{I}	The identity matrix $\mathbf{I} \in \mathbb{R}^{N \times N}$
\mathbf{f}	a vector containing elements $[f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N$
$\boldsymbol{\nu}$	weights vector for multiple annotators, $\boldsymbol{\nu} = [\nu_1, \dots, \nu_R] \in \mathbb{R}^R$
\mathbf{a}	sensitivities vector for multiple annotators $\mathbf{a} = [a_1, \dots, a_R] \in \mathbb{R}^R$
\mathbf{b}	specificities vector for multiple annotators $\mathbf{b} = [b_1, \dots, b_R] \in \mathbb{R}^R$
$\boldsymbol{\omega}, \gamma$	vector of parameters for a logistic regression model $\boldsymbol{\omega} = [\omega_1, \dots, \omega_P]^P$
\mathbf{K}_{ff}	Covariance matrix for function \mathbf{f} formed by the evaluation of $\kappa_f(\cdot, \cdot)$ over \mathbf{X} .
$\mathbf{K}_{f_j f_j}$	Covariance matrix belonging to \mathbf{f}_j .
$\hat{\mathbf{f}}$	vector containing $[\mathbf{f}_1, \dots, \mathbf{f}_J]^\top \in \mathbb{R}^{NJ}$
$\boldsymbol{\theta}$	parameters vector containing $[\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_J]^\top \in \mathbb{R}^{NJ}$
$\boldsymbol{\theta}_j$	vector related to the j -th parameter $[\theta_j(\mathbf{x}_1), \dots, \theta_j(\mathbf{x}_N)]^\top \in \mathbb{R}^N$
\mathbf{u}_q	“pseudo variables” vector $[\mu_q(\mathbf{z}_1^{(q)}), \dots, \mu_q(\mathbf{z}_M^{(q)})]^\top \in \mathbb{R}^M$ for the function μ_q
\mathbf{Z}_q	inducing points vector $[\mathbf{z}_1^{(q)}, \dots, \mathbf{z}_M^{(q)}] \in \mathbb{R}^{M \times P}$ for the function μ_q
\mathbf{u}	vector containing $[\mathbf{u}_1^\top, \dots, \mathbf{u}_Q^\top]^\top \in \mathbb{R}^{QM}$
\mathbf{K}_{uu}	A block-diagonal matrix with blocks $\mathbf{K}_{\mathbf{u}_q \mathbf{u}_q} \in \mathbb{R}^{M \times M}$
\mathbf{m}_q	the mean vector of the variational distribution $q(\mathbf{u}_q)$
\mathbf{V}_q	Covariance of the variational distribution $q(\mathbf{u}_q)$
\mathbf{m}	Mean vector containing $[\mathbf{m}_1^\top, \dots, \mathbf{m}_Q^\top]^\top \in \mathbb{R}^{QM}$
\mathbf{V}	Block-diagonal matrix with blocks given by the covariance matrices \mathbf{V}_q . holds elements

Abbreviations

RKHS	Reproductive kernel Hilbert space
FM	Models using a frequentist perspective
PM	Probabilistic or Bayesian Models
GP	Gaussian Process
CGP	Chained Gaussian Processes
ELBO	Evidence Lower Bound
MAP	Maximum a posteriori
EM	Expectation Maximization
CKA	Centered Kernel Alignment
GP	Gaussian processes
CCGPMA	Correlated chained Gaussian processes
KAAR	Kernel alignment-based annotator relevance
LKAAR	Localized kernel alignment-based annotator relevance
LF	Latent function
DNN	Deep neural networks
RDNN	Regularized deep neural networks
RCDNN	Regularized chained deep neural networks
SE	Squared exponential kernel
MV	Majority voting
GPR	Gaussian Processes Regression
LR	logistic regression
Av	average
MA	multiple annotators
DL	Deep learning
LFCR	Learning from crowds for regression
LRC	logistic regression classifier
MV	majority voting
MAE	Modelling annotators expertise
DGRL	Distinguishing good from random labelers
VGPCR	Variational Gaussian processes for crowdsourcing
SVGPCR	Scalable variational Gaussian processes for crowdsourcing
GRBAS	G: grade of disphony. R: Roughness. B: Breathiness. A: Anthenia. S: Strain.

Chapter 1

Preliminaries

1.1 Motivation

The accelerated use of information technologies in different domains is changing the way datasets are built. For instance, dedicated crowdsourcing platforms like Amazon Mechanical Turk (AMT), LabelMe, Crowdfunder allow extracting the “wisdom of crowds” (Welinder et al., 2010) to obtain large datasets labeled by multiple annotations. However, the concept of crowdsourcing goes beyond dedicated platforms; for example, the web, through its social nature, allows collecting a large volume of labeled datasets, where such labels commonly correspond to web-based ratings (products rating, translation rating, product tag). The attractiveness of crowdsourcing lies in the possibility of getting suitable quality labels at a low-cost (Albarqouni et al., 2016; Zhang et al., 2016).

The aforementioned poses a new challenge in the supervised learning context. Instead of having datasets labeled by one source (which is supposed to be an expert who provided the absolute gold standard), we have datasets labeled by multiple annotators with different and unknown expertise (Sung et al., 2017). Hence, typical supervised learning algorithms need to be adapted to face multi-labeler datasets.

Besides, in a local context, the research group in automatics (RGA) has successfully developed several research works related to the analysis of different medical images (ultrasound, X-rays, MRI, DMRI) aiming to support the identification of nerves structures and the diagnosis of various pathologies, including Parkinson, breast cancer, and brain cancer (Jimenez et al., 2018a,b; Vargas-Cardona et al., 2019). Most of these works include a supervised learning step, where it is necessary to hire experts to obtain the required labels for training the learning algorithms. However, the experts are scarce, and their time is costly; moreover, there exists disagreement among the annotations given

by the multiple annotators, especially if they have different levels of expertise.

Accordingly, from the global and local scenarios, it is necessary to continue developing methodologies that allow facing supervised learning problems in the context of multiple annotators.

1.2 Problem Statement

The main aim of a supervised learning task is to learn a function that maps from the input features to the output space, which is estimated by using a training set, where it is supposed that an expert provides the actual label (termed gold standard) for each instance (Bishop, 2006). Nonetheless, in many real-world applications, such a gold standard is not available since the experts are scarce, their time is expensive, and the labeling problem is tedious, and time-consuming (Liu et al., 2020; Zhu et al., 2019). Similarly, other labeling problems correspond to a subjective evaluation (e.g., sentimental analysis, products rating), and thus, the gold standard is not clear (Huang et al., 2018). Instead of the gold standard, we may have access to several noisy annotations provided by R heterogeneous annotators (also named as sources or labelers), where each source gives its version of the unknown gold standard (Kara et al., 2015). Those annotations can be collected in several ways. For instance, in the model proposed in Raykar et al. (2010), the sources correspond to physicians who make a diagnosis about the presence of cancer using medical images. Likewise, in the approach proposed by the authors in Zhu et al. (2019), the labels are provided by algorithms, which are used to measure the QT interval in an electrocardiogram signal (the QT interval is a measurement used to assess some electrical properties of the heart).

Accordingly, conversely to typical supervised learning settings, in multi-labeler scenarios, each instance is linked with multiple annotations provided by multiple annotators. Commonly, not all the labelers give an output for each input in the dataset. Thereby, it is not straightforward to apply traditional supervised learning algorithms in the presence of data from multiple annotators (Tao et al., 2018). In this sense, *learning from crowds* has been introduced as a general framework from two main perspectives: to fit the labels from multiple annotators or to adapt the supervised learning algorithms (Rizos and Schuller, 2020).

The first approach is known in the literature as “label aggregation” or “truth inference”, comprising the computation of a single hard label per sample as an estimation of the ground truth. The hard labels are then used to feed a standard supervised learning

algorithm (Morales-Alvarez et al., 2020). The straightforward method is the so-called majority voting–(MV), and it has been used in different multi-labeler problems due to its simplicity (Zhang et al., 2014). Still, MV assumes homogeneity in annotators’ reliability, which is hardly feasible in real applications, e.g., experts vs. spammers. Furthermore, the consensus is profoundly impacted by incorrect labels and outliers (Kara et al., 2015). Conversely, more elaborated models have been considered to improve the estimation of the correct tag through the well-known Expectation-Maximization–(EM) framework and by facing the imbalanced labeling issue (Dawid and Skene, 1979; Zhang et al., 2014).

The second approach jointly trains the supervised learning algorithm and models the annotators’ behavior. It has been shown that such strategies lead to better performance than those belonging to label aggregation. Thus, the features used to train the learning algorithm provide valuable information to puzzle out the ground truth (Ruiz et al., 2019). The most representative work in this area is exposed in Raykar et al. (2010), which offers an EM-based framework to learn the parameters of a logistic regression classifier and model the annotators’ behavior by computing their sensitivities and specificities. In fact, such a technique has inspired several models in the context of multi-labeler scenarios, including binary classification (Rodrigues et al., 2014b; Ruiz et al., 2019), multi-class discrimination (Gil et al., 2015; Morales-Alvarez et al., 2020), regression (Groot et al., 2011; Rodrigues et al., 2017), and sequence labeling (Rodrigues et al., 2017, 2014a). Furthermore, some works have addressed the multi-labeler problem using deep learning approaches, typically including an extra layer that codes the annotators’ information (Albarqouni et al., 2016; Guan et al., 2018; Rodrigues and Pereira, 2018a).

According to the above, among the works developed in this area, we recognize two main groups: the approaches based on frequentist models–(FMs) and the schemes based on probabilistic or Bayesian models–(BMs). Furthermore, we note that most works on the topic of learning from crowds are mainly based on two assumptions: *i) The performance of the annotators does not depend on the feature space* and *ii) independence among the annotators*; however, these assumptions are not valid in practice. Hence, the primary motivation of this proposal is to address these two problems.

The performance of the annotators does not depend on the feature space

In most multi-labeler approaches, it is necessary to estimate the performance of the annotators, which is usually measured in terms of the accuracy (Rodrigues et al., 2013), or sensitivity and specificity (Albarqouni et al., 2016; Raykar et al., 2010; Ruiz et al.,

2019) in classification settings; similarly, in regression approaches, such performance is measured in terms of the error variance (Groot et al., 2011). However, a restriction commonly seen in these algorithms is that they assume that the performance is consistent across the input space. This assumption is an impractical restriction since the expertise of the annotators may vary depending on the instance they label (Xie et al., 2019; Yan et al., 2014). For example, if we consider online annotators assessing some documents, they may have different labeling accuracy. Such differences may rely on whether they are more familiar with some specific topics related to studied documents (Wang and Bi, 2016).

Independence among the annotators Another assumption commonly seen is to consider independence among the annotators. This assumption is used to reduce the complexity of the model (Rodrigues, 2016; Venanzi et al., 2014), or based on the fact that it is plausible to guarantee that each labeler performs the annotation process individually (Tang et al., 2019). Nevertheless, this assumption cannot fit most real-world applications. For example, if the sources are humans, the independence assumption is hardly feasible because knowledge is a social construction; hence, people decisions will be correlated since they share information, they communicate with each other, or because they belong to a particular school of thought (Hahn et al., 2018; Surowiecki, 2005). On the other hand, if we consider that the sources are algorithms, where some of them are based on the same math principle, there likely exists a correlation among their labels (Zhu et al., 2019). Accordingly, the relaxation of this restriction could be used to improve the ground truth estimation (G. Rodrigo et al., 2019).

Therefore, some of the problems related to supervised learning with multiple annotators remain unsolved. For this reason, the following research question arises. How to develop a multiple annotators framework based on FMs and BMs aiming to code the labelers' performance as a function of the input space and considering dependencies among them to improve the multi-labeler representation ability for classification and regression tasks?

1.3 Mathematical preliminaries

In this section, we introduce the mathematical formulation for the problem of learning from multiple annotators. First, in Section 1.3.1 we recall the mathematical formulation of a traditional supervised learning scenario. We analyze two perspectives to face

such scenarios: The frequentist and the Bayesian point of view. Finally, Section 1.3.2, describes the composition of a dataset from multiple labelers and the main aims of algorithms dealing with this kind of data.

1.3.1 Methods for supervised learning

A classical supervised learning problem comprises the estimation of a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, where \mathcal{X} , and \mathcal{Y} are respectively the input space and the output space (Bishop, 2006). Such a function f is computed from of a training set $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, where $\mathbf{X} = \{\mathbf{x}_n \in \mathcal{X} \subseteq \mathbb{R}^P\}_{n=1}^N$ and $\mathbf{y} = \{y_n \in \mathcal{Y}\}_{n=1}^N$. Depending on the output space \mathcal{Y} nature, we recognize different supervised learning settings: *i) binary classification* holding $\mathcal{Y} \in \{-1, 1\}$ or $\mathcal{Y} \in \{0, 1\}$; *ii) multi-class classification*, where $\mathcal{Y} \in \{1, \dots, K\}$, being K the number of classes; and *iii) regression* where $\mathcal{Y} \in \mathbb{R}$ (Raykar et al., 2010).

Notice that in most supervised learning cases, we are not interested in knowing the exact form f , but we are focused on computing the output in a specific value $f(\mathbf{x}_*)$, where $\mathbf{x}_* \in \mathbb{R}^P$ (Rasmussen and Williams, 2006).

One of the most basic estimators is the well-known generalized linear estimator given by:

$$f(\mathbf{x}_*) = \vartheta_0 + \sum_{d=1}^D \vartheta_d \xi_d(\mathbf{x}_*), \quad (1.1)$$

where $\{\vartheta_d\}_{d=0}^D$ are the model parameters, and $\xi_d : \mathbb{R}^P \rightarrow \mathbb{R}$ is the d -th component of a pre-selected set of non-linear functions, named basis functions, aiming to deal with non-linearities in the data's structure (Theodoridis, 2015). The main advantage of this type of estimators lies in the fact the related models remains linear regarding the set of parameters $\{\vartheta_d\}_{d=0}^D$. Nevertheless, due to this estimator is parametric, for high-dimensional spaces, the number of basis functions D has to be large to obtain a proper performance, leading to overfitting (Barron, 1993). Such drawback has been addressed in the literature; for example, Support Vector Machines (SVMs) defines a set of basis functions centered in the training data. Then, a subset of them is selected during the training, where the number of selected functions is generally smaller than the training set's size. Another alternative propound to use parametric forms for the basis functions (the number of functions is given in advance), in which the parameters are estimated during the training stage. The most successful approach in this context is the well-known

artificial neural networks—(ANN) (Bishop, 2006).

On the other hand, in contrast of previous approaches, that are parametrized in terms of a set of basis functions, it is possible to introduce a symmetric positive bivariate function, termed *kernel*, which enables to build non-parametric estimators less prone to overfitting. Thus, a way to estimate $f(\mathbf{x}_*)$ is given as follows (Álvarez et al., 2012):

$$f(\mathbf{x}_*) = \mathbf{k}_*^\top (\mathbf{K} + \psi N \mathbf{I})^{-1} \mathbf{y}, \quad (1.2)$$

where $\mathbf{k}_* = [\kappa(\mathbf{x}_1, \mathbf{x}_*), \dots, \kappa(\mathbf{x}_N, \mathbf{x}_*)]^\top \in \mathbb{R}^{N \times 1}$, and $\kappa: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a kernel function. Besides, $\mathbf{K} \in \mathbb{R}^{N \times N}$ is formed by the evaluation of κ over the input set \mathbf{X} , and $\mathbf{I} \in \mathbb{R}^{N \times N}$ is the identity matrix. Finally, $\psi \in \mathbb{R}^+$ is a regularization hyper-parameter, which will be discussed in the next sections. Hence, such estimator in equation (1.2) can be derived from two different but related perspectives; namely, a frequentist (regularization), and a Bayesian perspective.

A frequentist perspective

From a frequentist perspective (regularization), the aim is to suppose the function of interest to belong to a reproducing kernel Hilbert space—(RKHS) \mathcal{H}_κ , which is generated by the chosen kernel κ ; thereby, we have $f \in \mathcal{H}_\kappa$. Accordingly, the estimator can be inferred by minimizing the following regularized functional

$$\hat{f} = \arg \min_{f \in \mathcal{H}_\kappa} \frac{1}{N} \sum_{n=1}^N (f(\mathbf{x}_n) - y_n)^2 + \psi \|f\|_{\mathcal{H}_\kappa}^2, \quad (1.3)$$

where $\|\cdot\|_{\mathcal{H}_\kappa}$ is the norm in \mathcal{H}_κ . Besides, we notice that the first term in equation (1.3) corresponds to the well-known empirical risk, which particularly is computed as the sum of the squared errors (Álvarez et al., 2012). Such a term penalizes the variation between predicted outputs $f(\mathbf{x}_n)$ and the corresponding true value y_n . The second term in equation (1.3) represents the regularization expression that controls the estimator's smoothness to avoid overfitting; hence, the larger ψ value, the smoother the estimation \hat{f} . The regularization concept in a RKHS plays a key role in the machine learning area, so it is necessary to review some notions about RKHS (Kanagawa et al., 2018). A RKHS is a Hilbert space of functions defined by a reproducing kernel κ , which is a symmetric positive definite function. In that sense, given a kernel function κ a RKHS is generated in such a way that the function $\kappa(\mathbf{x}, \cdot)$ is linked to \mathcal{H}_κ for all $\mathbf{x} \in \mathcal{X}$, and (Schölkopf

et al., 2002)

$$f(\mathbf{x}) = \langle f, \kappa(\mathbf{x}, \cdot) \rangle_{\mathcal{H}_\kappa}; \quad \forall f \in \mathcal{H}_\kappa,$$

where $\langle \cdot, \cdot \rangle_{\mathcal{H}_\kappa}$ is the inner product in \mathcal{H}_κ . The above expression corresponds to the *reproducing* property. Further, we recognize two additional properties, which are essential for the regularization perspective. First, the *representer* theorem indicates that in a RKHS the functions are built as a linear combination of the kernel in given points (Schölkopf et al., 2001)

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \kappa(\mathbf{x}, \mathbf{x}_n) = \mathbf{k}^\top \boldsymbol{\alpha}, \quad (1.4)$$

where $\mathbf{k} = [\kappa(\mathbf{x}_1, \mathbf{x}), \dots, \kappa(\mathbf{x}_N, \mathbf{x})]^\top \in \mathbb{R}^{N \times 1}$ and $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_N]^\top \in \mathbb{R}^{N \times 1}$ are the combination parameters. Second, the norm of a function in a given RKHS can be written as:

$$\|f\|_{\mathcal{H}_\kappa}^2 = \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}. \quad (1.5)$$

Then, replacing equations (1.4) and (1.5) in equation (1.3), the optimization problem becomes:

$$\hat{\boldsymbol{\alpha}} = \arg \min_{\boldsymbol{\alpha}} \frac{1}{N} \left(\boldsymbol{\alpha}^\top \mathbf{K} \mathbf{K} \boldsymbol{\alpha} - 2 \boldsymbol{\alpha}^\top \mathbf{K} \mathbf{y} + \mathbf{y}^\top \mathbf{y} \right) + \psi \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}. \quad (1.6)$$

We remark that the latter objective function is convex with respect to $\boldsymbol{\alpha}$; thereby, differentiating such an objective function with respect to $\boldsymbol{\alpha}$, and equaling zero, the minimizer for equation (1.6) yields:

$$\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \psi N \mathbf{I})^{-1} \mathbf{y}.$$

Accordingly, the output estimation for a new instance $f(\mathbf{x}_*)$ is given as follows

$$f(\mathbf{x}_*) = \sum_{n=1}^N \alpha_n \kappa(\mathbf{x}_*, \mathbf{x}_n) = \mathbf{k}_*^\top \hat{\boldsymbol{\alpha}}.$$

A Bayesian perspective

On the other hand, there exists another perspective to solve supervised learning problems, named Bayesian. We identify controversy in the exact definition of Bayesian approaches; for instance, Bishop (2006) establishes that the Bayesian formulation stands for considering f to be random; nevertheless, Kay and Kay (1993) sets that Bayesian estimation can also be used for deterministic computation of f . In this thesis, we follow the definition in Bishop (2006); thus, we consider as a Bayesian treatment only the models that fix f to be random (or the model parameters if it is the case).

For specific demonstration, we use the well-known Gaussian processes (GPs), one of the most common Bayesian approaches in the supervised learning context. Namely, a GP is a random process, where any finite set of samples follows a Gaussian distribution (Rasmussen and Williams, 2006). Typically, a GP is employed as a prior over functions. Hence, let be $f(\mathbf{x})$ a function that follows a GP, $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), \kappa_f(\mathbf{x}, \mathbf{x}'))$, where $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$ is the mean function (usually $m(\mathbf{x}) = 0$), and $\kappa_f(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$ is the covariance function with $\kappa_f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ being a given kernel function ($\mathbf{x}' \in \mathcal{X}$).

If we consider the finite set of inputs in \mathbf{X} , then $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N$ is drawn for a multivariate Gaussian distribution $\mathbf{f} \sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{ff})$, where $\mathbf{K}_{ff} \in \mathbb{R}^{N \times N}$ is the covariance matrix formed by the evaluation of κ_f over the input set \mathbf{X} .

From a Bayesian point of view, the GPs priors represent our beliefs about the properties of the function we are modeling (Álvarez et al., 2012). Let us consider the squared exponential (SE) kernel, which is given as:

$$\kappa_{\text{SE}}(\mathbf{x}, \mathbf{x}') = s^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}\right),$$

with the parameter $l \in \mathbb{R}^+$ defining the characteristic length-scale and $s \in \mathbb{R}^+$ specifying an output-scale amplitude. We notice in Figure 1.1 that a GP prior using a SE kernel function prefers smooth functions (Rodrigues and Pereira, 2018b). Then, the beliefs captured by the GP prior are updated in the presence of data through the *likelihood* function. This leads to an updated distribution, named the *posterior distribution*, that can be used to make predictions over new samples \mathbf{x}_* .

In a regression setting, the likelihood function is usually Gaussian and codes the linear relation between the observations and a given model corrupted with zero mean

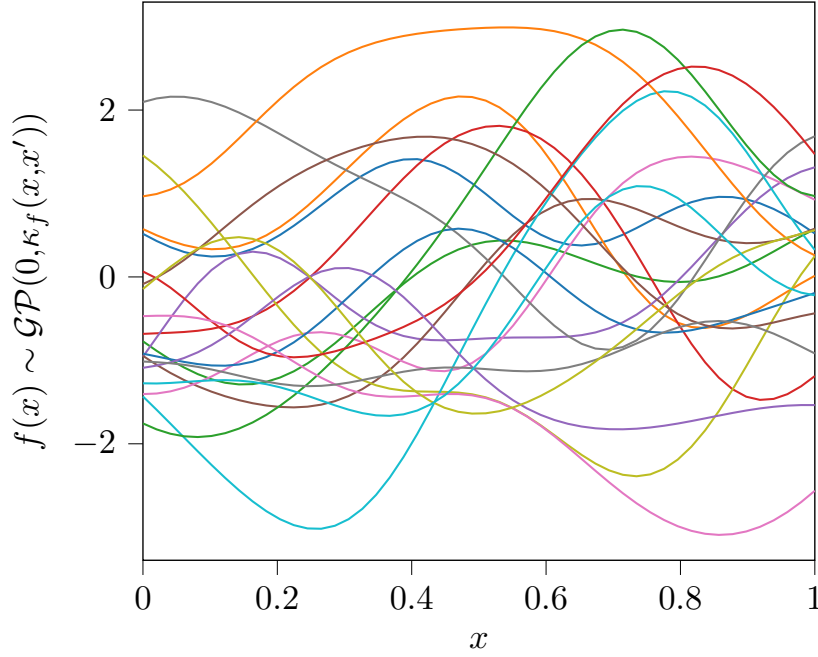


Figure 1.1: Samples from a GP with SE covariance function over an 1- D problem with $s = 1$ and $l = 0.2$ (Rasmussen and Williams, 2006).

Gaussian noise,

$$p(\mathbf{y}|\mathbf{f}, \mathbf{X}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(y_n|f(\mathbf{x}_n), \sigma^2),$$

where $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N$, and $\sigma^2 \in \mathbb{R}^+$ is the noise variance. Notice that the likelihood function factorizes over the data points, which assumes that the noise is independent and identically distributed. In this case, due to a Gaussian likelihood, the posterior distribution has an analytic solution. Given a new test point \mathbf{x}_* and the training data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, the predicted distribution is computed as:

$$p(f(\mathbf{x}_*)|\mathcal{D}, \mathbf{x}_*, \phi) = \mathcal{N}(\bar{f}(\mathbf{x}_*), k(\mathbf{x}_*, \mathbf{x}_*)), \quad (1.7)$$

where ϕ indicates the model hyper-parameters including σ^2 and the hyper-parameters related to the kernel function. Moreover, we have

$$\bar{f}(\mathbf{x}_*) = \mathbf{k}_{ff*}^\top (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \quad (1.8)$$

$$k(\mathbf{x}_*, \mathbf{x}_*) = k_{f_*f_*} - \mathbf{k}_{ff*}^\top (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{ff*}, \quad (1.9)$$

where $\mathbf{k}_{ff_*} \in \mathbb{R}^N$ is the cross-covariance function between \mathbf{f} and $f_* = f(\mathbf{x}_*)$; besides, $k_{f_*f_*} \in \mathbb{R}$ is the covariance function for f_* . Conversely, the model hyper-parameters ϕ are estimated by minimizing the minus logarithm of the marginalized likelihood function, as follows

$$\hat{\phi} = \arg \min_{\phi} (-\log p(\mathbf{y}|\mathbf{X})) = \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_{ff} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} + \frac{1}{2} \log |\mathbf{K}_{ff} + \sigma^2 \mathbf{I}|$$

where $|\cdot|$ indicates the determinant of a matrix.

Connections between Frequentist and Bayesian point of views

In general terms, the dimension of a RKHS can be infinity (Álvarez et al., 2012). However, if we limit the analysis to finite-dimensional RKHS, it is possible to determine that every RKHS can be described in terms of a feature map $\Phi : \mathcal{X} \rightarrow \mathbb{R}^D$, such that

$$\kappa(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D \Phi_d(\mathbf{x}) \Phi_d(\mathbf{x}'). \quad (1.10)$$

Moreover, the functions in a RKHS with kernel κ are formulated based on a set of parameters $\boldsymbol{\varpi} = [\varpi_1, \dots, \varpi_D]^\top \in \mathbb{R}^D$ as follows

$$f(\mathbf{x}) = \sum_{j=1}^D \varpi_j \Phi_j(\mathbf{x}) = \langle \boldsymbol{\varpi}, \boldsymbol{\Phi}(\mathbf{x}) \rangle, \quad \text{with } \|f_{\boldsymbol{\varpi}}\|_{\mathcal{H}_\kappa} = \|\boldsymbol{\varpi}\|_2, \quad (1.11)$$

where $\langle \cdot, \cdot \rangle$ and $\|\cdot\|_2$ are respectively the the Euclidean inner product and norm. Besides, $\boldsymbol{\Phi}(\mathbf{x}) = [\Phi_1(\mathbf{x}), \dots, \Phi_D(\mathbf{x})]^\top \in \mathbb{R}^D$. According to the above, the assumption $f \sim \mathcal{GP}(0, \kappa(\mathbf{x}, \mathbf{x}'))$ becomes

$$\boldsymbol{\varpi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D) \propto \exp\left(-\frac{1}{2} \|\boldsymbol{\varpi}\|_2^2\right),$$

where $\mathbf{I}_D \in \mathbb{R}^{D \times D}$ is an identity matrix. If we assume a Gaussian likelihood, we have

$$p(\mathbf{y}|\mathbf{f}, \mathbf{X}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(y_n|f(\mathbf{x}_n), \sigma^2) \propto \exp\left(-\frac{1}{2\sigma^2} \sum_{n=1}^N (\langle \boldsymbol{\varpi}, \boldsymbol{\Phi}(\mathbf{x}_n) \rangle - y_n)^2\right).$$

Then, the posterior distribution is proportional to

$$\exp\left(-\frac{1}{2\sigma^2} \sum_{n=1}^N (\langle \boldsymbol{\varpi}, \boldsymbol{\Phi}(\mathbf{x}_n) \rangle - y_n)^2 - \frac{1}{2} \|\boldsymbol{\varpi}\|_2^2\right).$$

We see that a maximum a posteriori–(MAP) estimation for the posterior will become in a minimization problem with Tikhonov regularization (Tikhonov and Arsenin, 1977), where the regularization parameter is related to the noise variance.

1.3.2 Learning from multiple annotators

A supervised learning scenario involves the estimation of a function $f : \mathcal{X} \rightarrow \mathcal{Y}$. Usually, each \mathbf{x}_n is assigned to a single y_n , i.e., the ground truth. Still, in several real-world problems instead of the ground truth, we have multiple labels provided by $R \in \mathbb{N}$ annotators with different levels of expertise (Raykar et al., 2010). It is common to find that each annotator r only labels $|N_r| \leq N$ samples, being $|N_r|$ the cardinality of the set $N_r \subseteq \{1, \dots, N\}$ that contains the indexes of samples labeled by the r -th annotator. Besides, we define the set $R_n \subseteq \{1, \dots, R\}$ holding the indexes of annotators that labeled the n -th instance. Thereby, it is possible to build a data set for the annotator $r \in \{1, 2, \dots, R\}$, $\mathcal{D}_r = \{\mathbf{X}_r, \mathbf{y}_r\}$, where $\mathbf{X}_r = \{\mathbf{x}_n\}_{n \in N_r} \in \mathbb{R}^{|N_r| \times P}$ and $\mathbf{y}_r = \{y_n^{(r)}\}_{n \in N_r} \in \mathcal{Y}_r$ (commonly $\mathcal{Y}_r = \mathcal{Y}, \forall r$) are the input feature matrix and the labels given by the r -th annotator, respectively. Besides, \mathbf{y}_r is composed by elements $y_n^{(r)}$, which is the r -th annotation of sample \mathbf{x}_n .

Now given the data set from multiple annotators $\mathcal{D} = \{\mathbf{X}, \mathbf{Y} = \{y_n^{(r)}\}_{n \in N_r, r \in R_n}\}$ the aims of a multi-labeler approach are: First, to estimate the unknown gold standard for the instances in the training set $\mathbf{y} = [y_1, \dots, y_N]$. Second, to annotators behavior. Finally, the third objective is to build a supervised learning model, which generalizes well on unseen data (Yan et al., 2014). A graphical comparison between a typical supervised learning dataset and a multi-labeler database is shown in Figure 1.2.

1.4 Literature review on supervised learning from multiple annotators

As we established previously, the area of *learning from crowds* faces supervised learning problems where the gold standard is not available. Among the developed works in this topic, we recognize two main groups: The approaches based on frequentist models (FM) and the schemes based on Bayesian models (BM).

Regarding FM, we recognize the model proposed in Tian and Zhu (2015). Here, the authors first perform a geometrical interpretation of the majority voting (MV) scheme, and they remark that the estimated true label corresponds to the annotation (which can

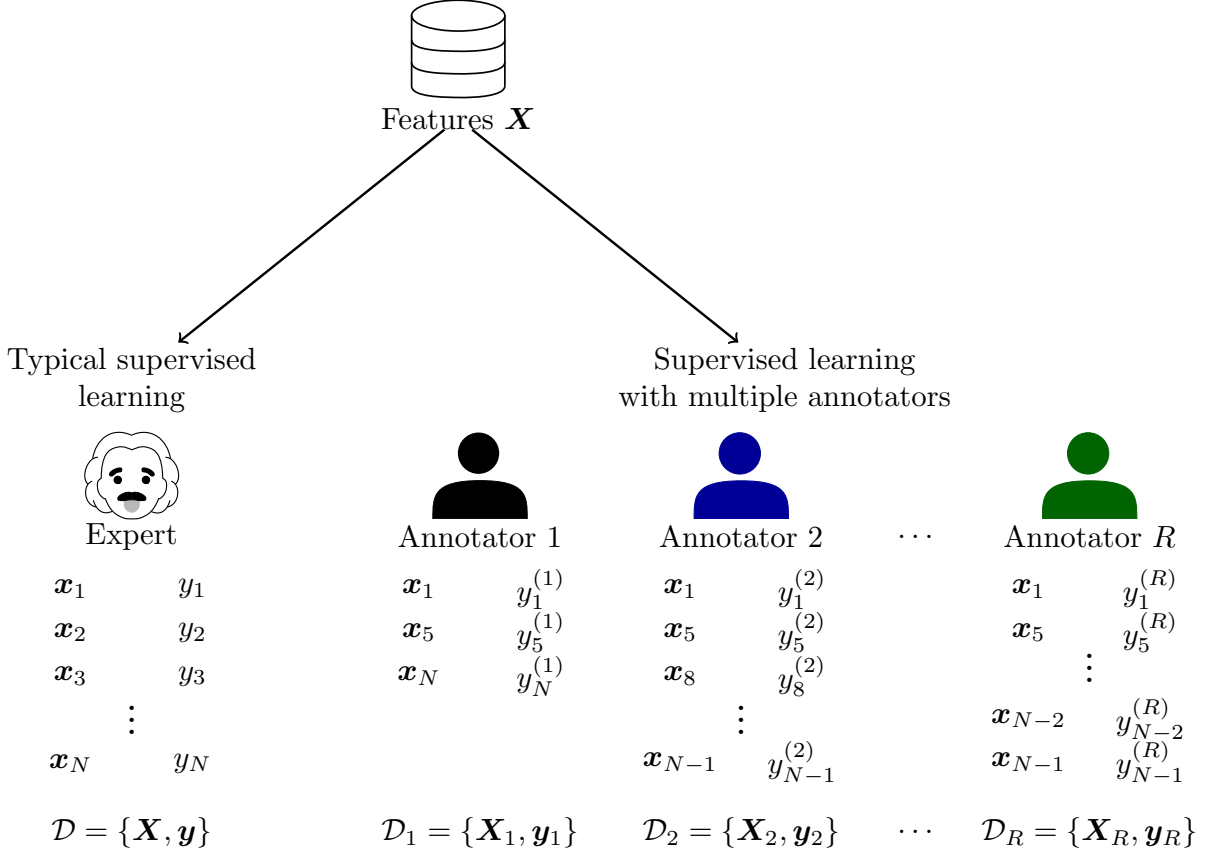


Figure 1.2: Graphical comparison between datasets for typical supervised learning settings and datasets from multiple annotators. We highlight that the annotators do not have access to the whole dataset; thus, it is common to face missing labels.

be seen as a point in a space generated by the MV approach) that has the largest distance to a decision hyperplane. Following this interpretation, they propose a weighted majority voting approach to deal with binary and multi-class classification, where the weights are related to the annotator reliability. Following the notion of margin in multi-class support vector machines (SVM), authors in [Tian and Zhu \(2015\)](#) define the *crowdsourcing margin*, which is the minimal difference between the aggregated score of the potential true label and the scores for other alternative labels. Accordingly, the annotators' reliability is estimated as the one that generates the largest margin between the potential true labels and other alternatives. However, this approach considers that the annotators' expertise is stationary across the input space; also, it does not take into account dependencies among the labelers. On the other hand, authors in [Imbajoa-Ruiz et al. \(2016\)](#) propose a mixture of classifiers aiming to deal with multi-labelers scenarios. The idea is to consider R classifiers $g_r(\mathbf{x})$, where each one of them is trained using the dataset

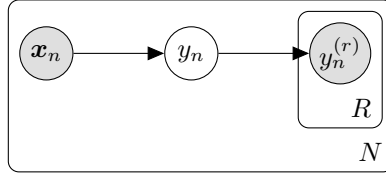


Figure 1.3: Graphical model for the approach proposed in Raykar et al. (2010). \mathbf{x}_n , y_n , and $y_n^{(r)}$ are respectively the input, the ground truth and the label given by the r -th annotator. Shaded nodes represent observed values, while unshaded nodes indicate latent variables.

for the r -th annotator $\mathcal{D}_r = \{\mathbf{X}_r, \mathbf{y}_r\}$, being \mathbf{X}_r and \mathbf{y}_r respectively the inputs and their corresponding outputs given by labeler r . Accordingly, the proposed classifier is given by $f(\mathbf{x}) = \sum_{r=1}^R \nu_r g_r(\mathbf{x})$, where $\boldsymbol{\nu} = [\nu_1, \dots, \nu_R] \in \mathbb{R}^R$ are the weighting factors. Each parameter $\nu_r \in \mathbb{R}$ is estimated using a version of the variable ranking approach over a kernel matrix, which is computed based on the labels from the r -th annotator and the input features. Still, this approach lacks interpretability since it is unclear the meaning of the parameters $\boldsymbol{\nu}$. On the other hand, authors in Rodrigues and Pereira (2018a) propose a Deep neural network-(DNN)-based approach to deal with multi-labeler data in both classification and regression tasks. The basic idea is to introduce an extra layer termed Crowdlayer, which allows training a DNN directly from the noisy labels from multiple annotators using backpropagation. The Crowdlayer is fed by what one usually define as the output layer of a DNN (e.g. softmax/sigmoid for classification, or linear for regression), and for each annotator r it learns a mapping from the output layer to the labels given by such labeler; hence, the Crowdlayer codes the annotators' reliabilities and biases.

So far, we have shown FM based on deterministic formulations such as SVM or DNN. Now, we focus on FM based on a probabilistic formulation. The most representative work is the proposed in Raykar et al. (2010), which corresponds to an extension of the early work in Dawid and Skene (1979) aiming to jointly learn a logistic regression-based classifier and the performance of the annotators in terms of sensitivity and specificity. Accordingly, the likelihood function (marginalizing with respect to \mathbf{y}) is given by

$$p(\mathbf{Y}|\mathbf{X}, \mathbf{a}, \mathbf{b}) = \prod_{n=1}^N \prod_{r=1}^R p_n p(y_n^{(r)} | y_n = 1, a_r) + (1 - p_n) p(y_n^{(r)} | y_n = 0, 1 - b_r), \quad (1.12)$$

where p_n is set as a logistic regression model with parameters $\boldsymbol{\omega} \in \mathbb{R}^P$

$$p_n = p(y_n = 1 | \mathbf{x}_n, \boldsymbol{\omega}) = \frac{1}{1 + \exp(-\boldsymbol{\omega}^\top \mathbf{x}_n)}. \quad (1.13)$$

In addition, $p(y_n^{(r)} | y_n = 1, a_r)$ and $p(y_n^{(r)} | y_n = 0, 1 - b_r)$ are modeled as Bernoulli distributions with parameter $a_r \in [0, 1]$ and $b_r \in [0, 1]$ respectively, indicating the annotator's sensitivity and specificity. Figure 1.3 shows the graphical model for the approach in Raykar et al. (2010). The model parameters $\mathbf{a} = [a_1, \dots, a_R]^\top \in \mathbb{R}^R$, $\mathbf{b} = [b_1, \dots, b_R] \in \mathbb{R}^R$, and $\boldsymbol{\omega}$ are estimated by maximizing the likelihood function in equation (1.12) and using the Expectation-Maximization-(EM) algorithm.

The model exposed below has a significant relevance since they have inspired many works in the area of *learning from crowds*, including binary classification (Morales-Álvarez et al., 2019; Rodrigues et al., 2014b; Ruiz et al., 2019; Sung et al., 2017), multiclass classification (Morales-Alvarez et al., 2020; Rodrigues et al., 2013), regression (Groot et al., 2011; Rodrigues et al., 2017), and sequence labeling (Rodrigues et al., 2014a).

According to the above, we notice some drawbacks. Firstly, such type of methods does not take into account the uncertainty in the predictions. Moreover, these approaches do not consider dependencies among the labelers nor that the relationship between the annotators' expertise and the input space, which does not fit real-world scenarios as we established in the problem statement.

On the other hand, few works have been focused on BM. In contrast to FM-based approaches, BM allows modeling the uncertainty in the labels, making them robust in the presence of noise. We recognize three strategies that use BM to face supervised learning settings with multiple annotators based on the graphical model in Figure 1.3. The first method is the proposed in Ruiz et al. (2019), termed Variational Gaussian processes for crowdsourcing-(VGPCR). The VGPCR's formulation is based on the work given in Rodrigues et al. (2014b). Both use a GPs-based framework to solve binary classification problems and uses a sensitivity-specificity-based model to measure the annotators' behavior. Nonetheless, unlike Rodrigues et al. (2014b), VGPCR treats the annotators' parameters (sensitivity and specificity) as random variables instead of fixed points. Also, it uses variational inference (instead of Expectation Propagation-(EP)) aiming to estimates the posterior distributions of the model random variables (Ruiz et al., 2019). In turn, Morales-Álvarez et al. (2019) proposed a scalable version of

VGPCR, named Scalable variational Gaussian processes for crowdsourcing–(SVGPCR). Authors in [Morales-Álvarez et al. \(2019\)](#) argue that a classical GP has a computational complexity of $\mathcal{O}(N^3)$ caused by the inversion of kernel matrix \mathbf{K}_{ff} , which is prohibited for large datasets. In contrast, they propose a GP sparse approximation via the so-called Variational Fourier features–(VFF) within a variational inference framework to learn a new kernel matrix to approximate \mathbf{K}_{ff} and to estimate the model’s random variables (the same as in VGPCR). In such way, the computational complexity is reduced to $\mathcal{O}(ND_f^2 + ND_fP)$, where $D_f \in \mathbb{Z}^+ \ll N$. Then, SVGPCR is extended by authors in [Morales-Alvarez et al. \(2020\)](#) to deal with multiclass classification problems in the context of multiple annotators. In this model, the behavior of each annotator is assessed in terms of her/his confusion matrix, which is assigned to a Dirichlet prior to perform fully Bayesian inference. Finally, regarding regression settings, authors in [Groot et al. \(2011\)](#) propose a GP-based model, where the labels are assumed to be a corrupted version of the hidden true labels by Gaussian noise; thus, $y_n^{(r)} = y_n + \mathcal{N}(0, \sigma_r)$, where $\sigma_r \in \mathbb{R}^+$ is the r -th annotator error-variance.

Similarly to FM, most BM are based on two assumptions: the outputs of the labelers do not depend on the input features and independence among the annotators. These assumptions have been widely discussed in the problem statement, where we have established that they do not fit real-world scenarios. Next, we describe the FM and BM that try to relax these assumptions.

Approaches to code the relationship between the annotators’ performance and the input space

Regarding the approaches that model the annotators’ behavior as a function of the input samples, we recognize that they are based on FM with probabilistic formulation; in fact, they are based on the graphical model shown in Figure 1.4. The first work is the proposed in [Yan et al. \(2010, 2014\)](#). They introduce an algorithm for binary classification, holding the joint conditional distribution as:

$$p(\mathbf{Y}, \mathbf{y} | \mathbf{X}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n) \prod_{r=1}^R p(y_n^{(r)} | \mathbf{x}_n, y_n), \quad (1.14)$$

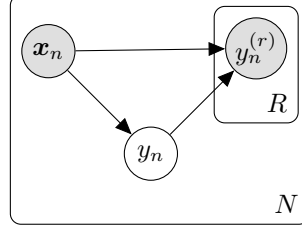


Figure 1.4: The graphical model used to code the labelers expertise as a function of the input space. \mathbf{x}_n , y_n , and $y_n^{(r)}$ are respectively the input, the ground truth and the label given by the r -th annotator. Shaded nodes represent observed values, while unshaded nodes indicate latent variables.

where $p(y_n|\mathbf{x}_n)$ is fixed as a logistic regression model (see equation (1.13)). Similarly, for the term $p(y_n^{(r)}|\mathbf{x}_n, y_n)$, they use two different models:

$$p(y_n^{(r)}|\mathbf{x}_n, y_n) = \mathcal{N}(y_n^{(r)}|y_n, \eta_r(\mathbf{x}_n)), \quad (1.15)$$

$$p(y_n^{(r)}|\mathbf{x}_n, y_n) = (1 - \eta_r(\mathbf{x}_n))^{|y_n^{(r)} - y_n|} \eta_r(\mathbf{x}_n)^{1 - |y_n^{(r)} - y_n|}, \quad (1.16)$$

where $\eta_r(\mathbf{x}_n)$ is fixed as a logistic regression with parameters $\boldsymbol{\gamma}_r \in \mathbb{R}^P$. Then, an EM algorithm is used to compute the model parameters.

On the other hand, the work in Zhang and Obradovic (2011) propose an extension of the work in Raykar et al. (2010) aiming to model the dependencies between the annotators and the input space. For doing so, they assume that the distribution of the input space \mathcal{X} can be approximated by using a Gaussian mixture model with L components. Hence, the author hypothesizes that each annotator r has a particular performance for each mixture component. Accordingly, they define $a_l^{(r)} \in [0, 1]$ and $b_l^{(r)} \in [0, 1]$ respectively as the sensitivity and specificity of the annotator r in the l -th mixture component, with $l \in \{1, \dots, L\}$. Therefore, the likelihood function for the incomplete data is given as

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{n=1}^N \prod_{l=1}^L \prod_{r=1}^R p_n p(y_n^{(r)}|y_n = 1, a_l^{(r)}) + (1 - p_n) p(y_n^{(r)}|y_n = 0, 1 - b_l^{(r)}), \quad (1.17)$$

where p_n is fixed to be a logistic regression model. Besides, $p(y_n^{(r)}|y_n = 1, a_l^{(r)})$ and $p(y_n^{(r)}|y_n = 0, 1 - b_l^{(r)})$ are modeled with Bernoulli distributions with parameters $a_l^{(r)}$ and $1 - b_l^{(r)}$ respectively. Similar to the previous approach, the model parameters are estimated by maximizing the likelihood and employing the EM algorithm.

So far, the works that we have discussed are intended to binary classification. Conversely, for regression settings, we only identify the work proposed by [Xiao et al. \(2013\)](#), which establishes the following generative model

$$y_n = f(\mathbf{x}_n) + \epsilon_f, \quad (1.18)$$

$$y_n^{(r)} = g_r(y_n, \mathbf{x}_n) + \epsilon_g, \quad (1.19)$$

where $\epsilon_f, \epsilon_g \in \mathbb{R}^+$ are modeled as Gaussian noise. Note that $f: \mathcal{X} \rightarrow \mathcal{Y}$ and $g_r: \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Y}_r$ characterize the regression function and the annotators model respectively. Here, the authors use a GPs-based approach for the regression function f ; thus:

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_{ff}), \quad (1.20)$$

where the covariance matrix $\mathbf{K} \in \mathbb{R}^{N \times N}$ is computed by using a kernel that involves a squared exponential term, a linear term, and a constant bias (see [Xiao et al. \(2013\)](#) for more details). Similarly, each g_r is also modeled as a GP; in this sense, the conditional distribution can be expressed as:

$$p(\mathbf{y}|\mathbf{X}) = \prod_{r=1}^R \mathcal{N}(\mathbf{y}_r|\mathbf{0}, \mathbf{S}_r), \quad (1.21)$$

where \mathcal{Y}_r is connected with \mathcal{X} and \mathcal{Y} via a covariance matrix $\mathbf{S}_r \in \mathbb{R}^{N \times N}$, where the element n, n' is computed using a particular kernel $\kappa_{\mathbf{S}}(\{y_n, \mathbf{x}_n\}, \{y_{n'}, \mathbf{x}_{n'}\})$ ([Xiao et al., 2013](#)). In this work, the latent functions f and $\{g_r\}_{r=1}^R$ are estimated by following the Maximum a posteriori-(MAP) solution.

Approaches to code Independence between annotators

One of the most common assumptions in learning from crowds is that the annotators make their decisions independently. However, as we remark in previous sections, that assumption is hardly plausible. Regarding this, we only recognize the work in [Zhu et al. \(2019\)](#) based on BM. Here, the authors introduce a covariance measure among the annotators; hence, the annotator model is expressed by

$$p(\mathbf{Y}|\mathbf{y}, \boldsymbol{\varphi}, \boldsymbol{\Sigma}) = \prod_{n=1}^N \mathcal{N}(\boldsymbol{\Upsilon}_n|y_n \mathbf{1}_R + \boldsymbol{\varphi}, \boldsymbol{\Sigma}), \quad (1.22)$$

where $\Sigma \in \mathbb{R}^{R \times R}$ is the covariance matrix of the R annotators, $\mathbf{1}_R \in \mathbb{R}^R$ is an all-ones vector, and $\varphi \in \mathbb{R}^R$ codes the bias of the annotators. Besides, $\Upsilon_n \in \mathbb{R}^R$ contains the labels provided for the n -th instance. Here, the model parameters are computed by using a Bayesian framework based on Gibbs sampling.

Figure 1.5 shows relevant state-of-the-art works. For summarizing, we arise the following observations:

- In general terms, we notice that most of the models for multiple annotators are based on FM to solve mainly binary classification tasks.
- We recognize only three relevant works (based on FM) to code the relationship between the input space and the annotators' behavior. Two of them are based on linear classifiers, two solve binary classification, and the remaining method employs a GP-based framework to solve a regression problem. Further, we observe that no one of these three approaches takes into account labelers' interdependencies.
- Concerning the codification of annotators' dependencies, we only identify a single algorithm based on BM to face regression settings. However, such a method does not code the labelers' performance as a function of the input space.
- Finally, we did not identify any model that codes the annotators' behavior under two assumptions: *i)* dependencies between the labelers; *ii)* the annotators' performance depends on the input space.

Accordingly, the principal aim of this thesis is to face the problems presented in Section 1.2 using both perspectives Frequentist and Bayesian approaches.

1.5 Objectives

1.5.1 General objective

To develop a supervised learning framework in the context of multiple annotators taking into account dependencies among the labelers and the fact that the annotators' performance is non-stationary across the input space, aiming to improve the representation of the labels in classification and regression tasks.

Annotators' Codification		
	<i>Frequentist models</i>	<i>Bayesian models</i>
	(1) + (2)	?
	?	?
Modeling interdependencies among the labelers (2)		- Gaussian distribution-Gibbs Sampling, Zhu et al. (2019)
Modeling the labelers' performance as function of the input space (1)	- Multiple Logistic Regression models. Yan et al. (2014) - Multiple GPs-Maximum a Posteriori, Xiao et al. (2013) .	
Using parameters to capture annotators' performance	- Deep Neural Networks. Rodriguez and Pereira (2018a) - Max-margin Majority Voting Tian and Zhu (2015)	- GPs Classification-Variational Inference. Morales-Álvarez et al. (2020) - GPs Regression-Exact Inference. Groot et al. (2011)

Figure 1.5: Relevant state-of-the-art works for supervised learning with multiple annotators.

1.5.2 Specific objectives

1. To develop a supervised learning approach (regression and classification) with multiple annotators, which uses a frequentist-based approach to code the labelers' expertise by considering dependencies among their decisions (labels).
2. To develop regression and classification approaches in multi-labelers scenarios, where the performance of the annotators is modeled using a frequentist-based strategy to code non-stationary and correlated labelers.
3. To develop a Bayesian-based method that jointly trains the supervised learning approach (regression or classification) and estimates the labelers' performance by taking into account dependencies among them and the relationship between their performance and the input space.

1.6 Outline and contributions

The main contributions are briefly introduced in the following sections, which are summarized in Figure 1.6.

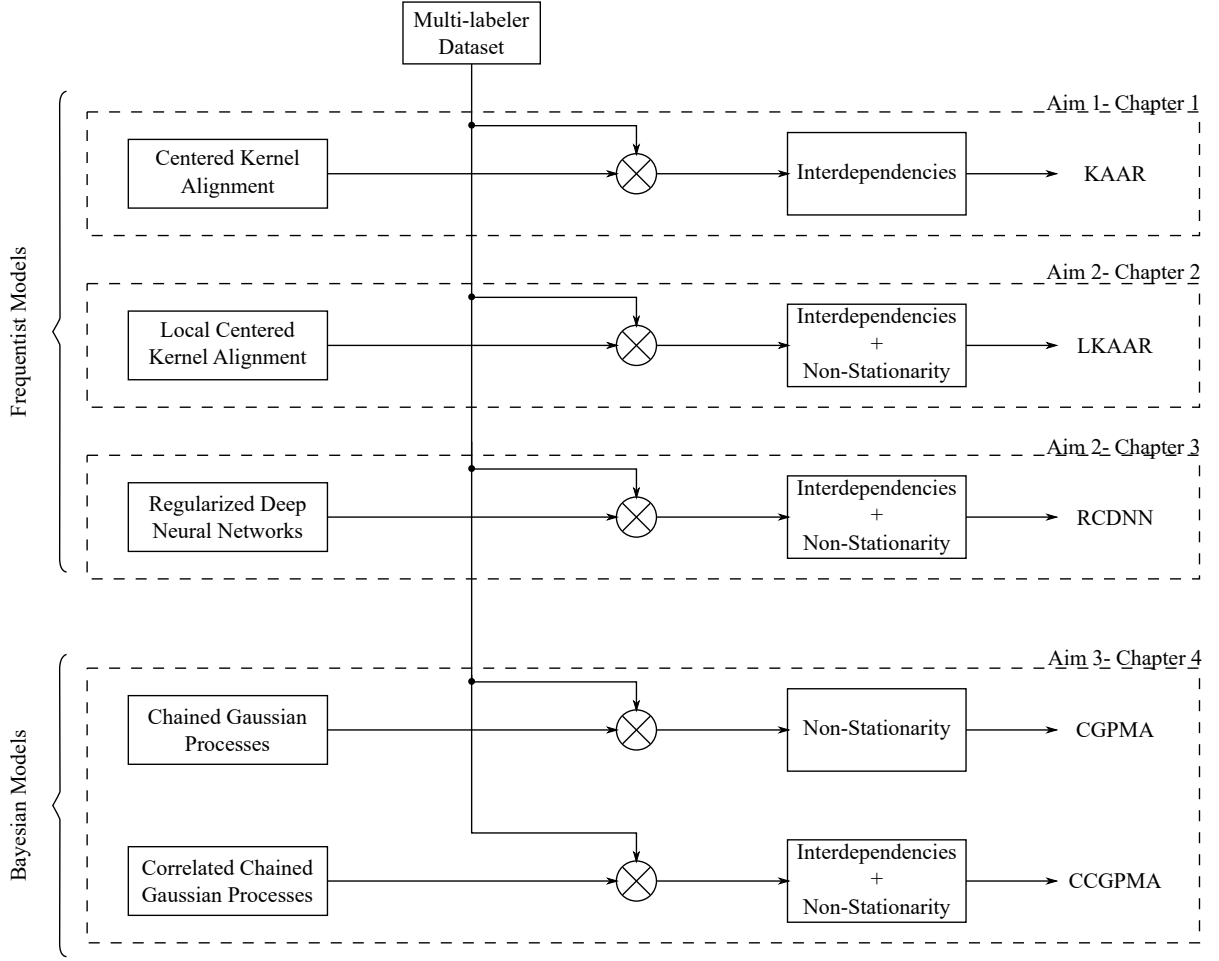


Figure 1.6: Relationship between the aims of this thesis and the developed methodologies. KAAR: Kernel alignment-based annotators relevance analysis. LKAAR: Localized kernel alignment-based annotator relevance analysis, RCDNN: Regularized chained deep neural network for Multiple Annotators, CGPMA: Chained Gaussian Processes for multiple annotators, and CCGPMA: Correlated chained Gaussian Processes for multiple annotators.

1.6.1 Kernel alignment-based annotator relevance analysis (KAAR)

The centered kernel alignment-(CKA) is commonly used for kernel selection in typical kernel-based learning models. The CKA approach comprises the computation of a convex combination of R kernels $\mathbf{K}_\nu = \sum_{r=1}^R \nu_r \mathbf{K}_r$, where each matrix \mathbf{K}_r is computed over the input features with a given kernel $\kappa_r(\mathbf{x}_n, \mathbf{x}_{n'})$; also, the weights $\nu = [\nu_1, \dots, \nu_R] \in [0, 1]^R$ are computed by quantifying the similarities between the combination of kernels and the target kernel \mathbf{F} , which in this case is a kernel computed over the labels (ground truth)

$\kappa_{\mathbf{y}}(y_n, y_{n'})$ (Cortes et al., 2012). Our first contribution comprises the use of CKA to measure the labelers' performance. For doing so, we define a convex combination of R kernels, one for annotator; $\mathbf{K}_{\boldsymbol{\nu}} = \sum_{r=1}^R \nu_r \mathbf{K}_r$, where \mathbf{K}_r is computed as $\kappa_r(y_n^{(r)}, y_{n'}^{(r)})$. Besides, we compute $\boldsymbol{\nu}$ as the average matching between the combination of kernels and the target kernel \mathbf{F} . We remark that for multi-labeler datasets, we do not have the actual labels; then, we choose the kernel over the features to be the target kernel $\kappa_X(\mathbf{x}_n, \mathbf{x}_{n'})$. We hypothesize that the input features code the main patterns of the unknown gold standard labels.

We highlight that in the case of multiple annotators, each parameter ν_r is proportional to the r -th annotator's performance. In turn, a new sample label is predicted as a convex combination of learners adopting the achieved KAAR-based coding $f(\cdot) = \sum_{r=1}^R \nu_r g_r(\cdot)$, where each $g_r(\cdot)$ is a supervised learning algorithm trained with the dataset for the r -th algorithm \mathcal{D}_r (see Figure 1.7). Our approach estimates the performance of the annotators using a non-parametric model, allowing it to be more flexible concerning the labels' distribution. Moreover, our methodology relaxes the assumption of independence between the annotators, which codes possible correlations in the annotators' opinions to model their expertise. This approach is related to the first specific aim, and it is described in Chapter 2.¹

1.6.2 Localized kernel alignment-based annotator relevance analysis (LKAAR)

As we exposed in Section 1.6.1, the first contribution, KAAR, uses the CKA approach to model the annotators' interdependencies. CKA assumes that the distribution of the input features is stationary; that is, the weight ν_r penalizes equally all the samples for the r -th kernel function. Thereby, KAAR does not consider the relationship between the input features and the labelers' behavior, which is not realistic as established in Section 1.2. To deal with this issue, we introduce a localized multiple kernel learning-based approach to compute $\mathbf{K}_q = \sum_{r=1}^R \mathbf{Q}_r \mathbf{K}_r \mathbf{Q}_r \in \mathbb{R}^{N \times N}$ (Gönen and Alpaydin, 2008; Zhang and Gao, 2018), where $\mathbf{Q}_r \in \mathbb{R}^{N \times N}$ is a diagonal matrix whose elements are defined by the vector $\mathbf{q}_r = [q_r(\mathbf{x}_1), \dots, q_r(\mathbf{x}_N)]^\top \in \mathbb{R}^N$. The combination factors $\mathbf{q} = [\mathbf{q}_1, \dots, \mathbf{q}_R]^\top \in \mathbb{R}^{N \times R}$ are estimated in such a way as to maximize the CKA between the kernel matrices \mathbf{K}_q and $\mathbf{F} \in \mathbb{R}^{N \times N}$, where we recall that for multiple annotators settings, \mathbf{F} holds elements

¹A MATLAB implementation of KAAR is available in <https://github.com/juliangilg/KAAR-Learning-from-multiple-annotators-using-Kernel-Alignment>

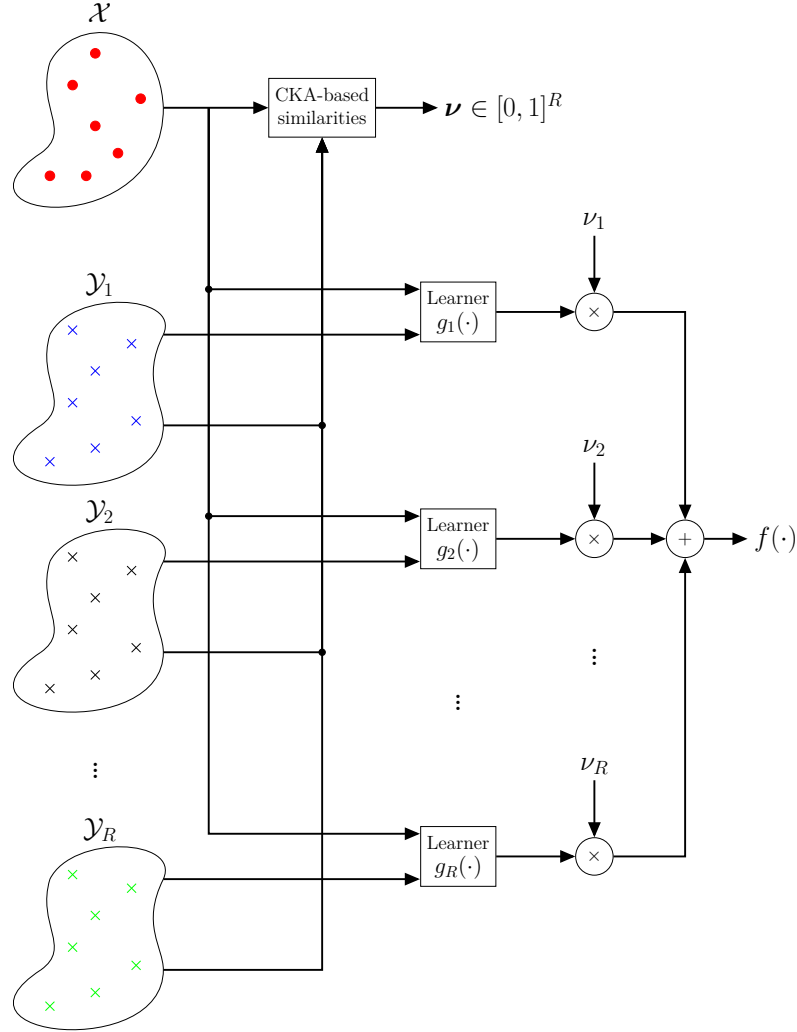


Figure 1.7: Kernel alignment-based annotator relevance analysis-(KAAR) pipeline. CKA stands for Centered Kernel Alignment. \mathcal{X} indicates the input space, and each \mathcal{Y}_r represents the output space for the r -th annotator. Moreover, $\boldsymbol{\nu} = [\nu_1, \dots, \nu_R]^\top$ is a vector containing the annotators' relevance parameters. Finally, the supervised learning mapping function f is computed as a convex combination of functions g_r that are trained from the r -th annotator's data.

$\kappa_X(\mathbf{x}_n, \mathbf{x}_{n'})$.

We remark that such combination factors \mathbf{q}_r estimates the annotators' performance considering it as a function of the input features and taking into account the interdependencies among the labelers. Besides, similar to KAAR, LKAAR is built as a convex combination of classifiers as shown in Figure 1.8; however, instead of having $\boldsymbol{\nu}$ as the

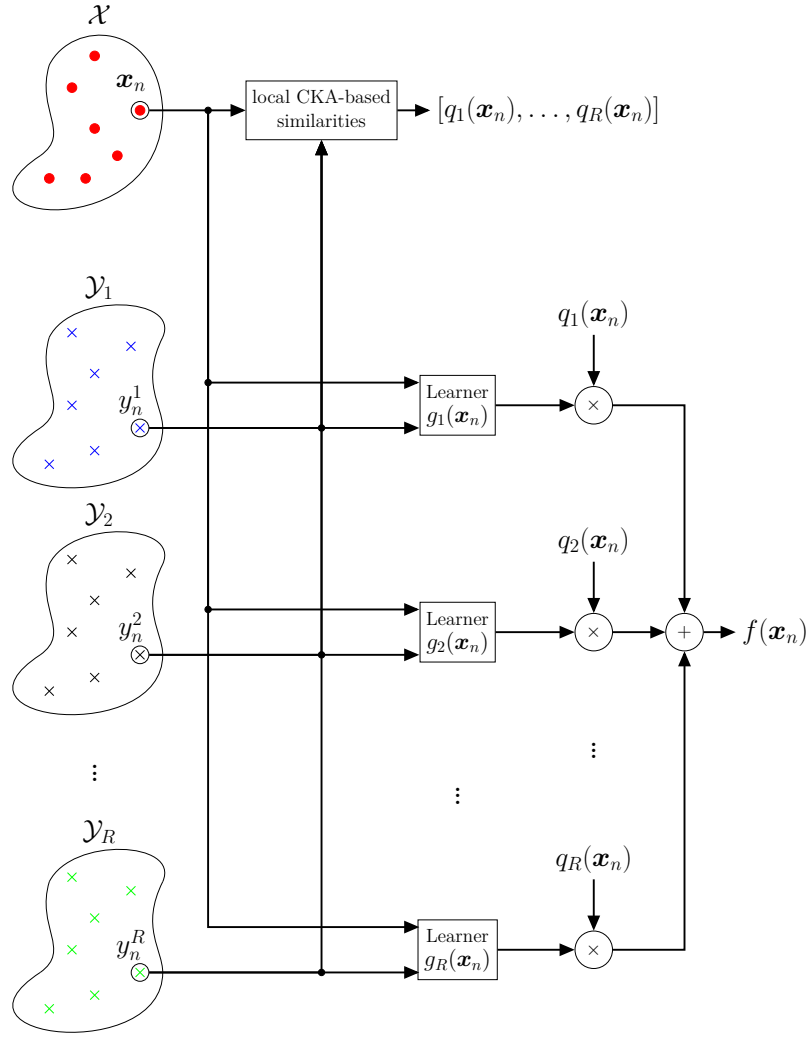


Figure 1.8: Localized kernel alignment-based annotator relevance analysis (KAAR) pipeline. CKA stands for Centered Kernel Alignment. \mathcal{X} indicates the input space, and each \mathcal{Y}_r represents the output space for the r -th annotator. Moreover, $[q_1(\mathbf{x}_n), \dots, q_R(\mathbf{x}_n)]^\top$ is a vector containing the annotators' relevance parameters for the n -th input sample. Finally, the supervised learning mapping function f is computed as a convex combination of functions g_r that are trained from the r -th annotator's data.

combination coefficients, which are constant across the input space, we use \mathbf{q} . LKAAR has three remarkable features: i) the performance of each annotator is a function of the input space; ii) the assumption of independence among the annotators is relaxed by modeling inter-annotators dependencies (Cohn and Specia, 2013); and iii) the performance of the annotators is estimated using a non-parametric model, allowing it to be more flexible to the distribution of the labels. This approach is related to the second

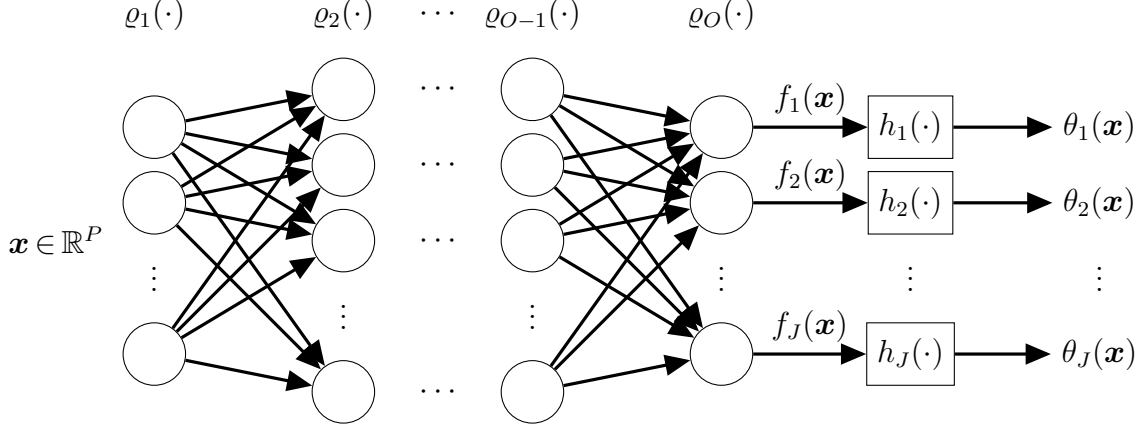


Figure 1.9: Regularized chained Deep Neural Network scheme. \mathbf{x} is an input vector, $\varrho_o(\cdot)$ is the o -th layer of a DNN, $f_j(\mathbf{x})$ is the j -th output neuron for the input \mathbf{x} . Moreover, each $\theta_j(\mathbf{x})$ is the j -th parameter of a given likelihood that is modeled as a function of the input space and J is the number of likelihood's parameters. Finally, $h_j(\cdot)$ is a deterministic function such that $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$.

specific aim, and it is described in Chapter 3.²

1.6.3 Regularized chained deep neural network for multiple annotators (RCDNN)

Our previous contributions KAAR and LKAAR are deal with a multi-labeler problem in two stages; the first comprises the estimation of the annotator's parameters, which are then used as weights factors for the combination of R learning models (second stage). In contrast, our third contribution is a Regularized deep neural network-based method that jointly estimates the annotators' performance and the supervised learning algorithm. This approach is inspired in the chained Gaussian Processes model-(CGP) (Saul et al., 2016), where the idea is to model each parameter $\theta_j(\mathbf{x})$, $j \in \{1, \dots, J\}$ in a given likelihood with multiple independent GPs priors (one GP prior per parameter). Unlike CGP, we consider that each neuron in the output layer of a DNN is linked to one of the parameters of a given likelihood through a deterministic function $h_j(\cdot)$. Thus, in a multi-labeler scenario, the annotators' parameters are modeled as a function of the input features. Moreover, we note that since each output in a DNN is computed as a linear combination of the outputs of a previous layer, our RCDNN codes interdependencies among the annotators. Besides, l1, l2, and Monte-Carlo Dropout-based regularizers are

²A MATLAB implementation of LKAAR is available in <https://github.com/juliangilg/LKAAR>

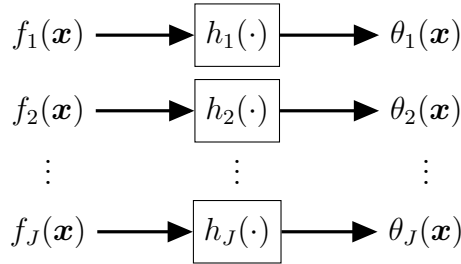


Figure 1.10: Chained Gaussian processes. Each $f_j(\mathbf{x})$ is a latent function that follows a GP. Moreover, each $\theta_j(\mathbf{x})$ is the j -th parameter of a given likelihood and J is the number of likelihood’s parameters. Finally, $h_j(\cdot)$ is a deterministic function such that $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$.

coupled within our method to deal with the overfitting issue in deep learning models. Like LKAAR, this approach is related to the second specific aim. RCDNN is described in Chapter 4 and summarized in Figure 1.9.³

1.6.4 Chained Gaussian processes for multiple annotators (CGPMA) and correlated chained Gaussian processes for multiple annotators (CCGPMA)

Up to this point, our contributions, KAAR, LKAAR, and RCDNN, solve the multi-labeler problem from a frequentist perspective. Conversely, the last two contributions are based on Bayesian approaches, specifically on GPs. First, we apply the CGP (Saul et al., 2016) model to a multi-labeler likelihood. As exposed in Section 1.6.3, CGP links each likelihood parameter $\theta_j(\mathbf{x})$ to a GP prior $f_j(\mathbf{x}) \sim \mathcal{GP}(0, \kappa_{f_j}(\cdot, \cdot))$, being $\kappa_{f_j}(\cdot, \cdot)$ the covariance function and $j \in \{1, \dots, J\}$. Such a connection between the parameter $\theta_j(\mathbf{x})$ and the $f(\mathbf{x})$ is performed via a deterministic function $h(\cdot) : \mathbb{R} \rightarrow \mathcal{M}_j$, where \mathcal{M}_j is the domain for $\theta(\mathbf{x})$ (see Figure 1.10). Accordingly, in a multi-labeler scenario, we are modeling the relationship between the annotators’ performance and the input features.

Unlike CGP, we consider that multiple correlated GPs model the likelihood’s parameters. For doing so, we take as a basis the ideas from a Multi-output GP–(MOGP) regression (Álvarez et al., 2012), where each output is coded as a weighted sum of shared latent functions via a semi-parametric latent factor model–(SLFM) (Teh et al., 2005). In contrast to the MOGP, we do not have multiple outputs but multiple functions chained

³A Python implementation of RCDNN is available in <https://github.com/juliangilg/RCDNN-MA>

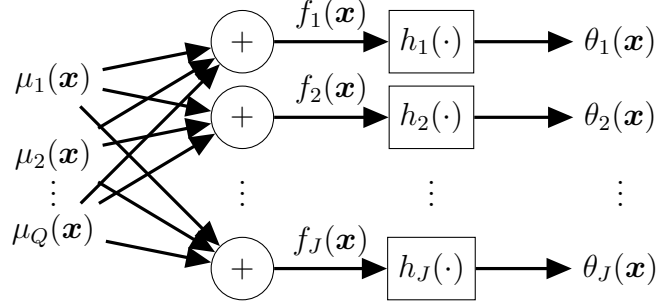


Figure 1.11: Correlated chained Gaussian processes. Each $f_j(\mathbf{x})$ is a latent function computed as a linear combination of functions $\{\mu_q(\mathbf{x})\}_{q=1}^Q$ that follow a GP. Moreover, each $\theta_j(\mathbf{x})$ is the j -th parameter of a given likelihood and J is the number of likelihood's parameters. Finally, $h_j(\cdot)$ is a deterministic function such that $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$.

to the given likelihood parameters. Thus, each latent function $f_j(\mathbf{x})$ is computed as

$$f_j(\mathbf{x}) = \sum_{q=1}^Q w_{j,q} \mu_q(\mathbf{x}), \quad (1.23)$$

where $\mu_q(\cdot) \sim \mathcal{GP}(0, k_q(\cdot, \cdot))$, with $k_q: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a kernel function, and $w_{j,q} \in \mathbb{R}$ is a combination coefficient ($Q \in \mathbb{N}$). Here, each LF is chained to the likelihood's parameters as shown in Figure 1.11. From the multiple annotators' point of view, the likelihood parameters are related to the labelers' behavior; thereby, CCGPMA models the labelers' behavior as a function of the input features while also taking into account annotators' interdependencies. CGPMA and CCGPMA are related to the third specific aim. Both models are described in Chapter 5.⁴

1.6.5 Thesis structure

This thesis is organized as follows. In Chapter 2, we introduce the KAAR model to code the annotators' performance taking into account dependencies among their labels, where the supervised learning algorithm is based on a linear combination of R learners. Chapter 3 discusses the LKAAR model, an extension of KAAR to model the relationship between the input features and the labelers' performance. In Chapter 4, we introduce a Regularized DNN-based approach to jointly training a supervised learning model and model the annotators' behavior considering both the dependencies among them and the relationship between their performance and the input space. Chapter 5 presents

⁴A Python implementation of CCGPMA is available in <https://github.com/juliangilg/CCGPMA>

the CGPMA and CCGPMA models, which addresses the multi-labeler problem from a Bayesian perspective. CGPMA considers the performance as a function of the inputs; while, CCGPMA extends CGPMA to capture dependencies among the labelers. Finally, in Chapter 6, we describe the conclusions, future works, and the academic discussion.

Chapter 2

Kernel alignment-based annotator relevance analysis

Most works in the context of learning from multiple annotators are based on parametric models, which lead to the following issues: i) if the quantity of the parameters is comparable with the number of samples available for training, the model can fall quickly in overfitting (Yan et al., 2014). ii) Linear or Gaussian constraints are commonly imposed to compute the optimal solutions analytically (Principe, 2010). However, real-world datasets cannot fulfill such assumptions. iii) Most of the works assume independence between the annotators; though it is suitable to consider that the labelers make their decisions independently, it is not true that these opinions are independent since there are possible correlations between the expert views (Esmaily et al., 2016).

In this chapter, we introduce a new kernel alignment-based annotator relevance analysis–(KAAR) approach to estimate the expertise of the labelers in scenarios where the gold standard is not available. KAAR computes each annotator’s relevance as an averaged matching between the input features and the expert labels. In turn, a new sample label is predicted as a convex combination of learners adopting the achieved KAAR-based coding. Unlike previous works, our approach estimates the annotators’ performance using a non-parametric model, allowing it to be more flexible concerning the distribution of the labels. Moreover, our methodology relaxes the assumption of independence between the annotators, highlighting possible correlations between the opinions to code their expertise. Our approach is similar to the proposed in Hua et al. (2018) in the sense that we perform the supervised learning task using a weighted combination of supervised learning algorithms. However, unlike such a work, we estimate the weights using a kernel alignment-based approach to quantify the matching between the

input features and the annotator expertise. We empirically show, using both simulated and real annotators for regression and classification tasks, that our methodology can be used to estimate the performance of multiple labelers even if the gold standard is not available, outperforming state-of-the-art techniques.

2.1 Centered kernel alignment fundamentals

In typical kernel-based learning models, the selection of a positive-definite function plays a crucial role. Multiple kernel learning–(MKL) deals with such an issue by defining a kernel \mathbf{K}_ν as a combination of R basis kernels; thus, we have (Gönen and Alpaydm, 2011)

$$\mathbf{K}_\nu = h_\nu \left(\{ \mathbf{K}_r \}_{r=1}^R \mid \nu \right), \quad (2.1)$$

where h_ν is the combination function, which can be linear (sum of kernels) or non-linear (e.g., product of kernels) (Varma and Babu, 2009). Moreover, $\mathbf{K}_r \in \mathbb{R}^{N \times N}$ is computed using a particular kernel function over the input features \mathbf{X} , $\kappa_r: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, and $\nu \in \mathbb{R}^R$ is a vector containing the combination parameters. In particular, we define the function h_ν as a convex combination; accordingly,

$$\mathbf{K}_\nu = h_\nu \left(\{ \mathbf{K}_r \}_{r=1}^R \mid \nu \right) = \sum_{r=1}^R \nu_r \mathbf{K}_r, \quad (2.2)$$

where $\nu \in [0, 1]^R$, and $\|\nu\|_2 = 1$ guarantees that \mathbf{K}_ν is positive definite ($\nu \in \mathbb{R}^R$ holds elements ν_r and $\|\cdot\|_2$ is the l2-norm). It is worth mentioning that adopting a convex sum favors the interpretability of parameters $\{\nu_r\}_{r=1}^R$, where the weight $\nu_r \in [0, 1]$ represents the r -th kernel relevance (Gönen and Alpaydm, 2011).

Centered Kernel Alignment–(CKA)-based approaches leverage a data-driven estimator of ν by quantifying the similarity among the kernels over input features and the kernel computed over the outputs $\mathbf{F} \in \mathbb{R}^{N \times N}$, which is also set as a positive definite function $\kappa_y: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. Thereby, the matching between \mathcal{Y} and \mathcal{X} spaces can be measured through an empirical estimate of the CKA value $\rho \in [0, 1]$ based on the \mathbf{K}_ν and \mathbf{F} kernels as follows (Cortes et al., 2012):

$$\rho(\mathbf{K}_\nu, \mathbf{F}) = \frac{\langle \bar{\mathbf{K}}_\nu, \bar{\mathbf{F}} \rangle_{\mathbf{F}}}{\|\bar{\mathbf{K}}_\nu\|_{\mathbf{F}} \|\bar{\mathbf{F}}\|_{\mathbf{F}}}, \quad (2.3)$$

where $\bar{\mathbf{F}}$ stands for a centered kernel matrix computed as: $\bar{\mathbf{F}} = \tilde{\mathbf{I}}\mathbf{F}\tilde{\mathbf{I}}$, being $\tilde{\mathbf{I}} = \mathbf{I} - \mathbf{1}\mathbf{1}^\top/N$ a centering matrix, $\mathbf{I} \in \mathbb{R}^{N \times N}$ is the identity, and $\mathbf{1} \in \mathbb{R}^N$ is an all-ones vector. Namely, the vector $\boldsymbol{\nu}$ can be inferred by minimizing equation (2.3); then, the higher the ν_r value the better the contribution of \mathbf{K}_r to match the target \mathbf{F} .

2.2 Kernel alignment-based annotator relevance analysis

As we pointed out previously, CKA quantifies the similarity between a pair of input-output spaces to perform a proper configuration of a convex set of kernel functions. Alike, our idea is to use CKA to estimate the annotators' expertise in scenarios where the gold standard is not available. For doing so, we assume that the input features code the main patterns of the unknown gold standard labels. Thus, we measure the CKA between a kernel extracted from the input features (our target) and R kernels computed from the annotations given by each expert. Accordingly, we rewrite the matrix \mathbf{F} as follows: $f_{nn'} = \kappa_X(\mathbf{x}_n, \mathbf{x}_{n'})$, with $n, n' \in \{1, 2, \dots, N\}$. For associating samples, several bivariate measures of similarity can be used, e.g., linear, Gaussian, polynomial, etc (Liu et al., 2011). Here, to favor the mathematical tractability and to avoid the influence of the free parameters regarding the annotator performance, we fix κ_X as a linear kernel, yielding: $f_{nn'} = \mathbf{x}_n \mathbf{x}_{n'}^\top$. In turn, the pairwise similarity within the r -th expert annotations is defined by rewriting κ_r . The form κ_r will depend on the label space nature (i.e., regression or classification), which is discussed in Section 2.2.1. On the other hand, the kernel matrix \mathbf{K}_ν , which codes the multiple annotation dependencies, is computed as a convex combination of the R kernels \mathbf{K}_r , holding elements: $k_{nn'}^{(r)} = \kappa_r(y_n^{(r)}, y_{n'}^{(r)})$. Correspondingly, to capture each annotator performance, we compute the ν_r value using the following CKA-based optimization:

$$\hat{\boldsymbol{\nu}} = \arg \max_{\boldsymbol{\nu}} \rho(\mathbf{K}_\nu, \mathbf{F}); \quad \text{s.t. } \|\boldsymbol{\nu}\|_2 = 1. \quad (2.4)$$

Note that we highlight the dependence of the kernel matrix \mathbf{K}_ν with respect to the vector $\boldsymbol{\nu}$. As a result, the weight $\hat{\nu}_r \in \hat{\boldsymbol{\nu}}$ in equation (2.4) explains the measured matching between the r -th expert, as stated in $\kappa_r(y_n^{(r)}, y_{n'}^{(r)})$, and the input features, as coded in $\kappa_X(\mathbf{x}_n, \mathbf{x}_{n'})$. Moreover, through an auxiliary vector $\mathbf{v} \in \mathbb{R}^R$, wherein the following equality is imposed $\boldsymbol{\nu} = \mathbf{v}/\|\mathbf{v}\|_2$, the optimization problem in equation (2.4) can be

solved through the minimization of a quadratic cost function:

$$\hat{\mathbf{v}} = \arg \max_{\mathbf{v}} \mathcal{L}(\mathbf{v}) = \arg \max_{\mathbf{v}} \mathbf{v}^\top \mathbf{\Gamma} \mathbf{v} - 2\mathbf{v}^\top \mathbf{a}; \quad \text{s.t.} \quad \|\mathbf{v}\|_2 \geq 0, \quad (2.5)$$

where the matrix $\mathbf{\Gamma} \in \mathbb{R}^{R \times R}$ collects inter-annotator dependencies as: $\Gamma_{rr'} = \langle \bar{\mathbf{K}}_r, \bar{\mathbf{K}}_{r'} \rangle_{\mathbf{F}}$; $r, r' \in \{1, 2, \dots, R\}$, and $\mathbf{a} \in \mathbb{R}^R$ quantifies the similarity between the r -th expert and the input feature space as: $a_r = \langle \bar{\mathbf{K}}_r, \bar{\mathbf{F}} \rangle_{\mathbf{F}}$. Accordingly, KAAR allows measuring the annotators' performance by taking into account dependencies among their behavior. A gradient descent-based approach is given to solve the optimization problem in equation (2.5) as follows:

$$\frac{d\mathcal{L}(\mathbf{v})}{d\mathbf{v}} = 2\mathbf{\Gamma}\mathbf{v} - 2\mathbf{a}. \quad (2.6)$$

After estimation of the relevance vector $\boldsymbol{\nu}$, we assess the output $y_{new} \in \mathcal{Y}$ of a new input $\mathbf{x}_{new} \in \mathbb{R}^P$ as the following convex combination of learners

$$\hat{y}_{new} = \frac{1}{\|\boldsymbol{\nu}\|_1} \sum_{r=1}^R \nu_r g_r(\mathbf{x}_{new}), \quad (2.7)$$

where $g_r: \mathbb{R}^P \rightarrow \mathcal{Y}$ is a supervised learning algorithm trained from the set \mathcal{D}_r , and $\|\cdot\|_1$ is the l1-norm.

Notably, our kernel alignment-based annotator relevance analysis approach—(KAAR) can deal with missing labels. Also, it is possible to use any supervised learning scheme to learn the function g_r . Summarizing, our KAAR (see Algorithm 1.) counts on the enhancement of learning from multiple annotators by its two central stages: i) Seeking a relevance vector $\boldsymbol{\nu}$, relying on the averaged matching between the annotator labels and the input data features, and ii) Predicting the output of a new sample as a convex combination of learners adopting the achieved CKA-based relevance vector that intends to enhance the data separability based on the explained discrimination of each provided expert.

2.2.1 KAAR for classification and regression

One of the advantages of our proposal is that it can be applied for regression and classification. The key factor comprises the election of a proper function for $\kappa_r(y_n^{(r)}, y_{n'}^{(r)})$. Following, we describe the kernel used for each setting (regression/classification).

Algorithm 1: KAAR description

-
- Data:** $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\}$, \mathbf{x}_{new}
- 1 Compute $\bar{\mathbf{K}}$ following equation (2.2).
 - 2 Compute $\bar{\mathbf{F}}$ from the kernel function $\kappa_X(\cdot, \cdot)$.
 - 3 Compute the annotators' parameters $\boldsymbol{\nu}$ by solving the optimization problem in equation (2.5).
 - 4 Learn a function $g_r(\cdot)$ for each annotator by using $\mathcal{D}_r = \{\mathbf{X}_r, \mathbf{y}_r\}$.
 - 5 Predict the output \hat{y}_{new} for \mathbf{x}_{new} following equation (2.7).
-

Classification

In this first case, $y_n^{(r)} \in \{1, \dots, K\}$, being K the number of classes. In turn, the pairwise similarity within the r -th expert annotations is defined by rewriting κ_r as follows:

$$\kappa_r(y_n^{(r)}, y_{n'}^{(r)}) = \begin{cases} 1, & \text{if } n, n' \in N_r \text{ and } y_n^{(r)} = y_{n'}^{(r)} \\ 0, & \text{Otherwise} \end{cases}, \quad (2.8)$$

where the condition $n, n' \in N_r$, inhibits the influence of missing labels.

Regression

For regression problems, we have that the labels $y_n^{(r)} \in \mathbb{R}$. Thus, as in Cortes et al. (2012), we choose a linear kernel as follows

$$\kappa_r(y_n^{(r)}, y_{n'}^{(r)}) = \begin{cases} y_n^{(r)} y_{n'}^{(r)}, & \text{if } n, n' \in N_r \\ 0, & \text{Otherwise} \end{cases}, \quad (2.9)$$

where the condition $n, n' \in N_r$, inhibits the influence of missing labels.

2.3 Experimental set-up

In this section, we describe the experiments' configurations to validate our KAAR approach in multiple annotators scenarios for classification and regression tasks.

Table 2.1: Datasets description.

	Name	Number of features	Number of instances	Number of classes
<i>semi-synthetic</i>	Breast	9	683	2
	Bupa	6	345	2
	Ionosphere	34	351	2
	Pima	8	768	2
	Tic-tac-toe	9	958	2
	Iris	4	150	3
	Wine	13	178	3
	Segmentation	18	2310	7
<i>fully real</i>	Voice	13	218	2
	Polarity	1200	10427	2
	Music	124	1000	10

2.3.1 Classification

Testing datasets

Aiming to test our KAAR approach as a classification tool from multiple annotators scenarios, we use eight datasets for binary classification of the well-known *UCI repository*¹. The chosen datasets include: Wisconsin Breast Cancer Database–(breast), BUPA liver disorders–(bupa), Johns Hopkins University Ionosphere database–(ionosphere), Pima Indians Diabetes Database–(pima), Tic-Tac-Toe Endgame database–(tic-tac-toe), Iris Plants Database–(iris), Wine Data set–(Wine), and Image Segmentation Data Set–(Segmentation).

Moreover, the proposed approach is also tested on real multiple annotators settings by applying it in two datasets. The first is a voice database, where the idea is to build a system for evaluating the voice quality. The Massachusetts Eye and Ear Infirmary Disordered Voice Database from the Kay Elemetrics company is used, which comprises voice records from healthy and different voice issues. Explicitly, a subset of $N = 218$ voice records is deemed: 51 healthy and 167 pathological. Four specialists assessed the quality following the GRBAS protocol, which comprises the evaluation of five qualitative scales: Grade of dysphonia–(G), Roughness–(R), Breathiness–(B), Aesthenia–(A), and Strain–(S). For each perceptual scale, the specialist assigns an integer tag ranging from 0 (healthy voice) to 3 (severe disease) (Arias-Londoño et al., 2011). Next, the well-known Mel-frequency cepstral coefficients–(MFCC) are computed for each voice signal to obtain an input space of 13 features ($P = 13$) (Gil et al., 2015). The automatic assessment of voice quality configures five independent multi-class classification problems. However,

¹<http://archive.ics.uci.edu/ml>

since we have information about each voice record’s diagnosis (pathological or normal), we map the labels to the set $\{1, 2\}$. The second dataset is named sentiment polarity, which corresponds to a collection of more than ten thousand sentences, where each sentence is labeled as positive or negative. From this collection, $N = 5000$ sentences were selected randomly and published in the AMT platform to obtain labels from multiples annotators about each sentence’s sentiment polarity. Besides, the remaining 5248 sentences were kept for testing. Each phrase is pre-processed to remove the stop words and is represented using a vocabulary with a size of 8919. A post-processing step based on Latent Semantic Analysis was carried out to reduce the dimensionality to $P = 1200$ features (for more details, see [Rodrigues et al. \(2013\)](#)). Further, the music genre data is analyzed, holding a collection of songs records labeled from one to ten depending on their music genre: classical, country, disco, hip-hop, jazz, rock, blues, reggae, pop, and metal. From this set, 700 samples were published randomly in the AMT platform to obtain labels from multiples sources (2946 labels were obtained from 44 workers). The feature extraction is performed by following the work by authors in [Rodrigues et al. \(2013\)](#), to obtain an input space with $P = 124$. Table 2.1 summarizes the tested datasets.

Provided and simulated annotations

Regarding the UCI repository datasets, which are mainly used for typical supervised learning approaches without considering multiple annotators, two different simulation methods are studied to avoid biased results: *i) Classifier disrupting* proposed in [Rodrigues et al. \(2013\)](#), where a logistic regression classifier is trained using the input features and the true labels $\{\mathbf{x}_n, y_n\}_{n=1}^N$, to obtain a classifier weight vector $\mathbf{w} \in \mathbb{R}^P$. Then, the label given by the r -th annotator is simulated as follows: $y_n^{(r)} = \mathbf{w}^\top \mathbf{x}_n + \mathcal{N}(0, \sigma_r^2)$, where $\mathcal{N}(0, \sigma_r^2)$ stands for an univariate normal distribution with zero mean and variance $\sigma_r^2 \in \mathbb{R}^+$. In turn, a sigmoid function is applied to map the disrupted label within the set $\{1, 2\}$. Accordingly, the higher the σ_r^2 value the lower the expertise level of the r -th labeler. *ii) Biased coin* used in [Yan et al. \(2010\)](#) that builds a binary number $\tau_n^{(r)} \in \{0, 1\}$ from a Bernoulli distribution ruled by the probability parameter $p_r \in [0, 1]$. Then, the simulated annotations of the r -th expert yields: $y_n^{(r)} = y_n$, if $\tau_n^{(r)} = 0$, otherwise, $y_n^{(r)} = \tilde{y}_n$, if $\tau_n^{(r)} = 1$; where \tilde{y}_n is the flipped version of y_n .

Concerning the voice quality dataset, the annotations from four experts are provided, $R = 4$. For concrete testing, we only take into account the R and B scales due to for scales A and S, the performance of the labelers is not satisfactory. For the scale G, the labelers expertise is quite similar (according to the analysis performed in [Gil et al.](#)

Table 2.2: A brief overview of the state-of-the-art methods tested. GPC: Gaussian Processes classifier, LRC: logistic regression classifier, MV: majority voting, MA: multiple annotators, MAE: Modelling annotators expertise, LFC: Learning from crowds, DGRL: Distinguishing good from random labelers.

Algorithm	Description
GPC-GOLD	A GPC using the real labels (upper bound).
GPC-MV	A GPC using the majority voting of the labels as the ground truth.
MA-LFC (Raykar et al., 2010)	A LRC with constant parameters across the input space.
MA-MAE (Yan et al., 2010)	A LRC where the sources parameters depend on the input space.
MA-DGRL (Rodrigues et al., 2013)	A multi-labeler approach that considers as latent variables the annotator performance.
MA-GPC (Rodrigues et al., 2014b)	A multi-labeler GPC, which is as an extension of MA-LFC by using a non-linear approach.

(2015)). Similarly, for the polarity sentiment dataset, we have labels from 203 workers. However, we only consider the annotators who labeled at least the 15% of the available instances; in this sense, we use the information from $R=7$ labelers. It is important to highlight that these 7 annotators do not label all the available instances. Further, for the music dataset, we only consider the annotators who labeled at least the 20% of the instances; thus, we use the information from $R=9$ labelers.

KAAR training and comparison methods

Owing to KAAR employs a linear kernel to define the input data similarities and a closed form to compare expert annotations (see equation (2.8)), it does not require any free parameter tuning to compute the relevance vector $\boldsymbol{\nu}$ in equations (2.4) and (2.5). Now, as the classification function g_r , we use a Gaussian processes-based model, whose hyperparameters are set by optimizing a marginal likelihood (Rasmussen and Williams, 2006). Also, the validation is assessed by estimating the classification performance as the Area Under the Curve (AUC) for both simulated and real data. A cross-validation scheme is carried out with 30 repetitions where 70% of the samples are utilized for training and the remaining 30% for testing (except for the sentiment polarity dataset since it clearly defines the training and testing sets). On the other hand, Table 2.2 presents the comparison methods

2.3.2 Regression

Testing datasets

We test our approach using three types of datasets: fully synthetic data, semi-synthetic data, and fully real datasets. First, We generate *fully synthetic data* as an one-dimensional regression problem, where the ground truth for the n -th sample corresponds to $y_n = \sin(2\pi x_n)$, where the input matrix \mathbf{X} is formed by randomly sampling 300 points within the range $[0, 1]$ from an uniform distribution. The test instances are obtained by extracting 300 equally spaced samples from the interval $[0, 1]$. Second, to control the label generation (Ruiz et al., 2019), we build *semi-synthetic data* from six datasets related to regression tasks from the well-known UCI repository. We selected the following datasets: Auto MPG Data Set–(Auto), Bike Sharing Dataset Data Set–(Bike), Concrete Compressive Strength Data Set–(Concrete), The Boston Housing Dataset–(Housing),² Yacht Hydrodynamics Data Set–(Yacht), and Relative location of CT slices on axial axis Data Set–(CT). Third, we evaluate our proposal on one *fully real dataset*. In particular, we use the music genre data³, holding a collection of songs records labeled from one to ten depending on their music genre: classical, country, disco, hip-hop, jazz, rock, blues, reggae, pop, and metal. From this set, 700 samples were published randomly in the AMT platform to obtain labels from multiples sources (2946 annotations from 44 workers). Notice that the music dataset configures a 10-class classification problem; however, in such an experiment, we are addressing a multiclass classification problem with a regression model. Such practice is not uncommon in machine learning, and it is usually known as “Least-square classification” (Bishop, 2006; Rasmussen and Williams, 2006). Table 2.3 summarizes the tested datasets for the regression case.

Simulated and provided labels

As we pointed out previously, *fully synthetic* and *semi-synthetic* datasets do not hold real annotations. Thus, it is necessary to generate these labels synthetically as a version of the gold standard corrupted by Gaussian noise, i.e., $y_n^{(r)} = y_n + \epsilon_r$, where $\epsilon_r \sim \mathcal{N}(0, \sigma_r)$, being σ_r the r -th annotator error-variance for the sample n . This simulation method has been used in several works such as Groot et al. (2011); Rodrigues et al. (2017). Further,

²See <https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html> for housing

³<http://fprodrigues.com/publications/learning-from-multiple-annotators-distinguishing-good-from-random-labelers/>

Table 2.3: Datasets for regression.

	Name	Number of features	Number of instances
<i>fully synthetic</i>	synthetic	1	100
<i>semi-synthetic</i>	Auto	8	398
	Bike	13	17389
	Concrete	9	1030
	Housing	13	506
	Yacht	6	308
	CT	384	53500
<i>fully real</i>	Music	124	1000

Table 2.4: A brief overview of state-of-the-art methods tested for regression tasks. GPR: Gaussian Processes Regression, LR: logistic regression, Av: average, MA: multiple annotators, DL: Deep learning, LFCR: Learning from crowds for regression.

Algorithm	Description
GPR-GOLD	A GPR using the real labels (upper bound).
GPR-Av	A GPR using the average of the labels as the ground truth.
MA-LFCR (Raykar et al., 2010)	A LR model for MA where the labelers' parameters are supposed to be constant across the input space.
MA-GPR (Rodrigues et al., 2014b)	A multi-labeler GPR, which is as an extension of MA-LFCR.
MA-DL (Rodrigues and Pereira, 2018a)	A Crowd Layer for DL, where the annotators' parameters are constant across the input space.

for the music dataset, we only consider the annotators who labeled at least the 20% of the instances; thus, we use the information from $R = 9$ labelers.

KAAR training and method comparison

Alike in the classification case, KAAR for regression settings employs a linear kernel to define the input data similarities and a closed-form to compare expert annotations (see equation (2.8)), it does not require any free parameter tuning to compute the relevance vector $\boldsymbol{\nu}$ in equations (2.4) and (2.5). Besides, as the regression function g_r , we use a Gaussian processes-based model, whose hyperparameters are set by optimizing a marginal likelihood (Rasmussen and Williams, 2006). The quality assessment is carried out by estimating the regression performance as the coefficient of determination- (R^2) . A cross-validation scheme is employed with 15 repetitions where 70% of the samples are utilized for training and the remaining 30% for testing (except for *fully synthetic dataset* and music dataset, since they clearly define the training and testing sets). Table 2.4 displays the employed methods of the state-of-the-art for comparison purposes. From

Table 2.4, we highlight that for the model MA-DL, the authors provided three different annotators’ codification: MA-DL-B, where the bias for the annotators is measured; MA-DL-S, where the labelers’ scale is computed; and measured; MA-DL-B+S, which is a version with both (Rodrigues and Pereira, 2018a).

2.4 Results and discussion

2.4.1 Classification

First, a representative experiment is carried out to verify the KAAR capability to code the annotator performance. Namely, samples belonging to the Setosa and the Versicolour classes in the iris dataset are considered for testing the KAAR-based coding on a linearly separable problem. Afterward, we simulate five annotators using the simulation methods described in Section 2.3 by using the parameters shown in Table 2.5. Moreover, Table 2.5 shows the KAAR-based coding results for the iris dataset regarding the matching weights in $\boldsymbol{\nu}$. We compute such weights as the CKA-based dependency between each expert kernel and the target for the two simulation methods. In particular, we analyze two different target kernels: from true labels and the input features ($\boldsymbol{\nu}_{TL}$ and $\boldsymbol{\nu}_{IF}$). Comparing the AUC for the simulated annotators and the KAAR-based results, overall, our approach can infer the annotators’ performance from the input features. Figure 2.1 shows a visual comparison among the following kernels computed: the kernel calculated over the ground truth labels, the kernel estimated from the input features (\mathbf{F}), and the kernels estimated for each expert $\{\mathbf{K}_r : r = 1, \dots, R\}$. Remarkably, the true labels and the input features kernels exhibit a significant coincidence; then, an appropriate parametrization reveals relevant information about the unknown gold standard. Besides, the annotation quality plays a vital role regarding the visual similarities between the true labels and the simulated annotations kernels; that is, the lower the quality, the lower the visual similarity. The latter statement can be corroborated in Table 2.5, which shows the parameters used in each method to simulate annotators with different levels of expertise. The AUC is computed between the true labels and the simulated annotations. On the other hand, we recall that one of the main aims of this first proposal is to measure the annotators’ performance by considering dependencies between them. Figure 2.2 we present a visual comparison between the five annotators’ estimated dependencies. In the left, we show the Pearson correlation coefficient between the annotations \mathbf{y}_r and $\mathbf{y}_{r'}$ for $r, r' \in \{1, \dots, R\}$ as the reference value. We remark that

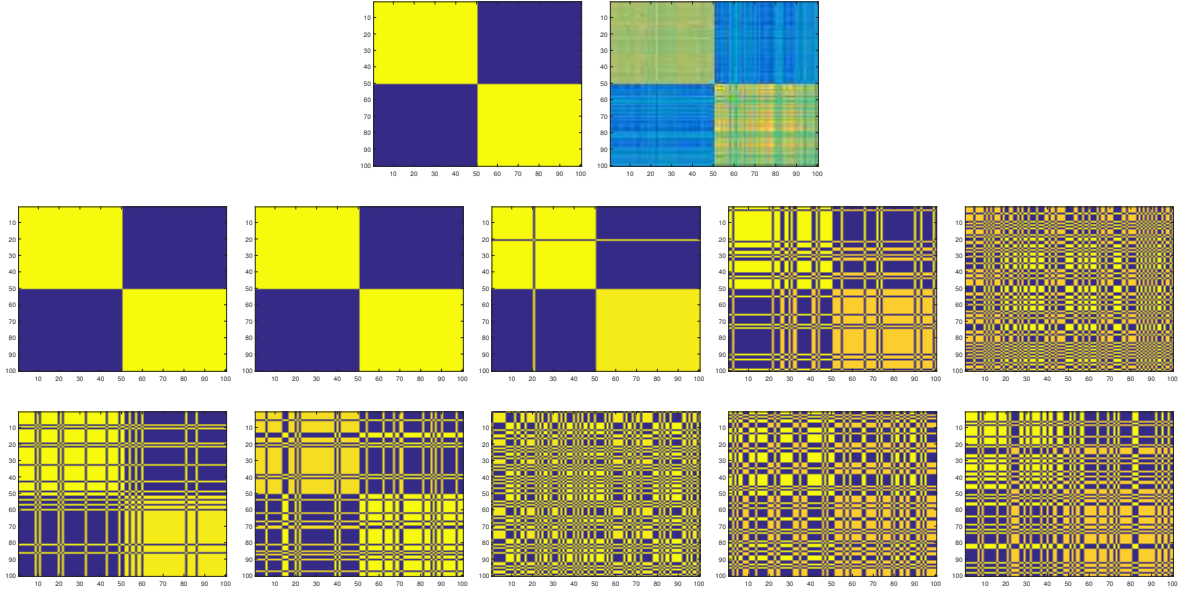


Figure 2.1: Iris dataset illustrative results. On the first row from left to right, we show the similarities among samples according to the true labels and the input features. On the second, and third rows, from left to right, we expose the kernels computed from the annotators labels concerning each simulation method (see Table 2.5). Vertical and horizontal axes display the sample index which is sorted regarding the value of the ground truth labels.

only positive values of the Pearson correlation coefficient are allowed; negative values are fixed as 0. Similarly, in the right, we estimate the dependencies $\Gamma_{rr'} \in [0, 1]$ between the r -th and r' -th annotator based on the CKA formulation as $\Gamma_{rr'} = \langle \bar{\mathbf{K}}_r, \bar{\mathbf{K}}_{r'} \rangle_{\mathbf{F}}$. Remarkably, comparing the reference and the estimated dependencies among the annotators, we can see that, although the exact behavior is not recovered, KAAR can identify the labelers' relevant relationship, which improves the annotators' representation.

Up to this point, we have verified that the relevance vector $\boldsymbol{\nu}$ captures the performance of multiple annotators even if the goal standard is not available and considering dependencies among the annotators. In turn, we perform classification experiments using datasets from the UCI repository, where we simulate five annotators with different levels of expertise using the simulation parameters exposed in the second column of Table 2.5. Here, the KAAR-based coding is employed within a *Gaussian Process-based classifier*-(GPC) to predict a new label. For concrete testing, the GPC kernel function is fixed as a squared exponential with automatic relevance determination (ARD), and the kernel parameters are computed by minimizing the marginal likelihood function. Ta-

Table 2.5: KAAR-based annotator coding results for the iris dataset. Targets: the true labels (ν_{TL}) and the input features (ν_{IF}).

Simulation method	Parameter value	AUC [%]	ν_{TL}	ν_{IF}
Classifier disrupting	$\sigma_1^2 = 10$	100	0.7071	0.6971
	$\sigma_2^2 = 10^2$	100	0.7071	0.6971
	$\sigma_3^2 = 10^4$	99.88	0	0.1613
	$\sigma_4^2 = 10^6$	87.28	0	0.0465
	$\sigma_5^2 = 10^{10}$	45.72	0	0
Biased coin	$p_1 = 0.1$	87.00	0.8423	0.8410
	$p_2 = 0.3$	84.00	0.5116	0.5186
	$p_3 = 0.5$	51.00	0	0
	$p_4 = 0.6$	35.00	0.0767	0.0843
	$p_5 = 0.7$	31.00	0.1511	0.1291

ble 3.3 reports the mean and the standard deviation for the predicted AUC. Moreover, the method with the highest performance is highlighted in bold, excluding the upper bound (GPC-GOLD), which is a GPC trained with the true labels. As seen, most of the classification methods from multiple annotators considered in this work outperform the majority voting baseline in most cases, which is not surprising since this baseline does not recognize different expertise. Remarkably, we can note that our approach is less prone to overfitting when compared with the parametric models (GPC-MV, MA-LFC, MA-MAE, MA-DGRL, MA-GPC). The above is a direct consequence of the number of parameters for modeling the performance of the annotators. For instance, MA-LFC and MA-GPC use $2 \times R$ parameters, MA-DGRL uses R parameters, and MA-MAE uses $P \times R$ parameters. Otherwise, our approach estimates the labelers expertise as closed-form similarities between the input space and the labelers annotations; hence, no free parameter tuning is required. However, it is important to highlight a trade-off between model overfitting and accuracy due to approaches that code the annotators' expertise using parametric representations, i.e., MA-GPC and MA-LFC, which can capture biased labelers.

The aforementioned is empirically demonstrated in Table 2.6; in some cases, MA-LFC and MA-GPC exhibit better performance than KAAR (mainly for the datasets bupa and segmentation). Nonetheless, we remark that our approach outperforms all the models considered for validation in general.

Now, Table 2.7 outlines the achieved performances for the datasets holding real-world annotators. We highlight in bold the best method, excluding the upper bound (GP-GOLD). As seen, for the voice quality dataset, the scale R allows exhibiting similar

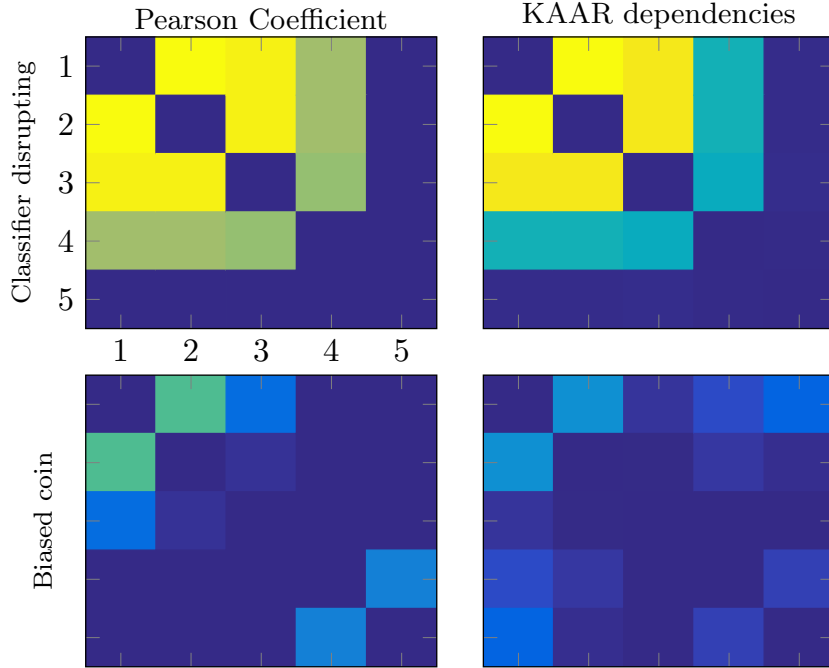


Figure 2.2: Annotators interdependencies. In the left column, from top to bottom, we show the Pearson coefficient computed over the labels from multiple annotators. In the right, we present the annotators’ dependencies estimated with our KAAR.

Table 2.6: UCI repository classification results. Bold: the highest performance excluding the upper bound (target) classifier GPC-GOLD

(a) Classifier disrupting.

Method	breast	bupa	ionosphere	pima	tic-tac-toe	iris	wine	segmentation	average
GPC-GOLD	99.04 ± 0.94	72.21 ± 3.69	95.02 ± 2.44	83.76 ± 1.98	99.97 ± 0.06	96.65 ± 3.81	99.40 ± 0.87	91.39 ± 2.05	92.18
GPC-MV	99.39 ± 0.41	65.15 ± 5.20	93.74 ± 3.04	82.59 ± 2.90	58.70 ± 5.11	99.75 ± 0.42	99.40 ± 0.87	94.13 ± 0.98	86.61
MA-LFC	99.51 ± 0.30	55.52 ± 17.25	86.97 ± 16.27	83.53 ± 1.73	62.89 ± 3.25	99.63 ± 0.39	99.88 ± 0.17	97.31 ± 0.42	85.66
MA-DGRL	99.61 ± 0.27	68.11 ± 7.62	90.51 ± 3.55	82.59 ± 2.82	61.10 ± 3.16	97.11 ± 1.87	99.78 ± 0.33	92.05 ± 2.36	86.36
MA-MAE	99.62 ± 0.31	66.37 ± 3.99	92.81 ± 2.83	83.10 ± 2.36	63.45 ± 2.80	99.55 ± 0.86	99.81 ± 0.27	95.59 ± 0.50	87.54
MA-GPC	95.08 ± 2.46	56.68 ± 10.39	94.79 ± 2.70	66.77 ± 4.30	60.07 ± 2.83	99.68 ± 0.34	99.81 ± 0.27	98.20 ± 0.22	83.88
KAAR	99.39 ± 0.39	65.33 ± 5.47	96.63 ± 1.57	83.91 ± 2.45	64.05 ± 3.58	99.64 ± 0.47	99.56 ± 0.81	94.50 ± 0.95	87.87

(b) Biased coin.

Method	breast	bupa	ionosphere	pima	tic-tac-toe	iris	wine	segmentation	average
GPC-GOLD	99.04 ± 0.94	72.21 ± 3.69	95.02 ± 2.44	83.76 ± 1.98	99.97 ± 0.06	96.65 ± 3.81	99.40 ± 0.87	91.39 ± 2.05	92.18
GPC-MV	94.13 ± 3.69	59.28 ± 5.95	65.48 ± 7.72	64.48 ± 8.57	62.54 ± 5.17	98.45 ± 1.84	97.83 ± 1.83	90.63 ± 2.73	79.10
MA-LFC	99.35 ± 0.43	71.43 ± 4.60	84.45 ± 4.27	83.02 ± 2.01	60.38 ± 3.09	99.56 ± 0.53	97.89 ± 1.78	99.41 ± 0.20	86.93
MA-DGRL	99.14 ± 1.26	68.89 ± 6.15	58.06 ± 8.45	82.23 ± 2.60	56.30 ± 4.63	94.01 ± 5.51	93.16 ± 5.77	96.24 ± 1.10	81.00
MA-MAE	96.57 ± 2.46	68.14 ± 4.58	67.34 ± 7.15	79.90 ± 3.49	50.43 ± 5.39	99.08 ± 0.92	95.25 ± 2.54	98.63 ± 0.33	81.92
MA-GPC	94.45 ± 2.51	43.19 ± 9.84	92.03 ± 6.05	68.96 ± 5.86	73.42 ± 4.15	99.45 ± 0.54	99.68 ± 0.35	99.41 ± 0.16	83.83
KAAR	99.04 ± 0.49	68.51 ± 5.42	85.51 ± 5.45	82.48 ± 2.72	90.13 ± 3.60	99.34 ± 0.91	99.33 ± 0.56	93.18 ± 1.78	89.69

AUC values for most of the approaches. The latter can be explained in the sense that for these features, the annotators share comparable performances in terms of the AUC. Besides, we note that our KAAR approach obtains the highest AUC values for the scale

Table 2.7: Real annotators datasets results. Bold: the method with the highest performance excluding the upper bound (target) classifier GPC-GOLD.

Method	AUC(%)				Average
	Voice Dataset R	Voice Dataset B	Polarity Dataset	Music Dataset	
GPC-GOLD	92.83	91.89	80.26	92.84	89.45
GPC-MV	88.91	82.05	77.62	88.79	84.34
MA-LFC	90.86	86.68	50.53	85.99	78.51
MA-MAE	89.95	82.07	48.73	81.92	75.66
MA-DGRL	93.11	80.93	56.13	88.32	79.62
MA-GPC	93.55	80.16	61.18	82.53	79.35
KAAR (proposal)	92.37	90.83	78.15	88.96	87.57

B, which is a remarkable result as long as in this scale the quality of the annotations decreases considerably for the R references. On the other hand, concerning the sentiment polarity dataset, the KAAR algorithm clearly outperforms all the validation methods; in fact, it achieves similar performance to the model trained with the gold standard (GPC-GOLD). The above can be explained in the sense that by taking into account the correlations between the labelers encourage a suitable coding about the performance of the annotators, which has a positive effect on the quality of the predictions.

2.4.2 Regression

First, we perform a controlled experiment aiming to verify the capability of our KAAR method to estimate the performance of multiple annotators taking into account their dependencies. For this first experiment, we use the *fully synthetic* dataset described in Section 2.3.2. We simulate five labelers ($R = 5$) with different expertise levels and using the parameters shown in Table 2.8. Besides, such a Table shows the KAAR-based coding results for this *fully synthetic* dataset regarding the matching weights in $\boldsymbol{\nu}$. From the classification result in Table 2.5, we note that our approach can measure the annotators' performance based on the similarities between the annotations from multiple labelers and the input features. Accordingly, for this first regression experiment, we compute the weights $\boldsymbol{\nu}$ as the CKA-based dependency between each expert kernel and the target kernel, which is computed over the matrix \mathbf{X} . From Table 2.8, we remark that the behavior of R^2 score for the simulated annotators and the KAAR-based results are proportional; hence, we identify that our approach can infer the annotator performance from the input features in scenarios of real-valued labels. Now, Figure 2.3 (column 1) shows the regression results generated by 5 different regression schemes based on GPs, where the r -th regressor g_r is trained with dataset \mathcal{D}_r . Remarkably, if we compared the

Table 2.8: KAAR-based annotator coding results the *fully synthetic* dataset.

Parameter value	R^2	ν
$v_1 = 0.2$	0.5698	0.5632
$v_2 = 0.33$	0.2423	0.3048
$v_3 = 0.5$	0.0412	0.1252
$v_4 = 2$	-2.8083	0.0025
$v_5 = 4$	-8.0806	0.0040

KAAR-based parameters ν in Table 2.8 and the regression results. We note that the higher the regression performance, the higher the KAAR-based performance estimation. Moreover, in the first two images (from top to bottom) of the second column of Figure 2.3, we expose the result of our GPR-GOLD (a GP-based regression model trained with the actual labels) and the regression results generated by the KAAR-based approach based on the combination of the regression results exposed in the first column. We notice that our method performs the regression task properly. In fact, the regression results of KAAR are similar to the ones from GPR-GOLD. On the other hand, we recall that our work’s main objective is to estimate the annotators’ performance by taking into account inter-dependencies. In the last two images of the second column in Figure 2.3, we expose a visual comparison between the five annotators’ estimated dependencies. In the left, we show the Pearson correlation coefficient between the annotations \mathbf{y}_r and $\mathbf{y}_{r'}$ for $r, r' \in \{1, \dots, R\}$ as the reference value. We remark that only positive values of the Pearson correlation coefficient are allowed; negative values are fixed as 0. Similarly, in the right, we estimate the dependencies $\Gamma_{rr'} \in [0, 1]$ between the r -th and r' -th annotator based on the CKA formulation as $\Gamma_{rr'} = \langle \bar{\mathbf{K}}_r, \bar{\mathbf{K}}_{r'} \rangle_{\mathbf{F}}$. Comparing the reference and the estimated dependencies among the annotators, we can see that, although the exact behavior is not recovered, KAAR can identify the labelers’ relevant relationship.

On the other hand, regarding the semy-synthetic datasets, Table 2.9 shows the re-

Table 2.9: Regression results in terms of R^2 score over *semi synthetic datasets*. Bold: the highest R^2 excluding the upper bound GPR-GOLD.

Method	auto	Bike	concrete	housing	yacht	ct	Average
GPR-GOLD	0.8698 \pm 0.0268	0.5612 \pm 0.0059	0.8101 \pm 0.0251	0.8288 \pm 0.0389	0.8043 \pm 0.0741	0.8652 \pm 0.0054	0.7899
GPR-Av	0.8518 \pm 0.0273	0.5432 \pm 0.0076	0.7749 \pm 0.0283	0.7984 \pm 0.0424	0.7667 \pm 0.0870	0.8223 \pm 0.0074	0.7595
MA-LFCR	0.7992 \pm 0.0231	0.3880 \pm 0.0073	0.5975 \pm 0.0388	0.7072 \pm 0.0538	0.6120 \pm 0.0775	0.6619 \pm 0.2492	0.6276
MA-GPR	0.8558 \pm 0.0269	0.4242 \pm 0.0175	0.7807 \pm 0.0312	0.7247 \pm 0.0628	0.7319 \pm 0.0965	0.0111 \pm 0.0038	0.5881
MA-DL-B	0.7890 \pm 0.0296	0.5926 \pm 0.0121	0.2281 \pm 0.0355	0.5323 \pm 0.0757	0.1897 \pm 0.0680	0.7590 \pm 0.2437	0.5151
MA-DL-S	0.7866 \pm 0.0274	0.5882 \pm 0.0122	0.2259 \pm 0.0312	0.5276 \pm 0.0852	0.1900 \pm 0.0722	0.8722 \pm 0.2388	0.5318
MA-DL-B+S	0.7826 \pm 0.0274	0.5874 \pm 0.0135	0.2216 \pm 0.0331	0.5326 \pm 0.0884	0.1815 \pm 0.1018	0.6782 \pm 0.2659	0.4973
KAAR	0.8541 \pm 0.0303	0.5667 \pm 0.0062	0.8034 \pm 0.0255	0.8195 \pm 0.0354	0.7729 \pm 0.0697	0.8884 \pm 0.0022	0.7841

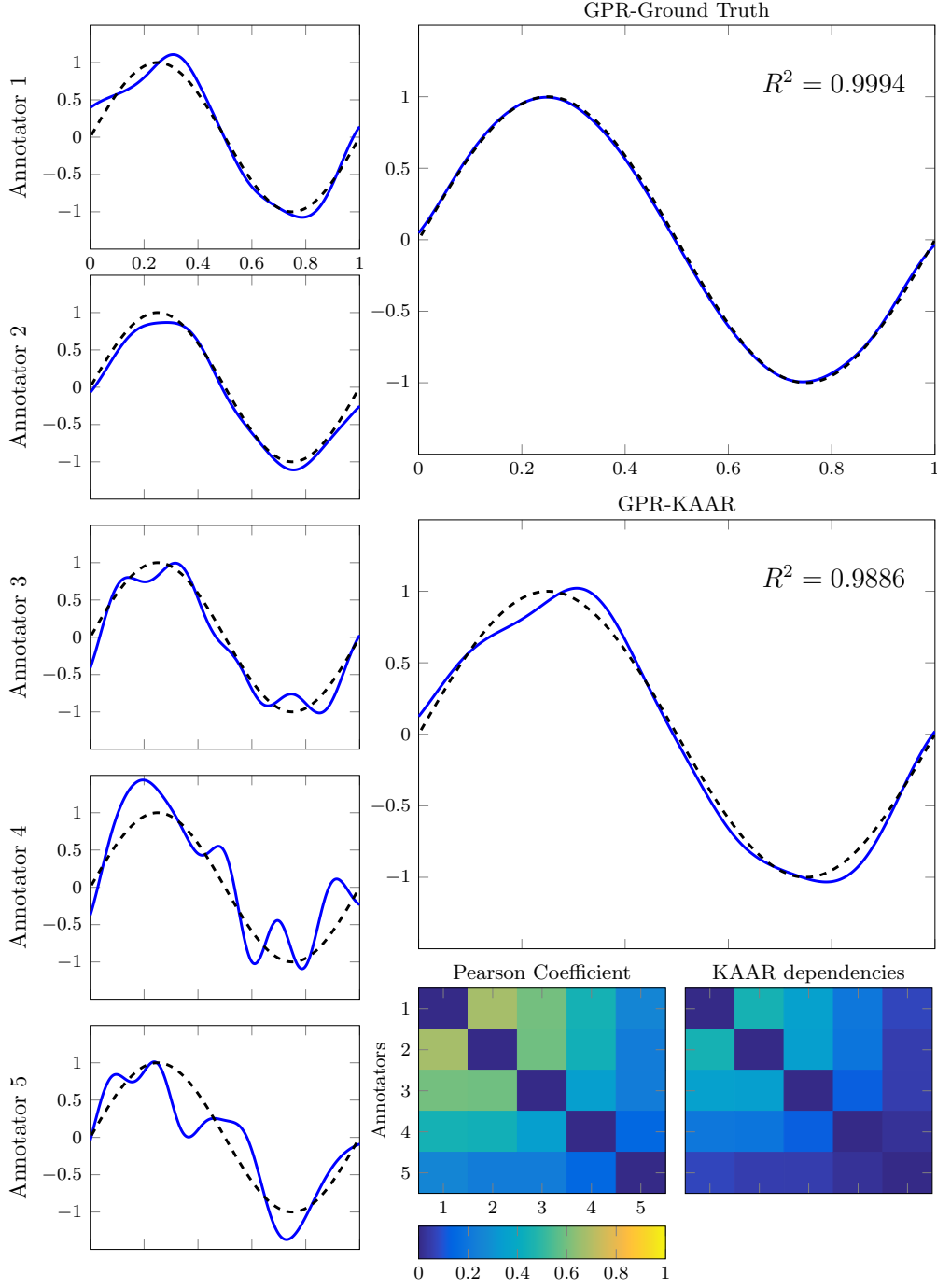


Figure 2.3: Regression illustrative results. On the first column, from top to bottom, we show the regression results for each of the annotators. On the second column, from top to bottom, we exhibit the result for GPR-GT, the result for our GPR-KAAR; finally, in the last row, we show the Pearson coefficient computed over the labels from multiple annotators and the dependencies estimated with our KAAR.

sults of the *semi synthetic datasets*. On average, our KAAR exhibits the best generalization performance in terms of the R^2 score. Now, regarding its GPs-based competitors (GPR-Av, MA-GPR), we note that the intuitive lower bound GPR-Av exhibits a worse prediction than our KAAR. However, we remark on MA-GPR’s behavior, which is lowest compared with its GPs-based competitors, even far worse than the supposed lower bound GPR-Av. The key to this particular outcome lies in its formulation; MA-GPR models the annotators’ behavior by assuming no dependencies among the labelers’ decisions.

Conversely, we analyze the results concerning the linear model MA-LFR; from the results, we notice that this approach’s prediction capacity is far lower than our KAAR. Such an outcome indicates that there may exist a non-linear structure in most databases. Finally, we analyze the DL-based models’ results; we remark a considerable low prediction capacity; in fact, they are even defeated by the linear model MA-LFR. We explain this behavior in the sense that the DL-based model uses an extra layer (termed Crowd-Layer), which is used to manage the data from multiple annotators. Such a layer does not offer a suitable codification of the labelers’ behavior.

Finally, the *fully real datasets* present the most challenging scenario, where both the input features and the labels come from real-world applications. Table 2.10 outlines

Table 2.10: Regression results in terms of R^2 score over *fully real dataset*. Bold: the highest R^2 excluding the upper bound GPR-GOLD.

Method	Music
GPR-GOLD	0.4889
GPR-Av	0.2744
MA-LFCR	0.1404
MA-GPR	0.0090
MA-DL-B	0.2339
MA-DL-S	0.2934
MA-DL-B+S	0.3519
KAAR	0.2816

the achieved performances. We notice that the DL approaches MA-DL-B+S and MA-DL-S obtain the best generalization performance in terms of the R^2 score, followed by our KAAR. Further, as theoretically expected, such performance lies between that of GPR-GOLD and GP-Av. Moreover, regarding the GPs-based competitor MA-GPR, we note that it exhibits the worst prediction capability with a R^2 close to zero. We argue that the above is a symptom of overfitting, which can be confirmed based on the fact that the training R^2 score is 0.4731, which is comparable with GPR-GOLD. Conversely,

the linear approach MA-LFCR performs worse than the theoretical lower bound GP-Av, which indicates a non-linear structure in the Music dataset. Also, we observe that all regression models presented a lower generalization performance than previous results over the same dataset. The above is a repercussion of solving a multi-class classification problem with regression models.

2.5 Summary

We propose a new kernel-alignment-based approach, termed KAAR, to support the enhancement of supervised learning in the context of multiple annotators data. To this end, KAAR computes each provided expert’s relevance through a CKA-based averaged matching between the annotator labels and the input data features. A convex combination of supervised learning algorithms is carried out by adopting the multiple annotator performances coded in the KAAR-based relevance analysis. Unlike previous works, our proposal estimates the performance of the annotators using a non-parametric model. It relaxes the assumption of independence between the labelers, which allows coding some biases and tendencies between the annotators’ opinions. We tested our approach in synthetic and real-world datasets for both classification and regression settings. For the synthetic experiments, we use some databases from the UCI repository. We simulate multiple annotators following the schemes in Section 2.3.1 and Section 2.3.2 for classification and regression, respectively. On the other hand, for the real-world datasets, we gather two classification tasks where the annotations are obtained from multiple experts’ opinions (voice quality problem) and using the crowd-sourcing platform AMT (polarity and music data). Moreover, for regression problems, we use the music dataset. The results show that the proposed method can deal with binary classification, multiclass classification, and regression problems with multiple labelers. In fact, in most cases, our approach achieves competitive or even better results when compared to different state-of-the-art models (Groot et al., 2011; Raykar et al., 2010; Rodrigues and Pereira, 2018a; Rodrigues et al., 2014b; Yan et al., 2014; Zhang et al., 2015). Furthermore, we experimentally demonstrate (using simulated annotators) that the performance of our approach does not depend on the model used for reproducing the annotations. Similarly, in the real-world datasets, it is evidenced that KAAR is not significantly affected when the labelers’ expertise decreases drastically. Remarkably, KAAR deals with scenarios with missing labels.

Still, KAAR assumes that the annotators’ performance only depends on the ground

truth labels, which is not entirely accurate in many real applications, as it was pointed out in [Yan et al. \(2014\)](#). Accordingly, future work must be oriented toward relaxing this assumption by considering that the annotators' expertise depends on the ground truth and the instance that they are labeling.

Chapter 3

Localized kernel alignment-based annotator relevance analysis

To model the annotators' behavior, it is necessary to learn some parameters related to their performance. Such parameters include the accuracy (Rodrigues et al., 2013), the confusion matrix (Gil et al., 2015), the error variance (Groot et al., 2011), and the bias (Rodrigues et al., 2017). In the literature, it is commonly founded that the parameters are modeled as fixed points (Rodrigues et al., 2014b) or as random variables (Morales-Alvarez et al., 2020), where it is considered that such parameters are homogeneous across the input data. The latter assumption is wrong since an expert makes decisions based not only on his/her expertise but also on the features observed from raw data (Yan et al., 2010, 2014).

On the other hand, the independence among the annotators is commonly used to reduce the complexity of the model (Venanzi et al., 2014), or based on the fact that it is plausible to guarantee that each labeler performs the annotation process individually (Tang et al., 2019). Nevertheless, this is not entirely correct due to there may exist correlations among the annotators (Zhang and Obradovic, 2011). For example, if the sources are humans, the independence assumption is hardly feasible because knowledge is a social construction; hence, people's decisions will be correlated when they share information, communicate with each other, or belong to a particular school of thought (Hahn et al., 2018; Surowiecki, 2005). Accordingly, the relaxation of this restriction could be used to improve the ground truth estimation (G. Rodrigo et al., 2019).

According to the related previously, in this chapter, we propose an approach to face supervised learning problems with multiple annotators. Our model is an extension of the KAAR model introduced in Chapter 2. Like KAAR, our LKAAR is built as a convex

combination of classifiers. The annotator's performance is computed from the matching between the input features and the labels. Similarly, both approaches take into account dependencies between the annotators. However, unlike KAAR, our LKAAR models the annotators' parameters as a function of the input features, which is an essential aspect in the behavior of the annotators, as has been established in the literature (Rodrigues et al., 2013; Xiao et al., 2013). Finally, We highlight that due to our approach is capable to model inconsistent annotators, it is more robust to outliers compared with such models that do not consider the relationship between the input features and the labelers' behavior. LKAAR estimates the annotators' performance for every region in the input space; meanwhile, the other approaches estimate such performance as an average of some parameters (Morales-Alvarez et al., 2020; Rodrigues et al., 2017; Ruiz et al., 2019). Consequently, it is known that the average operator suffers under the presence of outliers (Kara et al., 2015).

3.1 Localized Kernel alignment fundamentals

As we exposed in Chapter 2 an usual approach to perform kernel selection comprises a convex combination of R basis kernels $\mathbf{K}_\nu = \sum_{r=1}^R \nu_r \mathbf{K}_r$, where $\mathbf{K}_r \in \mathbb{R}^{N \times N}$ is a matrix holding elements $\kappa_r(\mathbf{x}_n, \mathbf{x}_{n'})$, $\kappa_r: \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$ is a particular kernel ($n, n' \in \{1, \dots, N\}$), and $\nu_r \in \mathbb{R}$ is a weighting factor.

The above combination assumes that the distribution of the input data is stationary, that is, the weight ν_r penalizes equally all the samples for the r -th kernel function, which cannot fulfill real-world scenarios. To deal with this issue, localized multiple kernel learning-based approaches compute \mathbf{K}_q as the following quadratic combination (Gönen and Alpaydin, 2008; Zhang and Gao, 2018):

$$\mathbf{K}_q = \sum_{r=1}^R \mathbf{Q}_r \mathbf{K}_r \mathbf{Q}_r, \quad (3.1)$$

where $\mathbf{Q}_r \in \mathbb{R}^{N \times N}$ is a diagonal matrix whose elements are defined by the vector $\mathbf{q}_r = [q_r(\mathbf{x}_1), \dots, q_r(\mathbf{x}_N)]^\top$. The combination factors $\{\mathbf{q}_r\}_{r=1}^R$ should be estimated in such a way as to maximize similarity between the kernel matrices \mathbf{K}_q and $\mathbf{F} \in \mathbb{R}^{N \times N}$, where \mathbf{F} holds elements $\kappa_y(y_n, y_{n'})$ and $\kappa_y: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. Then, the Centered Kernel Alignment–

(CKA) between \mathbf{K}_q and \mathbf{F} is computed as follows (Cortes et al., 2012):

$$\rho(\mathbf{K}_q, \mathbf{F}) = \frac{\langle \bar{\mathbf{K}}_q, \bar{\mathbf{F}} \rangle_{\text{F}}}{\|\bar{\mathbf{K}}_q\|_{\text{F}} \|\bar{\mathbf{F}}\|_{\text{F}}}, \quad (3.2)$$

where $\bar{\mathbf{F}}$ stands for a centered kernel computed as: $\bar{\mathbf{F}} = \mathbf{H}\mathbf{F}\mathbf{H}$, being $\mathbf{H} = \mathbf{I} - \mathbf{1}\mathbf{1}^\top/N$ a centering matrix, $\mathbf{I} \in \mathbb{R}^{N \times N}$ is the identity, $\mathbf{1} \in \mathbb{R}^N$ is an all-ones vector, and $\|\cdot\|_{\text{F}}$ and $\langle \cdot, \cdot \rangle_{\text{F}}$ stand for the Frobenius norm and inner product, respectively. Moreover, the centralized version of \mathbf{K}_q is estimated as: $\bar{\mathbf{K}}_q = \sum_{r=1}^R \mathbf{Q}_r \bar{\mathbf{K}}_r \mathbf{Q}_r$, with $\bar{\mathbf{K}}_r = \mathbf{H}\mathbf{K}_r\mathbf{H}$.

3.2 Localized kernel alignment-based annotator relevance analysis

We enhance our KAAR presented in Chapter 2, where a CKA-based approach is applied to code each labeler's expertise in scenarios where the gold standard is not available. However, unlike KAAR, we consider the annotator's expertise as a function of the input space. Accordingly, we assume that the input features in \mathbf{X} contain the relevant information regarding the ground truth label and we compute the CKA between a kernel extracted over the input features and a local combination of R labelers kernels. Consequently, the elements of the target matrix \mathbf{F} in equation (3.2) are computed through the kernel function: $\kappa_{\mathbf{X}} : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$. Further, to compute the local relevance of annotator r concerning the input \mathbf{x}_n , we infer the diagonal of \mathbf{Q}_r in equation (3.1) as follows:

$$q_r(\mathbf{x}_n) = \begin{cases} \beta_0^{(r)} + \sum_{n'=1}^N \beta_{n'}^{(r)} \kappa_{\beta}(\mathbf{x}_n, \mathbf{x}_{n'}), & \text{if } n \in N_r \\ 0, & \text{Otherwise} \end{cases}. \quad (3.3)$$

The condition $n \in N_r$ inhibits the influence of missing labels, $\{\beta_j^{(r)} \in \mathbb{R}\}_{j=0}^N$ are combination coefficients, and $\kappa_{\beta} : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$ is a kernel function. Now, to compute $\beta_j^{(r)}$, we employ the following CKA-based optimization problem:

$$\hat{\beta}_j^{(r)} = \arg \max_{\beta_j^{(r)}} \hat{\rho}(\mathbf{K}_q, \mathbf{F}) = \arg \max_{\beta_j^{(r)}} \frac{\langle \bar{\mathbf{K}}_q, \bar{\mathbf{F}} \rangle_{\text{F}}}{\|\bar{\mathbf{K}}\|_{\text{F}}}, \quad (3.4)$$

where:

$$\|\bar{\mathbf{K}}_q\|_F^2 = \text{tr}(\bar{\mathbf{K}}_q \bar{\mathbf{K}}_q^\top) = \text{tr}\left(\sum_{r,l=1}^R \langle \mathbf{Q}_r \bar{\mathbf{K}}_r \mathbf{Q}_r, \mathbf{Q}_l \bar{\mathbf{K}}_l \mathbf{Q}_l \rangle_F\right) \quad (3.5)$$

codes the inter-annotator dependencies, being $\text{tr}(\cdot)$ the trace operator. We have omitted the term $\|\bar{\mathbf{F}}\|_F$ in equation (3.2) since it does not depend on the combination parameters $\beta_j^{(r)}$. A gradient descent-based approach is provided to solve the optimization problem in equation (3.4) as follows (Petersen et al., 2008):

$$\frac{\partial \hat{\rho}(\mathbf{K}_q, \mathbf{F})}{\partial \beta_j^{(r)}} = \text{tr}\left(\left(\frac{\partial \hat{\rho}(\mathbf{K}_q, \mathbf{F})}{\partial \mathbf{Q}_r}\right)^\top \frac{\partial \mathbf{Q}_r}{\partial \beta_j^{(r)}}\right), \quad (3.6)$$

where $\frac{\partial \hat{\rho}(\mathbf{K}_q, \mathbf{F})}{\partial \mathbf{Q}_r} \in \mathbb{R}^{N \times N}$, yields:

$$\frac{\partial \hat{\rho}(\mathbf{K}_q, \mathbf{F})}{\partial \mathbf{Q}_r} = \frac{2 \|\bar{\mathbf{K}}_q\|_F \bar{\mathbf{K}}_r \mathbf{Q}_r \bar{\mathbf{F}} - 2 \langle \bar{\mathbf{K}}_q, \bar{\mathbf{F}} \rangle_F \|\bar{\mathbf{K}}_q\|_F^{-1} \bar{\mathbf{K}}_r \mathbf{Q}_r \bar{\mathbf{K}}_q}{\|\bar{\mathbf{K}}_q\|_F^2}, \quad (3.7)$$

$$\frac{\partial \hat{\rho}(\mathbf{K}_q, \mathbf{F})}{\partial \mathbf{Q}_r} = 2 \bar{\mathbf{K}}_r \mathbf{Q}_r \left(\frac{\bar{\mathbf{F}}}{\|\bar{\mathbf{K}}_q\|_F} - \frac{\langle \bar{\mathbf{K}}_q, \bar{\mathbf{F}} \rangle_F \bar{\mathbf{K}}_q}{\|\bar{\mathbf{K}}_q\|_F^3} \right), \quad (3.8)$$

and $\frac{\partial \mathbf{Q}_r}{\partial \beta_j^{(r)}} \in \mathbb{R}^{N \times N}$ is a diagonal matrix holding elements:

$$\frac{\partial [\mathbf{Q}_r]_{nn}}{\partial \beta_j^{(r)}} = \begin{cases} 1, & \text{if } j = 0 \\ \kappa_\beta(\mathbf{x}_n, \mathbf{x}_j), & \text{if } j \neq 0 \text{ and } n \in N_r \\ 0, & \text{if } j \neq 0 \text{ and } n \notin N_r \end{cases} \quad (3.9)$$

Next, to predict the output y_n , we propose the following convex combination of supervised learning models:

$$\hat{y}_n = \hat{g}(\mathbf{x}_n) = \frac{\sum_{r=1}^R q_r^2(\mathbf{x}_n) g_r(\mathbf{x}_n)}{\sum_{r=1}^R q_r^2(\mathbf{x}_n)}, \quad (3.10)$$

where $g_r: \mathbb{R}^P \rightarrow \mathbb{N}$ is a classification function learned from the r -th annotator's dataset

\mathcal{D}_r . In turn, to compute each labeler's expertise $q_r(\mathbf{x}_{new})$ for a new sample \mathbf{x}_{new} , we hypothesize that an annotator exhibits comparable performance for similar instances:

$$q_r^2(\mathbf{x}_{new}) = \frac{\sum_{n \in N_r} q_r^2(\mathbf{x}_n) \exp(-d(\mathbf{x}_{new}, \mathbf{x}_n))}{\sum_{n \in N_r} \exp(-d(\mathbf{x}_{new}, \mathbf{x}_n))}, \quad (3.11)$$

where $d(\cdot, \cdot) : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$ is the Euclidean distance. Finally, our localized kernel alignment-based annotator relevance analysis (LKAAR) can be summarized as in Algorithm 2.

Algorithm 2: LKAAR description

- Data:** $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\}$, \mathbf{x}_{new}
- 1 Compute $\bar{\mathbf{K}}_q$ following equation (3.1).
 - 2 Compute $\bar{\mathbf{F}}$ from the kernel function $\kappa_X(\cdot, \cdot)$.
 - 3 Solve the optimization problem: $\hat{\beta}_j^{(r)} = \arg \max_{\beta_j^{(r)}} \hat{\rho}(\mathbf{K}_q, \mathbf{F})$ in equation (3.4) from the gradient in equation (3.6).
 - 4 Compute the annotators' parameters $q_r^2(\mathbf{x}_n)$, using equation (3.3) and $\hat{\beta}_j^{(r)}$.
 - 5 Learn a classification function $g_r(\cdot)$ for each annotator by using \mathcal{D}_r .
 - 6 Given a new sample \mathbf{x}_{new} compute $q_r^2(\mathbf{x}_{new})$ using equation (3.11).
 - 7 Predict the output \hat{y}_{new} for \mathbf{x}_{new} following equation (3.10).
-

3.2.1 LKAAR for classification and regression

Similar to KAAR, this second proposal has the advantage of being flexible to be applied to regression and classification scenarios. The key factor comprises the election of a proper function for $\kappa_r(y_n^{(r)}, y_{n'}^{(r)})$. For classification, we use the kernel function in equation (2.8); conversely, for regression the kernel function in equation (2.9).

3.3 Experimental set-up

In this section, we describe the experiments' configurations to validate our LKAAR in multiple annotators scenarios for classification and regression tasks.

3.3.1 Classification

Testing datasets

To test our LKAAR approach, we use three kinds of datasets. First, we generate *fully synthetic data* as two multivariate Gaussians in 2D ($P = 2$). The data for the first class is sampled from two multivariate Gaussian distributions $\mathcal{N}([-3, -3]^\top, \Sigma)$, and $\mathcal{N}([3, 3]^\top, \Sigma)$; conversely, for the second class we use $\mathcal{N}([3, -3]^\top, \Sigma)$, and $\mathcal{N}([-3, 3]^\top, \Sigma)$. The covariance matrix is fixed as: $\Sigma = [0.8 \ 0.1; 0.1 \ 0.9]$.

Second, to control the label generation (Ruiz et al., 2019), we build *semi-synthetic data* from eight datasets devoted to binary and multi class-classification of the well-known UCI repository. Third, we evaluate our proposal on three *fully real datasets*, where both, the input features and the annotations are captured from real-world problems. The description for these two dataset groups is found in Section 2.3.1; besides, Table 2.1 summarizes the tested datasets.

Simulated and provided labels

As we pointed out previously, *fully synthetic* and *semi-synthetic* datasets do not hold real annotations. In this sense, it is necessary to generate these labels synthetically as a corrupted version of the gold standard. For this purpose, we use three different simulation methods: *i) Biased coin* that assumes a constant annotator’s performance across the input space; this method was already introduced in Section 2.3.1. Conversely, aiming to simulate annotators whose performance varies regarding the input samples, we use the next two methods. *ii) Non-homogeneous labels*, which assumes that the source performance depends on the input space (Yan et al., 2010). In that sense, it represents the input space by R clusters (for concrete testing, we use the K-means algorithm to define each cluster). Then, the r -th annotator is fixed as the “expert”, i.e., his/her labels correspond to the ground truth in samples belonging to the cluster r . For the rest of the samples, the annotator makes mistakes in the 35% of the cases, which are selected randomly. *iii) Biased coin (Non-homogeneous)* is an extension of *Biased coin* (Rodrigues et al., 2013), it also divides the input space by R clusters. In each cluster, it samples a number $\tau_{n \sim c}^{(r)}$ ($n \sim k$ stands for the sample n belonging to the cluster $c \in \{1, 2, \dots, R\}$) from a Bernoulli distribution with parameter $p_{c,r}$. Then, the simulated annotations of the r -th expert yields: $y_n^{(r)} = y_n$, if $\tau_{n \sim c}^{(r)} = 0$; otherwise, $y_n^{(r)} = \tilde{y}_n$, if $\tau_{n \sim c}^{(r)} = 1$. For the *fully real* data, Polarity and Music sets in provide the crowd labels from the AMT platform. Regarding the voice quality dataset, we have annotations from four experts. However,

for concrete testing, we only take into account the G, R, and B scales, as suggested by authors in Gil et al. (2015). The annotators' expertise is not satisfactory for A and S scales.

LKAAR training and method comparison

Overall, the Gaussian kernel is preferred in pattern classification because of its universal approximating ability and mathematical tractability (?). Hence, we fix the κ_X and κ_β kernel functions in equations (3.3) and (3.4) as:

$$\kappa_X(\mathbf{x}_n, \mathbf{x}_{n'}) = \kappa_\beta(\mathbf{x}_n, \mathbf{x}_{n'}) = \exp\left(\frac{-\|\mathbf{x}_n - \mathbf{x}_{n'}\|_2^2}{2l^2}\right), \quad (3.12)$$

where $l \in \mathbb{R}^+$ is the bandwidth and $\|\cdot\|_2$ is the L2 norm. For concrete testing, we fix the term l as the median of the input distances (Schölkopf et al., 2002). Next, to model the classification function $g_r(\cdot)$ in equation (3.10), we use three different approaches: multi-class linear classifier based on Logistic regression (LRC), a classifier based on support vector machines (SVMs) using a Gaussian kernel, and a Gaussian Processes classifier (GPC). The hyperparameters related to SVC are estimated by using a cross-validation scheme. Moreover, the GPC covariance function is computed via a squared exponential kernel, fixing the hyperparameters by optimizing the marginal likelihood (Rasmussen and Williams, 2006). The *one-vs-all* scheme is utilized in SVC and GPC to deal with multi-class problems.

The quality assessment is carried out by estimating the classification performance as the Area Under the Curve (AUC) and the overall accuracy (Acc). Further, the AUC is extended for multi-class settings, as discussed by authors in Fawcett (2006). A cross-validation scheme is employed with 30 repetitions where 70% of the samples are utilized for training and the remaining 30% for testing (except for the music dataset since it clearly defines the training and testing sets). Table 3.1 displays the employed methods of the state-of-the-art for comparison purposes. The Matlab codes for our LKAAR the state-of-the-art methods studied are publicly available.¹ We remark that the GPC-Gold is used only to provide an upper bound for our LKAAR.

¹GPC-MV, MA-LFC, MA-MAE, MA-DGRL, GPC-GTIC, KAAR, and LKAAR codes: <https://github.com/juliangilg>. MA-GPC codes: <http://www.fprodrigues.com/>

Table 3.1: A brief overview of the state-of-the-art methods tested. GPC: Gaussian Processes classifier, LRC: logistic regression classifier, MV: majority voting, MA: multiple annotators, MAE: Modelling annotators expertise, LFC: Learning from crowds, DGRL: Distinguishing good from random labelers, KAAR: kernel alignment-based annotator relevance analysis.

Algorithm	Description
GPC-GOLD	A GPC using the real labels (upper bound).
GPC-MV	A GPC using the majority voting of the labels as the ground truth.
MA-LFC (Raykar et al., 2010)	A LRC with constant parameters across the input space.
MA-DGRL (Rodrigues et al., 2013)	A multi-labeler approach that considers as latent variables the annotator performance.
MA-MAE (Yan et al., 2010)	A LRC where the sources parameters depend on the input space.
MA-GPC (Rodrigues et al., 2014b)	A multi-labeler GPC, which is as an extension of MA-LFC by using a non-linear approach.
KAAR (Gil-Gonzalez et al., 2018)	A kernel-based approach that employs a convex combination of classifiers and codes labelers dependencies.

3.3.2 Regression

Testing datasets

Aiming to test our LKAAR in regression scenarios, we use the synthetic, semi-synthetic, and real dataset presented in Section 2.3.2.

Simulated and provided labels

We recall that *fully synthetic* and *semi-synthetic* datasets do not hold real annotations from multiple labelers. In this sense, it is necessary to generate such labels synthetically as corrupted versions of the true labels and considering that the labelers' performance is a function of the input samples. Accordingly, we assume that the input space is represented by R clusters, where such representation is carried out from the K-means algorithm. Then, in cluster $c \in \{1, \dots, R\}$ the simulated labels follows, $y_n^{(r)} = y_n + \mathcal{N}(0, \sigma_{c,r})$, where $\sigma_{c,r}$ represent the error variance for the r -th labeler in cluster c . For real datasets, we perform a similar procedure as in Section 2.3.2.

LKAAR training and method comparison

Like in the classification problem, our LKAAR is built based on Gaussian kernels; in that sense, we define κ_X and κ_β as in equation (3.12). Next, to estimate each function g_r , we use a Gaussian processes-based model, whose hyperparameters are estimated by minimizing a marginal likelihood (Rasmussen and Williams, 2006). The regression performance is measured using the coefficient of determination—(R^2). A cross-validation

Table 3.2: A brief overview of state-of-the-art methods tested for regression tasks. GPR: Gaussian Processes Regression, LR: logistic regression, Av: average, MA: multiple annotators, DL: Deep learning, LFCR: Learning from crowds for regression.

Algorithm	Description
GPR-GOLD	A GPR using the real labels (upper bound).
GPR-Av	A GPR using the average of the labels as the ground truth.
MA-LFCR (Raykar et al., 2010)	A LR model for MA where the labelers' parameters are supposed to be constant across the input space.
MA-GPR (Rodrigues et al., 2014b)	A multi-labeler GPR, which is as an extension of MA-LFCR.
MA-DL (Rodrigues and Pereira, 2018a)	A Crowd Layer for DL, where the annotators' parameters are constant across the input space.
KAAR (Gil-Gonzalez et al., 2018)	A kernel-based approach that uses a convex combination of regression approaches and codes the labelers dependencies

scheme is employed with 30 repetitions where 70% of the samples are utilized for training and the remaining 30% for testing (except for music dataset, since it defines the training and testing sets). Table 3.2 displays the employed methods of the state-of-the-art for comparison; we remark that for the model MA-DL, the authors provided three different annotators' codification: MA-DL-B, where the bias for the annotators is measured; MA-DL-S, where the labelers' scale is computed; and measured; MA-DL-B+S, which is a version with both (Rodrigues and Pereira, 2018a).

3.4 Results and discussion

3.4.1 Classification

We perform a controlled experiment aiming to verify the LKAAR capability to estimate the performance of inconsistent annotators as a function of the input space and taking into account their dependencies. For this first experiment, we use the *fully synthetic* dataset described in Section 3.3.1. We simulate five labelers ($R = 5$) with different levels of expertise. For the *Biased coin* we fix $\mathbf{p} = [0.1, 0.3, 0.5, 0.6, 0.7]$, where the r -th component codes the corresponding labeler's performance. For the *Biased coin (Non-homogeneous)* approach, we divide the input space into five regions and define the following performance matrix $\mathbf{P} \in [0, 1]^{R \times R}$:

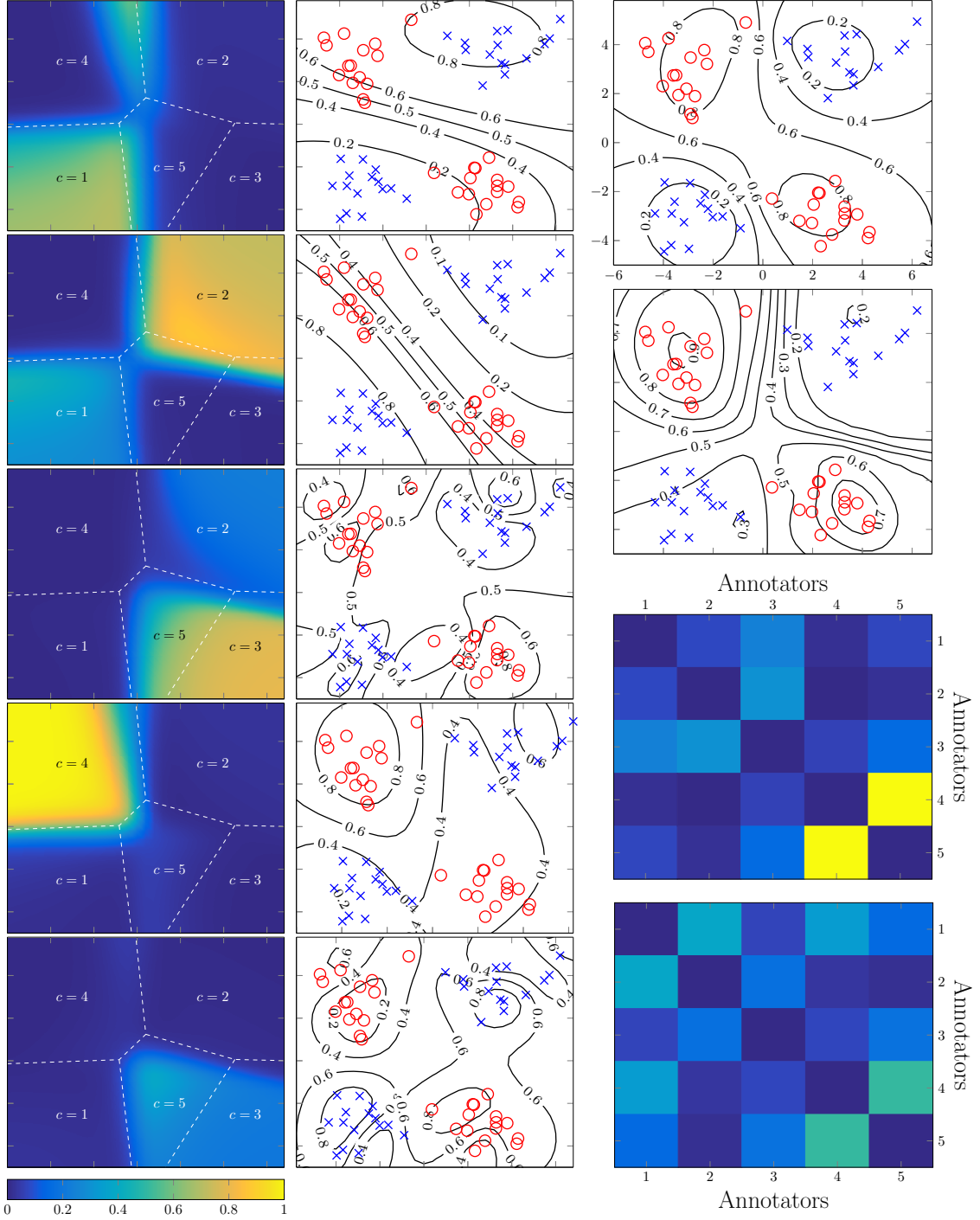


Figure 3.1: Fully synthetic dataset results. First column (top to bottom): relevance values $q_r^2(\mathbf{x}_n)$. Second column (top to bottom): decision boundaries produced by a GPC trained over each annotator dataset. Third column (top to bottom): decision boundaries generated by both the GPC-GOLD (gold standard) and the LKAAR. Also, it displays the dependencies among the annotators estimated by LKAAR (from the input samples) vs. the Pearson correlation coefficients (absolute value) from the labelers' annotations.

$$\mathbf{P} = \begin{bmatrix} 0.00 & 0.90 & 0.50 & 0.15 & 0.60 \\ 0.90 & 0.00 & 0.30 & 0.40 & 0.75 \\ 0.50 & 0.30 & 0.00 & 0.60 & 0.30 \\ 0.15 & 0.40 & 0.60 & 0.00 & 0.80 \\ 0.60 & 0.75 & 0.30 & 0.80 & 0.00 \end{bmatrix}, \quad (3.13)$$

holding elements $P_{c,r}$, which are related to the expertise of the r -th annotator in the region c . Accordingly, we note that the r -th annotator is an expert (its labels correspond to the ground truth) in the region $k = r$. Figure 3.1 (column 1) shows the LKAAR-based matching weights $q_r^2(\mathbf{x}_n)$. These weights are computed as the CKA-based similarities between each kernel expert and the target kernel from the input features. Ideally, such as target kernel would be calculated over the true labels; however, we are dealing with scenarios where these correct labels are not available. Comparing the parameters used for the simulation method *Biased coin (Non-homogeneous)* (see equation (3.13)) and the LKAAR-based results, we can note that our approach can infer the labeler performance. Thus, an appropriate parametrization allows capturing relevant information about the hidden ground truth. Now, Figure 3.1 (column 2) shows the decision boundaries generated by different Gaussian processes classifiers, where each one of them was trained with the labels from one of the annotators. Also, we can elucidate that the classification results generated by each of the simulated annotators ($g_r(\cdot)$) vary depending on the input space. Remarkably, we also note that the higher the classifier discrimination, the higher the LKAAR-based performance estimation. Moreover, in the first two images (from top to bottom) of the third column of Figure 3.1, we expose the result of our GPC-GOLD (a Gaussian processes classifier trained with the actual labels) and the decision boundaries generated by our LKAAR-based approach as in equation (3.10). We can observe that our method performs the classification task properly. Its decision boundaries are similar to those produced by the GPC-GOLD. As seen, the $q(\mathbf{x}_n)^{(r)}$ weights allow enhancing the decision boundaries around regions where the r -th expert exhibits high performances. Besides, one of the main aims of our work is to estimate the annotators' performance by taking into account inter-dependencies. So, in the last two images of the third column in Figure 3.1 we expose a visual comparison between the estimated dependencies among the five annotators. We estimate the dependencies $\Gamma_{rl} \in [0, 1]$ between the r -th and l -th

annotator from (3.5) as:

$$\Gamma_{rl} = \langle \mathbf{Q}_r \bar{\mathbf{K}}_r \mathbf{Q}_r, \mathbf{Q}_l \bar{\mathbf{K}}_l \mathbf{Q}_l \rangle_{\mathbf{F}} \quad (3.14)$$

Similarly, we consider the absolute value of the Pearson correlation coefficient between the annotations \mathbf{y}_r and \mathbf{y}_l for $r, l \in \{1, \dots, R\}$ as the reference value. Comparing the reference and the estimated dependencies among the annotators, we can see that, although the exact behavior is not recovered, LKAAR is able to identify the critical relationship among the labelers from the input features in \mathbf{X} .

On the other, regarding Semi-synthetic data results, Table 3.3(a) shows the results regarding the “Biased coin” simulated labels, where the performance of the experts is constant across the input space. On average, KAAR exhibits the best generalization performance in both AUC and Acc. However, we also note that our LKAAR-GPC gets the second-best performance in both metrics. The above is a satisfying result due to both approaches are based on the combination of R classifiers, where combinations coefficients are proportional to the labeler’s performance. Note that KAAR was designed under the assumption that the labelers’ performance is constant for every region in the input space. Now, concerning other competitors based on non-linear classifiers, we note that the approach based on majority voting gets a low performance, which is not unexpected since this method is the most naive to deal with multi-labelers scenarios. Besides, we note that the behavior of MA-GPC is surprising since its performance is a bit poor compared with other non-linear classifiers (KAAR, LKAAR-GPC, LKAAR-SVM). Linear classifiers such as MA-LFC and MA-DGRL obtains a better performance. This result indicates that MA-GPC is more prone to overfitting, which was empirically analyzed by authors in Ruiz et al. (2019). Next, we analyze the results concerning linear classifiers. Notably, simple classifiers as MA-LFC and MA-DGRL obtain competitive results compared with their non-linear competitors. The above suggests that there may exist a linear structure in some of the considered datasets. To confirm this supposition, we perform an additional experiment by training an LR-based classifier with the actual labels over all the datasets (we follow the same scheme as for GPC-GOLD). We obtain an AUC equal to 87.10 (on average), which is close to the results obtained by MA-LFC and MA-DGRL. We can elucidate that there exists a linear structure in some of the datasets, and under such settings, MA-LFC, and MA-DGRL sets an attractive option. Nonetheless, MA-MAE and LKAAR-LR obtain the worst generalization performance (MA-MAE performs worse than GPC-MV), which is a clear consequence of their feature-dependent model for the annotators. We recall that the process of generating

Table 3.3: UCI repository classification results. Bold: the highest AUC excluding the upper bound (target) classifier GPC-GOLD. Marked with *: the highest accuracy (Acc) except the upper bound.

(a) Biased coin labels.

Method		breast	bupa	ionosphere	pima	tic-tac-toe	iris	wine	segmentation	average
GPC-GOLD	AUC[%]	99.04 \pm 0.94	72.21 \pm 3.69	95.02 \pm 2.44	83.76 \pm 1.98	99.97 \pm 0.06	97.65 \pm 2.71	99.22 \pm 0.67	90.08 \pm 1.94	92.12
	Acc[%]	96.44 \pm 1.54	68.48 \pm 4.43	91.08 \pm 2.41	76.71 \pm 1.96	99.16 \pm 0.85	95.85 \pm 3.29	96.92 \pm 1.44	70.68 \pm 6.81	86.91
GPC-MV	AUC[%]	89.57 \pm 7.77	55.15 \pm 9.03	67.71 \pm 7.56	58.33 \pm 5.13	61.94 \pm 4.65	98.82 \pm 1.15	97.01 \pm 2.29	89.93 \pm 2.27	77.31
	Acc[%]	84.31 \pm 9.85	53.78 \pm 7.34	63.78 \pm 6.72	59.26 \pm 4.43	59.72 \pm 4.13	85.78 \pm 5.59	86.29 \pm 5.22	77.43 \pm 5.99	71.29
MA-LFC	AUC[%]	98.83 \pm 1.12	69.92 \pm 4.63	82.91 \pm 4.63	82.73 \pm 2.40	59.67 \pm 2.98	99.34 \pm 0.77	97.63 \pm 1.87	99.50 \pm 0.22	86.32
	Acc[%]	95.80 \pm 1.27	65.96 \pm 5.28*	83.43 \pm 4.12	76.35 \pm 2.10	66.82 \pm 2.19	93.85 \pm 4.66	90.13 \pm 4.82	94.89 \pm 0.75	83.40
MA-DGRL	AUC[%]	99.23 \pm 0.45	64.03 \pm 8.48	79.83 \pm 5.36	82.89 \pm 2.42	58.52 \pm 3.04	98.88 \pm 1.18	95.27 \pm 2.81	98.42 \pm 0.36	84.63
	Acc[%]	95.35 \pm 1.45	60.80 \pm 7.78	82.63 \pm 3.72	77.32 \pm 2.48	63.69 \pm 2.54	92.74 \pm 4.21	84.28 \pm 6.20	89.47 \pm 1.10	80.79
MA-MAE	AUC[%]	97.49 \pm 3.57	65.89 \pm 7.30	64.03 \pm 9.02	83.49 \pm 2.25	50.58 \pm 6.34	89.12 \pm 5.91	93.14 \pm 5.78	95.69 \pm 1.03	79.93
	Acc[%]	94.03 \pm 3.07	64.01 \pm 7.50	70.79 \pm 9.96	77.00 \pm 2.67*	51.56 \pm 8.00	78.74 \pm 11.07	81.70 \pm 8.82	90.99 \pm 2.92	76.10
MA-GPC	AUC[%]	94.45 \pm 2.51	43.19 \pm 9.84	92.03 \pm 6.05	68.06 \pm 5.86	73.42 \pm 4.15	99.09 \pm 0.81	99.73 \pm 0.38	99.33 \pm 0.14	83.78
	Acc[%]	96.24 \pm 1.02	54.23 \pm 4.34	80.57 \pm 6.34	73.33 \pm 2.47	71.39 \pm 3.32	93.85 \pm 3.07	93.84 \pm 3.17	93.95 \pm 0.92	82.18
KAAR	AUC[%]	99.05 \pm 0.60	67.14 \pm 6.50	93.73 \pm 2.36	81.34 \pm 2.49	91.32 \pm 3.60	99.59 \pm 0.64	99.16 \pm 0.60	92.10 \pm 1.26	90.43
	Acc[%]	95.89 \pm 1.01	63.91 \pm 5.12	89.46 \pm 2.55*	76.00 \pm 2.48	84.73 \pm 3.07*	95.11 \pm 3.16*	96.16 \pm 2.60	82.25 \pm 5.16	85.44*
LKAAR-LR	AUC[%]	97.35 \pm 2.28	60.71 \pm 7.25	69.36 \pm 8.28	76.18 \pm 3.22	54.02 \pm 3.53	95.53 \pm 2.91	98.05 \pm 1.30	98.23 \pm 0.28	81.18
	Acc[%]	89.53 \pm 5.34	57.21 \pm 6.64	69.43 \pm 7.78	69.99 \pm 3.32	55.02 \pm 3.54	85.85 \pm 5.76	92.70 \pm 3.96	90.93 \pm 1.02	76.33
LKAAR-SVM	AUC[%]	98.90 \pm 0.52	56.43 \pm 3.84	92.79 \pm 4.64	71.92 \pm 3.69	76.98 \pm 3.73	99.13 \pm 1.16	99.50 \pm 0.61	98.54 \pm 0.26	86.77
	Acc[%]	96.50 \pm 1.10*	55.10 \pm 2.97	88.63 \pm 4.11	68.81 \pm 3.19	71.71 \pm 3.84	90.89 \pm 7.62	95.91 \pm 2.33	92.67 \pm 1.12	82.53
LKAAR-GPC	AUC[%]	99.11 \pm 0.69	59.54 \pm 5.93	91.88 \pm 3.08	79.57 \pm 3.17	89.07 \pm 2.93	99.57 \pm 0.55	99.36 \pm 0.57	92.76 \pm 1.71	88.86
	Acc[%]	95.84 \pm 1.46	56.60 \pm 4.50	87.59 \pm 2.80	75.06 \pm 2.42	82.85 \pm 3.13	95.26 \pm 3.07	96.48 \pm 2.61*	82.58 \pm 4.68	84.03

(b) Non-homogeneous labels.

Method		breast	bupa	ionosphere	pima	tic-tac-toe	iris	wine	segmentation	average
GPC-GOLD	AUC[%]	99.04 \pm 0.94	72.21 \pm 3.69	95.02 \pm 2.44	83.76 \pm 1.98	99.97 \pm 0.06	97.65 \pm 2.71	99.22 \pm 0.67	90.08 \pm 1.94	92.12
	Acc[%]	96.44 \pm 1.54	68.48 \pm 4.43	91.08 \pm 2.41	76.71 \pm 1.96	99.16 \pm 0.85	95.85 \pm 3.29	96.92 \pm 1.44	70.68 \pm 6.81	86.91
GPC-MV	AUC[%]	99.11 \pm 0.58	70.95 \pm 3.90	93.14 \pm 3.49	81.21 \pm 2.57	87.83 \pm 4.11	99.63 \pm 0.39	98.41 \pm 1.38	91.48 \pm 1.48	90.22
	Acc[%]	96.29 \pm 1.48	66.60 \pm 4.31	87.90 \pm 3.26	74.87 \pm 2.32	81.96 \pm 3.46	95.33 \pm 3.05	93.96 \pm 3.34	82.68 \pm 5.30	84.95
MA-LFC	AUC[%]	98.72 \pm 0.93	71.53 \pm 4.18	82.08 \pm 4.79	82.29 \pm 2.22	61.13 \pm 3.28	98.75 \pm 1.44	96.83 \pm 1.75	99.58 \pm 0.11	86.36
	Acc[%]	95.63 \pm 1.79	69.68 \pm 4.20*	81.43 \pm 4.44	76.52 \pm 1.91*	64.88 \pm 2.86	94.44 \pm 4.62	87.74 \pm 4.67	95.40 \pm 0.71*	83.21
MA-DGRL	AUC[%]	99.30 \pm 0.39	68.00 \pm 4.09	77.60 \pm 7.50	81.72 \pm 2.57	61.83 \pm 2.80	98.78 \pm 1.34	95.33 \pm 3.35	98.31 \pm 0.32	85.11
	Acc[%]	94.63 \pm 1.77	65.77 \pm 3.47	81.94 \pm 3.42	76.45 \pm 2.81	66.45 \pm 2.24	94.59 \pm 2.96	84.91 \pm 6.50	89.58 \pm 0.99	81.79
MA-MAE	AUC[%]	99.28 \pm 0.60	70.82 \pm 3.90	78.91 \pm 6.01	81.80 \pm 2.57	60.35 \pm 3.28	85.97 \pm 2.39	98.20 \pm 1.33	97.27 \pm 0.28	85.32
	Acc[%]	96.31 \pm 1.38	66.92 \pm 3.35	82.25 \pm 3.99	76.12 \pm 2.77	65.64 \pm 2.37	94.81 \pm 4.14	89.31 \pm 5.79	92.94 \pm 0.76	83.04
MA-GPC	AUC[%]	95.81 \pm 2.94	49.81 \pm 11.72	94.46 \pm 3.09	67.83 \pm 4.24	81.44 \pm 3.81	99.15 \pm 1.03	99.85 \pm 0.24	99.42 \pm 0.14	85.97
	Acc[%]	96.70 \pm 1.37*	59.52 \pm 4.71	82.13 \pm 3.32	72.77 \pm 2.71	76.39 \pm 2.85	94.30 \pm 2.90	94.34 \pm 2.80	94.74 \pm 0.68	83.86
KAAR	AUC[%]	98.81 \pm 0.66	70.20 \pm 5.70	93.88 \pm 3.53	81.18 \pm 2.93	89.55 \pm 2.84	99.56 \pm 0.52	99.53 \pm 0.36	92.34 \pm 1.38	90.63
	Acc[%]	96.02 \pm 1.14	65.99 \pm 5.44	87.52 \pm 4.24	75.10 \pm 2.98	81.68 \pm 2.41	95.56 \pm 2.92*	96.54 \pm 2.11*	81.11 \pm 4.15	84.94
LKAAR-LR	AUC[%]	99.34 \pm 0.44	68.86 \pm 5.16	87.14 \pm 3.38	82.04 \pm 2.44	65.40 \pm 3.13	96.00 \pm 2.50	99.21 \pm 0.82	97.97 \pm 0.27	87.00
	Acc[%]	96.00 \pm 1.46	64.17 \pm 4.22	84.10 \pm 3.20	75.67 \pm 2.15	66.96 \pm 2.74	82.59 \pm 6.07	94.28 \pm 3.19	90.02 \pm 0.93	81.72
LKAAR-SVM	AUC[%]	98.29 \pm 0.80	64.37 \pm 3.36	96.98 \pm 2.01	77.80 \pm 2.28	89.82 \pm 2.14	98.05 \pm 1.90	99.53 \pm 0.47	97.89 \pm 0.32	90.34
	Acc[%]	96.36 \pm 1.02	63.14 \pm 3.67	92.19 \pm 2.43*	72.52 \pm 2.22	80.99 \pm 2.81	84.44 \pm 6.76	96.48 \pm 2.26	91.28 \pm 0.93	84.68
LKAAR-GPC	AUC[%]	99.00 \pm 0.75	71.07 \pm 5.05	93.37 \pm 2.91	81.23 \pm 2.21	91.97 \pm 2.01	99.57 \pm 0.61	99.64 \pm 0.34	92.61 \pm 1.73	91.06
	Acc[%]	96.03 \pm 1.32	66.92 \pm 4.79	87.75 \pm 3.90	75.10 \pm 2.65	84.09 \pm 2.43*	95.26 \pm 3.29	96.54 \pm 2.16*	80.98 \pm 3.91	85.33*

(c) Biased coin (Non-homogeneous) labels.

Method		breast	bupa	ionosphere	pima	tic-tac-toe	iris	wine	segmentation	average
GPC-GOLD	AUC[%]	99.04 ± 0.94	72.21 ± 3.69	95.02 ± 2.44	83.76 ± 1.98	99.97 ± 0.06	97.65 ± 2.71	99.22 ± 0.67	90.08 ± 1.94	92.12
	Acc[%]	96.44 ± 1.54	68.48 ± 4.43	91.08 ± 2.41	76.71 ± 1.96	99.16 ± 0.85	95.85 ± 3.29	96.92 ± 1.44	70.68 ± 6.81	86.91
GPC-MV	AUC[%]	90.78 ± 4.28	50.47 ± 6.19	82.91 ± 6.03	70.18 ± 6.29	65.91 ± 6.72	98.55 ± 1.38	97.75 ± 2.04	90.18 ± 1.71	80.84
	Acc[%]	86.63 ± 2.06	48.27 ± 4.84	75.65 ± 6.45	66.52 ± 5.16	64.66 ± 3.64	88.81 ± 5.00	86.92 ± 5.76	79.24 ± 4.99	74.59
MA-LFC	AUC[%]	97.99 ± 0.99	59.64 ± 8.08	72.66 ± 9.98	72.73 ± 3.43	52.88 ± 3.13	96.72 ± 8.98	96.47 ± 2.13	99.50 ± 0.15	81.07
	Acc[%]	96.00 ± 1.70*	56.41 ± 8.12	69.17 ± 12.53	58.10 ± 4.53	46.27 ± 3.03	92.30 ± 5.18	87.55 ± 4.97	95.06 ± 0.80*	75.11
MA-DGRL	AUC[%]	99.31 ± 0.42	61.77 ± 6.17	77.83 ± 7.02	81.66 ± 2.65	55.70 ± 3.95	98.76 ± 1.33	95.26 ± 3.30	98.32 ± 0.34	83.58
	Acc[%]	78.08 ± 2.22	55.64 ± 4.52	71.43 ± 5.15	76.90 ± 1.99*	60.64 ± 2.33	94.37 ± 2.66	84.84 ± 6.32	89.63 ± 0.89	76.44
MA-MAE	AUC[%]	95.22 ± 1.70	64.63 ± 9.77	64.18 ± 9.17	79.94 ± 2.64	52.36 ± 4.78	93.16 ± 5.08	96.25 ± 2.40	94.40 ± 1.26	80.02
	Acc[%]	87.15 ± 1.85	62.34 ± 8.46*	67.94 ± 7.19	75.94 ± 2.69	53.33 ± 6.42	81.70 ± 11.68	86.67 ± 5.15	88.38 ± 2.00	75.43
MA-GPC	AUC[%]	85.37 ± 5.90	40.79 ± 12.30	74.52 ± 4.57	73.17 ± 3.34	61.82 ± 4.51	98.71 ± 1.14	99.60 ± 0.41	99.35 ± 0.14	79.17
	Acc[%]	92.55 ± 2.17	52.82 ± 6.38	69.87 ± 4.41	62.42 ± 3.00	62.33 ± 2.98	93.85 ± 3.49	95.09 ± 2.65*	93.46 ± 0.83	77.80
KAAR	AUC[%]	97.81 ± 0.99	56.52 ± 9.13	82.20 ± 4.93	67.90 ± 3.16	75.34 ± 4.70	98.75 ± 1.10	97.91 ± 1.36	91.75 ± 1.41	83.52
	Acc[%]	77.19 ± 3.14	52.44 ± 7.79	72.60 ± 4.80	61.20 ± 2.95	70.69 ± 3.63	90.44 ± 5.48	91.45 ± 4.28	76.38 ± 5.05	74.05
LKAAR-LR	AUC[%]	99.52 ± 0.30	66.07 ± 6.14	82.99 ± 5.01	80.57 ± 3.31	52.32 ± 3.38	96.83 ± 1.24	99.27 ± 0.68	98.87 ± 0.30	84.43
	Acc[%]	92.47 ± 2.24	60.22 ± 6.67	78.92 ± 4.32	75.07 ± 2.65	55.64 ± 2.77	83.41 ± 6.92	94.59 ± 3.12	89.77 ± 0.99	78.76
LKAAR-SVM	AUC[%]	98.37 ± 1.00	52.35 ± 6.40	88.28 ± 5.13	66.84 ± 3.66	73.85 ± 2.43	96.22 ± 2.50	98.88 ± 0.80	97.59 ± 0.34	84.05
	Acc[%]	87.72 ± 5.17	50.96 ± 6.81	84.73 ± 6.66*	64.81 ± 3.11	70.02 ± 2.74	74.15 ± 5.98	91.82 ± 4.33	90.37 ± 1.24	76.82
LKAAR-GPC	AUC[%]	98.14 ± 1.04	58.36 ± 7.24	86.23 ± 4.47	73.80 ± 2.83	80.02 ± 1.45	99.61 ± 0.61	98.74 ± 0.93	92.24 ± 1.80	85.89
	Acc[%]	86.76 ± 4.33	54.52 ± 5.27	78.25 ± 5.51	69.64 ± 3.01	74.90 ± 2.99*	95.93 ± 3.15*	98.84 ± 3.57	87.71 ± 4.18	79.07

the input features; hence, this scenario is most convenient for MA-LFC and MA-DGRL.

On the other hand, Table 3.3(b) shows the results regarding the “Non-homogeneous labels”, where the performance of the experts depends on the input features. First, we observe that all classification models exposed better generalization performance in both AUC and Acc compared with the previous simulation method (see Table 3.3(a)). The above suggests an increase in the quality of the labels; in fact, we can confirm this suggestion by analyzing the behavior of GPC-MV, which gets quite competitive results. Besides, we highlight that our LKAAR-GPC exhibits, on average, the best performance among the nine multi-annotators classifiers (in terms of AUC and Acc). The behavior of the non-linear based approaches is, in principle, a bit unexpected in the sense that the most straightforward method (GPC-MV) reaches a comparable performance with more sophisticated strategies such as KAAR and LKAAR-SVM. However, as we already pointed out, this behavior is caused because, in this experimental set-up, the simulation method (“Non-homogeneous labels”) generates suitable quality labels, which favors the estimation of the unknown ground truth via majority voting. Besides, MA-GPC again gets a considerably low performance when compared with its other non-linear based competitors; still, this can be the result of a lack of generalization (overfitting). Now, let us analyze the results for the linear models. Remarkably, we notice that LKAAR-LR and MA-MAE get more competitive results when compared with their competitors MA-LFC and MA-DGRL. Since the labelers’ expertise is non-stationary across the input space, both LKAAR-LR and MA-MAE are suitable. Furthermore, we highlight that our LKAAR-LR outperforms all of its linear competitors, and its generalization performance is quite better than its natural competitor MA-MAE.

Finally, Table 3.3(c) shows the results as regards the experiment where the labels were simulated by using the method “Biased coin (Non-homogeneous)”. At first glance, this experiment is more challenging since there exists more difference between the performance of GPC-GOLD (the upper bound) and the classification scheme with the best performance (our LKAAR-GPC). Analyzing the non-linear classifiers, we note that again, MA-GPC obtains an insufficient performance, which is even worse than the lower bound (GPC-MV). Besides, we remark the performance of our LKAAR-SVM and LKAAR-GPC, which outperform all of its competitors (GPC-MV, KAAR, and MA-GPC). The latter is an exciting result, especially when we compare the KAAR and LKAAR-GPC performances (both approaches are based on the combination of GPCs). Concerning the linear classifiers, we remark that our LKAAR-LR surpass all of its competitors (MA-LFC, MA-DGRL, MA-MAE). Indeed, our approach obtained the second-best results,

Table 3.4: Fully real datasets results. Bold: the method with the highest performance excluding the upper bound (target) classifier GPC-GOLD.

Method	AUC([%])					
	Voice Dataset			Polarity Dataset	Music	Average
	G	R	B			
GPC-GOLD	93.66	93.66	93.66	80.26	92.84	90.81
GPC-MV	90.17	84.73	84.04	71.14	89.20	83.97
MA-LFC	89.99	90.59	87.27	72.06	89.63	85.90
MA-DGRL	85.45	90.14	79.33	74.35	89.20	83.69
MA-MAE	91.08	89.12	80.74	50.00	84.16	79.02
MA-GPC	91.50	91.16	80.81	77.18	79.18	83.96
KAAR	89.85	93.50	89.20	64.68	86.74	84.79
LKAAR-LR	90.39	92.92	88.94	76.99	89.36	87.72
LKAAR-SVM	92.06	93.02	86.98	75.68	90.79	87.70
LKAAR-GPC	90.78	93.60	89.79	77.16	90.69	88.40

also defeating non-linear classifiers such as GPC-MV, KAAR, and MA-GPC. The latter is evidence that our approach offers a better representation of the annotators' behavior. Besides, we notice that the performance of MA-MAE is quite low compared with our LKAAR-LR, which is, in principle, unexpected since both approaches compute the annotators' performance as a function of the input space. This result can be explained in two regards. First, unlike our LKAAR-LR, MA-MAE uses a logistic regression-based model to code the non-stationaries in the labelers' performance, which does not fit the labels generated in this experiment. Moreover, MA-MAE assumes independence between the annotators; the labelers make their decision independently, which decreases the modeling of the labeler's behavior.

Summarizing, we tested our approach in controlled scenarios by using three different strategies. First, we simulate annotators with homogeneity in their performance named *Biased coin labels*. The remaining two strategies named *Non-homogeneous labels* and *Biased coin (Non-homogeneous) labels* simulates inconsistent annotators, i.e., labelers, whose performance varies depending on the input features. Attained to the results (Table 3.3(a), Table 3.3(b), and Table 3.3(c)), we note that for consistent annotators KAAR offers the best performance. On the other hand, for inconsistent labelers, our LKAAR is the best option. However, we highlight that by considering inconsistent annotators, a more realistic scenario is configured (Raykar et al., 2010; Xiao et al., 2013; Yan et al., 2010); hence, our approach is presented as the most suitable option among the state-of-the-art models considered.

Until now, we have empirically demonstrated that our approach offers a better representation of the labelers' behavior since we compute the annotators' performance taking

into account dependencies among their decisions and considering that such a performance as a function of the input features. Nevertheless, the previous experiments configure a very controlled scenario due to the labels are simulated. Accordingly, these results could be biased by the simulation method. In this sense, the *fully real datasets* present the the most challenging scenario, where both the input samples and the labels come from real-world applications.

Table 3.4 outlines the achieved performances. First, we observe that for the voice data, the scales **G** and **R** exhibit a similar performance for all studied approaches. The latter can be explained in the sense that for these scales, the annotators exhibit a proper performance. On the other hand, we notice that for the scale **B**, there exists a considerable performance reduction, which indicates that the annotators' performance is lower when compared with scales **G** and **R**, as was empirically demonstrated in (Gil et al., 2015).

However, we highlight that our approach (LKAAR-LR, LKAAR-SVM, LKAAR-GPC) and KAAR exhibit the best generalization performances. Remarkably, LKAAR-LR outperforms all of its linear competitors. This is an outstanding outcome because it reflects that our approach offers a better representation of the labelers' behavior. We perform an additional experiment over the voice dataset for the scale B. For visualizations purposes; we reduce the dimensions to 2D by using the well-known t-student stochastic neighbor embedding algorithm (t-sne). Figure 3.2 (row 1) from left to right shows each expert's labels on the 2D t-sne space. Likewise, in the second row, from left to right, we show the LKAAR estimation of the dependencies among the annotators vs. the absolute value of the Pearson correlation coefficient between \mathbf{y}_r and \mathbf{y}_l , for $r, l \in \{1, \dots, R\}$. As seen, the labels provided by the first and third annotator are quite similar, which indicates the presence of a strong dependence between them. Remarkably, this structure in the annotation process is captured by LKAAR.

Next, regarding the Polarity dataset, we note that MA-GPC exhibits the best generalization performance followed by our LKAAR-GPC; actually, the difference between these two methods is not significant. Besides, the performance of MA-MAE and KAAR is quite low, which is, in principle, unexpected; however, we give some reasons to explain this anomaly. First, MA-MAE uses P parameters for the classification model and $P \times R$ parameters to model the labelers' performance. Moreover, from Table 2.1, we can note that for this dataset, $P = 1200$ and the number of training samples is 5000. The above is a problem because, considering that we have labels from 203 annotators, we need to estimate over 240000 parameters, which is leading MA-MAE to overfit. Now, the

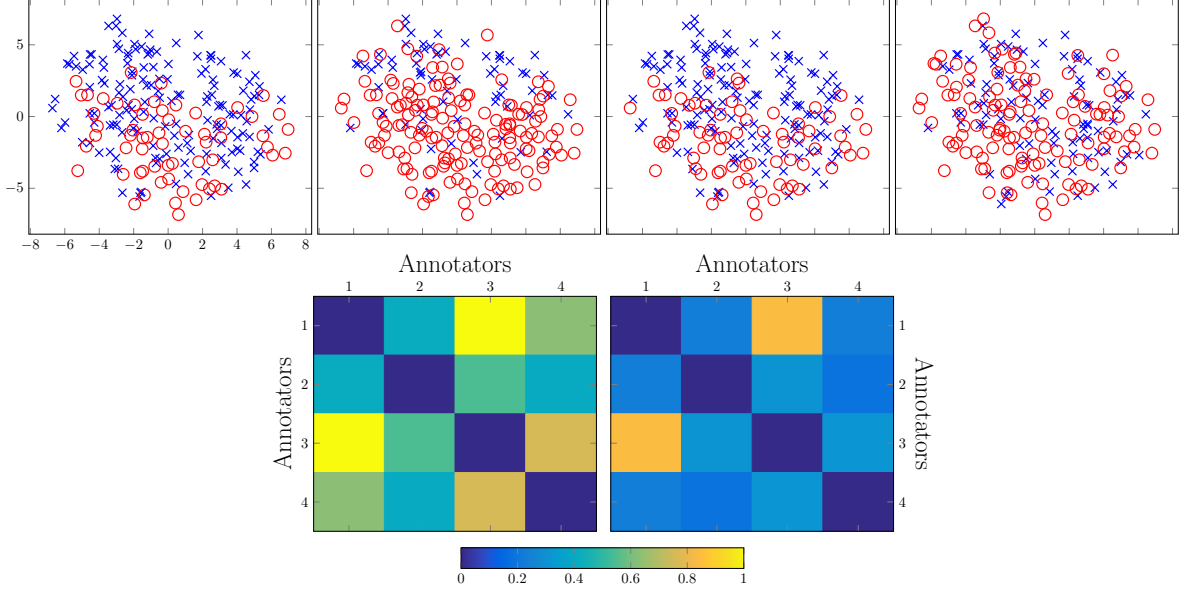


Figure 3.2: Annotator dependencies analysis for the Voice dataset (Scale B). First row (left to right): t-sne-based 2D projections holding the annotators labels. red circles: positive class, blue crosses: negative class. Second row (left to right): dependencies among annotators estimated by LKAAR (from the input samples) vs. the Pearson correlation coefficients (absolute value) from the experts’ labels.

behavior of KAAR is more unexpected than MA-MAE since it is based on a non-linear classifier, and it exhibited a competitive performance in previous experiments. We argue that the missing labels cause this unusual conduct because more than the 50% of the annotators labeled less than 30 instances. The above can represent a problem for KAAR due to it estimates the labelers’ performance as an average matching between the input features and the labels given by each annotator. Accordingly, the estimation of the labelers’ performance suffers a negative impact in scenarios with very few labels per annotator. We will expand this analysis in the final experiment.

Lastly, we analyze the results concerning Music dataset. We remark that our LKAAR-SVM and LKAAR-GPC obtains the best generalization performance in terms of AUC. Moreover, we highlight the achievement of our LKAAR-LR, given that it gets a competitive result when compared with the other linear classifiers, even defeat some non-linear classifiers (KAAR, MA-GPC, and GPC-MV). On the other hand, we note that MA-MAE and MA-GPC exhibit a significantly low performance, even lower than their intuitive lower bound (GPC-MV). This behavior is not uncommon, given that it has been repeated in the previous experiment; we argue that this outcome is the result of overfitting. However, we also note an additional issue, this dataset configures a multi-

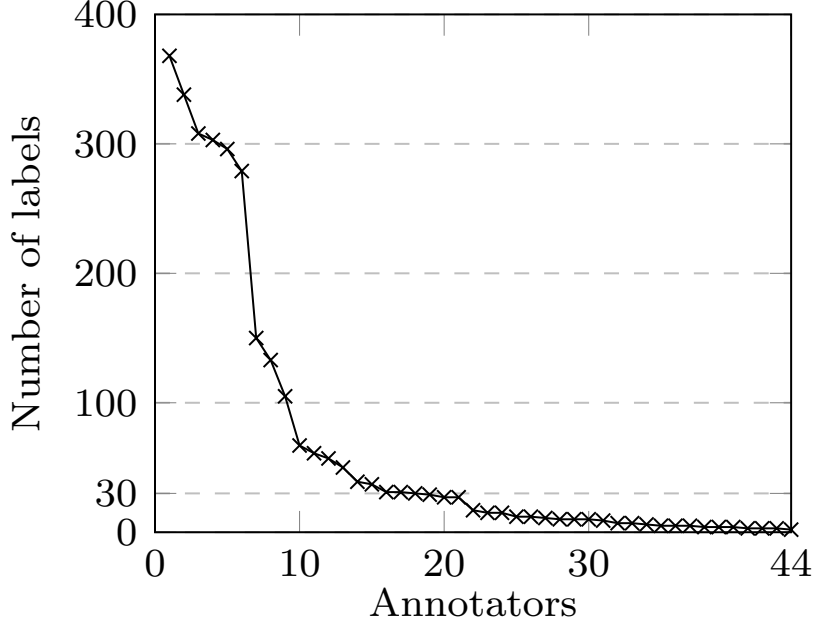


Figure 3.3: Number of labeled instances per annotator.

class classification problem; accordingly, we use a *one-vs-all* scheme for all of the binary classification (including MA-MAE and MA-GP) to deal with this. Nevertheless, such a scheme to deal with multi-class classification can lead to regions on the input space that are ambiguously classified (Bishop, 2006).

As a final experiment, we wish to evaluate the impact of the number of annotators on the performance of the multi-labeler classifiers. For concrete testing, we use the Music dataset, which holds a total of 2946 labeled instances from 44 real annotators. We sort the experts in a descending way regarding the number of cases marked. Figure 3.3 shows the number of samples labeled by each annotator. We can note that only a few annotators labeled more than 300 samples; besides, more than 50% of the labelers annotated less than 30 examples.

Figure 3.4 shows the classifiers' performance in terms of AUC as a function of the number of annotators. First, we notice that our LKAAR-SVM and LKAAR-GPC exhibit, on average, the best performances in terms of AUC (89.53, 87.60, respectively). Both approaches have no significant effects due to the number of annotators. In turn, we note an unusual behavior in MA-GPC and MA-MAE. Such action is caused by overfitting. On the other hand, we analyze that the linear models MA-LFC, LKAAR-LR, MA-GRL, and GPC-MV expose quite similar behavior, where their performance is not affected by the number of annotators. KAAR exhibits a suitable performance; however,

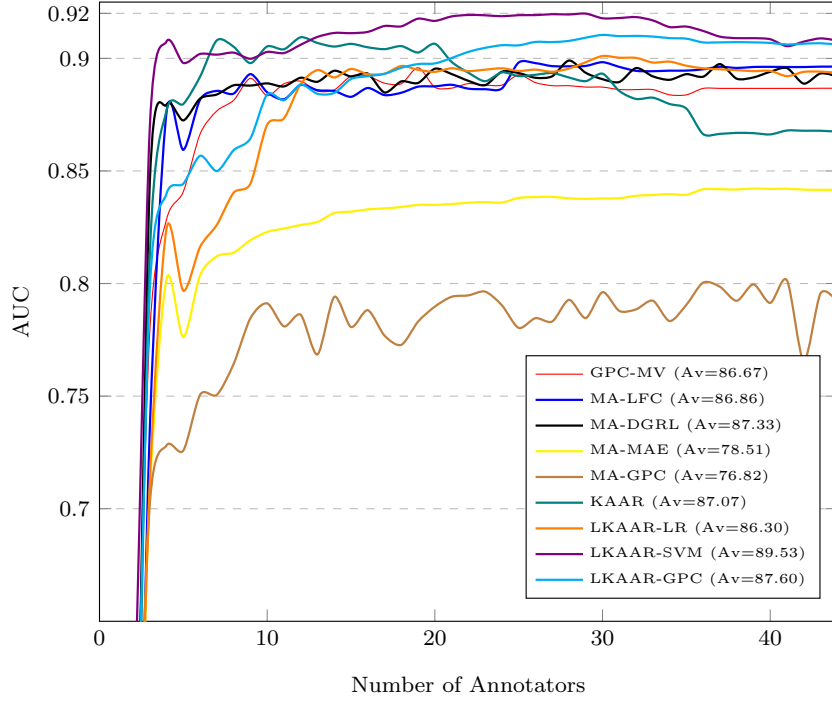


Figure 3.4: Generalization performance in terms of AUC as a function of the number of annotators.

after 20 annotators, we evidence a significant decrease in its execution; indeed, for 44 labelers, KAAR performs worse than linear classifiers. We can explain this result as follows: The point where the performance of KAAR starts to decrease coincides with the point where the number of labeled instances per annotator decreases. Accordingly, from the annotator 20, KAAR has to estimate the experts' performance with very few labels, and this could lead to miss-estimation, which harms the predictions.

3.4.2 Regression

Following a similar structure to the used for the classification results, we first carry out a controlled experiment to verify the LKAAR capabilities in the context of real-valued labels. For this first experiment, we use the *fully synthetic* dataset described in Section 3.3.1. We simulate five labelers ($R = 5$) with different levels of expertise and

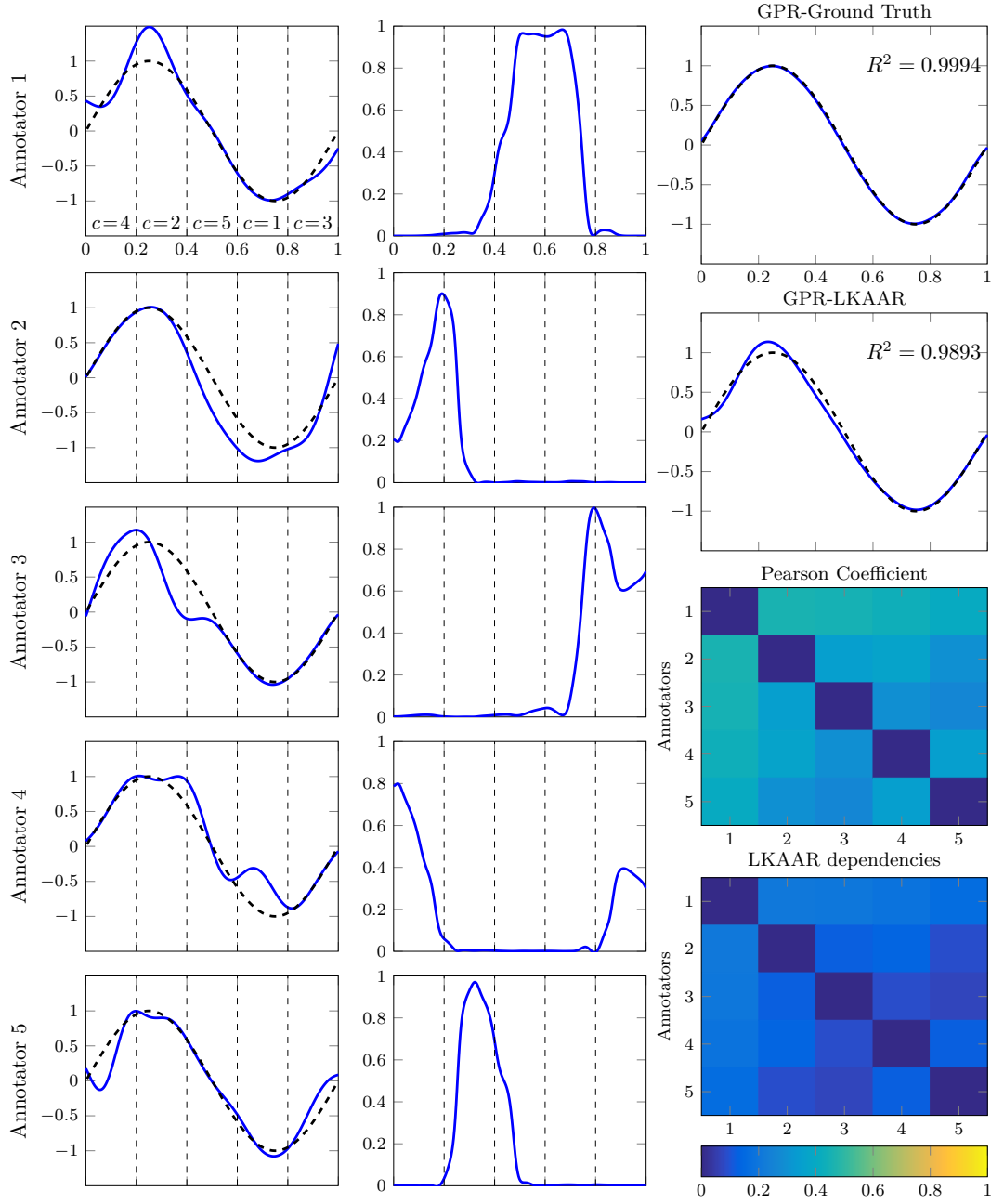


Figure 3.5: Regression illustrative results. On the first column, from top to bottom, we show the regression results g_r for each of the annotators. On the second column, from top to bottom, we present the LKAAR-based relevance q_r . In the third column, we exhibit the result for GPR-GOLD, the result for our GPR-LKAAR; finally, in the last row, we show the Pearson coefficient computed over the labels from multiple annotators and the dependencies estimated with our LKAAR.

using the following variance matrix $\mathbf{V} \in \mathbb{R}^{R \times R}$:

$$\mathbf{V} = \begin{bmatrix} 0.00 & 3.00 & 0.02 & 2.00 & 1.50 \\ 1.00 & 0.00 & 2.00 & 1.50 & 0.02 \\ 0.75 & 1.50 & 0.00 & 0.02 & 1.00 \\ 0.50 & 0.02 & 1.00 & 0.00 & 2.00 \\ 0.02 & 2.00 & 3.00 & 1.00 & 0.00 \end{bmatrix}, \quad (3.15)$$

holding elements $\sigma_{c,r}$, which are related to the expertise of the r -th annotator in the region c . We notice that the r -th annotator is an expert (its labels correspond to the ground truth) in the region $k=r$. Figure 3.5 column 1 and 2 present respectively the regression functions for each annotator g_r and the LKAAR-based weights $q_r^2(\mathbf{x}_n)$, which are related with the labelers' performance. Contrasting the parameters for the simulation process (see equation (3.15)) and the LKAAR-based results in the second column, we can note that our approach can detect the zones where the parameter $\sigma_{c,r}$ presents the lowest value and remarkably such zones match with regions where the regression results are closer to the gold standard (the above from a visual inspection). Further, in the first two images (from top to bottom) of the third column of Figure 3.1, we show the result of GPC-GOLD (a Gaussian processes regressor trained with the gold standard) and the regression results generated by our LKAAR-based approach as in equation (3.10). We notice that compared with the ground truth, our approach exhibit a proper performance; in fact, the result for our LKAAR is close to its theoretical upper bound GPC-GOLD. As seen, by estimating the $q_r(\mathbf{x}_n)$ weights as a function of the input space, we can enhance the zones where the regression function for r -th labeler, g_r , exhibits high performances. Besides, one of the main aims of this proposal is to estimate the annotators' performance by taking into account inter-dependencies. Therefore, in the last two images of the third column in Figure 3.5 we expose a visual comparison between the estimated dependencies among the five annotators. We estimate the dependencies $\Gamma_{rl} \in [0, 1]$ between the r -th and l -th annotator from (3.5). Likewise, we consider the absolute value of the Pearson correlation coefficient between the annotations \mathbf{y}_r and \mathbf{y}_l for $r, l \in \{1, \dots, R\}$ as the reference value. Comparing the reference and the estimated dependencies with LKAAR, we remark that our approach identifies critical relationships among the labelers.

On the other hand, Table 3.5 presents the results concerning semi-synthetic datasets. We highlight that our LKAAR exhibits the best regression performance in terms of the R^2 score. Now, analyzing the behavior of GPs-based competitors, we notice that GPR-

Av exhibits a lower performance when compared with our KAAR, which is not unexpected since GPR-Av corresponds to the most naive approach to deal with multi-labeler data. Besides, we remark a considerably low R^2 score (on average) for MA-GPR, even lower than the supposed lower bound GPR-Av. Our explanation for such an outcome lies in the MA-GPR’s formulation; MA-GPR models the annotators’ behavior without considering relationships between the input space and the annotators’ performance nor dependencies among them.

Regarding the linear model, MA-LFR, we note that its prediction performance is considerably lower than our method. The above indicates the presence of a non-linear structure in most considered datasets. Finally, analyzing the DL-based methods, we remark that their performance the lowest concerning the average of the R^2 score. Our explanation for this result is that the Crowdlayer used for the DL models to manage multi-labeler data does not offer a suitable representation of the annotators’ behavior.

As a final experiment, we use the Music dataset corresponding to *fully real datasets*, which configures the most challenging scenario; here, both the input samples and the labels from multiple annotators come from real-world applications. Table 3.6 outlines

Table 3.5: Regression results in terms of R^2 score over *semi synthetic datasets*. Bold: the highest R^2 excluding the upper bound GPR-GOLD.

Method	auto	Bike	concrete	housing	yacht	ct	Average
GPR-GOLD	0.8698 ± 0.0268	0.5612 ± 0.0059	0.8101 ± 0.0251	0.8288 ± 0.0389	0.8043 ± 0.0741	0.8652 ± 0.0054	0.7899
GPR-Av	0.8480 ± 0.0348	0.5479 ± 0.0069	0.7471 ± 0.1452	0.4972 ± 0.4034	0.7368 ± 0.1665	0.8261 ± 0.0079	0.7005
MA-LFCR	0.7975 ± 0.0230	0.3880 ± 0.0072	0.5965 ± 0.0369	0.7033 ± 0.0570	0.6065 ± 0.0834	0.6605 ± 0.2503	0.6254
MA-GPR	0.8558 ± 0.0269	0.4242 ± 0.0175	0.7807 ± 0.0312	0.7247 ± 0.0628	0.7319 ± 0.0965	0.0111 ± 0.0038	0.5881
MA-DL-B	0.7814 ± 0.0326	0.5914 ± 0.0131	0.2228 ± 0.0361	0.5351 ± 0.0781	0.1924 ± 0.0767	0.5477 ± 0.2991	0.4785
MA-DL-S	0.7835 ± 0.0302	0.5871 ± 0.0127	0.2186 ± 0.0328	0.5284 ± 0.0867	0.1809 ± 0.0976	0.9402 ± 0.0284	0.5398
MA-DL-B+S	0.7876 ± 0.0285	0.5877 ± 0.0116	0.2270 ± 0.0344	0.5339 ± 0.0881	0.1921 ± 0.0830	0.6431 ± 0.2999	0.4952
KAAR	0.8499 ± 0.0365	0.5672 ± 0.0056	0.7666 ± 0.0290	0.8037 ± 0.0448	0.5875 ± 0.0999	0.9051 ± 0.0023	0.7467
LKAAR	0.8484 ± 0.0377	0.5688 ± 0.0062	0.7986 ± 0.0246	0.8047 ± 0.0443	0.7511 ± 0.0875	0.8944 ± 0.0025	0.7777

Table 3.6: Regression results in terms of R^2 score over *fully real dataset*. Bold: the highest R^2 excluding the upper bound GPR-GOLD.

Method	Music
GPR-GOLD	0.4889
GPR-Av	0.2744
MA-LFCR	0.1404
MA-GPR	0.0090
MA-DL-B	0.2339
MA-DL-S	0.2934
MA-DL-B+S	0.3519
KAAR	0.2816
LKAAR	0.2865

the obtained performances in terms of the R^2 score. We remark that the DL-based approaches, MA-DL-B+S and MA-DL-S obtain the best generalization performance in terms of the R^2 score, followed by our LKAAR. Further, as theoretically expected, such performance lies between that of GPR-GOLD and GP-Av. Moreover, regarding the GPs-based competitor MA-GPR, we note that it exhibits the worst prediction capability with a R^2 close to zero. We argue that the above is a symptom of overfitting, as confirmed in Chapter 2. Besides, MA-LFCR presents the second-lowest performance, even worsen than the theoretical lower bound GPR-Av that suggests a non-linear structure in the music dataset. Finally, we notice a considerable reduction in the R^2 score compared with previous results (see Table 3.5), which is due to we are using regression schemes to solve a multi-class problem.

3.5 Summary

In this second proposal, we expose a localized kernel-alignment-based annotator relevance approach, named LKAAR, to support binary and multi-class classification problems in the presence of multiple annotators. Our LKAAR computes the relevance of each provided expert through a centered kernel alignment-based matching between the annotator labels and the input features, taking into account dependencies among the annotators and considering the labelers' performance as a function of the input data. Then, a combination of classifiers/regressors is carried out by adopting the multiple annotator performances coded in the LKAAR-based relevance analysis. We tested our approach in synthetic, semi-synthetic, and real-world datasets. For the synthetic experiments, we generate both the input data and the labels from multiple annotators. Similarly, for the semi-synthetic scenario, we use some databases from the UCI repository. Regarding the real-world datasets, we include three datasets, where the annotations are obtained from the opinion of multiple experts (voice quality problem) and using the crowd-sourcing platform AMT (Polarity data and Music). The results show that the proposed method deals with binary, multi-class classification, and regression problems, where the ground truth is not available. In fact, in most cases, our LKAAR achieves competitive or even better results when compared to different state-of-the-art classification and regression models (Gil-Gonzalez et al., 2018; Groot et al., 2011; Raykar et al., 2010; Rodrigues et al., 2013, 2014b; Yan et al., 2014). Besides, we empirically demonstrated that LKAAR offers a better coding of each annotator's behavior. The latter is also preserved for non-stationary labels across the input space (see Figure 3.1). More-

over, from Table 3.3 and Table 3.4, we can observe that our approach is not severely affected by low labelers performances. Besides, from Figure 3.3, we found that LKAAR exhibits a suitable generalization performance even after varying the number of annotators and the labeled samples, which is a remarkable output since it indicates that our proposal extracts relevant information from the labels, even if they are scarce.

Notwithstanding, like KAAR, our LKAAR solves the supervised learning problem as a convex combination of classifiers/regressors (one per labeler), which can be problematic if the number of labelers grows. Thereby, future work must be oriented to develop an approach that jointly estimates the annotators' performance while training the supervised learning algorithm.

Chapter 4

Regularized chained deep neural network for multiple annotators

This chapter introduces a Regularized Chained Deep Neural Network for Multiple Annotators, termed RCDNN, to jointly estimate the ground truth label and the annotators' performance. RCDNN is inspired in the Chained Gaussian Processes model-(CGP) (Saul et al., 2016), where each parameter in a given likelihood is coded with multiple independent Gaussian Processes-(GPs) priors (one GP prior per parameter). Unlike CGP, our method considers that the last layer models the parameters of an arbitrary likelihood. Thus, in a multi-labeler scenario, the annotators' parameters are coded as a function of the input space. Moreover, since each output in a deep model is computed as a linear combination of previous layers' outputs, our RCDNN can code interdependencies among the annotators. Besides, l1, l2, and Monte-Carlo Dropout-based regularizers are coupled within our method to deal with the overfitting issue in deep learning models. Our proposal follows the line of the works in Albarqouni et al. (2016); Rodrigues and Pereira (2018a) in the sense that RCDNN uses a deep-based approach to build a supervised learning model in the context of multiple annotators. However, while such approaches code the annotators' parameters as fixed points, we model them as functions to consider dependencies between the input features and the labelers' behavior. RCDNN is also similar to the LKAAR model introduced in Chapter 3. Both approaches model the annotators' performance as a function of the input instances and consider interdependencies among the labelers. Nonetheless, unlike LKAAR, where it is necessary to use as many classifiers as annotators, our approach only needs to train a single classifier from a deep learning representation, which is advantageous for a large number of labelers.

4.1 Chained Deep Neural Network

Let us consider an input-output dataset $\mathcal{D} = \{\mathbf{X} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}\}$, where $\mathbf{X} = \{\mathbf{x}_n \in \mathcal{X} \subseteq \mathbb{R}^P\}_{n=1}^N$ and $\mathbf{y} = \{y_n \in \mathcal{Y}\}_{n=1}^N$ hold the input and output spaces, respectively (with N instances and P features). Inspired by the Chained Gaussian Processes model–(CGP) (Saul et al., 2016), a likelihood function with J parameters is written as:

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N p(y_n|\theta_1(\mathbf{x}_n), \dots, \theta_J(\mathbf{x}_n)), \quad (4.1)$$

where $\boldsymbol{\theta} = [\boldsymbol{\theta}_1 \dots, \boldsymbol{\theta}_J]^\top \in \mathbb{R}^{NJ}$ is a parameter vector, and $\boldsymbol{\theta}_j = [\theta_j(\mathbf{x}_1) \dots \theta_j(\mathbf{x}_N)]^\top \in \mathbb{R}^N$. Here, each $\theta_j(\mathbf{x}) \in \mathcal{M}_j$ maps an input sample to the parameter space, being \mathcal{M}_j the domain for the j -th parameter ($j \in \{1, 2, \dots, J\}$). A Chained Deep Neural Network–(CDNN) can be introduced linking each likelihood parameter $\theta_j(\mathbf{x})$ to one of the J outputs of a deep neural network comprising S hidden layers. Accordingly, let $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_J(\mathbf{x})]^\top \in \mathbb{R}^J$ be a vector containing the J outputs of a deep network:

$$\mathbf{f}(\mathbf{x}) = (\varrho_S \circ \varrho_{S-1} \circ \dots \circ \varrho_1)(\mathbf{x}), \quad (4.2)$$

where \circ stands for function composition. Then, each parameter is computed as: $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$, where $h_j: \mathbb{R} \rightarrow \mathcal{M}_j$ is a deterministic function that maps each output $f_j(\mathbf{x})$ to the appropriate domain \mathcal{M}_j . Besides, each layer ϱ_s , with $s \in \{1, 2, \dots, S\}$, depends on a set of variables (neural network weights and bias) $\boldsymbol{\phi} = [\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_S]^\top$, which can be estimated by minimizing the following log likelihood cost (for i.i.d samples):

$$-\log(p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\phi})) = -\sum_{n=1}^N \log(p(y_n|\theta_1(\mathbf{x}_n), \dots, \theta_J(\mathbf{x}_n), \boldsymbol{\phi})). \quad (4.3)$$

Remarkably, the deep model in equation (4.2) allows exploiting the representation learning capability of neural networks within a chained framework through the likelihood in equation (4.3).

4.2 Regularized Chained Deep Neural Network for classification with multiple annotators

We follow the model proposed by authors in [Rodrigues et al. \(2013\)](#); here, a Regularized Chained Deep Neural Network-(RCDNN) is introduced for classification tasks from multiple annotations. Concerning this, let $\lambda_n^{(r)} \in \{0, 1\}$ be a binary variable representing the r -th annotator reliability: $\lambda_n^{(r)} = 1$ if $y_n^{(r)} = y_n$, and $\lambda_n^{(r)} = 0$ in other case. If $\lambda_n^{(r)} = 1$, the label $y_n^{(r)}$ is modeled by means of a categorical distribution; otherwise, if $\lambda_n^{(r)} = 0$, $y_n^{(r)}$ is supposed to follow an uniform distribution. In consequence, the likelihood function in equation (4.3) is rewritten as:

$$p(\mathbf{Y}|\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{r \in R_n} \left(\prod_{k=1}^K \zeta_{n,k}^{\delta(y_n^{(r)} - k)} \right)^{\lambda_n^{(r)}} \left(\frac{1}{K} \right)^{1 - \lambda_n^{(r)}}, \quad (4.4)$$

where $\delta(y_n^{(r)} - k) = 1$, if $y_n^{(r)} = k$, and $\delta(y_n^{(r)} - k) = 0$, otherwise. Moreover, $\zeta_{n,k} = p(y_n^{(r)} = k | \lambda_n^{(r)} = 1)$ is the estimation of the hidden ground truth for the n -th instance in class k .

Accordingly, an architecture holding $J = K + R$ outputs is required within our RCDNN for modeling the likelihood parameters $\boldsymbol{\theta}$ in equation (4.4). In particular, K output layers are fixed to estimate the hidden ground truth $\zeta_{n,k}$ based on a softmax function as follows:

$$\zeta_{n,k} = \frac{\exp(f_k(\mathbf{x}_n))}{\sum_{i=1}^K \exp(f_i(\mathbf{x}_n))}. \quad (4.5)$$

Furthermore, a step function can be used to compute the annotator reliability. Yet, the step function is approximated through R output layers $\{\varsigma_r(\cdot)\}_{r=1}^R$, fixing the well-known sigmoid activation to avoid discontinuities and favor the RCDNN implementation:

$$\lambda_n^{(r)} = \varsigma(f_{l_r}(\mathbf{x}_n)) = \frac{1}{1 + \exp(f_{l_r}(\mathbf{x}_n))}, \quad (4.6)$$

where $l_r = K + r \in \{K + 1, \dots, J\}$ is the index of the output linked to the estimation of the reliability of r -th expert. Afterward, the log likelihood of equation (4.4) is used to compute the RCDNN weights and bias in $\boldsymbol{\phi}$, as follows:

$$\phi^* = \arg \min_{\phi} - \sum_{n=1}^N \sum_{r \in R_n} \left[\lambda_n^{(r)}(\phi) \left(\sum_{k=1}^K \delta(y_n^{(r)} - k) \log(\zeta_{n,k}(\phi)) \right) - (1 - \lambda_n^{(r)}(\phi)) \log(K) \right], \quad (4.7)$$

where $\lambda_n^{(r)}(\phi)$ and $\zeta_{n,k}(\phi)$ highlight the dependency between the annotator reliability/ground truth estimation and the RCDNN weights and bias.

In turn, to avoid overfitting and favor the RCDNN generalization capability, l1 and l2 norm-based regularizers are used for dense layers; besides, Dropout layers are also added. We notice that both regularization schemes are implemented through the function composition presented in equation (4.2). Lastly, to exploit the RCDNN generalization, the well-known Monte-Carlo Dropout prediction strategy is used to estimate the expert's reliability $\hat{\lambda}_n^{(r)}$ and the ground truth label $\hat{\zeta}_{n,k}$, as follows (Géron, 2019):

$$\hat{\lambda}_n^{(r)} = \frac{1}{E} \sum_{e=1}^E \lambda_n^{(r)}(\phi^*, \Delta_e), \quad (4.8)$$

$$\hat{\zeta}_{n,k}^{(r)} = \frac{1}{E} \sum_{e=1}^E \zeta_{n,k}^{(r)}(\phi^*, \Delta_e); \quad (4.9)$$

where notation $\lambda_n^{(r)}(\phi^*, \Delta_e)$ and $\zeta_{n,k}^{(r)}(\phi^*, \Delta_e)$ stands for the dependency between the estimated output, the trained RCDNN weights and bias based on equation (4.7), and the set Δ_e holding Dropout layers. As seen, the Monte-Carlo Dropout-based predictions in equations (4.8) and (4.9) compute the RCDNN outputs as the sample mean over a stack of E predictions; each of them activates randomly the Dropout layers in Δ_e for the u -th iteration within a Monte-Carlo scheme. For RCDNN's implementation details see Section 4.3.4.

4.3 Experimental set-up

4.3.1 Tested datasets

The introduced RCDNN classifier for multiple annotators scenarios is tested in three different kind of datasets. The first category, termed *2D-PCA iris dataset*, is intended to show graphically how our method works. The Principal Component Analysis-(PCA) algorithm is applied to reduce the well-known Iris dataset dimension from four to

two (Géron, 2019), aiming to easily observe some preliminary results in a cartesian plane and illustrate how multiple annotations can be simulated. Second, we build *semi-synthetic data* from eight datasets devoted to binary and multi class-classification of the well-known UCI repository (see Section 2.3.1). Moreover, we test the publicly available bearing data collected by the Case Western Reserve University–(Western). The aim is to build a system to diagnose an electric motor’s status based on two accelerometers. The feature extraction was performed as in Hernández-Muriel et al. (2020) ($P = 7$, $N = 3413$, $K = 4$). Finally, our proposal is tested on three *fully real datasets* (Voice, Polarity, and Music), where both, the input features and the annotations are captured from real-world problems. The description for three dataset is found in Section 2.3.1; besides, Table 2.1 summarizes the tested datasets.

4.3.2 Provided and simulated annotations

Since the *semi-synthetic* datasets do not provide annotations from multiple labelers, to test our RCDNN classifier, it is necessary to simulate those annotations based on the ground truth, which is available for this kind of experiments. Taking into account that our approach is built under the consideration that the annotators’ performance depends on the input space, we use the schemes *Non-homogeneous labels*, and *Biased coin (Non-homogeneous)* presented in Section 3.3.1.

Regarding the voice quality dataset, the annotations from four experts are provided, $R = 4$. However, for concrete testing, only the G, R, and B scales are studied. Indeed, for scales A and S, the sources’ expertise are not satisfactory (Gil et al., 2015). Similarly, for the polarity sentiment dataset, labels from 203 workers are available. Annotators who labeled at least 15% of the available instances are kept, yielding $R = 7$ labelers. Finally, concerning the music dataset, 2946 labels were obtained from 44 instances. Nevertheless, in our experiments, the sources that labeled at least the 15% of the available instances are studied ($R = 9$).

4.3.3 Method comparison and quality assessment

Our model’s validation is carried out by estimating the classification performance as the Area Under the Curve (AUC) and the overall accuracy (Acc). The AUC is extended for multi-class scenarios, as discussed in Fawcett (2006). A cross-validation scheme is used with 30 repetitions, where 70% of the samples are utilized for training and the remaining 30% for testing, except for the music and polarity dataset since they clearly define

Table 4.1: A short overview of the tested state-of-the-art approaches. GPC: Gaussian Processes classifier, LRC: logistic regression classifier, MV: majority voting, MA: multiple annotators, MAE: Modelling annotators expertise, LFC: Learning from crowds, DGRL: Distinguishing good from random labelers, KAAR: kernel alignment-based annotator relevance analysis, LKAAR: localized kernel alignment-based annotator relevance analysis.

Approach	Brief description
GPC-GOLD	A GPC using the real labels (upper bound).
GPC-MV	A GPC using the majority voting of the labels as the ground truth.
MA-LFC (Raykar et al., 2010)	A LRC with constant parameters across the input space.
MA-DGRL (Rodrigues et al., 2013)	A multi-labeler approach that considers as latent variables the annotator performance.
MA-MAE (Yan et al., 2010)	A LRC where the sources parameters depend on the input space.
MA-GPC (Rodrigues et al., 2014b)	A multi-labeler GPC, which is an extension of MA-LFC by using a non-linear approach.
KAAR (Gil-Gonzalez et al., 2018)	A kernel-based approach that employs a convex combination of GPC, it codes the labelers dependencies.
LKAAR-(LR,SVM,GPC) (Gil-Gonzalez et al., 2018)	A localized kernel alignment-based annotator relevance analysis using a combination of LRC, SVM, GPC respectively. It models both the annotators dependencies and the relationship between the labelers' behavior and the input features.

the training and testing sets. Table 4.1 displays the state-of-the-art models that are considered for comparison purposes. The Matlab codes for the state-of-the-art methods studied are publicly available ¹. We highlight that the GPC-Gold is used only to provide an upper bound for our approach.

4.3.4 RCDNN detailed architecture and training

The proposed RCDNN architecture for multiple annotators comprises:

- **IN**: An input layer fed by the input samples $\mathbf{X} \in \mathbb{R}^{N \times P}$.
- $\varrho_1(\cdot)$: A dense layer coding relevant patterns from input features to perform. The number of neurons is set as $h = \lfloor \rho P \rfloor$, where $\rho \in \{0.5, 1, 1.5\}$ is chosen empirically; a linear-based activation function is used to code input data linear dependencies.
- $\varrho_2(\cdot)$: A dense layer fixing a *tanh*-based activation function with $J = K + R$ neurons to reveal non-linear relationships.
- $\varrho_3(\cdot)$: A fully-connected layer with K neurons and a *softmax*-based activation function, which is employed to estimate the hidden ground truth $\zeta_{k,n}$.
- $\varrho_4(\cdot)$: A dense layer with R neurons and a *sigmoid*-based activation function, which is used to compute the annotators' reliability in $\lambda_n^{(r)}$.

¹GPC-MV MA-LFC, MA-MAE, MA-DGRL, GPC-GTIC, KAAR, and LKAAR codes: <https://github.com/juliangilg>. MA-GPC codes: <http://www.fprodrigues.com/>

- For all provided ϱ_s layers l1 plus l2-based regularization strategy is used, searching the regularization weights within the range $\{1e-3, 1e-2, 1e-1\}$.
- BatchNormalization and Dropout layers are included between layers to avoid vanishing and exploding gradient issues. Also, it favors the RCDNN’s generalization capability as exposed in Section 4.2. See Figure 4.1 for details.
- The optimization problem in equation (4.7) is solved by using a Back-propagation algorithm as usual. Moreover, to favor scalability, we utilize a mini-batch-based gradient descent approach with automatic differentiation (RMSprop-based optimizer is fixed).

We clarify that our RCDNN is flexible, and it admits different deep structures such as Recurrent or Convolutional layers aiming to deal with complex tasks (e.g., computer vision or natural language processing). Moreover, our approach can build from different activation functions (RELU, ELU, sigmoid, softmax). However, the last layers (in this particular case ϵ_3 and ϵ_4) need to be designed to code each annotator’s behavior and the hidden ground truth. For example, the parameter $\lambda_n^{(r)}$ represents an estimation for the annotators’ reliability; accordingly, we need to use an activation function whose output belongs to the range $[0, 1]$.

4.4 Results and discussion

A controlled experiment is performed to estimate the performance of inconsistent labelers as a function of the input space while highlighting their dependencies. For this first experiment, the 2D PCA Iris dataset is employed (see Section 4.3.1). Besides, the data is divided into five clusters using the K-means technique to emulate five annotators using the approach “Biased coin (Non-homogeneous)”. A matrix $\mathbf{P} \in [0, 1]^{R \times R}$ is used to set a different score (annotator reliability) for each pair annotator-cluster, as follows:

$$\mathbf{P} = \begin{pmatrix} 0 & 0.9 & 0.5 & 0.15 & 0.6 \\ 0.9 & 0 & 0.3 & 0.4 & 0.75 \\ 0.5 & 0.3 & 0 & 0.6 & 0.3 \\ 0.15 & 0.4 & 0.6 & 0 & 0.8 \\ 0.6 & 0.75 & 0.3 & 0.8 & 0 \end{pmatrix}. \quad (4.10)$$

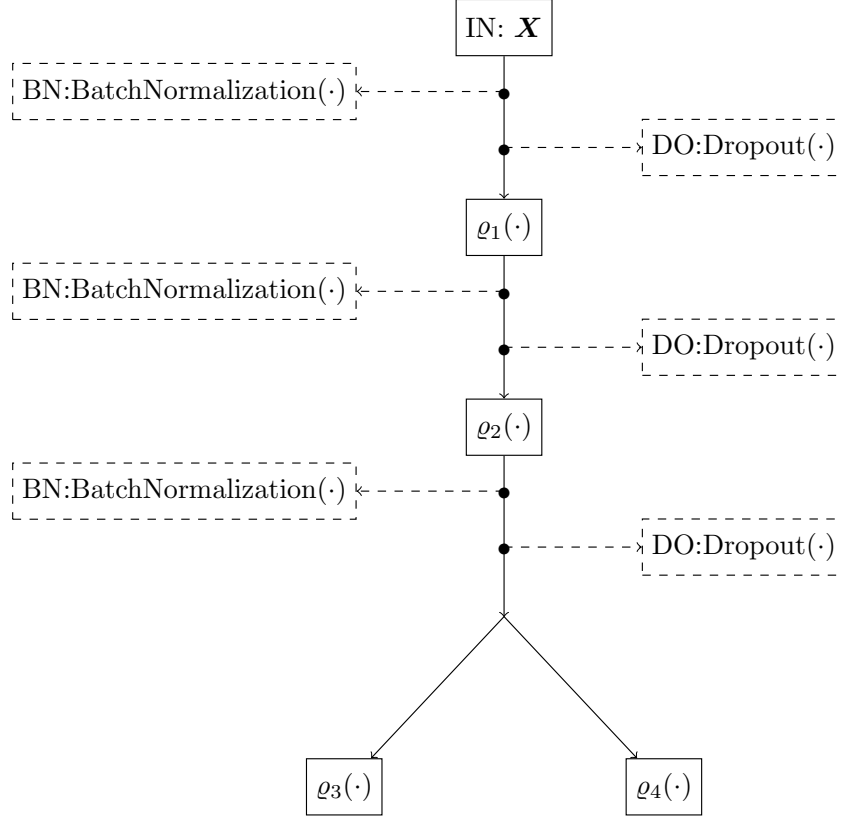


Figure 4.1: RCDNN architecture details. q_s stands for dense layer. q_1 holds a linear activation, q_2 includes a tanh-based activation, and q_3 and q_4 output the hidden ground truth label and the annotator’s reliability fixing a softmax and a sigmoid activation, respectively.

Note that the value $p_{c,r}$ refers to the probability that the annotator r fails labeling a sample that belongs to the cluster c ; thus, a zero-value means a perfect annotator for the correspondent cluster. The r -th annotator is an expert (its labels correspond to the ground truth) in the region $c = r$.

Figure 4.2 shows the decision boundaries generated by our approach for the first experiment. As shown, RCDNN offers a suitable representation for the multi-class classification problem; an AUC score of 0.9837 is achieved, which demonstrates its generalization capability, even in cases where the ground truth is unknown. Indeed, RCDNN codes both the relationship between the input space and the annotator’s behavior and the dependencies among their labels, which improves the quality of the expert codification (Gil-Gonzalez et al., 2018; Zhu et al., 2019). To empirically support the above statement, Figure 4.3 shows each annotator’s simulated accuracy and the reliability estimated by our RCDNN. The latter elucidates how our method performs a satisfactory

identification of the zones where the labelers have the best accuracy. The above is not unexpected because the annotators' accuracy (simulated) is compared with their reliability (estimated); hence, the regions where a specific labeler obtains the higher accuracy should match the regions where the estimated reliability is closer to 1.

In addition, Figure 4.4 shows a comparison between the Pearson correlation coefficients (absolute value) from the labelers' performance in equation (4.10), configuring the simulated dependencies among the annotators, and the Pearson correlation coefficients (absolute value) from the weight matrix $\Phi_{\varrho_4} \in \mathbb{R}^{(K+R) \times R}$ of the layer $\varrho_4(\cdot)$ (RCDNN annotators' dependencies estimation). Comparing the real and the estimated dependencies, it is noticeable that, even though the exact matrix is not recovered, our approach efficiently finds tendencies between annotators' performances. Thereby, the learned representation from hidden layers (see Figure 4.1) allows coding both linear and non-linear patterns that recover the expert dependencies from data. Then, our deep model estimates the unknown ground truth and the relationships between annotators.

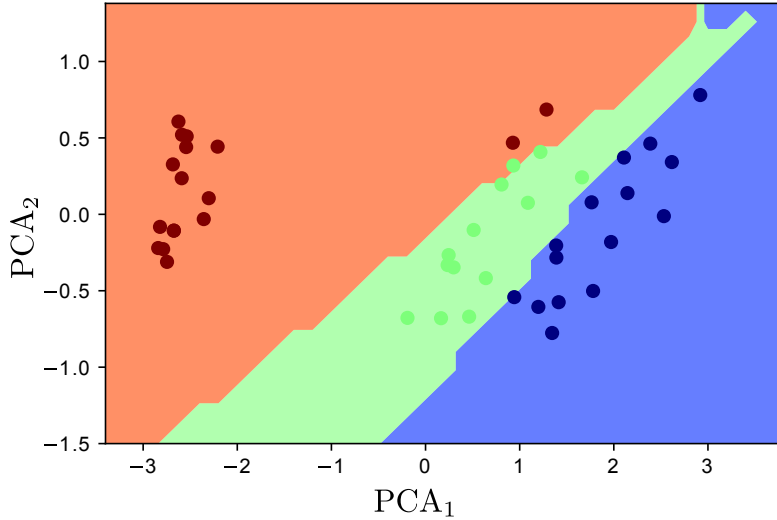


Figure 4.2: RCDNN's decision boundaries for the 2D-PCA Iris dataset (synthetic scenario). AUC= 0.9837. The point's color stands for the Iris dataset classes. PCA_1 and PCA_2 stand for the first and second PCA-based projections.

Table 4.2(a) shows the results concerning the “Non-homogeneous labels”, where it is supposed that the labelers' performance depends on the input space \mathcal{X} . First, we notice that most of the classification schemes present a considerably high performance for both AUC and Acc; in fact, the average AUC and Acc for all methods (except MA-DGRL and MA-MAE) are similar compared to the upper bound GPC-GOLD. The above behavior demonstrates high-quality labels, which is confirmed considering the performance of the

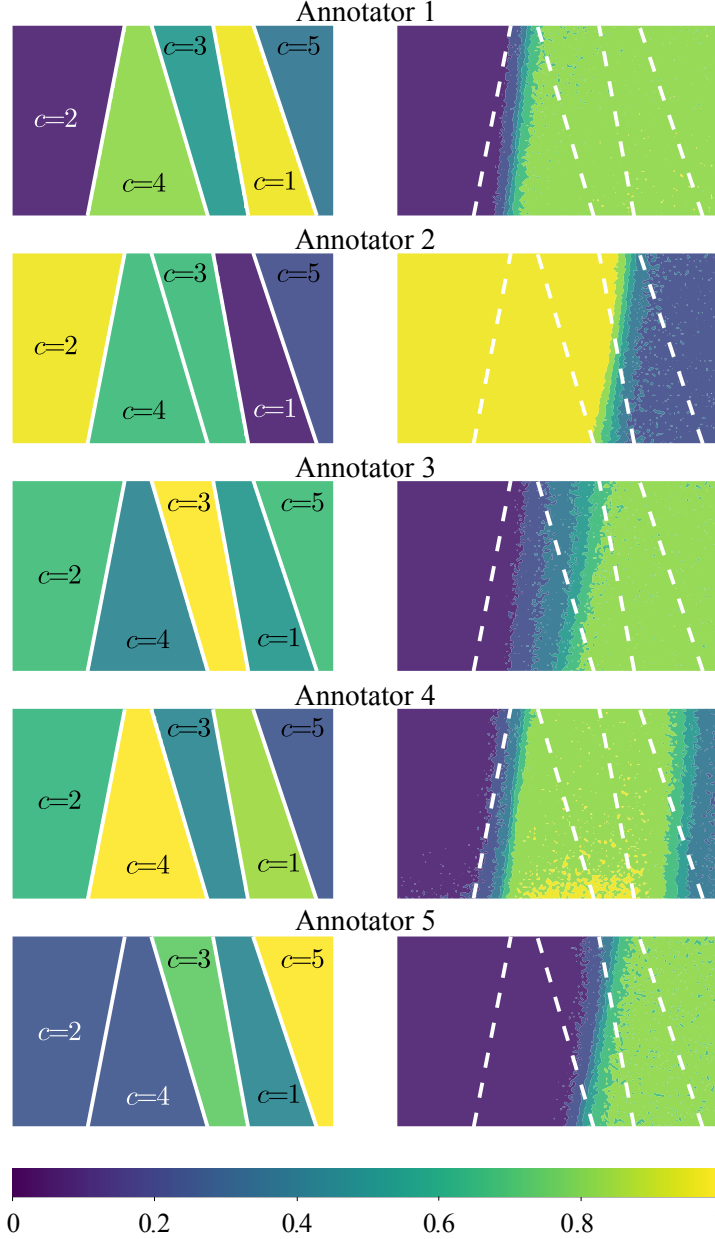


Figure 4.3: RCDNN-based annotators performance (reliability) estimation for the synthetic experiments (2D PCA Iris data). In the first column (from top to bottom), the simulated accuracy for each annotator is presented based on equation (4.10). The second column shows (from top to bottom) the estimated annotators' reliability (λ_r).

most naive approach GPC-MV. Furthermore, we highlight that our RCDNN presents the best average ranking and the second AUC and Acc scores. Then, from non-linear-based approaches, we notice that a naive approach, as GPC-MV, obtains similar performance compared with sophisticated ones, such as KAAR, LKAAR-SVM, and LKAAR-GPC.

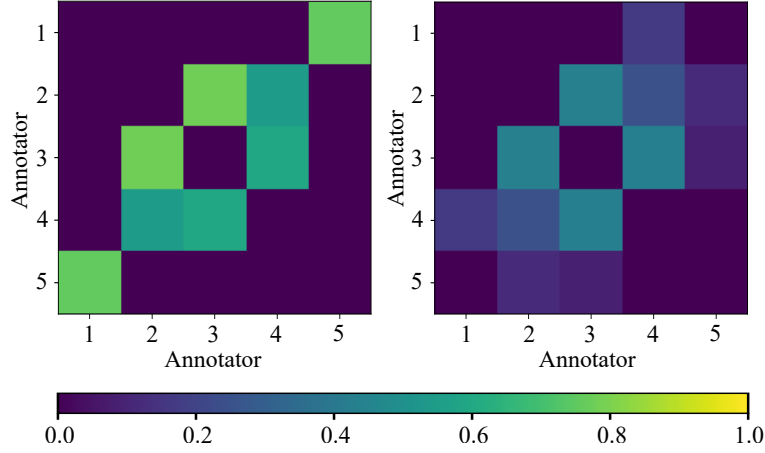


Figure 4.4: Target and estimated annotators dependencies for the synthetic 2D PCA Iris dataset. In the left, the Pearson correlation coefficients (absolute value) from simulated accuracies (experts reliability) in matrix \mathbf{P} of equation (4.10) are shown. In the right, the dependencies among the annotators estimated from the RCDNN q_4 layer’ weights are displayed.

Nevertheless, as we already comment, such an outcome is a consequence of simulating annotators with suitable quality, which favors the majority voting method. Besides, MA-GPC presents the lowest average ranking compared with its other non-linear methods, resulting from a lack of generalization (overfitting). Regarding the results for the linear models, they achieve lower performance than non-linear ones.

On the other hand, Table 4.2(b) shows the results concerning the simulation method “Biased coin (Non-homogeneous)”. At first sight, there exists a generalized lower performance compared with previous results in Table 4.2(a). To explain such an outcome, we recall the stimulation parameters \mathbf{P} in equation (4.10), where the element $1 - p_{c,r}$ (column r , row c) indicates the r -th annotator’s performance in region c . Accordingly, taking the average by column to the matrix $1 - \mathbf{P}$, we obtain the annotators accuracy [0.57, 0.53, 0.66, 0.61, 0.51]. We remark that the labelers’ accuracy is considerably low for this experiment, which impacts the algorithms’ performance. RCDNN achieves the best predictive performance in both the overall accuracy and the AUC score; RCDNN also obtains the best average ranking. Moreover, the non-linear competitors KAAR, LKAAR-GPC, and LKAAR-SVM achieve competitive results. However, GPC-MV and MA-GPC offer the lowest classification scores. Regarding GPC-MV, the result is explained because GPC-MV corresponds to the most naive approach. After all, it considers that the whole annotators achieve similar performance. On the other hand, the

MA-GPC achieves a similar performance compared with GPC-MV; such a behavior proves that MA-GPC is more prone to overfitting (Ruiz et al., 2019). Remarkably, simple classifiers as MA-LFC, MA-DGRL, and LKAAR-LR obtain competitive outcomes compared to the non-linear competitors; in fact, all the linear models excepting MA-MAE outperform GPC-MV and MA-GPC. An additional experiment is conducted: an LR-based classifier using the ground truth (following a similar scheme for GPC-GOLD) is trained overall datasets, obtaining an average AUC equal to 87.21 (close enough to the MA-DGRL and LKAAR-LR performances). Accordingly, a linear structure is presented in some of the studied datasets. In turn, MA-MAE obtains the worst generalization performance (even worse than GPC-MV). Such an outcome is a consequence of overfitting, empirically demonstrated in Ruiz et al. (2019). It is noteworthy that RCDNN and LKAAR-GP obtain similar results, which is expected since both approaches compute the annotators’ performance as a function of the input space while taking into account dependencies between the labelers. However, an unexpected result regarding the “tic-tac-toe” dataset arises, where LKAAR-GP far exceeds our approach’s performance. The above outcome is caused by the categorical features in such a dataset, which cannot be modeled with the chosen DNN architecture Figure 4.1. Still, our method can be easily adapted by setting different layers and activation functions. It is worth noting that the previous experiments were done under controlled scenarios using simulated annotations aiming to stress our method and compare its performance with recently developed approaches. In short, RCDNN offers the best advantages among the state-of-the-art models considered in AUC, overall accuracy, and average ranking.

Up to this point, RCDNN unravels the information hidden in noisy annotations (simulated) to estimate the unknown ground truth considering experts’ performance as a function of the input space and dependencies among labelers. However, the following experiments aim to demonstrate how our approach can outperform state-of-the-art methods even for real labelers, e.g., the challenge is higher as the input data and the annotations are obtained from real-world applications. Table 4.3 describes the results achieved using AUC as the metric to compare the state-of-the-art methods in five different real-world datasets.

First, analyzing the voice data, for the scales **G** and **R**, all the approaches give similar AUC values. In fact, for the scale **G**, the GPC-MV attains competitive performance. The latter can be explained in the sense that the annotators exhibit similar conduct for these scales (Gil et al., 2015). Conversely, for **B** scale, a generalized reduction is presented. Looking at RCDNN results for this database, it is noticeable that the achievement is

Table 4.2: Semi-synthetic datasets results. Bold: the highest AUC excluding the upper bound (target) classifier GPC-GOLD. Marked with *: the highest accuracy (Acc) except the upper bound. The last column presents the average ranking for both the AUC score and the overall accuracy (GPC-GOLD is not considered), the best average ranking for AUC is highlighted in bold, and the accuracy is marked with *.

(a) Non-homogeneous labels.

Method		breast	bupa	ionosphere	pima	tic-tac-toe	iris	wine	segmentation	western	Average
GPC-GOLD	AUC[%] Acc[%]	99.04 ± 0.94 96.44 ± 1.54	72.21 ± 3.69 68.48 ± 4.43	95.02 ± 2.44 91.08 ± 2.41	83.76 ± 1.98 76.71 ± 1.96	99.97 ± 0.06 99.16 ± 0.85	97.65 ± 2.71 95.85 ± 3.29	99.22 ± 0.67 96.92 ± 1.44	90.08 ± 1.94 70.68 ± 6.81	94.52 ± 0.57 79.75 ± 0.57	92.39 86.12
GPC-MV	AUC[%] Acc[%]	99.11 ± 0.58 96.29 ± 1.48	70.95 ± 3.90 66.60 ± 4.31	93.14 ± 3.49 87.90 ± 3.26	81.21 ± 2.57 74.87 ± 2.32	87.83 ± 4.11 81.96 ± 3.46	99.63 ± 0.39 95.33 ± 3.05	98.41 ± 1.38 93.96 ± 3.34	91.48 ± 1.48 82.68 ± 5.30	78.14 ± 4.15 63.35 ± 1.68	88.87 82.54
MA-LFC	AUC[%] Acc[%]	98.72 ± 0.93 95.63 ± 1.79	71.53 ± 4.18 69.68 ± 4.20*	82.08 ± 4.79 81.43 ± 4.44	82.29 ± 2.22 76.52 ± 1.91*	61.13 ± 3.28 64.88 ± 2.86	98.75 ± 1.44 94.44 ± 4.62	96.83 ± 1.75 87.74 ± 4.67	99.58 ± 0.11 95.40 ± 0.71*	87.77 ± 0.79 57.21 ± 1.32	86.72 80.32
MA-DGRL	AUC[%] Acc[%]	99.30 ± 0.39 94.63 ± 1.77	68.00 ± 4.09 65.77 ± 3.47	77.60 ± 7.50 81.94 ± 3.42	81.72 ± 2.57 76.45 ± 2.81	61.83 ± 2.80 66.45 ± 2.24	98.78 ± 1.34 94.59 ± 2.96	95.33 ± 3.35 84.91 ± 6.50	98.31 ± 0.32 89.58 ± 0.99	87.67 ± 0.85 60.55 ± 1.26	85.39 79.43
MA-MAE	AUC[%] Acc[%]	99.28 ± 0.60 96.31 ± 1.38	70.82 ± 3.90 66.92 ± 3.35	78.91 ± 6.01 82.25 ± 3.99	81.80 ± 2.57 76.12 ± 2.77	60.35 ± 3.28 65.64 ± 2.37	85.97 ± 2.39 94.81 ± 4.14	98.20 ± 1.33 89.31 ± 5.79	97.27 ± 0.28 92.94 ± 0.76	72.83 ± 0.80 52.41 ± 1.56	82.76 79.63
MA-GPC	AUC[%] Acc[%]	95.81 ± 2.94 96.70 ± 1.37*	49.81 ± 11.72 59.52 ± 4.71	94.46 ± 3.09 82.13 ± 3.32	67.83 ± 4.24 72.77 ± 2.71	81.44 ± 3.81 76.39 ± 2.85	99.15 ± 1.03 94.30 ± 2.90	99.85 ± 0.24 94.34 ± 2.80	99.42 ± 0.14 94.74 ± 0.68	94.14 ± 0.52 78.52 ± 1.11*	86.87 83.26*
KAAR	AUC[%] Acc[%]	98.81 ± 0.66 96.02 ± 1.14	70.20 ± 5.70 65.99 ± 5.44	93.88 ± 3.53 87.52 ± 4.24	81.18 ± 2.93 75.10 ± 2.98	89.55 ± 2.84 81.68 ± 2.41	99.56 ± 0.52 95.56 ± 2.92*	99.53 ± 0.36 96.54 ± 2.11	92.34 ± 1.38 81.11 ± 4.15	81.77 ± 1.02 64.58 ± 1.47	89.64 82.67
LKAAR-LR	AUC[%] Acc[%]	99.34 ± 0.44 96.00 ± 1.46	68.86 ± 5.16 64.17 ± 4.22	87.14 ± 3.38 84.10 ± 3.20	82.04 ± 2.44 75.67 ± 2.15	65.40 ± 3.13 66.96 ± 2.74	96.00 ± 2.50 82.59 ± 6.07	99.21 ± 0.82 94.28 ± 3.19	97.97 ± 0.27 90.02 ± 0.93	83.25 ± 1.22 51.49 ± 2.05	86.57 78.36
LKAAR-SVM	AUC[%] Acc[%]	98.29 ± 0.80 96.36 ± 1.02	64.37 ± 3.36 63.14 ± 3.67	96.98 ± 2.01 92.19 ± 2.43*	77.80 ± 2.28 72.52 ± 2.22	89.82 ± 2.14 80.99 ± 2.81	98.05 ± 1.90 84.44 ± 6.76	99.53 ± 0.47 96.48 ± 2.26	97.89 ± 0.32 91.28 ± 0.93	79.08 ± 0.95 53.73 ± 2.06	89.09 89.79
LKAAR-GPC	AUC[%] Acc[%]	99.00 ± 0.75 96.03 ± 1.32	71.07 ± 5.05 66.92 ± 4.79	93.37 ± 2.91 87.75 ± 3.90	81.23 ± 2.21 75.10 ± 2.65	91.97 ± 2.01 84.09 ± 2.43*	99.57 ± 0.61 95.26 ± 3.29	99.64 ± 0.34 96.54 ± 2.16	92.61 ± 1.73 80.98 ± 3.91	81.37 ± 1.37 65.20 ± 1.72	89.98 82.65
RCDNN (ours)	AUC[%] Acc[%]	99.47 ± 0.33 97.06 ± 1.19*	69.80 ± 6.07 63.69 ± 4.26	92.60 ± 2.80 86.79 ± 2.37	83.25 ± 3.13 76.00 ± 3.10	71.17 ± 3.76 68.06 ± 3.02	99.74 ± 0.26 95.33 ± 2.46	99.90 ± 0.13 97.84 ± 1.86*	99.15 ± 0.19 92.96 ± 1.06	89.61 ± 0.71 66.46 ± 1.82	89.41 82.68

(b) Biased coin (Non-homogeneous) labels.

Method		breast	bupa	ionosphere	pima	tic-tac-toe	iris	wine	segmentation	western	Average
GPC-GOLD	AUC[%] Acc[%]	99.04 ± 0.94 96.44 ± 1.54	72.21 ± 3.69 68.48 ± 4.43	95.02 ± 2.44 91.08 ± 2.41	83.76 ± 1.98 76.71 ± 1.96	99.97 ± 0.06 99.16 ± 0.85	97.65 ± 2.71 95.85 ± 3.29	99.22 ± 0.67 96.92 ± 1.44	90.08 ± 1.94 70.68 ± 6.81	94.52 ± 0.57 79.75 ± 1.28	92.39 86.12
GPC-MV	AUC[%] Acc[%]	90.78 ± 4.28 86.63 ± 2.06	50.47 ± 6.19 48.27 ± 4.84	82.91 ± 6.03 75.65 ± 6.45	70.18 ± 6.29 66.52 ± 5.16	65.91 ± 6.72 64.66 ± 3.64	98.55 ± 1.38 88.81 ± 5.00	97.75 ± 2.04 86.92 ± 5.76	90.18 ± 1.71 79.24 ± 4.99	74.40 ± 4.94 65.04 ± 1.52	80.13 73.53
MA-LFC	AUC[%] Acc[%]	97.99 ± 0.99 96.00 ± 1.70*	59.64 ± 8.08 56.41 ± 8.12	72.66 ± 9.98 69.17 ± 12.53	72.73 ± 3.43 58.10 ± 4.53	52.88 ± 3.13 46.27 ± 3.03	96.72 ± 8.98 92.30 ± 5.18	96.47 ± 2.13 87.55 ± 4.97	99.50 ± 0.15 95.06 ± 0.80*	84.97 ± 0.84 55.17 ± 1.33	81.51 72.89
MA-DGRL	AUC[%] Acc[%]	99.31 ± 0.42 78.08 ± 2.22	61.77 ± 6.17 55.64 ± 4.52	77.83 ± 7.02 71.43 ± 5.15	81.66 ± 2.65 76.90 ± 1.99*	55.70 ± 3.95 60.64 ± 2.33	98.76 ± 1.33 94.37 ± 2.66	95.26 ± 3.30 84.84 ± 6.32	98.32 ± 0.34 89.63 ± 0.89	86.61 ± 1.10 65.61 ± 1.28	83.91 75.24
MA-MAE	AUC[%] Acc[%]	95.22 ± 1.70 87.15 ± 1.85	64.63 ± 9.77 62.34 ± 8.46	64.18 ± 9.17 67.94 ± 7.19	79.94 ± 4.78 75.94 ± 2.69	52.36 ± 4.78 53.33 ± 6.42	93.16 ± 5.08 81.70 ± 11.68	96.25 ± 2.40 86.67 ± 5.15	94.40 ± 1.26 88.38 ± 2.00	61.40 ± 0.93 49.34 ± 4.15	77.95 72.53
MA-GPC	AUC[%] Acc[%]	85.37 ± 5.90 92.55 ± 2.17	40.79 ± 12.30 52.82 ± 6.38	74.52 ± 4.57 69.87 ± 4.41	73.17 ± 3.34 62.42 ± 3.00	61.82 ± 4.51 62.33 ± 2.98	98.71 ± 1.14 93.85 ± 3.49	99.60 ± 0.41 95.09 ± 2.65	99.35 ± 0.14 93.46 ± 0.83	93.09 ± 0.58 76.88 ± 1.19*	80.71 77.70
KAAR	AUC[%] Acc[%]	97.81 ± 0.99 77.19 ± 3.14	56.52 ± 9.13 52.44 ± 7.79	82.20 ± 4.93 72.60 ± 4.80	67.90 ± 3.16 61.20 ± 2.95	75.34 ± 4.70 70.69 ± 3.63	98.75 ± 1.10 90.44 ± 5.48	97.91 ± 1.36 91.45 ± 4.28	91.75 ± 1.41 76.38 ± 5.05	82.30 ± 0.73 64.61 ± 1.36	83.39 73.00
LKAAR-LR	AUC[%] Acc[%]	99.52 ± 0.30 92.47 ± 2.24	66.07 ± 6.14 60.22 ± 5.67*	82.99 ± 5.01 78.92 ± 4.32	80.57 ± 3.31 75.07 ± 2.65	52.32 ± 3.38 55.64 ± 2.77	96.83 ± 2.14 83.41 ± 6.92	99.27 ± 0.68 94.59 ± 3.12	97.87 ± 0.30 89.77 ± 0.99	81.03 ± 0.80 54.80 ± 2.05	84.05 76.10
LKAAR-SVM	AUC[%] Acc[%]	98.37 ± 1.00 87.72 ± 5.17	52.35 ± 6.40 50.96 ± 6.81	88.28 ± 5.13 84.73 ± 4.66*	66.84 ± 3.66 64.81 ± 3.11	73.85 ± 3.43 70.02 ± 2.74	96.22 ± 2.50 74.15 ± 7.90	98.88 ± 0.80 91.82 ± 4.33	97.59 ± 0.34 90.37 ± 1.24	79.19 ± 1.46 55.39 ± 3.03	82.39 74.44
LKAAR-GPC	AUC[%] Acc[%]	98.14 ± 1.04 86.76 ± 4.33	58.36 ± 7.24 54.52 ± 5.27	86.23 ± 4.47 78.26 ± 5.51	73.80 ± 2.83 69.64 ± 3.01	80.02 ± 4.15 74.90 ± 2.99*	99.61 ± 0.61 95.93 ± 3.15*	98.74 ± 0.93 93.84 ± 3.57	92.24 ± 1.80 78.71 ± 4.18	83.35 ± 0.75 66.58 ± 1.19	85.61 77.68
RCDNN (ours)	AUC[%] Acc[%]	99.26 ± 0.42 94.07 ± 2.00	64.16 ± 3.87 58.24 ± 5.13	83.41 ± 6.28 76.70 ± 6.19	82.08 ± 3.27 74.91 ± 3.77	65.31 ± 3.87 65.07 ± 1.17	99.51 ± 0.53 93.33 ± 3.30	99.77 ± 0.22 96.17 ± 2.57*	99.06 ± 0.20 91.28 ± 0.99	87.94 ± 1.03 61.56 ± 5.13	86.72 79.04*

similar among all the scales, which is an exceptional outcome that shows our method's capabilities to detect regions where annotators have superior execution.

In Polarity Dataset, an acceptable RCDNN's performance is attained compared to others. Our approach requires defining several layer weights in the deep model (Figure 4.1) concerning the number of features (P), labelers (R), and classes (K). For this particular dataset, those values are considerably higher: $P = 1200$, $R = 7$, and $K = 2$. Nevertheless, the introduced regularization strategy (l1, l2, plus Monte-Carlo Dropout)

Table 4.3: Fully real-world datasets results. Bold: the method with the highest performance excluding the upper bound (target) classifier GPC-GOLD.

Method	AUC([%])					
	Voice Dataset			Polarity Dataset	Music	Average
	G	R	B			
GPC-GOLD	93.66	93.66	93.66	80.26	92.84	90.81
GPC-MV	90.17	84.73	84.04	71.14	88.79	83.77
MA-LFC	89.99	90.59	87.27	72.06	85.99	85.18
MA-DGRL	85.45	90.14	79.33	56.13	88.32	79.86
MA-MAE	91.08	89.12	80.74	48.73	81.92	78.31
MA-GPC	91.50	91.16	80.81	61.18	82.53	81.43
KAAR	89.85	93.50	89.20	77.46	88.96	87.79
LKAAR-LR	90.39	92.92	88.94	68.28	84.43	84.99
LKAAR-SVM	92.06	93.02	86.98	72.70	89.98	87.70
LKAAR-GPC	90.78	93.60	89.79	76.50	86.44	87.42
RCDNN	92.24	94.19	92.57	76.04	93.29	89.66

allows computing an acceptable AUC performance of 76.04 in comparison with the best achieved by the KAAR method 77.46.

Lastly, in the case of Music data, our RCDNN obtains the best classification performance. On the other hand, MA-MAE and MA-GPC exhibit a significantly low performance, even lower than the intuitive lower bound (GPC-MV). This behavior has been repeated in the previous experiments because of the over-fitting issue. Nevertheless, an additional challenge is presented for the music dataset regarding the multi-class classification setting. Accordingly, a *one-vs-all* scheme is fixed for all of the binary classification methods (including MA-MAE and MA-GPC). Such a scheme to deal with multi-class classification can lead regions on the input space that are ambiguously classified (Bishop, 2006).

As a final experiment, we wish to analyze the impact of spammers and malicious annotators on the performance of our multi-labeler classifier. For concrete testing, we use the pima dataset, which holds 768 instances; from this dataset, we use 538 samples for training and the remaining 230 for testing. We create synthetic labels from 5 annotators generated from the Biased coin (Non-homogeneous) procedure (see Section 4.3.2 and equation (4.10)). According to Figure 4.5 (blue dots), we notice that from the 5 labelers, two are categorized as suitable labelers, one as Spammers and the remaining as Malicious. Then, we add R_e additional annotators aiming to test our approach in extreme scenarios, where the number of malicious or spammers annotators increases. The labels are simulated as follows: a random number $\tau_n^{(r)}$ is sampled from a Bernoulli distribution with parameter p_r ; then if $\tau_n^{(r)} = 0$, $y_n^{(r)} = y_n$, and $y_n^{(r)} = \tilde{y}_n$ otherwise. For

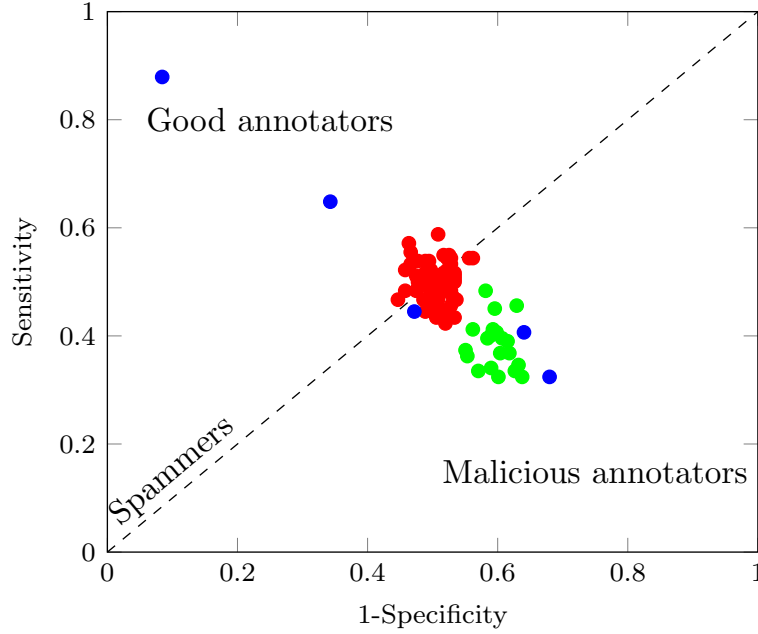


Figure 4.5: Receiver Operating Characteristic (ROC) plot for the annotators simulated within the spammers and malicious scenario. Blue dots indicate the basis annotators. Red dots show extra annotators with parameters $R_e = 65$, and $p_r = 0.5$. Green dots specify extra labelers with $R_e = 20$, and $p_r = 0.6$. We notice that annotators located in dashed line vicinity are considered Spammers. Similarly, labelers above the dashed line are regarded as good annotators; conversely, labelers located below such a line are malicious annotators.

Spammers, we use $R_e = 65$, and $p_r = 0.5$ (see red dots in Figure 4.5); alike, for malicious labelers, we fix $R_e = 20$, and $p_r = 0.6$ (see green dots in Figure 4.5).

Figure 4.6 presents the classifiers' performance as a function of the number of spammers (left in Figure 4.6) and malicious annotators (in right Figure 4.6). First, we analyze the effect of Spammers annotators on the RCDNN's performance. From the results in Figure 4.6 (left), we remark that when the number of Spammers is less than 40, the performance of our approach is not affected. However, when the number of Spammers exceeds 40, the RCDNN' AUC becomes unstable, oscillating between 0.6 and 0.8. Accordingly, we highlight that the critical point is presented when the percentages of good, spammers and malicious labelers are respectively 4.65%, 90.70%, and 4.65%; which shows that our RCDNN is robust in the presence of a high number of Spammers. Now, we compare our RCDNN with two state-of-the-art models, MA-LFC (linear model with the more competitive performance according to Table 4.3) and LKAAR-GPC (Non-linear model with the more competitive AUC in Table 4.3). We notice that the LKAAR-GPC

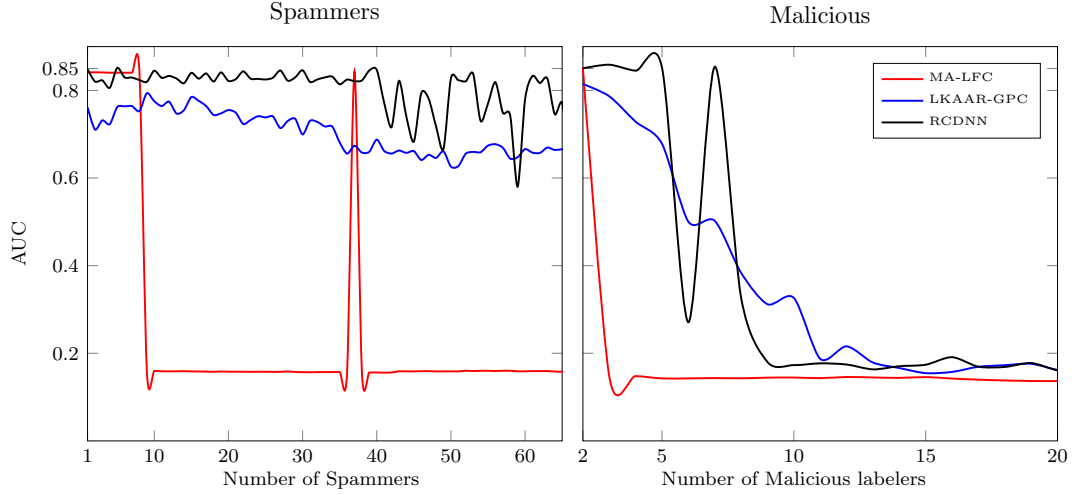


Figure 4.6: MA-LFC, LKAAR-GPC, and RCDNN performance (AUC) as a function of the number of labelers (spammers and malicious annotators).

behavior is similar to our approach when the number of spammers is greater than 35, the AUC starts to descend gradually. Conversely, we note that the MA-LFC’s performance is drastically affected by the spammers; in fact, for more than 8 spammers, the AUC is close to 0.2.

Second, we inspect the consequences when malicious labelers are added. From the results in Figure 4.6 (right), we note that our RCDNN is significantly affected when we have more than 5 malicious annotators; in that case, the AUC decreases from 0.85 approximately to a value near to 0.2. Thereby, we notice that the critical point is presented when the percentages of good, spammers, and malicious labelers are respectively 25%, 12.5%, 62.5%. In such a sense, for this experiment, we can affirm that our approach can deal with malicious labelers if the percentage of them is below 62.5%. Finally, studying the results related to LKAAR-GPC, we notice that LKAAR again performs similar to our RCDNN due to for more than 5 malicious labelers, LKAAR-GPC achieves AUC scores lower than 0.5; on the other hand, MA-LFC is susceptible since, for more than 2 malicious labelers, the AUC decreases to a value near to 0.2.

4.5 Summary

This paper introduces a novel Regularized Chained Deep Neural Network classifier, termed RCDNN, to deal with multiple annotator scenarios. Our method is built based on the ideas of the chained Gaussian Processes (Saul et al., 2016), where each pa-

parameter in a multi-labeler likelihood is modeled by using the outputs of a deep neural network. In such a way, RCDNN codes the annotators' expertise as a function of the input data and the dependencies among the labelers from the last hidden layer's weights. Besides, l_1 , l_2 , and Monte-Carlo Dropout regularization strategies are coupled within our RCDNN architecture and predictor to contract the overfitting challenge of deep models. The proposal is tested using different scenarios concerning the provided annotations: synthetic, semi-synthetic, and real-world experts. According to the results, RCDNN achieves robust predictive properties for the studied datasets, outperforming state-of-the-art methods while providing an estimation of each labeler's reliability and the dependencies among annotators.

However, we notice that RCDNN is based on a frequentist approach; thus, the estimations for the annotators' parameters and the supervised learning framework are deterministic, and it is not possible to capture the uncertainty related to them.

Chapter 5

Correlated chained Gaussian processes for multiple annotators

In this chapter, we propose a probabilistic model, named the correlated chained Gaussian Processes for multiple annotators–(CCGPMA), to jointly build a prediction algorithm applicable to classification and regression tasks. CCGPMA is based on the chained GPs model–(CGP) (Saul et al., 2016), which is a Multi-GPs framework where the parameters of an arbitrary likelihood function are modeled with multiple independent GPs (one GP prior per parameter). Unlike CGP, we consider that multiple correlated GPs model the likelihood’s parameters. For doing so, we take as a basis the ideas from a Multi-output GP–(MOGP) regression (Álvarez et al., 2012), where each output is coded as a weighted sum of shared latent functions via a semi-parametric latent factor model–(SLFM) (Teh et al., 2005). In contrast to the MOGP, we do not have multiple outputs but multiple functions chained to the given likelihood parameters. From the multiple annotators’ point of view, the likelihood parameters are related to the labelers’ behavior; thereby, CCGPMA models the labelers’ behavior as a function of the input features while also taking into account annotators’ interdependencies. Moreover, our proposal is based on the so-called inducing variables framework (Álvarez et al., 2010), in combination with stochastic variational inference (Hoffman et al., 2013). To the best of our knowledge, this is the first attempt to build a probabilistic approach to model the labelers’ behavior as a function of the input features while also considering annotators’ interdependencies. Achieved results, using both simulated and real-world data, show how our method can deal with both regression and classification problems from multi-labelers data.

5.1 Chained Gaussian processes

Let us consider an input-output dataset $\mathcal{D} = \{\mathbf{X} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}\}$, where $\mathbf{X} = \{\mathbf{x}_n \in \mathcal{X} \subseteq \mathbb{R}^P\}_{n=1}^N$ and $\mathbf{y} = \{y_n \in \mathcal{Y}\}_{n=1}^N$. In turn, let a GP be a collection of random variables $f(\mathbf{x})$ indexed by the input samples $\mathbf{x} \in \mathcal{X}$ holding a joint multivariate Gaussian distribution (Rasmussen and Williams, 2006). A GP is defined by its mean $m(\mathbf{x}) = \mathbb{E} f(\mathbf{x})$ (we consider $m(\mathbf{x}) = 0$) and covariance function $\kappa_f(\mathbf{x}, \mathbf{x}') = \mathbb{E} (f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))$, where $\kappa_f: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a given kernel function and $\mathbf{x}' \in \mathcal{X}$, yielding:

$$f(\mathbf{x}) \sim \mathcal{GP}(0, \kappa_f(\mathbf{x}, \mathbf{x}')). \quad (5.1)$$

If we consider the finite set of inputs in \mathbf{X} , then $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top \in \mathbb{R}^N$ is drawn for a multivariate Gaussian distribution $\mathbf{f} \sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_{ff})$, where $\mathbf{K}_{ff} \in \mathbb{R}^{N \times N}$ is the covariance matrix formed by the evaluation of $\kappa_f(\cdot, \cdot)$ over the input set \mathbf{X} .

Accordingly, using GPs for modeling the input-output data collection \mathcal{D} consists of constructing a joint distribution between a given likelihood function and one or multiple GP-based priors. To code each likelihood parameter as a random process, we employ the so-called chained GP-(CGP) that attaches such parameters to multiple independent GP priors, as follows (Saul et al., 2016):

$$p(\mathbf{y}, \hat{\mathbf{f}}|\mathbf{X}) = \prod_{n=1}^N p(y_n|\theta_1(\mathbf{x}_n), \dots, \theta_J(\mathbf{x}_n)) \times \prod_{j=1}^J \mathcal{N}(\mathbf{f}_j|\mathbf{0}, \mathbf{K}_{f_j f_j}), \quad (5.2)$$

where each $\{\theta_j(\mathbf{x}) \in \mathcal{M}_j\}_{j=1}^J$ represents the likelihood's parameters, being $J \in \mathbb{N}$ the number of parameters to represent the likelihood. Besides, each $\theta_j(\mathbf{x})$ holds a non-linear mapping from a GP prior, e.g., $\theta_j(\mathbf{x}) = h_j(f_j(\mathbf{x}))$, where $h_j: \mathbb{R} \rightarrow \mathcal{M}_j$ is a deterministic function that maps each latent function-(LF) $f_j(\mathbf{x})$, to the appropriate domain \mathcal{M}_j . Moreover, $\mathbf{f}_j = [f_j(\mathbf{x}_1), \dots, f_j(\mathbf{x}_N)]^\top \in \mathbb{R}^N$ is a LF vector that follows a GP prior, and $\hat{\mathbf{f}} = [\mathbf{f}_1, \dots, \mathbf{f}_J]^\top \in \mathbb{R}^{NJ}$. $\mathbf{K}_{f_j f_j} \in \mathbb{R}^{N \times N}$ is the covariance matrix belonging to the j -th GP prior, which is computed based on the kernel function $\kappa_j: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. The non-parametric formulation of a GP introduces computational loads through the inference process. For instance, considering that the dataset \mathcal{D} configures a regression problem, a GP modeling involves a computational complexity of $\mathcal{O}(N^3)$ to invert the matrix $\mathbf{K}_{f_j f_j}$ (Rasmussen and Williams, 2006). A common approach to reduce such computational complexity is to augment the GP prior with a set of $M \ll N$ inducing

variables (Hensman et al., 2015) $\mathbf{u}_j = [f_j(\mathbf{z}_1^{(j)}), \dots, f_j(\mathbf{z}_M^{(j)})]^\top \in \mathbb{R}^M$ through additional evaluations of $f_j(\cdot)$ at unknown locations $\mathbf{Z}_j = [\mathbf{z}_1^{(j)}, \dots, \mathbf{z}_M^{(j)}] \in \mathbb{R}^{M \times P}$, which decreases the GP's computational complexity to $\mathcal{O}(NM^2)$. Further, the following augmented GP prior arises:

$$p(\mathbf{f}_j, \mathbf{u}_j) = \mathcal{N} \left(\begin{bmatrix} \mathbf{f}_j \\ \mathbf{u}_j \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_{\mathbf{f}_j \mathbf{f}_j} & \mathbf{K}_{\mathbf{f}_j \mathbf{u}_j} \\ \mathbf{K}_{\mathbf{u}_j \mathbf{f}_j} & \mathbf{K}_{\mathbf{u}_j \mathbf{u}_j} \end{bmatrix} \right), \quad (5.3)$$

where $\mathbf{K}_{\mathbf{f}_j \mathbf{u}_j} \in \mathbb{R}^{N \times M}$ is the cross-covariance matrix formed by the evaluation of the kernel function $\kappa_j(\cdot, \cdot)$ between \mathbf{X} and \mathbf{Z}_j . Likewise, $\mathbf{K}_{\mathbf{u}_j \mathbf{u}_j} \in \mathbb{R}^{M \times M}$ is the inducing points-based covariance matrix. Then, the distribution of \mathbf{f}_j conditioned to the inducing points \mathbf{u}_j can be written as:

$$p(\mathbf{f}_j | \mathbf{u}_j) = \mathcal{N} \left(\mathbf{f}_j | \mathbf{K}_{\mathbf{f}_j \mathbf{u}_j} \mathbf{K}_{\mathbf{u}_j \mathbf{u}_j}^{-1} \mathbf{u}_j, \mathbf{K}_{\mathbf{f}_j \mathbf{f}_j} - \mathbf{K}_{\mathbf{f}_j \mathbf{u}_j} \mathbf{K}_{\mathbf{u}_j \mathbf{u}_j}^{-1} \mathbf{K}_{\mathbf{u}_j \mathbf{f}_j} \right), \quad (5.4)$$

$$p(\mathbf{u}_j) = \mathcal{N} \left(\mathbf{u}_j | \mathbf{0}, \mathbf{K}_{\mathbf{u}_j \mathbf{u}_j} \right). \quad (5.5)$$

In most cases equations (5.4) and (5.5) are non-conjugate to the likelihood, finding the posterior distribution $p(\mathbf{f}, \mathbf{u} | \mathbf{y})$ is not tractable analytically; therefore, we resort to a deterministic approximation of the posterior distribution using variational inference. Hence, the actual posterior can be approximated by a parametrized variational distribution $p(\hat{\mathbf{f}}, \mathbf{u} | \mathbf{y}) \approx q(\hat{\mathbf{f}}, \mathbf{u})$ (Titsias, 2009), as:

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) = \prod_{j=1}^J p(\mathbf{f}_j | \mathbf{u}_j) q(\mathbf{u}_j), \quad (5.6)$$

where $\mathbf{u} = [\mathbf{u}_1^\top, \dots, \mathbf{u}_J^\top]^\top \in \mathbb{R}^{MJ}$; moreover, $p(\mathbf{f}_j | \mathbf{u}_j)$ is defined in equation (5.4), and $q(\mathbf{u})$ is the posterior approximation over the inducing variables:

$$q(\mathbf{u}) = \prod_{j=1}^J q(\mathbf{u}_j) = \prod_{j=1}^J \mathcal{N}(\mathbf{u}_j | \mathbf{m}_j, \mathbf{V}_j). \quad (5.7)$$

The approximation for the posterior distribution comprises the estimation of the following variational parameters: the mean vectors $\mathbf{m}_j \in \mathbb{R}^M$ and the covariance matrices

$\mathbf{V}_j \in \mathbb{R}^{M \times M}$. Such an assessment is carried out by maximizing an evidence lower bound–(ELBO). Due to space restrictions, the ELBO derivation and details are included in the supplementary material.

5.1.1 Correlated chained Gaussian processes

From Section 5.1, we note that the CGP model assumes independence between priors, thereby lacking a correlation structure between GPs. As mentioned before, we consider that the annotators are correlated. We will enable this aspect of the model by assuming dependencies among the latent parameters of the chained GP. In particular, we introduce the correlated chained GPs–(CCGP) to model correlations between the GP latent functions, which are supposed to be generated from a semi-parametric latent factor model–(SLFM) (Teh et al., 2005):

$$f_j(\mathbf{x}_n) = \sum_{q=1}^Q w_{j,q} \mu_q(\mathbf{x}_n), \quad (5.8)$$

where $f_j : \mathcal{X} \rightarrow \mathbb{R}$ is an LF, $\mu_q(\cdot) \sim \mathcal{GP}(0, k_q(\cdot, \cdot))$ with $k_q : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ being a kernel function, and $w_{j,q} \in \mathbb{R}$ is a combination coefficient ($Q \in \mathbb{N}$). Here, each LF is chained to the likelihood’s parameters to extend the joint distribution in equation (5.2) as follows:

$$p(\mathbf{y}, \hat{\mathbf{f}}, \mathbf{u} | \mathbf{X}) = p(\mathbf{y} | \boldsymbol{\theta}) \prod_{j=1}^J p(\mathbf{f}_j | \mathbf{u}) p(\mathbf{u}), \quad (5.9)$$

where $\boldsymbol{\theta} = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_J]^\top \in \mathbb{R}^{NJ}$ holds the model’s parameters and $\boldsymbol{\theta}_j = [\theta_j(\mathbf{x}_1), \dots, \theta_j(\mathbf{x}_N)]^\top \in \mathbb{R}^N$ relates the j -th parameter with the input space. Our CCGP employs the inducing variables-based method for sparse approximations of GPs (Hensman et al., 2015). For each $\mu_q(\cdot)$, we introduce a set of $M \leq N$ “pseudo variables” $\mathbf{u}_q = [\mu_q(\mathbf{z}_1^{(q)}), \dots, \mu_q(\mathbf{z}_M^{(q)})]^\top \in \mathbb{R}^M$ through evaluations of $\mu_q(\cdot)$ at unknown locations $\mathbf{Z}_q = [\mathbf{z}_1^{(q)}, \dots, \mathbf{z}_M^{(q)}] \in \mathbb{R}^{M \times P}$. Note that $\mathbf{u} = [\mathbf{u}_1^\top, \dots, \mathbf{u}_Q^\top]^\top \in \mathbb{R}^{QM}$, yielding:

$$p(\mathbf{f}_j | \mathbf{u}) = \mathcal{N}(\mathbf{f}_j | \mathbf{K}_{f_j u} \mathbf{K}_{uu}^{-1} \mathbf{u}, \mathbf{K}_{f_j f_j} - \mathbf{K}_{f_j u} \mathbf{K}_{uu}^{-1} \mathbf{K}_{u f_j}), \quad (5.10)$$

$$p(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{0}, \mathbf{K}_{uu}) = \prod_{q=1}^Q \mathcal{N}(\mathbf{u}_q | \mathbf{0}, \mathbf{K}_{u_q u_q}), \quad (5.11)$$

where $\mathbf{K}_{uu} \in \mathbb{R}^{QM \times QM}$ is a block-diagonal matrix with blocks $\mathbf{K}_{u_q u_q} \in \mathbb{R}^{M \times M}$, based on the kernel function $\kappa_q(\cdot, \cdot)$. The covariance matrix $\mathbf{K}_{f_j f_j} \in \mathbb{R}^{N \times N}$ holds elements $\sum_{q=1}^Q w_{j,q} w_{j',q} \kappa_q(\mathbf{x}_n, \mathbf{x}_{n'})$, with $\mathbf{x}_n, \mathbf{x}_{n'} \in \mathbf{X}$. Likewise, $\mathbf{K}_{f_j u} = [\mathbf{K}_{f_j u_1}, \dots, \mathbf{K}_{f_j u_Q}] \in \mathbb{R}^{N \times QM}$, where $\mathbf{K}_{f_j u_q} \in \mathbb{R}^{N \times M}$ gathers elements $w_{j,q} \kappa_q(\mathbf{x}_n, \mathbf{z}_m^{(q)})$, $m \in \{1, \dots, M\}$. Alike CGP, in most cases, the CCGP posterior distribution $p(\hat{\mathbf{f}}, \mathbf{u} | \mathbf{y})$ has not an analytical solution, so the actual posterior can be approximated by a parametrized variational distribution $p(\hat{\mathbf{f}}, \mathbf{u} | \mathbf{y}) \approx q(\hat{\mathbf{f}}, \mathbf{u})$, as:

$$q(\hat{\mathbf{f}}, \mathbf{u}) = p(\hat{\mathbf{f}} | \mathbf{u}) q(\mathbf{u}) = \prod_{j=1}^J p(\mathbf{f}_j | \mathbf{u}) \prod_{q=1}^Q q(\mathbf{u}_q), \quad (5.12)$$

where $p(\mathbf{f}_j | \mathbf{u})$ is given by equation (5.10), $q(\mathbf{u}_q) = \mathcal{N}(\mathbf{u}_q | \mathbf{m}_q, \mathbf{V}_q)$, and $q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{V})$. Also, $\mathbf{m}_q \in \mathbb{R}^M$, and $\mathbf{V}_q \in \mathbb{R}^{M \times M}$ are respectively the mean and covariance of variational distribution $q(\mathbf{u}_q)$; similarly, $\mathbf{m} = [\mathbf{m}_1^\top, \dots, \mathbf{m}_Q^\top]^\top \in \mathbb{R}^{QM}$, and $\mathbf{V} \in \mathbb{R}^{QM \times QM}$ is a block-diagonal matrix with blocks given by the covariance matrices \mathbf{V}_q . We remark that the variational approximation given by equation (5.12) is well known, and it has been used in several GPs models, including [Moreno-Muñoz et al. \(2018\)](#); [Saul et al. \(2016\)](#). The approximation for the posterior distribution comprises the computation of the following variational parameters: the mean vectors $\{\mathbf{m}_q\}_{q=1}^Q$ and the covariance matrices $\{\mathbf{V}_q\}_{q=1}^Q$. Such an estimation is carried out by maximizing an evidence lower bound–(ELBO), which is given as:

$$\mathcal{L} = \sum_{n=1}^N \mathbb{E}_{q(\mathbf{f}_1), \dots, q(\mathbf{f}_J)} [\log p(y_n | \theta_{1,n}, \dots, \theta_{J,n}) - \sum_{q=1}^Q \mathbb{D}_{KL}(q(\mathbf{u}_q) || p(\mathbf{u}_q))], \quad (5.13)$$

where $\theta_{j,n} = \theta_j(\mathbf{x}_n)$, with $j \in \{1, \dots, J\}$, and $\mathbb{D}_{KL}(\cdot || \cdot)$ is the Kullback-Leibler divergence and $q(\mathbf{f}_j)$ is defined as follows:

$$q(\mathbf{f}_j) = \mathcal{N}(\mathbf{f}_j | \mathbf{K}_{f_j u} \mathbf{K}_{uu}^{-1} \mathbf{m}, \mathbf{K}_{f_j f_j} + \mathbf{K}_{f_j u} \mathbf{K}_{uu}^{-1} (\mathbf{V} - \mathbf{K}_{uu}) \mathbf{K}_{uu}^{-1} \mathbf{K}_{u f_j}). \quad (5.14)$$

Yet, in presence of non-Gaussian likelihoods, the variational expectations–(VEs) computation in equation (5.13) cannot be solved analytically ([Moreno-Muñoz et al., 2018](#); [Saul et al., 2016](#)). Hence, aiming to model different data types, i.e., classification and regression tasks, we need to find a generic alternative to solve the integrals related

to these expectations. In that sense, we use the Gaussian-Hermite quadratures approach as in Hensman et al. (2015); Saul et al. (2016). We remark such ELBO is used to infer the model’s hyperparameters such as the inducing points, the kernel hyperparameters, and the combination factors $w_{j,q}$ equation (5.8). It is worth mentioning that the CCGPs objective functions exhibit an ELBO that allows Stochastic Variational Inference–(SVI) (Blei et al., 2017). Hence, the optimization is solved through a *mini-batch*-based approach from noisy estimates of the global objective gradient, which allows dealing with large scale datasets (Hensman et al., 2015; Moreno-Muñoz et al., 2018; Saul et al., 2016). Finally, we notice that the computational complexity for our CCGP is dominated by the inversion of \mathbf{K}_{uu} with $\mathcal{O}(QM^3)$ and \mathbf{K}_{fu} with $\mathcal{O}(JNQM^2)$.

5.2 Correlated chained GP for multiple annotators-CCGPMA

5.2.1 Classification

To model categorical data from multiple annotators using our CCGPMA, we use the framework proposed in Rodrigues et al. (2013), which introduces a binary variable $\lambda_n^{(r)} \in \{0, 1\}$ representing the r -th labeler’s reliability as a function of each sample \mathbf{x}_n . If $\lambda_n^{(r)} = 1$, the r -th annotator is supposed to provide the actual label, yielding to a categorical distribution. Conversely, $\lambda_n^{(r)} = 0$ indicates that the r -th annotator gives an incorrect output, which is modeled by a uniform distribution. Therefore, the likelihood function is given as:

$$p(\mathbf{Y}|\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{r \in R_n} \left(\prod_{k=1}^K \zeta_{k,n}^{\delta(y_n^{(r)}, k)} \right)^{\lambda_n^{(r)}} \left(\frac{1}{K} \right)^{(1-\lambda_n^{(r)})}, \quad (5.15)$$

where $\delta(y_n^{(r)}, k) = 1$, if $y_n^{(r)} = k$, otherwise $\delta(y_n^{(r)}, k) = 0$. Besides, $\zeta_{k,n} = p(y_n^{(r)} = k | \lambda_n^{(r)} = 1)$ is an estimation of the unknown ground truth. Accordingly, $J = K + R$ LFs are required within our CCGPMA approach, aiming to model the likelihood’s parameters $\boldsymbol{\theta}$. In particular, K LFs are used to model $\zeta_{k,n}$ based on a softmax function as:

$$\zeta_{k,n} = \Xi(f_k(\mathbf{x}_n)) = \frac{\exp(f_k(\mathbf{x}_n))}{\sum_{j=1}^K \exp(f_j(\mathbf{x}_n))}. \quad (5.16)$$

Besides, R LFs are utilized to compute each $\lambda_n^{(r)}$ from a step function; therefore, $\lambda_n^{(r)} = 1$ if $f_{l_r}(\mathbf{x}_n) \geq 0$, otherwise, $\lambda_n^{(r)} = 0$ ($r \in \{1, \dots, R\}$). $l_r = K + r \in \{K + 1, \dots, J\}$ indexes the r -th annotator's LF. It is worth mentioning that we approximate the step function through the well-known sigmoid function ς to avoid discontinuities and favor the CCGPMA implementation. Alike to CCGP, we use variational inference to approximate the posterior distribution of our CCGPMA. In consequence, the actual posterior $p(\hat{\mathbf{f}}, \mathbf{u}|\mathbf{Y})$ is approximated following equation (5.12). Besides, we can derive a CCGPMA ELBO, yielding

$$\begin{aligned} \mathcal{L} = & \sum_{n=1}^N \sum_{r \in R_n} \mathbb{E}_{q(f_1), \dots, q(f_J)} \left[\log \left(\prod_{k=1}^K \zeta_{k,n}^{\delta(y_n^{(r)}, k)} \right)^{\lambda_n^{(r)}} \left(\frac{1}{K} \right)^{(1-\lambda_n^{(r)})} \right] - \dots \\ & \dots - \sum_{q=1}^Q \mathbb{D}_{KL}(q(\mathbf{u}_q) || p(\mathbf{u}_q)). \end{aligned} \quad (5.17)$$

Finally, given a new sample \mathbf{x}_* , we are interested in the mean and variance for predictive distributions related to the ground truth $\zeta_{k,*} = p(y_* = k)$, and the labelers' reliabilities $\lambda_*^{(r)}$. Accordingly, for $\zeta_{k,*}$ we obtain

$$\mathbb{E}[\zeta_{k,*} | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] \approx \int \Xi(f_k(\mathbf{x}_*)) q(\mathbf{f}_*) d\mathbf{f}_*, \quad (5.18)$$

where $q(\mathbf{f}_*) = \int p(\mathbf{f}_* | \mathbf{u}) q(\mathbf{u}) d\mathbf{u}$. Similarly, for the predictive variance of $\zeta_{k,*}$, we use the expression $\text{Var}[\zeta_{k,*}] = \mathbb{E}[\zeta_{k,*}^2] - \mathbb{E}[\zeta_{k,*}]^2$; hence, we need to compute $\mathbb{E}[\zeta_{k,*}^2]$ as

$$\mathbb{E}[\zeta_{k,*}^2 | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] \approx \int \Xi(f_k(\mathbf{x}_*))^2 q(\mathbf{f}_*) d\mathbf{f}_*. \quad (5.19)$$

On the other hand, regarding the predictive mean and variance for $\lambda_*^{(r)}$, we have

$$\mathbb{E}[\lambda_*^{(r)} | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] = \int \varsigma(f_{l_r,*}) q(\mathbf{f}_*) d\mathbf{f}_*. \quad (5.20)$$

For the variance of $\lambda_*^{(r)}$, we use the expression $\text{Var}[\lambda_*^{(r)}] = \mathbb{E}[(\lambda_*^{(r)})^2] - \mathbb{E}[\lambda_*^{(r)}]^2$; hence, we need to compute

$$\mathbb{E}[(\lambda_*^{(r)})^2 | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] = \int \varsigma(f_{l_r,*})^2 q(\mathbf{f}_*) d\mathbf{f}_*. \quad (5.21)$$

In this case, integrals in equations (5.18) to (5.21) have not closed solution; hence, we approximate them using the Gaussian-Hermite quadrature.

5.2.2 Regression

On the other hand, For real-valued outputs, e.g., $\mathcal{Y} \in \mathbb{R}$, we follow the multi-annotator model used in Groot et al. (2011); Raykar et al. (2010); Rodrigues et al. (2017); Xiao et al. (2013), where each output $y_n^{(r)}$ is considered to be a corrupted version of the hidden ground truth y_n . Then:

$$p(\mathbf{Y} | \boldsymbol{\theta}) = \prod_{n=1}^N \prod_{r \in R_n} \mathcal{N}(y_n^{(r)} | y_n, v_n^{(r)}), \quad (5.22)$$

where $v_n^{(r)} \in \mathbb{R}^+$ is the r -th annotator error-variance for the instance n . In turn, to model this likelihood function with CCGPMA, it is necessary to chain each likelihood's parameter to a latent function f_j . Thus, we require $J = R + 1$ LFs; one to model the hidden ground truth, such that $y_n = f_1(\mathbf{x}_n)$, and R LFs to model each error-variance $v_n^{(r)} = \exp(f_{l_r}(\mathbf{x}_n))$, with $r \in \{1, \dots, R\}$, and $l_r = r + 1 \in \{2, \dots, J\}$. Note that we use an exponential function to map from f_{l_r} to $v_n^{(r)}$, aiming to guarantee $v_n^{(r)} > 0$ ($f_{l_r} \in \mathbb{R}$).

Similar to the classification problem, the actual posterior $p(\hat{\mathbf{f}}, \mathbf{u} | \mathbf{Y})$ is approximated following equation (5.12). Further, the CCGPMA ELBO in regression settings is given by

$$\mathcal{L} = \sum_{n=1}^N \sum_{r \in R_n} \mathbb{E}_{q(\mathbf{f}_1), \dots, q(\mathbf{f}_J)} \left[\log \mathcal{N}(y_n^{(r)} | y_n, v_n^{(r)}) \right] - \sum_{q=1}^Q \mathbb{D}_{KL}(q(\mathbf{u}_q) || p(\mathbf{u}_q)). \quad (5.23)$$

Now, given a new sample \mathbf{x}_* , we are interested in the mean and variances for predictive distributions concerning the ground truth y_* , and the labelers' error-variances $v_*^{(r)}$. First, for y_* we have that since $\mathbf{y} = \mathbf{f}_1$, the posterior distribution for y_* corresponds to $q(f_{1*})$, yielding:

$$\mathbb{E}[y_* | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] = \mu_{1,*} \quad (5.24)$$

$$\text{Var}[y_* | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] = s_{1,*}, \quad (5.25)$$

where $\mu_{1,*}$, and $s_{1,*}$ are respectively the mean and variance of $q(f_{1*})$. Then, for $v_*^{(r)}$, we note that due to $\mathbf{v}_r = \exp(\mathbf{f}_{l_r})$, the posterior distribution for $v_*^{(r)}$ follows a log-normal distribution with parameters $\mu_{l_r,*}$ and $s_{l_r,*}$, which respectively correspond to the mean and variance of $q(f_{l_r,*})$. In this sense, the mean and variance of $v_*^{(r)}$ are given as:

$$\mathbb{E}[v_*^{(r)} | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] = \exp\left(\mu_{l_r,*} + \frac{s_{l_r,*}}{2}\right). \quad (5.26)$$

$$\text{Var}[v_*^{(r)} | \mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u}] = \exp(2\mu_{l_r,*} + s_{l_r,*}) (\exp(s_{l_r,*}) - 1). \quad (5.27)$$

5.3 Experimental set-up

In this section we introduce the experiments' configurations to validate our CCGPMA in both classification and regression settings.

5.3.1 Classification

Testing datasets

We test our approach using three types of datasets: *fully synthetic data*, *semi-synthetic data*, and *fully real datasets*.

First, we generate *fully synthetic data* as one-dimensional ($P = 1$) multi-class classification problem ($K = 3$). The input feature matrix \mathbf{X} is built by randomly sampling $N = 100$ points from an uniform distribution within the interval $[0, 1]$. The true label for the n -th sample is generated by taking the $\arg \max_i \{t_{n,i} : i \in \{1, 2, 3\}\}$, where $t_{n,1} = \sin(2\pi x_n)$, $t_{n,2} = -\sin(2\pi x_n)$, and $t_{n,3} = -\sin(2\pi(x_n + 0.25)) + 0.5$. Besides, the

test instances are obtained by extracting 200 equally spaced samples from the interval $[0, 1]$. Second, to control the label generation, we build *semi-synthetic data* from seven datasets of the UCI repository¹ focused on binary and multi class-classification, which are explained in Section 2.3.1 and the Western dataset introduced in Section 4.3.1. Besides, we use additional datasets from UCI repository: Occupancy Detection Data Set–(Occupancy) ($N = 20560$, $P = 7$, $K = 2$), and Skin Segmentation Data Set–(Skin) ($N = 245057$, $P = 4$, $K = 2$). Finally, we evaluate our proposal on two *fully real datasets*, where both the input features and the annotations are captured from real-world problems. Namely, the Voice and Music dataset presented in Section 2.3.1.

Simulated and provided labels

Note that the *fully synthetic* and the *semi-synthetic* datasets do not hold real annotations. Therefore, it is necessary to simulate those labels as corrupted versions of the hidden ground truth. Here, the simulations are performed by assuming: i) dependencies among annotators, and ii) the labelers' performance is modeled as a function of the input features. In turn, the generative model of our approach based on SLFM (termed SFLM-C) is used to build the labels, as follows:

- Define Q deterministic functions $\hat{\mu}_q : \mathcal{X} \rightarrow \mathbb{R}$, and their combination parameters $\hat{w}_{l,r,q} \in \mathbb{R}$, $\forall r \in R, n \in N$.
- Compute $\hat{f}_{l,r,n} = \sum_{q=1}^Q \hat{w}_{l,r,q} \hat{\mu}_q(\hat{x}_n)$, where $\hat{x}_n \in \mathbb{R}$ is the n -th component of $\hat{\mathbf{x}} \in \mathbb{R}^N$, being $\hat{\mathbf{x}}$ the 1-D representation of the input features in \mathbf{X} by using the well-known t -distributed Stochastic Neighbor Embedding approach (Maaten and Hinton, 2008).
- Calculate $\hat{\lambda}_n^{(r)} = \varsigma(\hat{f}_{l,r,n})$, where $\varsigma(\cdot) \in [0, 1]$ is the sigmoid function.
- Finally, find the r -th label as $y_n^{(r)} = \begin{cases} y_n, & \text{if } \lambda_n^{(r)} \geq 0.5 \\ \tilde{y}_n, & \text{if } \lambda_n^{(r)} < 0.5 \end{cases}$, where \tilde{y}_n is a flipped version of the actual label y_n .

Moreover, aiming to validate our approach under different labels' distributions, we employ the simulation method *Biased coin (Non-homogeneous accuracy)*, which assumes that the annotators' performance depends on the input space (see Section 3.3.1).

¹<http://archive.ics.uci.edu/ml>

CCGPMA training and method comparison

The classification performance is assessed as the Area Under the Curve (AUC) and the overall accuracy (Acc). Further, the AUC is extended for multi-class settings, as discussed by authors in [Fawcett \(2006\)](#). We use a cross-validation scheme with 15 repetitions where 70% of the samples are utilized for training and the remaining 30% for testing (except for the music dataset training and testing sets are clearly defined). Table 5.1 displays the employed methods of the state-of-the-art for comparison purposes. The abbreviations are fixed as: Gaussian Processes classifier (GPC), logistic regression classifier (LRC), majority voting (MV), multiple annotators (MA), Modelling annotators expertise (MAE), Learning from crowds (LFC), Distinguishing good from random labelers (DGRL), kernel alignment-based annotator relevance analysis (KAAR).

On the other hand, the Radial basis function (RBF) kernel is preferred in pattern classification because of its universal approximating ability and mathematical tractability. Hence, for all GP-based approaches, the kernel functions are fixed as:

$$\kappa(\mathbf{x}_n, \mathbf{x}_{n'}) = s^2 \exp \left(\frac{-\|\mathbf{x}_n - \mathbf{x}_{n'}\|_2^2}{2l^2} \right), \quad (5.28)$$

where $\|\cdot\|^2$ stands for the L2 norm, $n, n' \in \{1, 2, \dots, N\}$, and $s, l \in \mathbb{R}^+$ are the kernel hyper-parameters. For concrete testing, we fix $s = 1$, while l is estimated by optimizing the corresponding ELBO in equation (5.17) (we use a gradient-based optimization. The required gradients, and the predictive distributions are presented in Section A.2.1). Moreover, for CGPMA, we fix $Q = R + K$, since each LF $f_j(\cdot)$ is linked to $u_q(\cdot)$. On the other hand, for CCGPMA, each $f_j(\cdot)$ is built as a convex combination of $\mu_q(\cdot)$ (see equation (5.8)); therefore, there is no restriction concerning Q . However, to make a fair comparison with CGPMA, we also fix $Q = R + K$ in CCGPMA. For the *fully synthetic datasets*, we use $M = 10$ inducing points per latent function, and for the remaining experiments, we test with $M = 40$, and $M = 80$. For all the experiments, we use stochastic inference with a mini-batch size of 100. The CCGPMA's Python code is publicly available.²

²<https://github.com/juliangilg/CCGPMA>

Table 5.1: A brief overview of the state-of-the-art methods tested.

Algorithm	Description
GPC-GOLD	A GPC using the real labels (upper bound).
GPC-MV	A GPC using the MV of the labels as the ground truth.
MA-LFC-C (Raykar et al., 2010)	A LRC with constant parameters across the input space.
MA-DGRL (Rodrigues et al., 2013)	A multi-labeler approach that considers as latent variables the annotator performance.
MA-GPC (Rodrigues et al., 2014b)	A multi-labeler GPC, which is as an extension of MA-LFC.
MA-GPCV (Morales-Álvarez et al., 2019)	An extension of MA-GPC that includes variational inference and priors over the labelers' parameters.
MA-DL (Rodrigues and Pereira, 2018a)	A Crowd Layer for DL, where the annotators' parameters are constant across the input space.
KAAR (Gil-Gonzalez et al., 2018)	A kernel-based approach that employs a convex combination of classifiers and codes labelers dependencies.
CCGPMA-C	A particular case of our CCGPMA for classification, where $Q = J$, and we fix $w_{j,q} = 1$, if $j = q$, otherwise $w_{j,q} = 0$.

5.3.2 Regression

Testing datasets

Aiming to test our CCGPMA in regression scenarios, we use the synthetic, semi-synthetic, and real dataset presented in Section 2.3.2.

Simulated and provided labels

For *fully synthetic* and *semi-synthetic* datasets do not hold real annotations. Thus, it is necessary to generate these labels synthetically as a version of the gold standard corrupted by Gaussian noise, i.e., $y_n^{(r)} = y_n + \epsilon_n^{(r)}$, where $\epsilon_n^{(r)} \sim \mathcal{N}(0, v_n^{(r)})$, being $v_n^{(r)}$ the r -th annotator error-variance for the sample n . Note that we are interested in modeling such an error-variance for the r -th annotator as a function of the input features, which is correlated with the other labelers' variances. In turn, an SLFM-based approach is used to build the labels, as follows:

- Define Q functions $\hat{\mu}_q : \mathcal{X} \rightarrow \mathbb{R}$, and the combination parameters $\hat{w}_{l_r,q} \in \mathbb{R}$, $\forall r, q$.
- Compute $\hat{f}_{l_r,n} = \sum_{q=1}^Q \hat{w}_{l_r,q} \hat{\mu}_q(\hat{x}_n)$, where \hat{x}_n is the n -th component of $\hat{\mathbf{x}} \in \mathbb{R}$, which is an 1-D representation of input features \mathbf{X} by using the t-distributed Stochastic Neighbor Embedding approach (Maaten and Hinton, 2008).
- Finally, determine $\hat{v}_n^{(r)} = \exp(\hat{f}_{l_r,n})$.

Besides, aiming to validate our approach using different labels' distributions, we use the simulation method described in Section 3.3.2 (termed *Non-homogeneous error-variance*), which assumes that the annotators' performance depends on the input space.

CCGPMA training and method comparison

The quality assessment is carried out by estimating the regression performance as the coefficient of determination—(R^2) and Pearson correlation coefficient—(Pear). A cross-validation scheme is employed with 15 repetitions where 70% of the samples are utilized for training and the remaining 30% for testing (except for *fully synthetic dataset*, since it clearly defines the training and testing sets). Table 2.4 displays the employed methods of the state-of-the-art for comparison purposes.

Besides, for CGPMA, we set $Q = R + 1$, where each LF $f_j(\cdot)$ is linked to $u_q(\cdot)$. On the other hand, for CCGPMA, each $f_j(\cdot)$ is built as a convex combination of $\mu_q(\cdot)$ (see equation (5.8)); therefore, there is no restriction concerning Q . However, to make a fair comparison with CGPMA, we also fix $Q = R + 1$ in CCGPMA. For the *fully synthetic datasets*, we use $M = 10$ inducing points per latent function, and for the remaining experiments, we test with $M = 40$, and $M = 80$. For all the experiments, we use stochastic inference with a mini-batch size of 100. The model parameters are estimated by optimizing the ELBO in equation (5.23). Such optimization is performed via a gradient-based algorithm (The required gradients, and the predictive distributions are presented in Section A.2.2)

5.4 Results and discussion

5.4.1 Classification

We first perform a controlled experiment to test the CCGPMA capability when dealing with binary and multi-class classification. We use the *fully synthetic* dataset described in Section 5.3.1. Besides, five labelers ($R = 5$) are simulated with different levels of expertise. To simulate the annotators' performance based on the method SLFM-C, we define $Q = 3 \hat{\mu}_q(\cdot)$ functions, yielding:

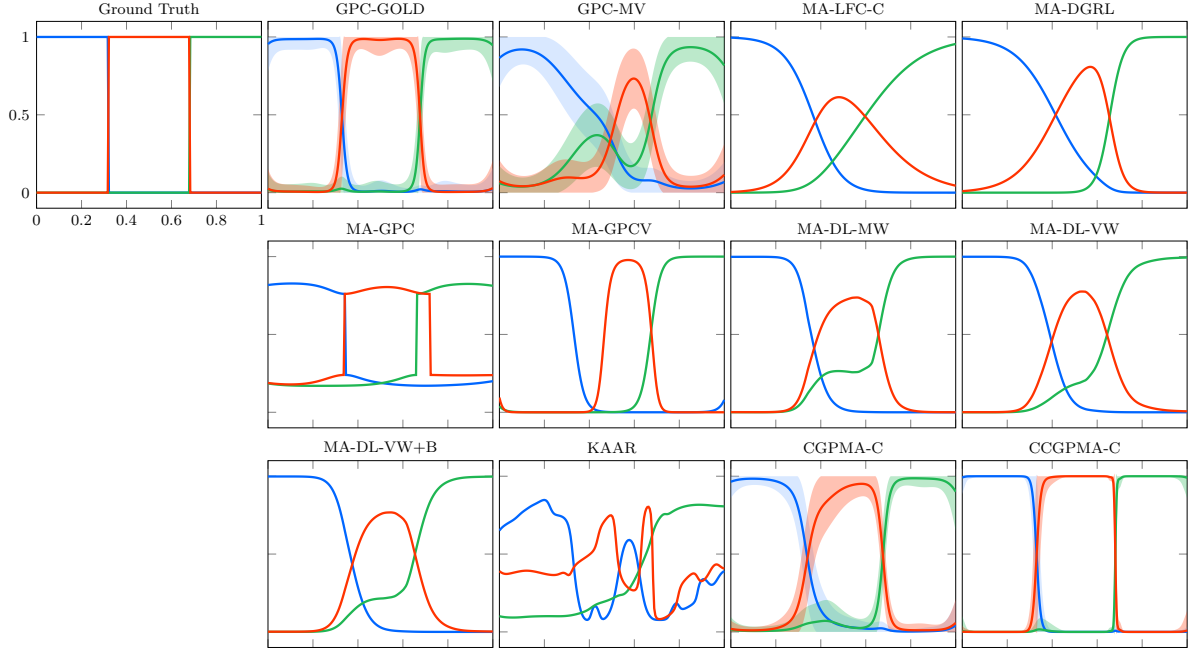


Figure 5.1: Fully synthetic dataset results. The PLP is shown, comparing the prediction of our CCGPMA-C($AUC = 1$) and CCGPMA-C($AUC = 0.9999$) against: the theoretical upper bound GPC-GOLD($AUC = 1.0$), the lower bound GPC-MV($AUC = 0.9809$), and the state-of-the-art approaches MA-LFC-C($AUC = 0.9993$), MA-DGRL($AUC = 0.9999$), MA-GPC($AUC = 0.9977$), MA-GPCV($AUC = 0.9515$), MA-DL-MW($AUC = 0.9989$), MA-DL-VW($AUC = 0.9972$), MA-DL-VW+B($AUC = 0.9994$), KAAR(0.9099). Note that the shaded region in GPC-MV, CGPMA-C, and CCGPMA-C indicates the area enclosed by the mean \pm two standard deviations. There is no shaded region for approaches lacking prediction uncertainty.

$$\hat{\mu}_1(x) = 4.5 \cos(2\pi x + 1.5\pi) - 3 \sin(4.3\pi x + 0.3\pi), \quad (5.29)$$

$$\hat{\mu}_2(x) = 4.5 \cos(1.5\pi x + 0.5\pi) + 5 \sin(3\pi x + 1.5\pi), \quad (5.30)$$

$$\hat{\mu}_3(x) = 1, \quad (5.31)$$

where $x \in [0, 1]$. Besides, the combination weights are gathered within the following combination matrix $\hat{\mathbf{W}} \in \mathbb{R}^{Q \times R}$:

$$\hat{\mathbf{W}} = \begin{bmatrix} 0.4 & 0.7 & -0.5 & 0.0 & -0.7 \\ 0.4 & -1.0 & -0.1 & -0.8 & 1.0 \\ 3.1 & -1.8 & -0.6 & -1.2 & 1.0 \end{bmatrix}, \quad (5.32)$$

holding elements $\hat{w}_{l,r,q}$. Similarly, For the *Non-homogeneous accuracy* approach, we divide the input space into five regions and define the performance matrix $\mathbf{P} \in [0, 1]^{R \times R}$ defined by equation (4.10).

For visual inspection purposes, we perform an initial experiment using the simulation method SLFM-C. Figure 5.1 shows the predictive label’s probability–(PLP), $p(y_* = k|\mathbf{x}_*)$, and the AUC for all studied approaches regarding the *fully synthetic* data. Notice that for methods MA-GPC, MA-GPCV, and KAAR, we use the *one-vs-all* scheme for this experiment (such methods were defined only for binary classification settings). Accordingly, for those models, the PLP corresponds to scores rather than probabilities. Besides, regarding the PLP of our CGPMA and CCGPMA, we provide the mean and variance for the predictive distribution $\zeta_{k,*} = p(y_* = k|\mathbf{x}_*, \hat{\mathbf{f}}, \mathbf{u})$, which are computed based on equations (5.18) and (5.19). As seen in Figure 5.1, KAAR, MA-GPC, and MA-GPCV presents a different shape than the ground truth; moreover, KAAR and MA-GPCV exhibit the worst AUC, even worse than the intuitive lower bound GPC-MV. We explain such conduct in the sense that these approaches are designed to deal with binary labels (Gil-Gonzalez et al., 2018; Rodrigues et al., 2014b; Ruiz et al., 2019). To face such a problem, we use the *one-vs-all* scheme; still, it can lead to ambiguously classified regions (Bishop, 2006). We note an akin predictive AUC concerning MA-DL methods and the linear approaches MA-LFC-C and MA-DGRL. Nonetheless, the linear techniques exhibit a PLP less similar to the Ground truth, which is due to MA-LFC-C and MA-DGRL only can deal with linearly separable data. Further, we analyze the results of our CGPMA-C and its particular enhancement CCGPMA-C. We remark that our methods’ predictive AUC is pretty close to deep learning and linear models. Unlike them, our CGPMA-C and CCGPMA-C show the most accurate PLP compared to the absolute gold standard. CCGPMA-C behaves quite similarly to GPC-GOLD, which is the theoretical upper bound. Finally, from the GPC-MV, we do not identify notable differences with the rest of the approaches (excluding KAAR and MA-GPCV).

From the above, we recognize that analyzing both the predictive AUC and the PLP, our CCGPMA-C exhibits the best performance obtaining similar results compared to the

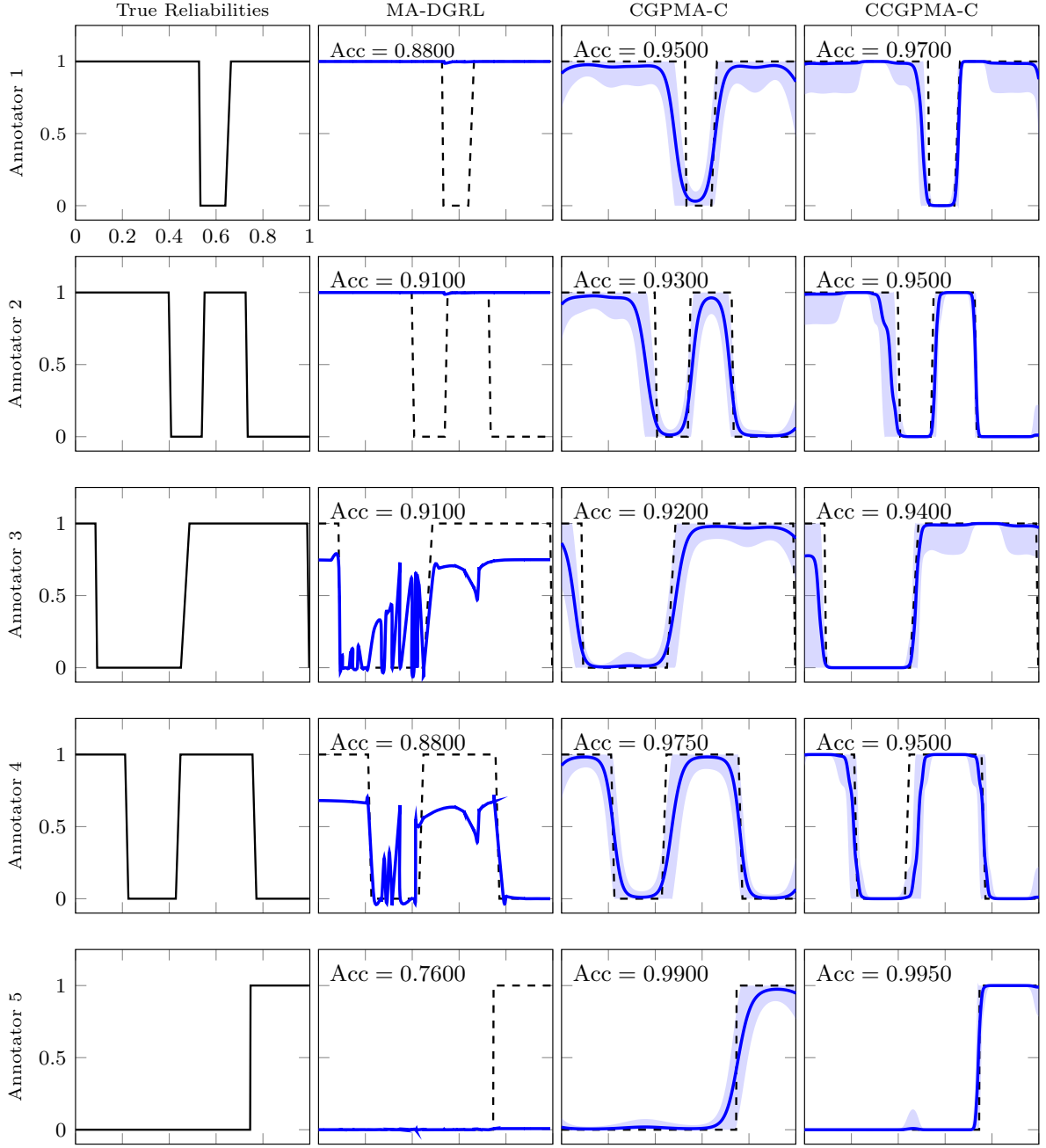


Figure 5.2: Fully synthetic data reliability results. From top to bottom, the first column exposes the true reliabilities (λ_r). The subsequent columns present the estimation of the reliabilities performed by state-of-the-art models, where the correct values are provided in dashed lines. The shaded region in CGPMA-C and CCGPMA-C indicates the area enclosed by the mean \pm two standard deviations. Also, the accuracy (Acc) is provided.

intuitive upper bound (GPC-GOLD). Accordingly, CCGPMA-C proffers a more suitable representation of the labelers’ behavior than its competitors. Indeed, CCGPMA-C codes both the annotators’ dependencies and the relationship between the input features and the annotators’ performance. To empirically support the above statement, Figure 5.2 shows the estimated per-annotator reliability, where we only take into account models that include such types of parameters (MA-DGRL, CGPMA-C, and CCGPMA-C). As seen, MA-DGRL (see column 2 in Figure 5.2) does not offer a proper representation of the annotators’ behavior. CGPMA-C and CCGPMA-C (columns 3 and 4 in Figure 5.2) outperforms MA-DGRL, which is a direct repercussion of modeling the labelers’ parameters as functions of the input features. We observe that CCGPMA-C exhibits the best performance in terms of accuracy; such an outcome is due to this method improves the quality of the annotators’ model by considering correlations among their decisions (Gil-Gonzalez et al., 2018; Zhu et al., 2019)).

It is worth mentioning that the Semi-synthetic experiments are a common practice in the *learning from crowds* area (Morales-Álvarez et al., 2019; Ruiz et al., 2019), where the input features comes from real-world datasets whilst the labels from multiple annotators are simulated following the *fully synthetic data* set-up (see equations (5.29) to (5.32)). Table 5.2 (a) shows the results concerning semi-synthetic datasets and using the simulation method SLFM-C. On average, our CCGPMA-C accomplishes the best predictive AUC; moreover, we note that CGPMA-C reaches the second-best performance. Furthermore, the GPs-based competitors achieve competitive results (GPC-MV, MA-GPC, MA-GPCV, and KAAR). On the other hand, the GPC-MV method obtains a significantly lower performance than our CCGPMA-C, which is explained because GPC-MV is the most naive approach since it considers that the whole annotators exhibit the same performance. Conversely, analyzing the results from MA-GPC, MA-GPCV, and KAAR, we note that they perform worse than GPC-MV. We explain such an outcome in two ways. First, these approaches do not model the relationship between the input features and the annotators’ performance. Second, as exposed in a previous experiment MA-GPC, MA-GPCV, and KAAR use a *one-vs-all* to deal with multi-class problems, which can lead to ambiguously classified regions (Bishop, 2006). The latter can be confirmed in the results for the multi-class dataset “Western” ($K = 4$), where the predictive AUC for such approaches are the lowest. Then, analyzing the results from the DL-based strategies, we note a slightly better performance compared to the GPs-based methods (excluding CGPMA-C and CCGPMA-C). However, the DL-based performs considerably worse than our proposal because the CrowdLayer provides straightforward

Table 5.2: Semi-synthetic classification results. Bold: the highest AUC excluding the upper bound (target) classifier GPC-GOLD. Marked with *: the highest accuracy (Acc) except the upper bound.

(a) SLFM-C.

Method		Breast	Bupa	Ionosphere	Pima	TicTacToe	Occupancy	Skin	Western	Wine	Segmentation	Average
GPC-GOLD($M=40$)	AUC[%]	99.07 \pm 0.45	69.75 \pm 4.66	94.90 \pm 2.35	83.78 \pm 3.02	84.29 \pm 3.34	99.56 \pm 0.06	99.97 \pm 0.01	91.85 \pm 0.61	99.87 \pm 0.15	95.96 \pm 1.96	91.90
	Acc[%]	96.62 \pm 1.00	65.38 \pm 4.03	85.91 \pm 3.61	77.92 \pm 2.48	79.61 \pm 2.55	98.83 \pm 0.09	99.62 \pm 0.13	72.09 \pm 1.01	96.67 \pm 2.24	75.26 \pm 8.69	84.79
GPC-GOLD($M=80$)	AUC[%]	99.03 \pm 0.46	69.97 \pm 4.83	95.13 \pm 2.25	83.74 \pm 2.97	84.91 \pm 3.23	99.56 \pm 0.06	99.97 \pm 0.01	92.50 \pm 0.57	99.88 \pm 0.16	97.81 \pm 0.41	92.25
	Acc[%]	96.62 \pm 1.07	66.09 \pm 3.94	86.16 \pm 3.36	77.84 \pm 2.51	80.07 \pm 2.47	98.83 \pm 0.10	99.58 \pm 0.11	73.00 \pm 1.30	96.79 \pm 2.26	86.57 \pm 1.87	86.15
GPC-MV($M=40$)	AUC[%]	98.97 \pm 0.45	53.66 \pm 5.16	75.66 \pm 5.72	53.99 \pm 7.60	66.20 \pm 3.57	75.85 \pm 19.16	84.58 \pm 0.90	86.58 \pm 3.31	81.79 \pm 2.12	95.62 \pm 2.28	77.29
	Acc[%]	94.93 \pm 1.27	49.42 \pm 7.21	64.91 \pm 3.93	66.81 \pm 2.80	62.36 \pm 4.52	80.40 \pm 1.84	82.64 \pm 0.71	60.81 \pm 5.05	59.51 \pm 6.22	75.62 \pm 9.54	69.74
GPC-MV($M=80$)	AUC[%]	98.92 \pm 0.48	56.98 \pm 5.29	77.79 \pm 5.50	53.02 \pm 6.74	67.44 \pm 2.57	63.12 \pm 19.68	84.20 \pm 0.80	84.46 \pm 0.89	83.23 \pm 4.87	97.49 \pm 0.47	76.66
	Acc[%]	94.96 \pm 1.45	54.42 \pm 5.51	64.91 \pm 3.93	67.73 \pm 1.46	63.12 \pm 2.17	78.41 \pm 2.38	83.06 \pm 0.33	57.64 \pm 1.80	61.73 \pm 11.56	84.00 \pm 1.67	71.00
MA-LFC-C	AUC[%]	87.89 \pm 5.10	45.93 \pm 14.44	73.58 \pm 9.01	81.19 \pm 3.13	60.04 \pm 2.61	89.42 \pm 0.79	94.40 \pm 0.08	84.00 \pm 2.11	96.92 \pm 3.57	98.92 \pm 0.31	81.23
	Acc[%]	67.71 \pm 2.85	46.54 \pm 11.32	62.96 \pm 6.01	76.62 \pm 2.21	60.39 \pm 2.48	70.83 \pm 0.78	93.52 \pm 0.29	59.71 \pm 1.41	90.25 \pm 5.12	92.17 \pm 0.78	72.07
MA-DGRL	AUC[%]	97.57 \pm 1.89	57.24 \pm 3.36	64.53 \pm 7.21	81.38 \pm 2.90	61.29 \pm 2.30	49.71 \pm 1.05	93.79 \pm 1.07	81.43 \pm 1.50	97.95 \pm 2.21	98.97 \pm 0.38	78.39
	Acc[%]	89.92 \pm 3.28	53.72 \pm 3.84	33.27 \pm 3.01	77.00 \pm 2.71	62.06 \pm 2.52	66.73 \pm 0.46	89.30 \pm 1.01	51.21 \pm 1.58	89.01 \pm 4.28	93.27 \pm 1.05	70.55
MA-GPC	AUC[%]	98.11 \pm 1.16	54.46 \pm 5.78	66.31 \pm 14.74	53.25 \pm 17.80	60.79 \pm 9.95	92.57 \pm 7.96	80.89 \pm 0.60	86.71 \pm 1.14	94.17 \pm 2.62	97.34 \pm 0.35	78.46
	Acc[%]	96.13 \pm 1.07	50.38 \pm 4.69	62.14 \pm 11.13	52.41 \pm 14.30	58.70 \pm 8.58	60.62 \pm 1.99	76.70 \pm 0.46	72.18 \pm 1.61	90.49 \pm 3.42	88.86 \pm 2.10	70.86
MA-GPCV	AUC[%]	82.70 \pm 5.47	55.67 \pm 6.83	62.38 \pm 8.71	62.17 \pm 5.90	61.04 \pm 10.03	60.22 \pm 2.66	76.29 \pm 3.74	84.51 \pm 1.47	97.35 \pm 1.72	99.24 \pm 0.27	74.16
	Acc[%]	68.72 \pm 3.24	52.37 \pm 6.67	60.94 \pm 5.64	61.21 \pm 5.83	65.69 \pm 2.33	56.25 \pm 3.01	76.67 \pm 2.30	69.45 \pm 1.48	94.20 \pm 3.35	95.95 \pm 1.19	70.15
MA-DL-MW	AUC[%]	94.70 \pm 1.73	52.37 \pm 5.68	75.35 \pm 5.43	61.78 \pm 2.67	68.27 \pm 2.96	64.09 \pm 2.26	86.36 \pm 0.57	90.92 \pm 0.56	97.28 \pm 1.09	99.50 \pm 0.17	79.06
	Acc[%]	90.24 \pm 4.52	51.67 \pm 6.17	70.25 \pm 3.90*	55.12 \pm 3.64	62.87 \pm 2.83	69.01 \pm 1.30	46.03 \pm 0.26	62.70 \pm 1.05	88.52 \pm 3.30	93.91 \pm 1.33	69.03
MA-DL-VW	AUC[%]	95.26 \pm 2.45	53.27 \pm 6.18	69.87 \pm 4.97	60.63 \pm 3.36	67.71 \pm 2.67	68.40 \pm 3.45	86.56 \pm 0.68	91.73 \pm 0.67	98.07 \pm 1.52	99.72 \pm 0.11	79.12
	Acc[%]	90.57 \pm 1.66	51.67 \pm 5.65	67.30 \pm 4.90	54.46 \pm 1.97	64.42 \pm 2.47	63.99 \pm 1.71	46.19 \pm 0.29	68.98 \pm 1.48	91.23 \pm 4.77	95.89 \pm 0.92*	69.47
MA-DL-VW+B	AUC[%]	94.65 \pm 2.42	52.81 \pm 6.31	71.96 \pm 4.53	61.23 \pm 3.78	67.80 \pm 3.42	67.82 \pm 3.86	86.68 \pm 0.67	91.64 \pm 0.85	98.17 \pm 1.55	99.72 \pm 0.09	79.25
	Acc[%]	90.37 \pm 1.88	52.31 \pm 5.53	67.92 \pm 4.76	56.68 \pm 3.70	62.52 \pm 3.28	64.15 \pm 2.36	46.04 \pm 0.28	64.16 \pm 1.85	90.74 \pm 4.75	95.65 \pm 1.10	69.05
KAAR	AUC[%]	80.58 \pm 2.74	59.20 \pm 6.63	70.46 \pm 7.39	58.02 \pm 4.06	63.81 \pm 5.45	69.16 \pm 2.06	51.58 \pm 4.74	85.88 \pm 1.20	99.43 \pm 1.05	92.17 \pm 1.90	73.03
	Acc[%]	70.18 \pm 3.58	56.22 \pm 5.20*	63.14 \pm 6.68	61.44 \pm 2.64	56.32 \pm 2.61	53.62 \pm 1.10	39.84 \pm 2.78	62.35 \pm 1.86	93.58 \pm 4.53	77.85 \pm 5.82	63.45
CGPMA-C($M=40$)	AUC[%]	99.20 \pm 0.38	57.13 \pm 4.68	83.56 \pm 10.02	82.01 \pm 3.14	70.56 \pm 3.04	82.20 \pm 2.73	92.62 \pm 1.20	91.78 \pm 0.66	99.82 \pm 0.18	96.79 \pm 0.65	85.56
	Acc[%]	96.10 \pm 1.26	47.95 \pm 7.36	65.66 \pm 3.70	76.83 \pm 2.81	67.62 \pm 2.51	61.15 \pm 1.56	90.62 \pm 1.34	71.13 \pm 1.83	95.31 \pm 2.41	82.22 \pm 2.20	75.46
CGPMA-C($M=80$)	AUC[%]	99.14 \pm 0.38	56.96 \pm 4.74	86.15 \pm 6.96	82.04 \pm 3.18	70.48 \pm 3.12	99.08 \pm 0.26	90.46 \pm 1.64	91.85 \pm 0.57	99.84 \pm 0.12	94.06 \pm 0.61	87.01
	Acc[%]	95.62 \pm 1.07	51.47 \pm 5.99	65.16 \pm 3.95	76.97 \pm 2.31*	66.97 \pm 2.93	85.73 \pm 0.88	85.32 \pm 1.50	70.66 \pm 1.27	95.56 \pm 1.95*	71.43 \pm 3.47	76.43
CCGPMA-C($M=40$)	AUC[%]	99.38 \pm 0.27	60.22 \pm 5.06	87.84 \pm 6.72	78.10 \pm 6.22	74.95 \pm 5.39	91.98 \pm 2.00	85.70 \pm 2.66	93.09 \pm 0.51	99.44 \pm 0.33	97.67 \pm 0.53	86.84
	Acc[%]	96.68 \pm 1.23	50.71 \pm 6.73	65.72 \pm 4.43	67.13 \pm 11.17	67.64 \pm 5.45*	87.08 \pm 1.79	82.39 \pm 0.96	73.63 \pm 1.10	93.46 \pm 3.04	84.61 \pm 3.21	76.90
CCGPMA-C($M=80$)	AUC[%]	99.33 \pm 0.30	59.19 \pm 5.65	90.55 \pm 6.29	80.45 \pm 5.10	73.12 \pm 3.23	97.75 \pm 2.00	89.42 \pm 2.20	93.15 \pm 0.50	99.43 \pm 0.33	97.58 \pm 0.43	88.00
	Acc[%]	97.17 \pm 0.87*	51.86 \pm 5.54	65.97 \pm 4.18	67.56 \pm 14.16	65.30 \pm 3.11	91.56 \pm 5.64*	86.84 \pm 1.99	73.98 \pm 1.55*	92.96 \pm 3.30	83.56 \pm 4.17	77.68*

(b) Non-homogeneous accuracy.

Method		Breast	Bupa	Ionosphere	Pima	TicTacToe	Occupancy	Skin	Western	Wine	Segmentation	Average
GPC-GOLD($M=40$)	AUC[%]	99.07 \pm 0.45	69.75 \pm 4.66	94.90 \pm 2.35	83.78 \pm 3.02	84.29 \pm 3.34	99.56 \pm 0.06	99.97 \pm 0.01	91.85 \pm 0.61	99.87 \pm 0.15	95.96 \pm 1.96	91.90
	Acc[%]	96.62 \pm 1.00	65.38 \pm 4.03	85.91 \pm 3.61	77.92 \pm 2.48	79.61 \pm 2.55	98.83 \pm 0.09	99.62 \pm 0.13	72.09 \pm 1.01	96.67 \pm 2.24	75.26 \pm 8.69	84.79
GPC-GOLD($M=80$)	AUC[%]	99.03 \pm 0.46	69.97 \pm 4.83	95.13 \pm 2.25	83.74 \pm 2.97	84.91 \pm 3.23	99.56 \pm 0.06	99.97 \pm 0.01	92.50 \pm 0.57	99.88 \pm 0.16	97.81 \pm 0.41	92.25
	Acc[%]	96.62 \pm 1.07	66.09 \pm 3.94	86.16 \pm 3.36	77.84 \pm 2.51	80.07 \pm 2.47	98.83 \pm 0.10	99.58 \pm 0.11	73.00 \pm 1.30	96.79 \pm 2.26	86.57 \pm 1.87	86.15
GPC-MV($M=40$)	AUC[%]	93.23 \pm 6.06	63.66 \pm 8.50	76.47 \pm 15.09	66.27 \pm 15.45	64.57 \pm 2.64	60.90 \pm 20.43	54.09 \pm 12.70	90.27 \pm 0.55	98.53 \pm 1.60	95.40 \pm 2.11	76.34
	Acc[%]	91.67 \pm 4.19	57.24 \pm 5.96	70.06 \pm 12.75	69.44 \pm 7.35	63.01 \pm 2.28	79.59 \pm 3.19	25.28 \pm 16.83	68.86 \pm 1.58	82.10 \pm 9.94	69.32 \pm 9.29	67.66
GPC-MV($M=80$)	AUC[%]	92.89 \pm 6.23	63.88 \pm 8.19	79.65 \pm 11.68	61.61 \pm 15.74	66.18 \pm 2.75	56.55 \pm 17.28	57.72 \pm 17.30	91.57 \pm 0.48	98.56 \pm 1.62	97.18 \pm 1.27	76.58
	Acc[%]	91.12 \pm 4.45	56.86 \pm 5.03	73.84 \pm 8.35*	68.31 \pm 11.19	63.50 \pm 2.20	78.30 \pm 3.52	29.50 \pm 22.09	71.11 \pm 1.45	83.21 \pm 9.51	78.54 \pm 9.38	69.43
MA-LFC-C	AUC[%]	77.96 \pm 3.82	64.90 \pm 6.59	74.70 \pm 4.40	81.90 \pm 2.36	54.56 \pm 2.92	85.01 \pm 4.96	82.45 \pm 2.51	85.24 \pm 0.97	98.86 \pm 1.03	99.05 \pm 0.37	80.50
	Acc[%]	72.26 \pm 6.65	62.88 \pm 6.05*	75.91 \pm 4.28	76.19 \pm 2.40*	56.85 \pm 2.92	69.44 \pm 4.08	55.25 \pm 1.97	55.64 \pm 1.44	94.07 \pm 2.73	93.58 \pm 1.20	71.21
MA-DGRL	AUC[%]	56.88 \pm 5.65	66.41 \pm 7.39	57.18 \pm 8.85	83.11 \pm 1.92	54.56 \pm 3.08	70.06 \pm 14.56	61.12 \pm 1.24	74.22 \pm 1.68	93.17 \pm 3.55	99.89 \pm 0.13	71.66
	Acc[%]	59.19 \pm 4.56	62.12 \pm 6.70	65.09 \pm 4.97	77.75 \pm 2.28	59.56 \pm 4.28	41.43 \pm 24.82	43.57 \pm 0.36	26.09 \pm 1.27	63.58 \pm 4.57	28.03 \pm 1.61	52.64
MA-GPC	AUC[%]	40.52 \pm 6.09	45.19 \pm 6.26	46.84 \pm 5.68	40.88 \pm 3.82	57.93 \pm 7.38	34.14 \pm 3.27	53.13 \pm 13.53	86.53 \pm 0.85	89.83 \pm 3.19	97.69 \pm 0.49	53.48
	Acc[%]	58.63 \pm 4.55	49.36 \pm 4.14	30.06 \pm 4.16	56.85 \pm 4.28	60.58 \pm 9.38	26.67 \pm 0.99	57.55 \pm 11.06	75.41 \pm 0.75	84.07 \pm 1.53	92.27 \pm 0.92	53.09
MA-GPCV	AUC[%]	62.95 \pm 3.77	59.88 \pm 5.32	64.12 \pm 5.34	65.27 \pm 3.44	55.37 \pm 5.35	27.93 \pm 2.26	55.48 \pm 9.45	89.34 \pm 0.96	88.91 \pm 7.83	98.94 \pm 0.35	66.82
	Acc[%]	59.22 \pm 2.44	57.37 \pm 4.31	61.19 \pm 4.51	64.99 \pm 2.84	53.06 \pm 3.64	25.22 \pm 4.43	49.12 \pm 1.53	72.12 \pm 0.67	83.95 \pm 3.70	95.06 \pm 1.15	62.13
MA-DL-WM	AUC[%]	81.95 \pm 9.55	53.06 \pm 4.27	75.77 \pm 4.64	68.72 \pm 2.81	68.94 \pm 2.57	72.50 \pm 1.59	44.23 \pm 1.79	93.32 \pm 0.54	96.07 \pm 1.79	99.63 \pm 0.24	75.42
	Acc[%]	70.31 \pm 5.58	52.05 \pm 3.58	71.38 \pm 5.96	65.92 \pm 3.04	63.68 \pm 2.99	75.05 \pm 4.48	52.82 \pm 2.27	73.34 \pm 1.72	85.93 \pm 0.41	95.59 \pm 1.24*	70.61
MA-DL-VW	AUC[%]	83.06 \pm 6.45	52.84 \pm 4.49	73.22 \pm 7.61	65.76 \pm 3.91	68.67 \pm 2.35	65.54 \pm 12.99	43.92 \pm 1.65	93.45 \pm 0.39	95.98 \pm 1.96	99.69 \pm 0.17	74.21
	Acc[%]	72.13 \pm 5.13	51.99 \pm 4.20	69.31 \pm 7.22	64.27 \pm 2.83	62.99 \pm 3.74	65.00 \pm 20.34	52.68 \pm 0.32	75.53 \pm 0.87*	86.67 \pm 4.60	95.57 \pm 1.22	69.61
MA-DL-VW+B	AUC[%]	83.55 \pm 5.23	53.29 \pm 4.88	74.51 \pm 7.66	66.45 \pm 3.56	68.81 \pm 3.14	54.14 \pm 19.10	43.90 \pm 1.61	93.49 \pm 0.46*	95.90 \pm 1.80	99.66 \pm 0.26	73.37
	Acc[%]	69.27 \pm 7.52	52.69 \pm 3.51	69.31 \pm 6.85	64.56 \pm 3.36	63.80 \pm 4.78	52.32 \pm 24.94	52.75 \pm 0.38	75.05 \pm 1.68	85.56 \pm 4.27	95.60 \pm 1.69	68.09
KAAR	AUC[%]	66.66 \pm 4.06	54.45 \pm 9.87	75.43 \pm 5.52	73.77 \pm 9.17	50.00 \pm 0.00	82.84 \pm 2.11	92.26 \pm 1.44	84.67 \pm 1.53	97.85 \pm 1.89	96.30 \pm 0.68	77.42
	Acc[%]	59.74 \pm 3.59	53.65 \pm 9.66	73.27 \pm 6.51	70.62 \pm 4.34	33.52 \pm 4.41	72.21 \pm 0.35	57.22 \pm 10.69	63.55 \pm 1.81	90.49 \pm 4.19	80.27 \pm 3.31	65.45
CGPMA-C($M=40$)	AUC[%]	81.02 \pm 9.70	61.31 \pm 5.68	81.70 \pm 3.87	81.88 \pm 1.90	69.01 \pm 2.21	70.00 \pm 6.05	96.77 \pm 0.64	89.55 \pm 0.67	99.61 \pm 0.40	63.20 \pm 5.94	79.41
	Acc[%]	71.58 \pm 8.86	58.78 \pm 4.42	67.11 \pm 4.94	75.96 \pm 2.63	64.98 \pm 2.95	29.01 \pm 0.67	87.66 \pm 0.36*	65.76 \pm 2.17	88.40 \pm 9.58	28.52 \pm 1.23	63.78
CGPMA-C($M=80$)	AUC[%]	81.42 \pm 8.07	61.71 \pm 5.09	83.38 \pm 3.49	81.36 \pm 1.79	71.66 \pm 2.16	67.74 \pm 4.01	84.82 \pm 1.60	89.48 \pm 0.65	99.32 \pm 0.53	67.59 \pm 1.28	78.85
	Acc[%]	69.72 \pm 3.89	59.29 \pm 4.23	70.00 \pm 5.06	75.67 \pm 2.78	67.38 \pm 2.12*	45.52 \pm 4.48	62.81 \pm 1.41	67.62 \pm 1.48	94.44 \pm 3.43	15.94 \pm 1.63	62.84
CCGPMA-C($M=40$)	AUC[%]	94.63 \pm 6.00	61.04 \pm 5.00	84.28 \pm 3.40	81.34 \pm 2.51	60.73 \pm 3.51	96.97 \pm 4.95	85.50 \pm 13.37	92.56 \pm 0.41	99.75 \pm 0.27	98.00 \pm 0.38	85.48
	Acc[%]	86.08 \pm 11.94	59.29 \pm 4.53	70.19 \pm 5.13	75.79 \pm 2.97	61.48 \pm 4.71	88.84 \pm 13.85*	67.77 \pm 17.47	73.13 \pm 1.40	96.30 \pm 0.21	85.06 \pm 0.37	76.39
CCGPMA-C($M=80$)	AUC[%]	97.00 \pm 6.08	61.32 \pm 4.43	86.59 \pm 3.01	80.94 \pm 2.29	59.50 \pm 2.67	85.21 \pm 4.68	89.20 \pm 13.78	92.85 \pm 0.46	99.74 \pm 0.27	97.43 \pm 2.61	84.98
	Acc[%]	91.74 \pm 8.74*	59.04 \pm 3.81	72.26 \pm 5.98	75.70 \pm 2.87	62.80 \pm 2.52	77.92 \pm 4.03	72.96 \pm 2.32	73.70 \pm 1.37	95.80 \pm 2.57	83.70 \pm 6.69	76.56*

lar, the simulated labels (see Section 5.3.1) follows the MA-DGRL model, favoring its performance. Though MA-LFC-C achieves competitive performance compared to the DL-based methods, it is considerably lower than our proposal. In fact, the MA-LFC-C formulation assumes that the annotators' behavior is homogeneous across the input space, which does not correspond to the labels simulation procedure.

On the other hand, Table 5.2 (b) shows the results concerning the simulation method "Non-homogeneous accuracy". At first sight, there exists a generalized lower performance compared to previous results in Table 5.2(a); the above indicates that the annotators' performance is lower, which impacts the algorithms' performance. Our CCGPMA achieves the best predictive performance in both the overall accuracy and the AUC score. Moreover, the non-linear competitors KAAR and DL-based models achieve competitive results. However, GPC-MV, MA-GPC, and MA-GPCV offer the lowest classification scores. Regarding GPC-MV, the result is explained because GPC-MV corresponds to the most naive approach. After all, it considers homogeneity in the labelers' performance. Conversely, MA-GPC and MA-GPCV achieve a worse performance than GPC-MV; such behavior indicates that those methods do not properly represent the labelers' behavior. Regarding linear approaches, we highlight the outcome for MA-LFC, which achieves competitive performance compared to the non-linear competitors.

To summarize, we tested our approach in controlled scenarios by using two different strategies. First, we simulate correlated annotators, where their performance is a function of the input space, named *SLFM-C*. The remaining strategy named *Non-homogeneous accuracy* simulates inconsistent annotators, i.e., labelers, whose performance varies depending on the input features. Attained to the results (Table 5.2 (a), and Table 5.2 (b)), we note that our CCGPMA outperforms all its competitors considering both metrics AUC and the overall accuracy.

Finally, we test the *fully real datasets*, which configure the most challenging scenario. The input features and the labels from multiple experts come from real-world applications. Table 5.3 outlines the achieved AUC. First, we observe that for the voice data, G and R scales exhibit a similar AUC for all considered approaches; in fact, GPC-MV obtains a result comparable with the upper bound GPC-GOLD. The latter can be explained in the sense that the annotators exhibit a suitable performance for these scales, i.e., the provided labels are similar to the ground truth. On the other hand, a reduction in the predictive AUC is observed for scale B, which is a consequence of diminishing the labelers' performance compared to scales G and R, as demonstrated in Gil et al. (2015). Our approaches exhibit the best generalization performances for the three scales

in the voice dataset. Remarkably, CGPMA-C and CCGPMA-C do not suffer significant changes in the scale B, which is an outstanding outcome because it reflects that our method offers a better representation of the labelers' behavior against low-quality annotations. Finally, we review the AUC for the Music dataset. Achieved results show a low performance for the MA-GPC, even lower than their intuitive lower bound (GPC-MV). Notably, our CCGPMA-C reaches the best predictive AUC, being comparable with the intuitive upper bound.

Table 5.3: AUC classification results for the fully real datasets. Bold: the highest performance excluding the GPC-GOLD bound.

Method	G	Voice R	B	Music	Average
GPC-GOLD($M = 40$)	0.9481	0.9481	0.9481	0.9358	0.9450
GPC-GOLD($M = 80$)	0.9484	0.9484	0.9484	0.9178	0.9407
GPC-MV($M = 40$)	0.8942	0.9373	0.8001	0.8871	0.8797
GPC-MV($M = 80$)	0.9301	0.9377	0.7962	0.8897	0.8884
MA-LFC-C	0.9122	0.9130	0.8406	0.8599	0.8814
MA-DGRL	0.9127	0.9164	0.8259	0.8832	0.8845
MA-GPC	0.8660	0.8597	0.4489	0.8253	0.7500
MA-GPCV	0.9283	0.9208	0.8835	0.8677	0.9001
MA-DL-MW	0.8957	0.8966	0.8123	0.8567	0.8653
MA-DL-VW	0.8942	0.8929	0.8092	0.9167	0.8782
MA-DL-VW+B	0.9030	0.8937	0.8218	0.8573	0.8689
KAAR	0.9109	0.9351	0.8969	0.8896	0.9081
CGPMA-C($M = 40$)	0.9324	0.9406	0.8696	0.9025	0.9113
CGPMA-C($M = 80$)	0.9324	0.9417	0.8708	0.8987	0.9109
CCGPMA-C($M = 40$)	0.9318	0.9422	0.9002	0.9446	0.9297
CCGPMA-C($M = 80$)	0.9243	0.9383	0.8907	0.9456	0.9247

5.4.2 Regression

We first perform a controlled experiment aiming to verify the capability of our CGPMA and CCGPMA to estimate the performance of inconsistent annotators as a function of the input space and taking into account their dependencies. For this first experiment, we use the *fully synthetic* dataset described in Section 5.3.2. We simulate five labelers ($R = 5$) with different levels of expertise. To simulate the error-variances based on the approach SLFM-R, we define $Q = 3$ functions $\hat{\mu}_q(\cdot)$, which are given as

$$\hat{\mu}_1(x) = 4.5 \cos(2\pi x + 1.5\pi) - 3 \sin(4.3\pi x + 0.3\pi) + 4 \cos(7\pi x + 2.4\pi), \quad (5.33)$$

$$\hat{\mu}_2(x) = 4.5 \cos(1.5\pi x + 0.5\pi) + 5 \sin(3\pi x + 1.5\pi) - 4.5 \cos(8\pi x + 0.25\pi), \quad (5.34)$$

$$\hat{\mu}_3(x) = 1, \quad (5.35)$$

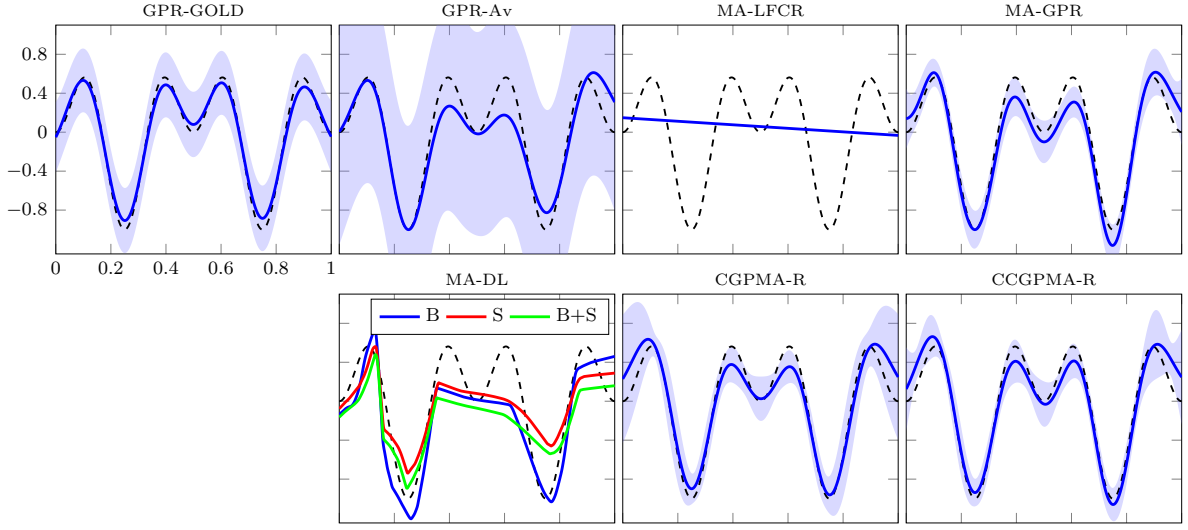


Figure 5.3: Fully synthetic dataset results. We compare the prediction of our CCGPMA-R ($R^2 = 0.9438$), and CGPMA-R ($R^2 = 0.9280$) with the theoretical upper bound GPR-GOLD ($R^2 = 0.9843$) and lower bound GPR-Av ($R^2 = 0.8718$), and state-of-the-art approaches, MA-LFCR ($R^2 = -0.0245$), MA-GPR ($R^2 = 0.9208$), MA-DL-B ($R^2 = 0.7020$), MA-DL-S ($R^2 = 0.6559$), MA-DL-B+S ($R^2 = 0.5997$). Note that we provided the Gold Standard in dashed lines. The shaded region in GPR-Av, MA-GPR, CGPMA-R, and CCGPMA-R indicates the area enclosed by the mean plus or minus two standard deviations. We remark that there is no shaded region for MA-LFCR, and DLMA since they do not provide information about the prediction uncertainty.

where $x \in [0, 1]$. Besides, we define the following combination matrix $\hat{\mathbf{W}} \in \mathbb{R}^{Q \times R}$, where

$$\hat{\mathbf{W}} = \begin{bmatrix} -0.10 & 0.01 & -0.05 & 0.01 & -0.01 \\ 0.10 & -0.01 & 0.01 & -0.05 & 0.05 \\ -2.3 & -1.77 & 0.54 & 0.9 & 1.42 \end{bmatrix}, \quad (5.36)$$

holding elements $w_{l_r, q}$. Likewise, For the *Non-homogeneous error-variance* approach, we divide the input space into five regions and define the performance matrix \mathbf{V} defined by equation (3.15).

Figure 5.3 shows the predictive performance of all methods in this first experiment. The results show two clear groups: those based on GPs (GPR-Av, MA-GPR, CGPMA-R, and CCGPMA-R), which expose the best performance in terms of the R^2 score, and those based on other types of approaches (MA-LFCR, and MA-DL), whose performance is not satisfactory. The behavior of MA-LFCR is low since it only can deal with linear

problems. Besides, concerning MA-DL and its three variations (S, B, and S+B), we note that this approach can deal with non-linear dynamics. However, MA-DL reaches a significantly low performance (even lower than the most naive approach, GPR-Av). To explain such an outcome, we remark that MA-DL comprises the introduction of an additional layer, the “CrowdLayer”, which allows the training of neural networks, directly from the noisy labels of multiple annotators (Rodrigues and Pereira, 2018a). Yet, such a CrowdLayer provides a very simple codification of the annotators’ performance to guarantee a low computational cost (Morales-Álvarez et al., 2019); therefore, MA-DL does not provide a proper codification of the annotators’ behavior. On the other hand, among the GP-based methods, the proposed CCGPMA-R achieves the best performance in terms of R^2 , followed closely by CGPMA-R and MA-GPR.

Besides, concerning the high performance of our CCGPMA-R (the best in terms of R^2 score), we hypothesize that such an outcome is a consequence of our method offers a better representation of the labelers’ behavior when compared to its competitors. To empirically support the above hypothesis, Figure 5.4 shows the estimated error-variances for this first experiment; here, we only take into account the models that include these parameters in their formulations. As seen in Figure 5.4, MA-LFCR and MA-GPR offer the worst representation for the annotator’s performance, which is expected due to such models do not take into account the relationship between the annotators and the input space. Conversely, CGPMA-R and CCGPMA-R outperform the models named previously. This outcome is a consequence that such two approaches compute the error-variance as a function of the input features, allowing for a better codification of the labelers’ behavior. Besides, by making a visual inspection and analyzing the R^2 scores, CCGPMA-R performs better than CGPMA-R, because the former codes properly the annotators’ interdependencies (Zhu et al., 2019). Finally, we remark that although our CCGPMA-R achieves the best representation of the annotators’ performance, the result for Annotator 4 exhibits a lower performance in terms of R^2 score compared to the other labelers. Such an outcome is caused by the quasi-periodic behavior in the error-variances for those labelers, which our approach cannot capture because we are using a RBF-based kernel.

Second, Table 5.4 (a) shows the results of the *semi synthetic datasets*. On average, our CCGPMA-R exhibits the best generalization performance in terms of the R^2 score. On the other hand, regarding its GPs-based competitors (GPR-Av, MA-GPR, and CGPMA-R), we first note that the performance of CGPMA-R exhibits a similar (but lower) performance than CCGPMA-R. The above is a consequence of that conversely to

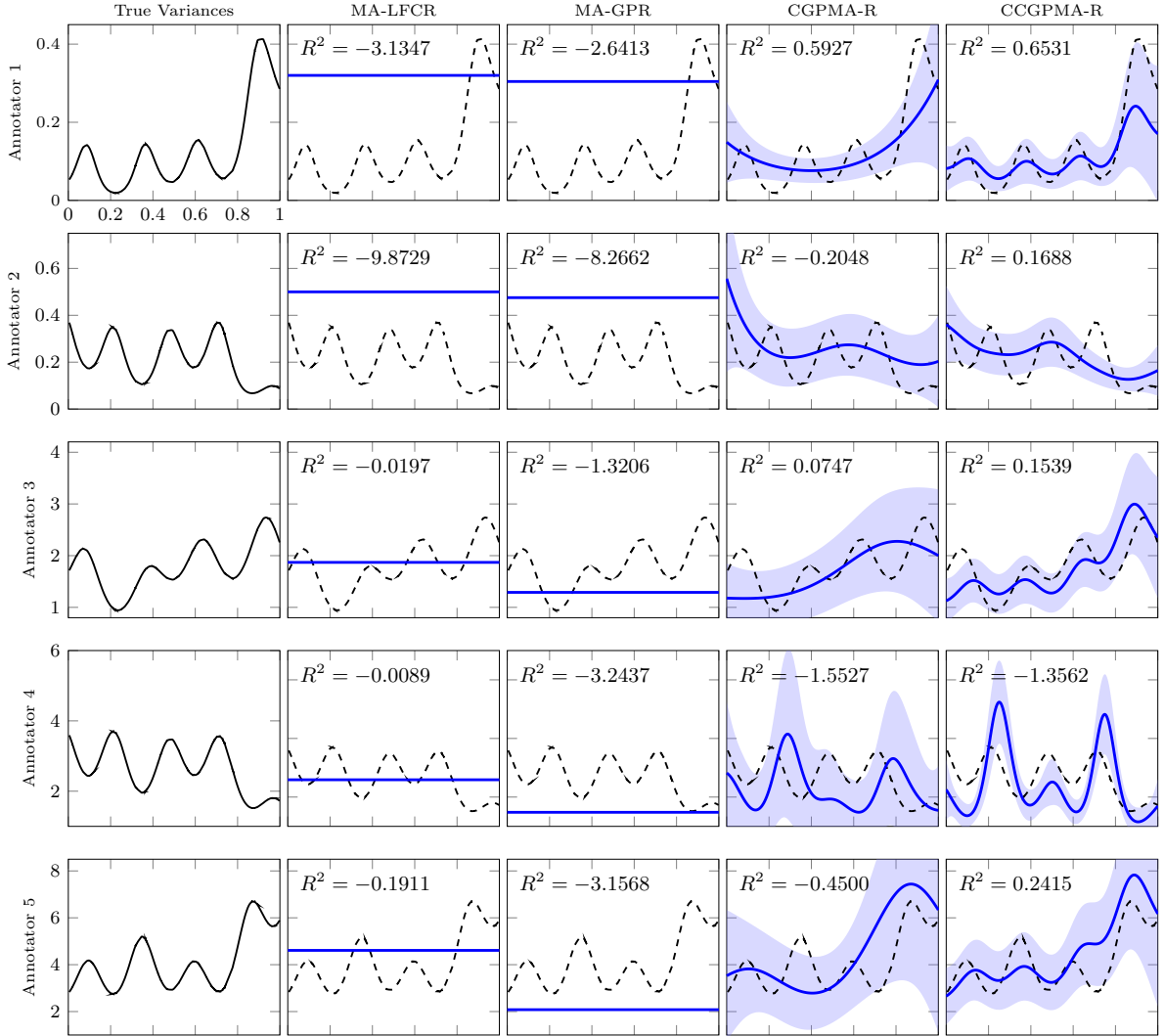


Figure 5.4: Estimated values of error-variance for the five annotators in the *fully synthetic* experiment. In the first column, from top to bottom, we expose the error-variances used to simulate the labels from each annotator. Furthermore, the subsequent columns from top to bottom present the estimation of such error-variances performed by state-of-the-art models that include these kinds of parameters in their formulation; moreover, the true error-variances are provided in dashed lines. The shaded region in CGPMA-R and CCGPMA-R indicates the area enclosed by the mean plus or minus two standard deviations. We remark that there is no shaded region for MA-LFCR, and MA-GPR since these approaches perform a fixed-point estimation for the annotators' parameters. Finally, we remark that the R^2 score between the true and estimated error variances are provided.

CGPMA-R, our CCGPMA-R models the annotators' interdependencies. Secondly, the lower bound GPR-Av exhibits a significantly worse prediction than our approaches. We

Table 5.4: Semi-synthetic regression results. Bold: the highest R^2 excluding the upper bound (target) classifier GPR-GOLD. Marked with *: the highest Pearson coefficient (Pear) except the upper bound.

(a) SLFM-R.

Method		Auto	Bike	Concrete	Housing	Yacht	CT	Average
GPC-GOLD($M=40$)	R^2	0.8604 ± 0.0271	0.5529 ± 0.0065	0.8037 ± 0.0254	0.8235 ± 0.0419	0.8354 ± 0.0412	0.8569 ± 0.0055	0.7888
	Pear	0.9359 ± 0.0132	0.7440 ± 0.0046	0.8997 ± 0.0142	0.9152 ± 0.0209	0.9350 ± 0.0110	0.9268 ± 0.0030	0.8928
GPC-GOLD($M=80$)	R^2	0.8612 ± 0.0279	0.5603 ± 0.0063	0.8271 ± 0.0230	0.8275 ± 0.0399	0.8240 ± 0.0339	0.8648 ± 0.0047	0.7942
	Pear	0.9361 ± 0.0139	0.7490 ± 0.0042	0.9124 ± 0.0126	0.9174 ± 0.0198	0.9281 ± 0.0141	0.9309 ± 0.0025	0.8956
GPR-MV($M=40$)	R^2	0.8425 ± 0.0286	0.5280 ± 0.0100	0.7589 ± 0.0279	0.7834 ± 0.0463	0.7588 ± 0.0498	0.8070 ± 0.0130	0.7464
	Pear	0.9260 ± 0.0155	0.7289 ± 0.0073	0.8765 ± 0.0168	0.8966 ± 0.0265	0.8913 ± 0.0162	0.9067 ± 0.0058	0.8710
GPR-MV($M=80$)	R^2	0.8406 ± 0.0304	0.5397 ± 0.0085	0.7765 ± 0.0274	0.7903 ± 0.0451	0.7676 ± 0.0535	0.8167 ± 0.0089	0.7552
	Pear	0.9246 ± 0.0168	0.7378 ± 0.0056	0.8852 ± 0.0162	0.8995 ± 0.0253	0.8964 ± 0.0184	0.9117 ± 0.0030	0.8759
MA-LFCR	R^2	0.7973 ± 0.0218	0.3385 ± 0.0051	0.6064 ± 0.0384	0.7122 ± 0.0509	0.6403 ± 0.0186	0.8400 ± 0.0014	0.6558
	Pear	0.8953 ± 0.0116	0.5892 ± 0.0058	0.7823 ± 0.0230	0.8497 ± 0.0311	0.8091 ± 0.0137	0.9172 ± 0.0008	0.8071
MA-GPR	R^2	0.8456 ± 0.0281	0.4448 ± 0.0187	0.7769 ± 0.0367	0.7685 ± 0.0632	0.7842 ± 0.1027	0.0105 ± 0.0045	0.6051
	Pear	0.9227 ± 0.0148	0.6694 ± 0.0140	0.8852 ± 0.0209	0.8809 ± 0.0352	0.9007 ± 0.0427	0.1410 ± 0.0072	0.7330
MA-DL-B	R^2	0.7766 ± 0.0253	0.5854 ± 0.0107	0.2319 ± 0.0328	0.5317 ± 0.1005	0.2089 ± 0.0783	0.6903 ± 0.2689	0.5041
	Pear	0.8951 ± 0.0190	0.7712 ± 0.0047	0.5168 ± 0.0327	0.7555 ± 0.0619	0.5123 ± 0.0892	0.9698 ± 0.0138	0.7368
MA-DL-S	R^2	0.7761 ± 0.0279	0.5828 ± 0.0149	0.2363 ± 0.0252	0.5352 ± 0.0948	0.1822 ± 0.0985	0.8418 ± 0.2368	0.5257
	Pear	0.8977 ± 0.0160	0.7736 ± 0.0060*	0.5175 ± 0.0322	0.7540 ± 0.0607	0.4883 ± 0.1189	0.9584 ± 0.0161*	0.7316
MA-DL-B+S	R^2	0.7717 ± 0.0239	0.5816 ± 0.0181	0.2369 ± 0.0322	0.5330 ± 0.0850	0.1974 ± 0.0895	0.5517 ± 0.2316	0.4787
	Pear	0.8936 ± 0.0182	0.7727 ± 0.0071	0.5172 ± 0.0429	0.7537 ± 0.0580	0.5048 ± 0.0969	0.9632 ± 0.0157	0.7342
CGPMA-R($M=40$)	R^2	0.8476 ± 0.0229	0.5464 ± 0.0069	0.8169 ± 0.0231	0.7244 ± 0.2973	0.8049 ± 0.0482	0.8236 ± 0.0132	0.7606
	Pear	0.9280 ± 0.0109	0.7395 ± 0.0047	0.9063 ± 0.0126	0.8234 ± 0.2705	0.9183 ± 0.0235	0.9117 ± 0.0066	0.8712
CGPMA-R($M=80$)	R^2	0.8342 ± 0.0217	0.5560 ± 0.0074	0.8190 ± 0.0254	0.7259 ± 0.3018	0.7928 ± 0.0884	0.8371 ± 0.0104	0.7608
	Pear	0.9212 ± 0.0119	0.7459 ± 0.0051	0.9078 ± 0.0137	0.8093 ± 0.3051	0.9206 ± 0.0348	0.9188 ± 0.0046	0.8706
CCGPMA-R($M=40$)	R^2	0.8558 ± 0.0248	0.5284 ± 0.0117	0.7976 ± 0.0270	0.8169 ± 0.0468	0.8409 ± 0.0548	0.8219 ± 0.0062	0.7769
	Pear	0.9323 ± 0.0115	0.7345 ± 0.0050	0.8960 ± 0.0153	0.9129 ± 0.0263	0.9454 ± 0.0109	0.9123 ± 0.0017	0.8889
CCGPMA-R($M=80$)	R^2	0.8534 ± 0.0243	0.5467 ± 0.0069	0.8220 ± 0.0259	0.8215 ± 0.0466	0.8691 ± 0.0473	0.8252 ± 0.0083	0.7897
	Pear	0.9305 ± 0.0109	0.7443 ± 0.0039	0.9092 ± 0.0140*	0.9150 ± 0.0255*	0.9583 ± 0.0102	0.9138 ± 0.0017	0.8952*

(b) Non-homogeneous error-variance.

Method		Auto	Bike	Concrete	Housing	Yacht	CT	Average
GPC-GOLD($M=40$)	R^2	0.8604 ± 0.0271	0.5529 ± 0.0065	0.8037 ± 0.0254	0.8235 ± 0.0419	0.8354 ± 0.0412	0.8569 ± 0.0055	0.7888
	Pear	0.9359 ± 0.0132	0.7440 ± 0.0046	0.8997 ± 0.0142	0.9152 ± 0.0209	0.9350 ± 0.0110	0.9268 ± 0.0030	0.8928
GPC-GOLD($M=80$)	R^2	0.8612 ± 0.0279	0.5603 ± 0.0063	0.8271 ± 0.0230	0.8275 ± 0.0399	0.8240 ± 0.0339	0.8648 ± 0.0047	0.7942
	Pear	0.9361 ± 0.0139	0.7490 ± 0.0042	0.9124 ± 0.0126	0.9174 ± 0.0198	0.9281 ± 0.0141	0.9309 ± 0.0025	0.8956
GPR-MV($M=40$)	R^2	0.8388 ± 0.0373	0.5348 ± 0.0065	0.7648 ± 0.0262	0.7795 ± 0.0508	0.7752 ± 0.0579	0.8169 ± 0.0103	0.7517
	Pear	0.9277 ± 0.0154	0.7335 ± 0.0055	0.8784 ± 0.0151	0.8921 ± 0.0280	0.9030 ± 0.0152	0.9109 ± 0.0043	0.8743
GPR-MV($M=80$)	R^2	0.8402 ± 0.0396	0.5471 ± 0.0056	0.7819 ± 0.0270	0.7892 ± 0.0468	0.7792 ± 0.0616	0.8268 ± 0.0061	0.7607
	Pear	0.9281 ± 0.0170	0.7421 ± 0.0045	0.8874 ± 0.0154	0.8968 ± 0.0259	0.9061 ± 0.0166	0.9153 ± 0.0030	0.8793
MA-LFCR	R^2	0.7909 ± 0.0210	0.3868 ± 0.0065	0.6050 ± 0.0333	0.7045 ± 0.0610	0.6230 ± 0.0290	0.8605 ± 0.0014	0.6618
	Pear	0.8992 ± 0.0098	0.6223 ± 0.0051	0.7817 ± 0.0198	0.8512 ± 0.0284	0.8024 ± 0.0119	0.9277 ± 0.0007	0.8141
MA-GPR	R^2	0.8597 ± 0.0355	0.4408 ± 0.0169	0.7921 ± 0.0269	0.7446 ± 0.0624	0.7800 ± 0.0565	0.0116 ± 0.0034	0.6048
	Pear	0.9300 ± 0.0180	0.6661 ± 0.0125	0.8915 ± 0.0153	0.8661 ± 0.0348	0.8899 ± 0.0309	0.1415 ± 0.0068	0.7309
MA-DL-B	R^2	0.7681 ± 0.0297	0.5996 ± 0.0105	0.2543 ± 0.0267	0.5264 ± 0.0973	0.1843 ± 0.0820	0.7612 ± 0.2679	0.5156
	Pear	0.8923 ± 0.0196	0.7804 ± 0.0055	0.5395 ± 0.0387	0.7466 ± 0.0572	0.4963 ± 0.0970	0.9806 ± 0.0061	0.7393
MA-DL-S	R^2	0.7687 ± 0.0316	0.5944 ± 0.0097	0.2495 ± 0.0332	0.5279 ± 0.1031	0.1809 ± 0.0963	0.8231 ± 0.3349	0.5241
	Pear	0.8919 ± 0.0171	0.7830 ± 0.0051	0.5312 ± 0.0407	0.7484 ± 0.0604	0.4840 ± 0.1281	0.9801 ± 0.0091	0.7364
MA-DL-B+S	R^2	0.7765 ± 0.0315	0.5909 ± 0.0113	0.2446 ± 0.0379	0.5454 ± 0.0933	0.1889 ± 0.0813	0.6294 ± 0.2498	0.4960
	Pear	0.8954 ± 0.0168	0.7813 ± 0.0051	0.5251 ± 0.0482	0.7598 ± 0.0571	0.4946 ± 0.1163	0.9747 ± 0.0122	0.7385
CGPMA-R($M=40$)	R^2	0.8570 ± 0.0335	0.5462 ± 0.0062	0.8154 ± 0.0265	0.8260 ± 0.0503	0.8025 ± 0.0748	0.8561 ± 0.0062	0.7839
	Pear	0.9338 ± 0.0154*	0.7399 ± 0.0044	0.9055 ± 0.0148	0.9166 ± 0.0253	0.9177 ± 0.0417	0.9264 ± 0.0029	0.8900
CGPMA-R($M=80$)	R^2	0.8137 ± 0.0667	0.5544 ± 0.0069	0.8250 ± 0.0260	0.8302 ± 0.0479	0.7458 ± 0.1670	0.8622 ± 0.0049	0.7719
	Pear	0.9133 ± 0.0298	0.7456 ± 0.0049	0.9107 ± 0.0143*	0.9185 ± 0.0243*	0.9064 ± 0.0748	0.9296 ± 0.0024*	0.8873
CCGPMA-R($M=40$)	R^2	0.8567 ± 0.0342	0.5265 ± 0.0068	0.7980 ± 0.0264	0.8172 ± 0.0548	0.8442 ± 0.0738	0.8006 ± 0.0093	0.7739
	Pear	0.9342 ± 0.0156	0.7340 ± 0.0033	0.8960 ± 0.0151	0.9116 ± 0.0282	0.9566 ± 0.0064	0.9035 ± 0.0024	0.8893
CCGPMA-R($M=80$)	R^2	0.8571 ± 0.0326	0.5428 ± 0.0073	0.8202 ± 0.0257	0.8288 ± 0.0515	0.8687 ± 0.0587	0.8403 ± 0.0099	0.7930
	Pear	0.9335 ± 0.0160	0.7429 ± 0.0045	0.9082 ± 0.0144	0.9183 ± 0.0258	0.9653 ± 0.0136*	0.9205 ± 0.0040	0.8981*

remark on MA-GPR's behavior, which is lowest compared to its GPs-based competitors,

even far worse than GPR-Av. The key to this outcome lies in the MA-GPR formulation; since it models the annotators' behavior by assuming that their performance does not depend on the input features and considering that the labelers make their decisions independently, which does fit the process that we use to simulate the labels. Next, we analyze the results concerning the linear model MA-LFR; attained to the results, we note that this approach's prediction capacity is far lower than ours. The above outcome suggests that there may exist a non-linear structure in most databases. However, we highlight a particular result for the dataset CT, where MA-LFCR exhibits the best performance defeating all its competitors based on non-linear models. From the above, we intuit that the CT dataset may have a linear structure. To confirm this supposition, we perform an additional experiment over CT by training a regression scheme based on LR with the actual labels (we follow the same scheme as for GPR-GOLD). We obtain an R^2 score equal to 0.8541 (on average), which is close to the results obtained by GPR-GOLD. Thus, we can elucidate that there exists a linear structure in the dataset CT. Finally, we analyze the results for the DL-based models. Similar to the experiments over *fully synthetic datasets*, we note a considerable low prediction capacity; in fact, they are even defeated by the linear model MA-LFR. Again, we attribute this behavior to the fact that the CrowdLayer (used to manage the data from multiple annotators) does not offer a suitable codification of the labelers' behavior. Nevertheless, taking the above into account, we observe a remarkable result in the Bike dataset, where the DL-based approaches offer the best performance, even defeating the supposed upper-bound GPR-GOLD. To explain that, it is necessary to analyze the meaning of the target variable in such a dataset. Regarding the description of this dataset,³ the target variables indicate the count of total rental bikes, including both casual and registered in a day. The above suggests that there may exist a quasi-periodic structure in the dataset, which cannot be captured by the GPR-GOLD since it uses a non-periodic kernel (RBF). To support our suppositions, an additional experiment was performed over this dataset by training the model GPR-GOLD with the following kernel:

$$\kappa(\mathbf{x}_n, \mathbf{x}_{n'}) = s^2 \exp \left[-\frac{1}{2} \sum_{p=1}^P \left(\frac{\sin(\frac{\pi}{T_p}(x_{p,n} - x_{p,n'}))}{l_p} \right)^2 \right], \quad (5.37)$$

where $s \in \mathbb{R}$ is the variance parameter, $l_p \in (\mathbb{R}^+)$ is the length-scale parameter for the p -th dimension, and $T_p \in (\mathbb{R}^+)$ is the period for the p -th dimension. Therefore, we

³Such description can be found in <https://archive.ics.uci.edu/ml/datasets/bike+sharing+dataset>

obtain an R^2 score equal to 0.5952 (on average), which is greater than the obtained by the DL-based approaches, indicating a quasi-periodic structure in the Bike dataset as we had supposed.

On the other hand, Table 5.4(b) presents the results concerning the simulation method “Non-homogeneous error-variance”. We highlight that our CCGPMA exhibits the best regression performance in both metrics R^2 score and the Pearson coefficient. Now, analyzing the behavior of GPs-based competitors, we notice that GPR-Av exhibits the best performance compared to MA-GPR. Our explanation for such an outcome lies in the formulation of MA-GPR since it models the annotators’ behavior without considering relationships between the input space and the annotators’ performance. Regarding the linear model, MA-LFR, we note that its prediction performance is considerably lower than our method. The above indicates the presence of a non-linear structure in most considered datasets. Finally, analyzing the DL-based methods, we remark that their performance is the lowest concerning the average of the R^2 score. Our explanation for this result is that the Crowdlayer used for the DL models to manage multi-labeler data does not offer a suitable representation of the annotators’ behavior.

Finally, we use the *fully real datasets*, which present the most challenging scenario, where both the input samples and the labels come from real-world applications. Ta-

Table 5.5: Regression results in terms of R^2 score over *fully real dataset*. Bold: the highest R^2 excluding the upper bound GPR-GOLD.

Method	Music
GPR-GOLD($M = 40$)	0.4704
GPR-GOLD($M = 80$)	0.4889
GPR-Av($M = 40$)	0.2572
GPR-Av($M = 80$)	0.2744
MA-LFCR	0.1404
MA-GPR	0.0090
MA-DL-B	0.2339
MA-DL-S	0.2934
MA-DL-B+S	0.3519
CGPMA-R($M = 40$)	0.3345
CGPMA-R($M = 80$)	0.3531
CCGPMA-R($M = 40$)	0.3337
CCGPMA-R($M = 80$)	0.3872

ble 5.5 outlines the achieved performances. We remark that our CCGPMA-R with $M = 80$ obtains the best generalization performance in terms of the R^2 score. Further, as theoretically expected, its performance lies between that of GPR-GOLD and GP-Av. Moreover, regarding the GPs-based competitors (MA-GPR and CGPMA-R), we note

that our CGPMA-R is just a bit lower than CCGPMA-R. On the other hand, MA-GPR exhibits the worst prediction capability with a R^2 close to zero. We suppose the above is a symptom of overfitting, which can be confirmed due to the training R^2 score for MA-GPR is 0.4731, which is comparable with GPR-GOLD. Conversely, the linear approach MA-LFCR exhibits the second-lowest performance and performs worse than the theoretical lower bound GP-Av, which indicates a non-linear structure in the Music dataset. Finally, analyzing the results from the deep learning approaches, we note that the variation MA-DL-B+S exhibits a similar performance compared to our CGPMA-R; however, it is slightly lower than our CCGPMA-R. We highlight that despite the capacities of deep learning, our approach CCGPMA-R offers a better representation of annotators' behavior, unlike the deep learning techniques, which measure such performance using a single parameter.

5.5 Summary

We propose a novel Gaussian Process-based approach to deal with Multiple Annotators scenarios, termed Correlated Chain Gaussian Process for Multiple Annotators (CCGPMA). Our method is built as an extension of the chained GP (Saul et al., 2016), introducing a semi-parametric latent factor model (SLFM) to exploit correlations between the GP latent functions that model the parameters of a given likelihood function. To the best of our knowledge, CCGPMA is the first attempt to build a probabilistic framework that codes the annotators' expertise as a function of the input data and exploits the correlations among the labelers' answers. Besides, we highlight that our approach can be used with different likelihood, which allows us to deal with both categorical data (classification) and real-valued (regression). We tested our approach for classification and regression tasks using different scenarios concerning the provided annotations: synthetic, semi-synthetic, real-world experts. According to the results, we remark that our CCGPMA can achieve robust predictive properties for the studied datasets, outperforming state-of-the-art methods.

On the other hand, as observed in Section 5.4.2, the kernel selection plays a key role in the performance of our CCGPMA. Hence, An improper kernel selection leads to unsatisfactory performance.

Chapter 6

Final Remarks

6.1 Conclusions

This dissertation presented several specific methodologies based on frequentist and Bayesian models to face supervised learning problems with multiple annotators. In that sense, five strategies were proposed to code the annotators' behavior, considering the relationship between the input features and the labelers' performance and the dependencies among their outputs. The introduced approaches naturally lead to data-dependent processing tuned to the particular samples constraints and to the considered learning scenario, including regression and classification (binary and multi-class) tasks. The introduced framework is tested in both controlled scenarios, where the input samples and the labels from multiple annotators come from real-world problems. Overall, attained results demonstrated that proposed approaches code the labelers' performance taking into account dependencies among the annotators and considering the labelers' performance as a function of the input data, favoring the learning performance in terms of task accuracy and data interpretability in comparison to state-of-the-art methods. Following, the main concluding remarks regarding each provided representation strategy are described:

- A kernel-alignment-based approach that allows coding the annotators' performance in the absence of the ground truth. The introduced strategy termed KAAR can compute each provided labeler's relevance through a CKA-based averaged matching between the given labels and the input data features; it also considers similarities between all annotators' labels. Hence, KAAR relaxes the assumption of independence between the labelers, which allows coding some biases and tenden-

cies between the annotators' opinions. Then, a convex combination of supervised learning algorithms is carried out by adopting the multiple annotator performances coded in the KAAR-based relevance analysis. For the sake of comparison, KAAR is tested in synthetic and real-world datasets for both classification and regression settings. Results show that KAAR can deal with binary classification, multiclass classification, and regression problems with multiple annotators.

- A localized kernel-alignment-based annotator relevance approach, named LKAAR, to support binary, multi-class classification, and regression problems in the presence of multiple annotators, which configures an extension of KAAR aiming to code the relationship between the input features and the labelers' behavior. Accordingly, our LKAAR computes the relevance of each provided expert through a centered kernel alignment-based matching between the annotator labels and the input features, taking into account dependencies among the annotators and considering the labelers' performance as a function of the input data. Like KAAR, the LKAAR-based relevance analysis is used as weights in a combination of classifiers/regressors (one for labeler), which are used to solve the supervised learning problem from multiple annotators. We tested our approach in synthetic, semi-synthetic, and real-world datasets. The results show that the proposed method deals with binary, multi-class classification, and regression problems, where the ground truth is not available. In fact, in most cases, our LKAAR achieves competitive or even better results when compared to different state-of-the-art classification and regression approaches.
- A regularized chained Deep neural network (RCDNN) to support binary and multi-class classification. RCDNN is built based on the concepts of the chained Gaussian processes (Saul et al., 2016), where the parameters of a given multi-labeler likelihood are linked to the output neurons of a deep neural network, where some of those parameters are related to the annotators' performance. Accordingly, RCDNN codes such performance as a function of the input space \mathcal{X} and the annotators' interdependencies. Besides, l1, l2, and Monte-Carlo Dropout regularization strategies are coupled within our RCDNN architecture and predictor to contract the overfitting challenge of deep models. The proposal is tested over binary and multiclass classification settings using different scenarios concerning the provided annotations: synthetic, semi-synthetic, and real-world experts. According to the results, RCDNN achieves robust predictive properties for the studied datasets, out-

performing state-of-the-art methods while providing an estimation of each labeler's reliability and the dependencies among annotators.

- A novel Gaussian Process-based approach to deal with Multiple labelers scenarios termed Correlated Chain Gaussian Process for Multiple Annotators (CCGPMA). Our method is built as an extension of the chained GP (Saul et al., 2016), introducing a semi-parametric latent factor model-(SLFM) to exploit correlations between the GP latent functions that model the parameters of a given likelihood function. Hence, our CCGPMA jointly solves the supervised learning problem and estimates the annotators' behavior. Notably, CCGPMA allows codifying the labelers' performance as a function of the input space while capturing correlations between them. We highlight that our approach can be used with different likelihood functions, which allows us to deal with categorical data (classification) and real-valued (regression). We tested our approach for classification and regression tasks using different scenarios concerning the provided annotations: synthetic, semi-synthetic, real-world experts. According to the results, we remark that our CCGPMA can achieve robust predictive properties for the studied datasets, outperforming state-of-the-art methods.

Then, for comparison purposes, in Table 6.1 we expose the results of our approaches and the state-of-the-art methods in Table 5.1 over three semi-synthetic datasets (Ionosphere, Western, and Segmentation), where the labels were simulated by following the procedure Section 5.3.1, which is the most challenging since it simulates correlated labelers, where their performance depends on the input features. We show the non-parametric Friedman test results to establish the statistical significance of such results. The null hypothesis settles that all algorithms perform equal (Demšar, 2006). Also, we fix the significance threshold as $p < 0.05$.

From Table 6.1, we note that analyzing both the average AUC and the average ranking, the approaches proposed in this thesis are competitive (KAAR, LKAAR) or even outperform (RCDNN, CGPMA, and CCGPMA) the considered state-of-the-art approaches. Besides focusing on our approaches, we notice two groups; the first includes the methods CCGPMA, CGPMA, and RCDNN (exhibit the best performance) that jointly estimates the annotators' performance and the supervised function f . On the other hand, the second group is conformed by the approaches that perform such estimation independently, KAAR and LKAAR (the lowest performance). We highlight that the first group offers better prediction performance in terms of the AUC score. The above is explained

in the sense that as it has been mentioned in different studies (Morales-Alvarez et al., 2020; Morales-Álvarez et al., 2019; Raykar et al., 2010; Ruiz et al., 2019), methods that jointly build the supervised learning algorithm and the annotators’ performance achieve superior performance since the interaction between such factors (labelers’ behavior and the supervised learning model) provides critical information to puzzle out the actual labels. We apply the Friedman test to verify the significance of results in Table 6.1. As seen, we obtain a Chi-square of 26.84 with $p - value = 0.0084$. Thus, we have enough statistical evidence to determine that our CCGPMA exhibits a better performance than state-of-the-art competitors.

Table 6.1: AUC classification results for the semi-synthetic and fully real datasets. Bold: the highest performance excluding the GPC-GOLD bound. The last column presents the average ranking for AUC score, in bold the best average ranking. The Friedman test returns a Chi-square value of 26.84 ($p - value = 0.0081$).

Method	Ionosphere	Western	Segmentation	G	Voice R	B	Music	Average AUC	Average ranking
GPC-GOLD	0.9513	0.9250	0.9781	0.9484	0.9484	0.9484	0.9358	0.9479	--
GPC-MV	0.7779	0.8658	0.9749	0.9301	0.9377	0.8001	0.8897	0.8823	6.1428
MA-LFC-C	0.7358	0.8400	0.9892	0.9122	0.9130	0.8406	0.8599	0.8701	8.0000
MA-DGRL	0.6453	0.8143	0.9897	0.9127	0.9164	0.8259	0.8832	0.8553	8.7142
MA-GPC	0.6631	0.8671	0.9734	0.8660	0.8597	0.4489	0.8253	0.7862	11.142
MA-GPCV	0.6238	0.8451	0.9924	0.9283	0.9208	0.8835	0.8677	0.8659	7.4285
MA-DL-MW	0.7535	0.9092	0.9950	0.8957	0.8966	0.8123	0.8567	0.8741	7.8571
MA-DL-VW	0.6987	0.9173	0.9972	0.8942	0.8929	0.8092	0.9167	0.8751	7.4285
MA-DL-VW+B	0.7196	0.9164	0.9972	0.9030	0.8937	0.8218	0.8573	0.8727	7.5714
KAAR	0.7046	0.8588	0.9217	0.9109	0.9351	0.8969	0.8896	0.8739	7.8571
LKAAR	0.7121	0.8359	0.9355	0.9206	0.9360	0.8979	0.8998	0.8768	7.1428
RCDNN	0.6818	0.8520	0.9634	0.9224	0.9419	0.9257	0.9329	0.8885	5.7142
CGPMA-C	0.8615	0.9185	0.9679	0.9324	0.9417	0.8708	0.9025	0.9136	4.0000
CCGPMA-C	0.9023	0.9307	0.9774	0.9318	0.9422	0.9002	0.9456	0.9328	2.1428

Finally, we intend to analyze the impact of malicious annotators on the performance of our approaches and compared them with some state-of-the-art methods. Similar to Section 4.4, we employ the Pima dataset with 768 instances; thus, we use 538 samples for training and the remaining 230 for testing. We create synthetic labels from 5 annotators generated from the Biased coin (Non-homogeneous) procedure (see Section 3.3.1 and equation (4.10)). According to Figure 6.1 (left), we notice that from the 5 labelers (red dots), two are categorized as suitable labelers, one as Spammers and the remaining as Malicious. Then, we simulate 20 additional annotators to analyze the behavior of our approaches in scenarios where the number of malicious annotators increases. The labels from malicious annotators are generated as follows (see green dots on left of Figure 6.1): a random number $\alpha_n^{(r)}$ is sampled from a Bernoulli distribution with parameter $p_r = 0.6$; then if $\alpha_n^{(r)} = 0$, $y_n^{(r)} = y_n$, and $y_n^{(r)} = \tilde{y}_n$ otherwise.

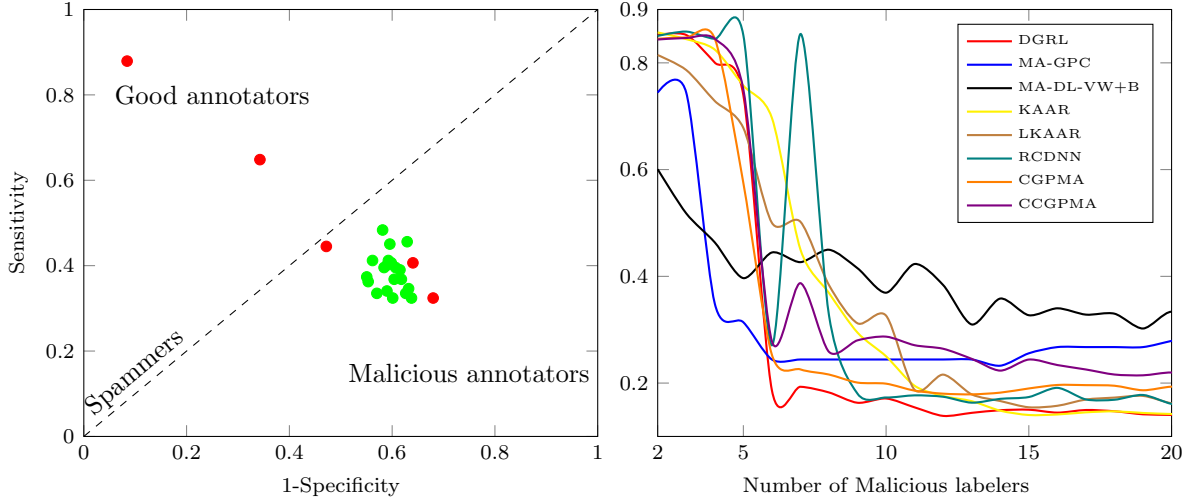


Figure 6.1: Malicious annotators. On the left, we show the Receiver Operating Characteristic (ROC) plot for the simulated annotators. Red dots indicate the basis annotators. Green dots show malicious annotators. On the right, we present the performance (AUC) of DGRL, MA-GPC, MA-DL-VW+B, KAAR, LKAAR, RCDNN, CGPMA, and CCGPMA as a function of the number of malicious annotators.

From the results in Figure 6.1 (right), we note that our approaches are significantly affected when the number of malicious annotators is great than 5 malicious annotators, which suggest the presence of a critical point is presented when the percentages of good, spammers and malicious labelers are respectively 25%, 12.5%, 62.5%. Accordingly, for this experiment, we can affirm that our approach can deal with malicious labelers if the percentage of them is below 62.5%. On the other hand, regarding the state-of-the-art competitors MA-DGRL, MA-GPC, and MA-DL-VW+B, we notice that all of them (excepting MA-DGRL) are more susceptible to the presence of malicious labelers; in fact, their AUC is affected whit three or more malicious annotators. Conversely, we highlight the performance of MA-DGRL, which is very close to our methods. To explain such an outcome, we remark that MA-GPC and MA-DL-VW+B estimate the annotators' performance as an average of some parameters. However, MA-DGRL estimates such performance for every region in the input space, which makes it more robust to the presence of malicious annotators.

6.2 Future work

We have presented a framework based on frequentist and Bayesian models aiming to face supervised learning problems with data from multiple annotators. However, from the achieved theoretical and experimental results, there are still many issues that can be addressed to improve the learning performance. In particular, the following remarks could be of interest for future work approaches:

- KAAR and LKAAR independently compute the annotators' behavior and the supervised learning function f . In fact, such a function f is estimated by using a convex combination of R classifiers/regressors, which could be inconvenient for a large number of labelers. In that sense, an interesting extension of such works could be to design a model that jointly relates the supervised learning function and the annotators' behavior based on LKAAR or KAAR.
- Concerning RCDNN, extending RCDNN for regression tasks is an exciting research line, i.e., based on the model introduced in [Groot et al. \(2011\)](#). Next, we plan to use other deep structures, i.e., Convolutional and Recurrent layers and different activation functions, to apply our approach in more complex tasks such as computer vision or natural language processing.
- Moreover, regarding the RCDNN model, we recall that it only was proposed to deal with classification (binary and multi-class) tasks. Accordingly, we plan to extend such an approach for real-valued data aiming to deal with interesting applications such as text regression [Rodrigues and Pereira \(2018a\)](#).
- Regarding CCGPMA, we can use convolution processes ([Álvarez and Lawrence, 2011](#)) instead of the SLFM, aiming to obtain a better representation of the correlations among the labelers. Also, our approach can be extended for multi-task learning in the context of multiple annotators ([Álvarez et al., 2012](#)). Accordingly, the likelihood for multi-task learning for multiple annotators can be given as:

$$p(\mathbf{Y}_1, \dots, \mathbf{Y}_E | \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_E) = \prod_{e=1}^E \prod_{n=1}^N \prod_{r_e \in R_n^e} p(y_{r_e, n}^e | \theta_1(\mathbf{x}_n), \dots, \theta_{J_e}(\mathbf{x}_n)), \quad (6.1)$$

$\mathbf{Y}_e \in \mathbb{R}^{N \times R_e}$ is a matrix containing the labels $y_{r_e, n}^e$ given by R_e annotators for the e -th task. Further, $\boldsymbol{\theta}_e = [\boldsymbol{\theta}_1^e, \dots, \boldsymbol{\theta}_{J_e}^e]^\top \in \mathbb{R}^{N J}$, and $\boldsymbol{\theta}_j^e = [\theta_j^e(\mathbf{x}_1), \dots, \theta_j^e(\mathbf{x}_N)]^\top \in \mathbb{R}^N$ is the j -th parameter corresponding to the e -th task.

- We recognize an increasing interest in multi-labeler approaches from the Bioinformatics community; particularly, in applications related to computer-aided diagnosis (CAD) systems (G. Rodrigo et al., 2019; Raykar et al., 2010), such as predicting malignancy of pulmonary nodules, mitosis detection in breast cancer Albarqouni et al. (2016), and the assessment of voice quality Gil et al. (2015). In this type of system, the medical diagnosis (labels) represents the health risk levels, which are inherently ordered. Accordingly, an ordinal-regression problem is configured (Liu et al., 2018), which cannot be tackled by using multi-class classification or regression methods. In that sense, it is interesting to extend our approaches to deal with ordinal data aiming to contribute to the Bioinformatics area.
- Following the previous idea, we remark that in this work, we used a dataset termed “Voice”, which comprises the evaluation of voice records following the perceptual scale GRBAS. Such a dataset configures an ordinal regression problem; however, we convert it into a binary classification problem for validation purposes (gold standard’s availability). In that sense, it is interesting to apply our approaches to deal with the Voice dataset and building a system for the automatic assessment of voice quality based on the GRBAS scale in the presence of multiple annotators.
- Finally, we recall that this dissertation is focused on dealing with the training of SL algorithms from noisy labels adopting the *learning from crowds* point of view. However, this is not the only perspective to cast the noisy labels drawback. Specifically, we recognize the *self-learning* approaches (Han et al., 2019), which correspond to a pre-training step within a DNN architecture to correct possibly noisy labels or the missing labels imputation. *Self-learning* approaches have been applied in areas such as computer vision (Zhai et al., 2019), and natural language processing (Artetxe et al., 2018). Thereby, we establish that it is interesting to develop a model that combines the concepts of *self-learning* and *learning from crowds* to face SL scenarios with noisy labels.

6.3 Academic discussion

6.3.1 Journal and conference papers

1. Gil-González, J., Valencia-Duque, A., Álvarez-Meza, A., Orozco-Gutiérrez, Á., & García-Moreno, A. (2021). *Regularized Chained Deep Neural Network Classifier*

- for Multiple Annotators*. Applied Sciences, 11(12), 5409.. (Q1-A1).
2. Gil-Gonzalez, J., Giraldo, J., Orozco-Gutierrez, A., Alvarez-Meza, A. & Álvarez, M.A. *Correlated Chained Gaussian Processes for Datasets with Multiple Annotators*. IEEE transactions on Neural networks and learning systems (submitted). (Q1-A1).
 3. Gil-Gonzalez, J., Orozco-Gutierrez, A., & Alvarez-Meza, A. (2021). *Learning from multiple inconsistent and dependent annotators to support classification tasks*. Neurocomputing, 423, 236-247. (Q1-A1).
 4. Gil-Gonzalez, J., Alvarez-Meza, A., & Orozco-Gutierrez, A. (2018). *Learning from multiple annotators using kernel alignment*. Pattern Recognition Letters, 116, 150-156. (Q1-A1).
 5. González, J. G., Álvarez, A. M., & Orozco, Á. A. (2018, November). *Gaussian Processes Regression with Multiple Annotators: When the Annotator Performance Is Not Homogeneous*. In Iberoamerican Congress on Pattern Recognition (pp. 91-99). Springer, Cham. (International conference).
 6. Gil-González, J., Álvarez-Meza, A., Echeverry-Correa, J., Orozco-Gutiérrez, A., & Álvarez-López, M. (2017, May). *Enhancement of nerve structure segmentation by a correntropy-based pre-image approach*. TecnoLógicas, 20(39), 199-210. (A2)
 7. González, J. G., Álvarez, A. M., Valencia, A. F., & Orozco, Á. A. (2017, November). *Automatic Peripheral Nerve Segmentation in Presence of Multiple Annotators*. In Iberoamerican Congress on Pattern Recognition (pp. 246-254). Springer, Cham. (International Conference)

6.3.2 Repositories

Model	Github repository
KAAR	https://github.com/juliangilg/KAAR
LKAAR	https://github.com/juliangilg/LKAAR
RCDNN	https://github.com/juliangilg/RCDNN_MA
CCGPMA	https://github.com/juliangilg/CCGPMA

Appendix A

CCGPMA supplementary material

A.1 Derivation of CCGPMA lower bounds

$$\begin{aligned}\log p(\mathbf{Y}) &= \log \int p(\mathbf{Y}|\hat{\mathbf{f}})p(\hat{\mathbf{f}}|\mathbf{u})p(\mathbf{u})d\hat{\mathbf{f}}d\mathbf{u}, \\ &= \log \int \prod_{j=1}^J p(\mathbf{f}_j|\mathbf{u})q(\mathbf{u}) \frac{p(\mathbf{Y}|\hat{\mathbf{f}}) \prod_{j=1}^J p(\mathbf{f}_j|\mathbf{u})p(\mathbf{u})}{\prod_{j=1}^J p(\mathbf{f}_j|\mathbf{u})q(\mathbf{u})} d\hat{\mathbf{f}}d\mathbf{u}.\end{aligned}\tag{A.1}$$

Using Jensen's inequality, we obtain

$$\log p(\mathbf{Y}) \geq \int \left(\prod_{j=1}^J p(\mathbf{f}_j|\mathbf{u})q(\mathbf{u}) \right) \log \left(\frac{p(\mathbf{Y}|\hat{\mathbf{f}}) \prod_{j=1}^J p(\mathbf{f}_j|\mathbf{u})p(\mathbf{u})}{\prod_{j=1}^J p(\mathbf{f}_j|\mathbf{u})q(\mathbf{u})} \right) d\hat{\mathbf{f}}d\mathbf{u},\tag{A.2}$$

$$\mathcal{L} = \prod_{j=1}^J \int q(\mathbf{f}_j) \log p(\mathbf{Y}|\hat{\mathbf{f}})d\hat{\mathbf{f}} - \sum_{q=1}^Q \mathbb{D}_{KL}(q(\mathbf{u}_q)||p(\mathbf{u}_q)),\tag{A.3}$$

where $\mathbb{D}_{KL}(\cdot||\cdot)$ is the Kullback-Leibler-(KL) divergence and $q(\mathbf{f}_j)$ is defined as follows:

$$q(\mathbf{f}_j) = \int p(\mathbf{f}_j|\mathbf{u})q(\mathbf{u})d\mathbf{u} = \mathcal{N}(\mathbf{f}_j|\boldsymbol{\mu}_j, \mathbf{S}_j),\tag{A.4}$$

and

$$\boldsymbol{\nu}_j = \mathbf{K}_{f_j u} \mathbf{K}_{uu}^{-1} \mathbf{m}, \quad (\text{A.5})$$

$$\mathbf{S}_j = \mathbf{K}_{f_j f_j} + \mathbf{K}_{f_j u} \mathbf{K}_{uu}^{-1} (\mathbf{V} - \mathbf{K}_{uu}) \mathbf{K}_{uu}^{-1} \mathbf{K}_{u f_j}. \quad (\text{A.6})$$

Moreover, note that the solution of the variational expectation (analytical or numerical), will depend on the form of the likelihood $p(\mathbf{Y}|\hat{\mathbf{f}})$:

$$\mathbb{E}_{\prod_{j=1}^J q(f_j(\mathbf{x}_n))} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] = \prod_{j=1}^J \int q(\mathbf{f}_j) \log p(\mathbf{Y}|\hat{\mathbf{f}}) d\hat{\mathbf{f}}, \quad (\text{A.7})$$

For numerical solution, we will use Gaussian-Hermite quadratures. Thus, considering for simplicity an univariate case, the expectations can be approximated as:

$$\mathbb{E}_{q(f_j)} [\log p(\mathbf{Y}|f_j)] \approx \frac{1}{\sqrt{\pi}} \sum_{s=1}^S o_s \log p(\mathbf{Y}|\sqrt{2\sigma_j}g_s + \nu_j), \quad (\text{A.8})$$

where ν_j and σ_j are respectively the mean and the variance of the distribution $q(f_j)$. Besides, the pair of values o_s and g_s are obtained from the Hermite polynomial $H_n(x) = (-1)^n e^{x^2} \frac{d}{dx^n} e^{-x^2}$.

A.1.1 Gradients w.r.t. the variational parameters

Here, we aim to maximize the lower bound in equation (A.3) with respect to the variational parameters \mathbf{m}_q and \mathbf{V}_q of each distribution $q(\mathbf{u}_q)$.

Useful Equalities: These equalities are useful to compute the gradients [Hensman et al. \(2015\)](#):

$$\frac{\partial}{\partial \mu} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [h(x)] = \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} \left[\frac{\partial}{\partial x} h(x) \right]. \quad (\text{A.9})$$

$$\frac{\partial}{\partial \sigma^2} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [h(x)] = \frac{1}{2} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} \left[\frac{\partial^2}{\partial x^2} h(x) \right]. \quad (\text{A.10})$$

First, the bound derivatives w.r.t \mathbf{m}_q , yields:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{m}_q} = \underbrace{\frac{\partial}{\partial \mathbf{m}_q} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]}_{\text{Variational Expectation-(VE) part}} - \underbrace{\frac{\partial}{\partial \mathbf{m}_q} \sum_{q=1}^Q \mathbb{D}_{KL}(q(\mathbf{u}_q)||p(\mathbf{u}_q))}_{\text{KL part}}, \quad (\text{A.11})$$

where the KL part w.r.t. \mathbf{m}_q is defined as:

$$\frac{\partial}{\partial \mathbf{m}_q} \mathbb{D}_{KL}(q(\mathbf{u}_q)||p(\mathbf{u}_q)) = \mathbf{K}_{\mathbf{u}_q \mathbf{u}_q}^{-1} \mathbf{m}_q. \quad (\text{A.12})$$

Now, the VE part w.r.t. \mathbf{m}_q is given as:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{m}_q} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] &= \frac{\partial}{\partial \mathbf{m}_q} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})], \\ &\text{Chain Rule: } \boldsymbol{\nu}_j = h(\mathbf{m}_q), \\ &= \underbrace{\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial \hat{\mathbf{f}}} \log p(\mathbf{Y}|\hat{\mathbf{f}}) \right]}_{\text{See Likelihoods}} \frac{\partial \boldsymbol{\nu}_j}{\partial \mathbf{m}_q}. \end{aligned} \quad (\text{A.13})$$

Second, the bound derivatives w.r.t \mathbf{V}_q are defined as:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{V}_q} = \underbrace{\frac{\partial}{\partial \mathbf{V}_q} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]}_{\text{VE part}} - \underbrace{\frac{\partial}{\partial \mathbf{V}_q} \sum_{q=1}^Q \mathbb{D}_{KL}(q(\mathbf{u}_q)||p(\mathbf{u}_q))}_{\text{KL part}}, \quad (\text{A.14})$$

where the KL part w.r.t. \mathbf{V}_q , yields:

$$\frac{\partial}{\partial \mathbf{V}_q} \mathbb{D}_{KL}(q(\mathbf{u}_q)||p(\mathbf{u}_q)) = \frac{1}{2} (\mathbf{K}_{\mathbf{u}_q \mathbf{u}_q}^{-1} - \mathbf{V}_q^{-1}).$$

Now, the VE part w.r.t. \mathbf{V}_q is given as:

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{V}_q} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] &= \frac{\partial}{\partial \mathbf{V}_q} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})], \\
&\text{Chain Rule: } \mathbf{s}_j = h(\mathbf{V}_q), \\
&= \frac{1}{2} \mathbb{E}_{q(\hat{\mathbf{f}})} \underbrace{\left[\frac{\partial^2}{\partial \hat{\mathbf{f}}^2} \log p(\mathbf{Y}|\hat{\mathbf{f}}) \right]}_{\text{See Likelihoods}} \frac{\partial \mathbf{s}_j}{\partial \mathbf{V}_q}, \tag{A.15}
\end{aligned}$$

where $\mathbf{s}_j \in \mathbb{R}^N$ is a vector containing the diagonal elements of matrix \mathbf{S}_j . On the other hand, the gradients w.r.t. the hyper-parameters are similar to the ones derived in [Moreno-Muñoz et al. \(2018\)](#).

A.2 Likelihood functions

The models exposed previously accept a wide variety of likelihood functions $p(\mathbf{Y}|\hat{\mathbf{f}})$ ([Saul et al., 2016](#)). To incorporate any new likelihood, we need to compute the following expressions:

1. The likelihood $p(\mathbf{Y}|\hat{\mathbf{f}})$.
2. The Log-likelihood $\log p(\mathbf{Y}|\hat{\mathbf{f}})$.
3. The variational expectation of the Log-likelihood $\mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$.
4. The first order derivatives $\frac{\partial}{\partial \hat{\mathbf{f}}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$.
5. The variational expectation of the first order derivatives $\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial \hat{\mathbf{f}}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right]$.
6. The second order derivatives $\frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$.
7. The variational expectation of the second order derivatives $\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right]$.
8. The predictive mean and variance.

A.2.1 Multiclass classification with multiple annotators

To model categorical data from multiple annotators using our CCGPMA, we use the framework proposed in [Rodrigues et al. \(2013\)](#), which introduces a binary variable $\lambda_n^{(r)} \in \{0, 1\}$ representing the r -th labeler's reliability as a function of each sample \mathbf{x}_n .

If $\lambda_n^{(r)} = 1$, the r -th annotator is supposed to provide the actual label, yielding to a categorical distribution. Conversely, $\lambda_n^{(r)} = 0$ indicates that the r -th annotator gives an incorrect output, which is modeled by a uniform distribution.

1. $p(\mathbf{Y}|\boldsymbol{\theta})$ is written as:

$$p(\mathbf{Y}|\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{r \in R_n} \left(\prod_{k=1}^K \zeta_{k,n}^{\delta(y_n^{(r)}, k)} \right)^{\lambda_n^{(r)}} \left(\frac{1}{K} \right)^{(1-\lambda_n^{(r)})},$$

where $\delta(y_n^{(r)}, k) = 1$, if $y_n^{(r)} = k$, and $\delta(y_n^{(r)}, k) = 0$ otherwise. Moreover, $\zeta_{k,n} = p(y_n^{(r)} = k | \lambda_n^{(r)} = 1)$; note that, since $\lambda_n^{(r)} = 1$, $\zeta_{k,n} = p(y_n = k)$ is an estimation of the unknown ground truth. Accordingly, to use our CCGPMA over this model, we need $J = K + R$ LFs; K LFs to model $\zeta_{k,n}$, yielding:

$$\zeta_{k,n} = \Xi(f_k(\mathbf{x}_n)) = \frac{\exp(f_k(\mathbf{x}_n))}{\sum_{j=1}^K \exp(f_j(\mathbf{x}_n))}. \quad (\text{A.16})$$

Besides, we need R LFs to perform a soft estimation of each $\lambda_n^{(r)}$; therefore, $\lambda_n^{(r)} = \varsigma(f_{l_r}(\mathbf{x}_n))$ with $r \in \{1, \dots, R\}$, and $l_r = K + r \in \{K + 1, \dots, J\}$. Here, $\varsigma(f_{l_r}(\mathbf{x}_n))$ is the logistic sigmoid function: $\varsigma(a) = 1/(1 + \exp(-a))$, where $a \in \mathbb{R}$.

2. $\log p(\mathbf{Y}|\hat{\mathbf{f}})$

$$\begin{aligned} \log p(\mathbf{Y}|\hat{\mathbf{f}}) &= \log \prod_{n=1}^N \prod_{r \in R_n} \left(\prod_{k=1}^K \zeta_{k,n}^{\delta(y_n^{(r)}, k)} \right)^{\lambda_n^{(r)}} \left(\frac{1}{K} \right)^{(1-\lambda_n^{(r)})}, \\ &= \sum_{n=1}^N \sum_{r \in R_n} \left\{ \lambda_n^{(r)} \sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + (1 - \lambda_n^{(r)}) \log\left(\frac{1}{K}\right) \right\}. \end{aligned} \quad (\text{A.17})$$

3. $\mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$, yields:

$$\begin{aligned} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] &= \int \prod_{j=1}^J q(\hat{\mathbf{f}}_j) \log p(\mathbf{Y}|\hat{\mathbf{f}}) d\hat{\mathbf{f}}, \\ &= \sum_{n=1}^N \sum_{r \in R_n} \mathbb{E}_{q(f_{l_r,n})} [\lambda_n^{(r)}] \sum_{k=1}^K \delta(y_n^{(r)}, k) \mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\log(\zeta_{k,n})] + \dots \\ &\quad \dots + (1 - \mathbb{E}_{q(f_{l_r,n})} [\lambda_n^{(r)}]) \log\left(\frac{1}{K}\right), \end{aligned} \quad (\text{A.18})$$

where

$$\mathbb{E}_{q(f_{l_r,n})}[\lambda_n^{(r)}] = \int q(f_{l_r,n}) \varsigma(f_{l_r,n}) df_{l_r,n}. \quad (\text{A.19})$$

$$\mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\log(\zeta_{k,n})] = \int q(f_{1,n}) \dots q(f_{K,n}) \log(\Xi(f_k(\mathbf{x}_n))) df_{1,n} \dots df_{K,n}. \quad (\text{A.20})$$

The previous integrals have not a close solution; hence, we approximate them by using the Gaussian-Hermite quadratures.

$$\mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\log(\zeta_{k,n})] \approx \left(\frac{1}{\sqrt{\pi}} \right)^K \sum_{s_K=1}^{S_1} \dots \sum_{s_1=1}^{S_K} o_{s_K} \dots o_{s_1} \log(\Xi(g_{s_K} \sqrt{2\sigma_{k,n}} + \iota_{k,n})), \quad (\text{A.21})$$

$$\mathbb{E}_{q(f_{l_r,n})}[\lambda_n^{(r)}] \approx \frac{1}{\sqrt{\pi}} \sum_{s=1}^S o_s \varsigma(g_s \sqrt{2\sigma_{l_r,n}} + \iota_{l_r,n}), \quad (\text{A.22})$$

where

$$\Xi(g_{s_K} \sqrt{2\sigma_{k,n}} + \iota_{k,n}) = \frac{\exp(g_{s_K} \sqrt{2\sigma_{k,n}} + \iota_{k,n})}{\sum_{i=1}^K \exp(g_{s_i} \sqrt{2\sigma_{i,n}} + \iota_{i,n})}, \quad (\text{A.23})$$

$$\varsigma(g_s \sqrt{2\sigma_{l_r,n}} + \iota_{l_r,n}) = \frac{1}{1 + \exp(-(g_s \sqrt{2\sigma_{l_r,n}} + \iota_{l_r,n}))}. \quad (\text{A.24})$$

$$4. \frac{\partial}{\partial \hat{\mathbf{f}}} [\log p(\mathbf{Y} | \hat{\mathbf{f}})]$$

$$\frac{\partial}{\partial \hat{\mathbf{f}}} [\log p(\mathbf{Y} | \hat{\mathbf{f}})] = \frac{\partial}{\partial \hat{\mathbf{f}}} \left[\sum_{n=1}^N \sum_{r \in R_n} \left\{ \lambda_n^{(r)} \sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + (1 - \lambda_n^{(r)}) \log\left(\frac{1}{K}\right) \right\} \right] \quad (\text{A.25})$$

Case 1

$$\begin{aligned} \frac{\partial}{\partial f_{k,n}} \left[\sum_{n=1}^N \sum_{r \in R_n} \left\{ \lambda_n^{(r)} \sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) \right\} \right] &= \sum_{r \in R_n} \lambda_n^{(r)} \frac{\partial}{\partial f_{k,n}} \left[\sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) \right], \\ &= \sum_{r \in R_n} \lambda_n^{(r)} \left(\delta(y_n^{(r)}, k) - \zeta_{k,n} \right). \end{aligned} \quad (\text{A.26})$$

Case 2

$$\frac{\partial}{\partial f_{l_r,n}} \left[\sum_{n=1}^N \sum_{r \in R_n} \left\{ \lambda_n^{(r)} \sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + (1 - \lambda_n^{(r)}) \log\left(\frac{1}{K}\right) \right\} \right] = \begin{cases} \Omega, & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases},$$

where

$$\begin{aligned}\Omega &= \frac{\partial \lambda_n^{(r)}}{\partial f_{l_r, n}} \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k, n}) \right) - \frac{\partial \lambda_n^{(r)}}{\partial f_{l_r, n}} \left(\log \left(\frac{1}{K} \right) \right), \\ &= \left(\lambda_n^{(r)} - (\lambda_n^{(r)})^2 \right) \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k, n}) + \log(K) \right).\end{aligned}\quad (\text{A.27})$$

$$5. \mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial \hat{\mathbf{f}}} \left[\log p(\mathbf{Y} | \hat{\mathbf{f}}) \right] \right]$$

Case 1

$$\begin{aligned}\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial f_{k, n}} \left[\log p(\mathbf{Y} | \hat{\mathbf{f}}) \right] \right] &= \sum_{r \in R_n} \mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n}) q(f_{l_r, n})} \left[\lambda_n^{(r)} \left(\delta(y_n^{(r)}, k) - \zeta_{k, n} \right) \right], \\ &= \sum_{r \in R_n} \mathbb{E}_{q(f_{l_r, n})} \left[\lambda_n^{(r)} \right] \left(\delta(y_n^{(r)}, k) - \mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n})} [\zeta_{k, n}] \right),\end{aligned}\quad (\text{A.28})$$

where

$$\mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n})} [\zeta_{k, n}] \approx \left(\frac{1}{\sqrt{\pi}} \right)^K \sum_{s_K=1}^{S_1} \dots \sum_{s_1=1}^{S_K} o_{s_K} \dots o_{s_1} \Xi(g_{s_K} \sqrt{2\sigma_{k, n}} + \iota_{k, n}), \quad (\text{A.29})$$

and $\mathbb{E}_{q(f_{l_r, n})} [\lambda_n^{(r)}]$ is approximated using equation (A.22),

Case 2

$$\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial f_{l_r, n}} \left[\log p(\mathbf{Y} | \hat{\mathbf{f}}) \right] \right] = \begin{cases} \mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n}) q(f_{l_r, n})} [\Omega], & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases},$$

where

$$\begin{aligned}\mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n}) q(f_{l_r, n})} [\Omega] &= \\ &= \mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n}) q(f_{l_r, n})} \left[\left(\lambda_n^{(r)} - (\lambda_n^{(r)})^2 \right) \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k, n}) + \log(K) \right) \right], \\ &= \left(\mathbb{E}_{q(f_{l_r, n})} [\lambda_n^{(r)}] - \mathbb{E}_{q(f_{l_r, n})} [(\lambda_n^{(r)})^2] \right) \times \dots \\ &\quad \dots \times \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n})} [\log(\zeta_{k, n})] + \log(K) \right),\end{aligned}\quad (\text{A.30})$$

where $\mathbb{E}_{q(f_{1, n}) \dots q(f_{K, n})} [\log(\zeta_{k, n})]$, and $\mathbb{E}_{q(f_{l_r, n})} [\lambda_n^{(r)}]$ are approximated by equa-

tions (A.21) and (A.22) respectively, and

$$\mathbb{E}_{q(f_{l_r,n})}[(\lambda_n^{(r)})^2] = \int q(f_{l_r,n}) (\sigma(f_{l_r,n}))^2 df_{l_r,n} \approx \frac{1}{\sqrt{\pi}} \sum_{s=1}^S o_s [\varsigma (g_s \sqrt{2\sigma_{l_r,n}} + \iota_{l_r,n})]^2. \quad (\text{A.31})$$

$$6. \frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$$

$$\frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] = \frac{\partial^2}{\partial \hat{\mathbf{f}}^2} \left[\sum_{n=1}^N \sum_{r \in R_n} \left\{ \lambda_n^{(r)} \sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + (1 - \lambda_n^{(r)}) \log\left(\frac{1}{K}\right) \right\} \right]. \quad (\text{A.32})$$

Case 1

$$\begin{aligned} \frac{\partial^2}{\partial f_{k,n}^2} \left[\sum_{n=1}^N \sum_{r \in R_n} \left\{ \lambda_n^{(r)} \sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + (1 - \lambda_n^{(r)}) \log\left(\frac{1}{K}\right) \right\} \right] &= \\ &= \sum_{r \in R_n} \lambda_n^{(r)} \frac{\partial}{\partial f_{k,n}} [\delta(y_n^{(r)}, k) - \zeta_{k,n}] , \\ &= - \sum_{r \in R_n} \lambda_n^{(r)} (\zeta_{k,n} - \zeta_{k,n}^2) \end{aligned} \quad (\text{A.33})$$

Case 2

$$\frac{\partial^2}{\partial f_{l_r,n}^2} \left[\sum_{n=1}^N \sum_{r \in R_n} \left\{ \lambda_n^{(r)} \sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + (1 - \lambda_n^{(r)}) \log\left(\frac{1}{K}\right) \right\} \right] = \begin{cases} \frac{\partial \Omega}{\partial f_{l_r,n}}, & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases},$$

where

$$\begin{aligned} \frac{\partial \Omega}{\partial f_{l_r,n}} &= \frac{\partial}{\partial f_{l_r,n}} [(\lambda_n^{(r)} - (\lambda_n^{(r)})^2)] \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + \log(K) \right), \\ &= \left(2(\lambda_n^{(r)})^3 - 3(\lambda_n^{(r)})^2 + \lambda_n^{(r)} \right) \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + \log(K) \right). \end{aligned} \quad (\text{A.34})$$

$$7. \mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right]$$

Case 1

$$\begin{aligned}
\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial^2}{\partial f_{k,n}^2} \left[\log p(\mathbf{Y}|\hat{\mathbf{f}}) \right] \right] &= - \sum_{r \in R_n} \mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n}) q(f_{l_r,n})} \left[\lambda_n^{(r)} (\zeta_{k,n} - \zeta_{k,n}^2) \right], \\
&= - \sum_{r \in R_n} \mathbb{E}_{q(f_{l_r,n})} \left[\lambda_n^{(r)} \right] \mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\zeta_{k,n}] + \dots \\
&\quad \dots + \sum_{r \in R_n} \mathbb{E}_{q(f_{l_r,n})} \left[\lambda_n^{(r)} \right] \mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\zeta_{k,n}^2],
\end{aligned} \tag{A.35}$$

where $\mathbb{E}_{q(f_{l_r,n})} [\lambda_n^{(r)}]$, and $\mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\zeta_{k,n}]$ are respectively approximated by equations (A.22) and (A.29), and

$$\mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\zeta_{k,n}^2] \approx \left(\frac{1}{\sqrt{\pi}} \right)^K \sum_{s_K=1}^{S_1} \dots \sum_{s_1=1}^{S_K} o_{s_K} \dots o_{s_1} [\Xi (g_{s_K} \sqrt{2\sigma_{k,n}} + \iota_{k,n})]^2, \tag{A.36}$$

Case 2

$$\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial^2}{\partial f_{l_r,n}^2} \left[\log p(\mathbf{Y}|\hat{\mathbf{f}}) \right] \right] = \begin{cases} \mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n}) q(f_{l_r,n})} \left[\frac{\partial \Omega}{\partial f_{l_r,n}} \right], & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases},$$

where

$$\begin{aligned}
&\mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n}) q(f_{l_r,n})} \left[\frac{\partial \Omega}{\partial f_{l_r,n}} \right] = \\
&= \mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n}) q(f_{l_r,n})} \left[\left(2 \left(\lambda_n^{(r)} \right)^3 - 3 \left(\lambda_n^{(r)} \right)^2 + \lambda_n^{(r)} \right) \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \log(\zeta_{k,n}) + \log(K) \right) \right], \\
&= \left(2 \mathbb{E}_{q(f_{l_r,n})} \left[\left(\lambda_n^{(r)} \right)^3 \right] - 3 \mathbb{E}_{q(f_{l_r,n})} \left[\left(\lambda_n^{(r)} \right)^2 \right] + \mathbb{E}_{q(f_{l_r,n})} \left[\lambda_n^{(r)} \right] \right) \\
&\quad \times \left(\sum_{k=1}^K \delta(y_n^{(r)}, k) \mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\log(\zeta_{k,n})] + \log(K) \right),
\end{aligned} \tag{A.37}$$

where $\mathbb{E}_{q(f_{l_r,n})} [\lambda_n^{(r)}]$, $\mathbb{E}_{q(f_{l_r,n})} \left[\left(\lambda_n^{(r)} \right)^2 \right]$, $\mathbb{E}_{q(f_{l_r,n})} \left[\left(\lambda_n^{(r)} \right)^3 \right]$, and $\mathbb{E}_{q(f_{1,n}) \dots q(f_{K,n})} [\log(\zeta_{k,n})]$ are approximated respectively by equations (A.22), (A.29) and (A.31); also,

$$\mathbb{E}_{q(f_{l_r,n})} [(\lambda_n^{(r)})^3] = \int q(f_{l_r,n}) (\sigma(f_{l_r,n}))^3 df_{l_r,n} \approx \frac{1}{\sqrt{\pi}} \sum_{s=1}^S o_s [g_s \sqrt{2\sigma_{l_r,n}} + \iota_{l_r,n}]^3. \tag{A.38}$$

8. The predictive distributions.

Given a new sample \mathbf{x}_* , we are interested in the predictive distributions for the

ground truth $\zeta_{k,*}$, and the labelers' reliabilities $\lambda_*^{(r)}$.

Predictive distribution for $\zeta_{k,*}$ Each $\zeta_{k,n}$ is computed by applying the Softmax function over each function $f_{k,n}$, $\forall k \in \{1, \dots, K\}$. Accordingly, the mean for $\zeta_{k,*}$ can be approximated as

$$\begin{aligned} \mathbb{E}[\zeta_{k,*}] &= \int \Xi(f_{k,*}) q(f_{1,*}) \dots q(f_{K,*}) df_{1,*} \dots df_{K,*}, \\ &\approx \left(\frac{1}{\sqrt{\pi}} \right)^K \sum_{s_K=1}^{S_1} \dots \sum_{s_1=1}^{S_K} o_{s_K} \dots o_{s_1} \Xi \left(g_{s_k} \sqrt{2\sigma_{k,*}} + \iota_{k,*} \right). \end{aligned} \quad (\text{A.39})$$

For the predictive variance of $\zeta_{k,*}$, we use the expression $\text{Var}[\zeta_{k,*}] = \mathbb{E}[\zeta_{k,*}^2] - \mathbb{E}[\zeta_{k,*}]^2$; hence, we need to compute $\mathbb{E}[\zeta_{k,*}^2]$, which is given as

$$\begin{aligned} \mathbb{E}[\zeta_{k,*}^2] &= \int \Xi(f_{k,*})^2 q(f_{1,*}) \dots q(f_{K,*}) df_{1,*} \dots df_{K,*}, \\ &\approx \left(\frac{1}{\sqrt{\pi}} \right)^K \sum_{s_K=1}^{S_1} \dots \sum_{s_1=1}^{S_K} o_{s_K} \dots o_{s_1} \left[\Xi \left(g_{s_k} \sqrt{2\sigma_{k,*}} + \iota_{k,*} \right) \right]^2. \end{aligned} \quad (\text{A.40})$$

Predictive distribution for $\lambda_*^{(r)}$ In this case, λ_r is computed by applying a Sigmoid function over the latent functions $\varsigma(\mathbf{f}_{l_r})$, $\forall r$. Thus, the mean for $\varsigma(f_*)$ is given as

$$\mathbb{E}[\varsigma(f_{l_r,*})] = \int \varsigma(f_{l_r,*}) q(f_{l_r,*}) df_{l_r,*} \approx \frac{1}{\sqrt{\pi}} \sum_{s=1}^S o_s \varsigma \left(g_s \sqrt{2\sigma_{l_r,*}} + \iota_{l_r,*} \right). \quad (\text{A.41})$$

For the variance of $\varsigma(f_*)$, we use the expression $\text{Var}[\varsigma(f_{l_r,*})] = \mathbb{E}[\varsigma(f_{l_r,*})^2] - \mathbb{E}[\varsigma(f_{l_r,*})]^2$; hence, we need to compute $\mathbb{E}[\varsigma(f_{l_r,*})^2]$, which is given as

$$\mathbb{E}[\varsigma(f_{l_r,*})^2] = \int \varsigma(f_{l_r,*})^2 q(f_{l_r,*}) df_{l_r,*} \approx \frac{1}{\sqrt{\pi}} \sum_{s=1}^S o_s \left[\varsigma \left(g_s \sqrt{2\sigma_{l_r,*}} + \iota_{l_r,*} \right) \right]^2. \quad (\text{A.42})$$

A.2.2 Gaussian distribution for regression with multiple annotators

For real-valued outputs, e.g., $\mathcal{Y} \subset \mathbb{R}$, we follow the multi-annotator model used in Groot et al. (2011); Raykar et al. (2010); Rodrigues et al. (2017); Xiao et al. (2013), where each output $y_n^{(r)}$ is considered to be a corrupted version of the hidden ground truth y_n .

1. $p(\mathbf{Y}|\hat{\mathbf{f}})$

$$p(\mathbf{Y}|\hat{\mathbf{f}}) = \prod_{n=1}^N \prod_{r \in R_n} \mathcal{N}(y_n^{(r)}|y_n, v_n^{(r)}), \quad (\text{A.43})$$

where, $y_n = f_1(\mathbf{x}_n)$, and $v_n^{(r)} = \exp(f_{l_r}(\mathbf{x}_n))$, with $r \in \{1, \dots, R\}$, and $l_r = r + 1 \in \{2, \dots, J\}$; hence, the number of required LFs are $J = R + 1$.

2. $\log p(\mathbf{Y}|\hat{\mathbf{f}})$

$$\begin{aligned} \log p(\mathbf{Y}|\hat{\mathbf{f}}) &= \log \prod_{n=1}^N \prod_{r \in R_n} \mathcal{N}(y_n^{(r)}|y_n, v_n^{(r)}), \\ &= -\frac{1}{2} \sum_{n=1}^N \sum_{r \in R_n} \left\{ \log(2\pi) + \log(v_n^{(r)}) + \frac{(y_n^{(r)} - y_n)^2}{v_n^{(r)}} \right\}, \end{aligned} \quad (\text{A.44})$$

where $\log v_n^{(r)} = \log(\exp(f_{l_r, n})) = f_{l_r, n}$

3. $\mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$

$$\begin{aligned} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] &= \int \prod_{j=1}^J q(\hat{\mathbf{f}}_j) \log p(\mathbf{Y}|\hat{\mathbf{f}}) d\hat{\mathbf{f}}, \\ &= \sum_{n=1}^N \sum_{r \in R_n} \mathbb{E}_{q(f_{1,n})q(f_{l_r,n})} [\log \mathcal{N}(y_n^{(r)}|y_n, v_n^{(r)})]. \end{aligned} \quad (\text{A.45})$$

Since $q(f_{1,n})$ and $q(f_{l_r,n})$ obey Gaussian distributions, the above integrals can be solved analytically as in [Lázaro-Gredilla and Titsias \(2011\)](#). Thus, we have

$$\begin{aligned} &\int q(f_{1,n})q(f_{l_r,n}) \log \mathcal{N}(y_n^{(r)}|y_n, v_n^{(r)}) df_{1,n} df_{l_r,n} \\ &= \log \mathcal{N}(y_n^{(r)}|\mu_{1,n}, \exp(\mu_{l_r,n} - \frac{\sigma_{l_r,n}}{2})) - \frac{\sigma_{l_r,n}}{4} - \frac{s_{1,n} \exp(-\mu_{l_r,n} + \frac{\sigma_{l_r,n}}{2})}{2}. \end{aligned} \quad (\text{A.46})$$

Thus,

$$\begin{aligned} \mathbb{E}_{q(\hat{\mathbf{f}})} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] &= \sum_{n=1}^N \sum_{r \in R_n} \log \mathcal{N}(y_n^{(r)}|\mu_{1,n}, \exp(\mu_{l_r,n} - \frac{\sigma_{l_r,n}}{2})) - \dots \\ &\dots - \frac{\sigma_{l_r,n}}{4} - \frac{s_{1,n} \exp(-\mu_{l_r,n} + \frac{\sigma_{l_r,n}}{2})}{2}, \end{aligned} \quad (\text{A.47})$$

where $f_{1,n}$ is the n -th position of vector \mathbf{f}_1 , $f_{l_r,n}$ is the n -th position of vector \mathbf{f}_{l_r} , $\mu_{1,n}$ is the n -th position of vector $\boldsymbol{\mu}_1$, $\mu_{l_r,n}$ is the n -th position of vector $\boldsymbol{\mu}_{l_r}$, $s_{1,n}$ is the n -th position of vector \mathbf{s}_1 , and $\sigma_{l_r,n}$ is the n -th position of vector \mathbf{s}_{l_r} .

4. $\frac{\partial}{\partial \hat{\mathbf{f}}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$

$$\frac{\partial}{\partial \hat{\mathbf{f}}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] = \frac{\partial}{\partial \hat{\mathbf{f}}} \left[-\frac{1}{2} \sum_{n=1}^N \sum_{r \in R_n} \left\{ \log(2\pi) + f_{l_r,n} + \frac{(y_n^{(r)} - f_{1,n})^2}{\exp(f_{l_r,n})} \right\} \right]. \quad (\text{A.48})$$

Case 1

$$\frac{\partial}{\partial f_{1,n}} \left[-\frac{1}{2} \sum_{n=1}^N \sum_{r \in R_n} \left\{ \frac{(y_n^{(r)} - f_{1,n})^2}{\exp(f_{l_r,n})} \right\} \right] = \sum_{r \in R_n} \frac{(y_n^{(r)} - f_{1,n})}{\exp(f_{l_r,n})}. \quad (\text{A.49})$$

Case 2

$$\frac{\partial}{\partial f_{l_r,n}} \left[-\frac{1}{2} \sum_{n=1}^N \sum_{r \in R_n} \left\{ f_{l_r,n} + \frac{(y_n^{(r)} - f_{1,n})^2}{\exp(f_{l_r,n})} \right\} \right] = \begin{cases} -\frac{1}{2} + \frac{(y_n^{(r)} - f_{1,n})^2}{2 \exp(f_{l_r,n})}, & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases}. \quad (\text{A.50})$$

5. $\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial \hat{\mathbf{f}}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right]$

Case 1

$$\begin{aligned} \mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial f_{1,n}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right] &= \sum_{r \in R_n} \mathbb{E}_{q(f_{1,n})q(f_{l_r,n})} \left[\frac{(y_n^{(r)} - f_{1,n})}{\exp(f_{l_r,n})} \right], \\ &= \sum_{r \in R_n} \left(y_n^{(r)} - \mathbb{E}_{q(f_{1,n})}[f_{1,n}] \right) \mathbb{E}_{q(f_{l_r,n})} [\exp(-f_{l_r,n})], \end{aligned} \quad (\text{A.51})$$

where

$$\mathbb{E}_{q(f_{1,n})}[f_{1,n}] = \iota_{1,n}, \quad (\text{A.52})$$

$$\mathbb{E}_{q(f_{l_r,n})} [\exp(-f_{l_r,n})] = \exp \left(-\iota_{l_r,n} + \frac{\sigma_{l_r,n}}{2} \right). \quad (\text{A.53})$$

Thus,

$$\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial f_{1,n}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right] = \sum_{r \in R_n} (y_n^{(r)} - \iota_{1,n}) \exp \left(-\iota_{l_r,n} + \frac{\sigma_{l_r,n}}{2} \right). \quad (\text{A.54})$$

Case 2

$$\begin{aligned} \mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial f_{l_r,n}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right] &= \\ &= \begin{cases} \frac{1}{2} \mathbb{E}_{q(f_{1,n})} [(y_n^{(r)} - f_{1,n})^2] \mathbb{E}_{q(f_{l_r,n})} [\exp(-f_{l_r,n})] - \frac{1}{2}, & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases}, \end{aligned} \quad (\text{A.55})$$

where

$$\begin{aligned}\mathbb{E}_{q(f_{1,n})} \left[(y_n^{(r)} - f_{1,n})^2 \right] &= \mathbb{E}_{q(f_{1,n})} \left[(y_n^{(r)})^2 - 2y_n^{(r)} f_{1,n} - (f_{1,n})^2 \right], \\ &= \left(y_n^{(r)} - \iota_{1,n} \right)^2 + s_{1,n}.\end{aligned}\quad (\text{A.56})$$

Furthermore, $\mathbb{E}_{q(f_{l_r,n})} [\exp(-f_{l_r,n})]$ is computed as in equation (A.53). Thus,

$$\begin{aligned}\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial}{\partial f_{l_r,n}} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right] &= \\ &= \begin{cases} \frac{1}{2} \left(\left(y_n^{(r)} - \iota_{1,n} \right)^2 + s_{1,n} \right) \exp \left(-\iota_{l_r,n} + \frac{\sigma_{l_r,n}}{2} \right) - \frac{1}{2}, & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases}.\end{aligned}\quad (\text{A.57})$$

$$6. \frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})]$$

$$\frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] = \frac{\partial^2}{\partial \hat{\mathbf{f}}^2} \left[-\frac{1}{2} \sum_{n=1}^N \sum_{r \in R_n} \left\{ \log(2\pi) + f_{l_r,n} + \frac{(y_n^{(r)} - f_{1,n})^2}{\exp(f_{l_r,n})} \right\} \right]. \quad (\text{A.58})$$

Case 1

$$\frac{\partial^2}{\partial f_{1,n}^2} \left[-\frac{1}{2} \sum_{n=1}^N \sum_{r \in R_n} \left\{ \log(2\pi) + f_{l_r,n} + \frac{(y_n^{(r)} - f_{1,n})^2}{\exp(f_{l_r,n})} \right\} \right] = - \sum_{r \in R_n} \exp(-f_{l_r,n}). \quad (\text{A.59})$$

Case 2

$$\begin{aligned}\frac{\partial^2}{\partial f_{l_r,n}^2} \left[-\frac{1}{2} \sum_{n=1}^N \sum_{r \in R_n} \left\{ \log(2\pi) + f_{l_r,n} + \frac{(y_n^{(r)} - f_{1,n})^2}{\exp(f_{l_r,n})} \right\} \right] &= \\ &= \begin{cases} \frac{-(y_n^{(r)} - f_{1,n})^2}{2 \exp(f_{l_r,n})}, & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases}.\end{aligned}\quad (\text{A.60})$$

$$7. \mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial^2}{\partial \hat{\mathbf{f}}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right]$$

Case 1

$$\begin{aligned}
\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial^2}{\partial f_{1,n}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right] &= \sum_{r \in R_n} \mathbb{E}_{q(f_{l_r,n})} [-\exp(-f_{l_r,n})], \\
&= - \sum_{r \in R_n} \exp \left(-\iota_{l_r,n} + \frac{\sigma_{l_r,n}}{2} \right). \quad (\text{A.61})
\end{aligned}$$

Case 2

$$\begin{aligned}
\mathbb{E}_{q(\hat{\mathbf{f}})} \left[\frac{\partial^2}{\partial f_{l_r,n}^2} [\log p(\mathbf{Y}|\hat{\mathbf{f}})] \right] &= \\
&= \begin{cases} -\frac{1}{2} \mathbb{E}_{q(f_{1,n})} [(y_n^{(r)} - f_{1,n})^2] \mathbb{E}_{q(f_{l_r,n})} [\exp(-f_{l_r,n})], & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases}, \\
&= \begin{cases} -\frac{1}{2} \left((y_n^{(r)} - \iota_{1,n})^2 + s_{1,n} \right) \exp \left(-\iota_{l_r,n} + \frac{\sigma_{l_r,n}}{2} \right), & \text{if } r \in R_n \\ 0, & \text{Otherwise} \end{cases}. \quad (\text{A.62})
\end{aligned}$$

8. The predictive distributions.

Given a new sample \mathbf{x}_* , we are interested in the predictive distributions for the ground truth y_* , and the labelers' error-variances $v_*^{(r)}$.

Predictive distribution for y_* Due to $\mathbf{y} = \mathbf{f}_1$, the posterior distribution for y_* corresponds to $q(f_{1,*})$. Accordingly,

$$\mathbb{E}[y_*] = \mathbb{E}[f_{1,*}] = \iota_{1,*} \quad (\text{A.63})$$

$$\text{Var}[y_*] = \text{Var}[f_{1,*}] = s_{1,*} \quad (\text{A.64})$$

Predictive distribution for $v_*^{(r)}$ In this case, since $\mathbf{v}_r = \exp(\mathbf{f}_{l_r})$, the posterior distribution for $v_*^{(r)}$ follows a log-normal distribution with parameters $\iota_{l_r,*}$ and $\sigma_{l_r,*}$, which respectively correspond to the mean and variance of $q(f_{l_r,*})$. In this sense, the mean and variance of $v_*^{(r)}$ are given as

$$\mathbb{E}[v_*^{(r)}] = \int \exp(f_{l_r,*}) q(f_{l_r,*}) df_{l_r,*} = \exp \left(\iota_{l_r,*} + \frac{\sigma_{l_r,*}}{2} \right). \quad (\text{A.65})$$

$$\begin{aligned}
\text{Var}[v_*^{(r)}] &= \int \left(\exp(f_{l_r,*}) - \mathbb{E}[v_*^{(r)}] \right)^2 q(f_{l_r,*}) df_{l_r,*} \\
&= \exp(2\iota_{l_r,*} + \sigma_{l_r,*}) (\exp(\sigma_{l_r,*}) - 1). \quad (\text{A.66})
\end{aligned}$$

References

- Albarqouni, S., Baur, C., Achilles, F., Belagiannis, V., Demirci, S., and Navab, N. (2016). Aggnet: deep learning from crowds for mitosis detection in breast cancer histology images. *IEEE transactions on medical imaging*, 35(5):1313–1321. (pages 1, 3, 72, and 122)
- Álvarez, M., Luengo, D., Titsias, M., and Lawrence, N. D. (2010). Efficient multioutput Gaussian processes through variational inducing kernels. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 25–32. (page 89)
- Álvarez, M. A. and Lawrence, N. D. (2011). Computationally efficient convolved multiple output Gaussian processes. *The Journal of Machine Learning Research*, 12:1459–1500. (page 121)
- Álvarez, M. A., Rosasco, L., Lawrence, N. D., et al. (2012). Kernels for vector-valued functions: A review. *Foundations and Trends® in Machine Learning*, 4(3):195–266. (pages 6, 8, 10, 25, 89, and 121)
- Arias-Londoño, J. D., Godino-Llorente, J. I., Gutiérrez-Arriola, J., Osma-Ruiz, V., and Sáenz-Lechón, N. (2011). Automatic grbas assessment using complexity measures and a multiclass gmm-based detector. *Models and Analysis of Vocal Emissions for Biomedical Applications*, pages 111–114. (page 33)
- Artetxe, M., Labaka, G., and Agirre, E. (2018). A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. *arXiv preprint arXiv:1805.06297*. (page 122)
- Barron, A. R. (1993). Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information theory*, 39(3):930–945. (page 5)
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. springer. (pages 2, 5, 6, 8, 36, 65, 85, 103, and 105)
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877. (page 94)
- Cohn, T. and Specia, L. (2013). Modelling annotator bias with multi-task gaussian processes: An application to machine translation quality estimation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 32–42. (page 23)

- Cortes, C., Mohri, M., and Rostamizadeh, A. (2012). Algorithms for learning kernels based on centered alignment. *J. Mach. Learn. Res.*, 13(Mar):795–828. (pages 21, 29, 32, and 50)
- Dawid, A. P. and Skene, A. M. (1979). Maximum likelihood estimation of observer error-rates using the em algorithm. *Applied statistics*, pages 20–28. (pages 3 and 13)
- Demšar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research*, 7:1–30. (page 118)
- Esmaily, M., Yazdi, H. S., Abbassi, S., and Monsefi, R. (2016). Hierarchical cooperation of experts in learning from crowds. In *ICCKE*, pages 211–217. IEEE. (page 28)
- Fawcett, T. (2006). An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874. (pages 54, 76, and 99)
- G. Rodrigo, E., Aledo, J. A., and Gámez, J. A. (2019). Machine learning from crowds: A systematic review of its applications. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(2):e1288. (pages 4, 48, and 122)
- Géron, A. (2019). *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems*. O’Reilly Media. (pages 75 and 76)
- Gil, J., Álvarez, M., and Orozco, Á. (2015). Automatic assessment of voice quality in the context of multiple annotations. In *EMBC*, pages 6236–6239. IEEE. (pages 3, 33, 34, 48, 54, 63, 76, 83, 107, and 122)
- Gil-Gonzalez, J., Alvarez-Meza, A., and Orozco-Gutierrez, A. (2018). Learning from multiple annotators using kernel alignment. *Pattern Recognition Letters*, 116:150–156. (pages 55, 56, 70, 77, 79, 100, 103, and 105)
- Gönen, M. and Alpaydin, E. (2008). Localized multiple kernel learning. In *Proceedings of the 25th international conference on Machine learning*, pages 352–359. ACM. (pages 21 and 49)
- Gönen, M. and Alpaydin, E. (2011). Multiple kernel learning algorithms. *The Journal of Machine Learning Research*, 12:2211–2268. (page 29)
- Groot, P., Birlutiu, A., and Heskes, T. (2011). Learning from multiple annotators with Gaussian processes. In *ICANN*, pages 159–164. Springer. (pages 3, 4, 14, 15, 36, 46, 48, 70, 96, 121, and 133)
- Guan, M. Y., Gulshan, V., Dai, A. M., and Hinton, G. E. (2018). Who said what: Modeling individual labelers improves classification. In *Thirty-Second AAAI Conference on Artificial Intelligence*. (page 3)
- Hahn, U., von Sydow, M., and Merdes, C. (2018). How communication can make voters choose less well. *Topics in cognitive science*. (pages 4 and 48)

- Han, J., Luo, P., and Wang, X. (2019). Deep self-learning from noisy labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5138–5147. (page 122)
- Hensman, J., Matthews, A. G., and Ghahramani, Z. (2015). Scalable variational Gaussian process classification. *Proceedings of Machine Learning Research*, 38:351–360. (pages 91, 92, 94, and 125)
- Hernández-Muriel, J. A., Bermeo-Ulloa, J. B., Holguin-Londoño, M., Álvarez-Meza, A. M., and Orozco-Gutiérrez, Á. A. (2020). Bearing health monitoring using relief-based feature relevance analysis and HMM. *Applied Sciences*, 10(15):5170. (page 76)
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. (2013). Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347. (page 89)
- Hua, G., Long, C., Yang, M., and Gao, Y. (2018). Collaborative active visual recognition from crowds: A distributed ensemble approach. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):582–594. (page 28)
- Huang, J., Li, Y., Tao, J., Lian, Z., Niu, M., and Yang, M. (2018). Deep learning for continuous multiple time series annotations. In *Proceedings of the 2018 on Audio/Visual Emotion Challenge and Workshop*, pages 91–98. ACM. (page 2)
- Imbajoa-Ruiz, D. E., Gustin, I., Bolaños-Ledezma, M., Arciniegas-Mejía, A. F., Guasmayan-Guasmayan, F., Bravo-Montenegro, M., Castro-Ospina, A. E., and Peluffo-Ordóñez, D. H. (2016). Multi-labeler classification using kernel representations and mixture of classifiers. In *Iberoamerican Congress on Pattern Recognition*, pages 343–351. Springer. (page 12)
- Jimenez, C., Diaz, D., Salazar, D., Alvarez, A., Orozco, A., and Henao, O. (2018a). Nerve structure segmentation from ultrasound images using random under-sampling and an svm classifier. In *International Conference Image Analysis and Recognition*, pages 571–578. Springer. (page 1)
- Jimenez, D. A., García, H. F., Álvarez, A. M., Orozco, Á. A., and Holguín, G. (2018b). A kernelized morphable model for 3d brain tumor analysis. In *International Conference Image Analysis and Recognition*, pages 529–537. Springer. (page 1)
- Kanagawa, M., Hennig, P., Sejdinovic, D., and Sriperumbudur, B. K. (2018). Gaussian processes and kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*. (page 6)
- Kara, Y. E., Genc, G., Aran, O., and Akarun, L. (2015). Modeling annotator behaviors for crowd labeling. *Neurocomputing*, 160:141–156. (pages 2, 3, and 49)
- Kay, S. M. and Kay, S. M. (1993). *Fundamentals of statistical signal processing: estimation theory*, volume 1. Prentice-hall Englewood Cliffs, NJ. (page 8)
- Lázaro-Gredilla, M. and Titsias, M. K. (2011). Variational heteroscedastic gaussian process regression. In *ICML*. (page 134)

- Liu, W., Principe, J. C., and Haykin, S. (2011). *Kernel adaptive filtering: a comprehensive introduction*, volume 57. John Wiley & Sons. (page 30)
- Liu, X., Zou, Y., Song, Y., Yang, C., You, J., and K Vijaya Kumar, B. (2018). Ordinal regression with neuron stick-breaking for medical diagnosis. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0. (page 122)
- Liu, Y., Zhang, W., Yu, Y., et al. (2020). Truth inference with a deep clustering-based aggregation model. *IEEE Access*, 8:16662–16675. (page 2)
- Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-SNE. *Journal of machine learning research*, 9(Nov):2579–2605. (pages 98 and 100)
- Morales-Alvarez, P., Ruiz, P., Coughlin, S., Soriano, R. M., and Katsaggelos, A. (2020). Scalable variational gaussian processes for crowdsourcing: Glitch detection in ligo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. (pages 3, 14, 15, 48, 49, and 119)
- Morales-Álvarez, P., Ruiz, P., Santos-Rodríguez, R., Molina, R., and Katsaggelos, A. K. (2019). Scalable and efficient learning from crowds with Gaussian processes. *Information Fusion*, 52:110–127. (pages 14, 15, 100, 105, 106, 110, and 119)
- Moreno-Muñoz, P., Artés, A., and Alvarez, M. (2018). Heterogeneous multi-output Gaussian process prediction. In *Advances in neural information processing systems*, pages 6711–6720. (pages 93, 94, and 127)
- Petersen, K. B., Pedersen, M. S., et al. (2008). The matrix cookbook. *Technical University of Denmark*, 7(15):510. (page 51)
- Principe, J. C. (2010). *Information theoretic learning: Renyi’s entropy and kernel perspectives*. Springer Science & Business Media. (page 28)
- Rasmussen, C. E. and Williams, C. K. (2006). *Gaussian processes for machine learning*, volume 1. MIT press Cambridge. (pages xi, 5, 8, 9, 35, 36, 37, 54, 55, and 90)
- Raykar, V. C., Yu, S., Zhao, L. H., Valadez, G. H., Florin, C., Bogoni, L., and Moy, L. (2010). Learning from crowds. *J. Speech Lang. Hear. Res.*, 11(Apr):1297–1322. (pages xi, 2, 3, 5, 11, 13, 14, 16, 35, 37, 46, 55, 56, 62, 70, 77, 96, 100, 119, 122, and 133)
- Rizos, G. and Schuller, B. W. (2020). Average jane, where art thou?—recent avenues in efficient machine learning under subjectivity uncertainty. In *International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems*, pages 42–55. Springer. (page 2)
- Rodrigues, F., Lourenco, M., Ribeiro, B., and Pereira, F. (2017). Learning supervised topic models for classification and regression from crowds. *IEEE transactions on PAMI*. (pages 3, 14, 36, 48, 49, 96, and 133)
- Rodrigues, F., Pereira, F., and Ribeiro, B. (2013). Learning from multiple annotators: Distinguishing good from random labelers. *Pattern Recognition Letters*, 34(12):1428–1436. (pages 3, 14, 34, 35, 48, 49, 53, 55, 70, 74, 77, 94, 100, and 127)

- Rodrigues, F., Pereira, F., and Ribeiro, B. (2014a). Sequence labeling with multiple annotators. *Machine learning*, 95(2):165–181. (pages 3 and 14)
- Rodrigues, F. and Pereira, F. C. (2018a). Deep learning from crowds. In *Thirty-Second AAAI Conference on Artificial Intelligence*. (pages 3, 13, 37, 38, 46, 56, 72, 100, 110, and 121)
- Rodrigues, F. and Pereira, F. C. (2018b). Heteroscedastic gaussian processes for uncertainty modeling in large-scale crowdsourced traffic data. *Transportation research part C: emerging technologies*, 95:636–651. (page 8)
- Rodrigues, F., Pereira, F. C., and Ribeiro, B. (2014b). Gaussian process classification and active learning with multiple annotators. In *ICML*, pages 433–441. (pages 3, 14, 35, 37, 46, 48, 55, 56, 70, 77, 100, and 103)
- Rodrigues, F. M. P. D. (2016). *Probabilistic models for learning from crowdsourced data*. PhD thesis. (page 4)
- Ruiz, P., Morales-Álvarez, P., Molina, R., and Katsaggelos, A. K. (2019). Learning from crowds with variational gaussian processes. *Pattern Recognition*, 88:298–311. (pages 3, 14, 36, 49, 53, 59, 83, 103, 105, and 119)
- Saul, A., Hensman, J., Vehtari, A., and Lawrence, N. (2016). Chained Gaussian processes. In *Artificial Intelligence and Statistics*, pages 1431–1440. (pages 24, 25, 72, 73, 87, 89, 90, 93, 94, 115, 117, 118, and 127)
- Schölkopf, B., Herbrich, R., and Smola, A. J. (2001). A generalized representer theorem. In *International conference on computational learning theory*, pages 416–426. Springer. (page 7)
- Schölkopf, B., Smola, A. J., Bach, F., et al. (2002). *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press. (pages 6 and 54)
- Sung, H.-E., Chen, C.-K., Xiao, H., and Lin, S.-D. (2017). A classification model for diverse and noisy labelers. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 58–69. Springer. (pages 1 and 14)
- Surowiecki, J. (2005). *The wisdom of crowds*. Anchor. (pages 4 and 48)
- Tang, W., Yin, M., and Ho, C.-J. (2019). Leveraging peer communication to enhance crowdsourcing. In *The World Wide Web Conference*, pages 1794–1805. ACM. (pages 4 and 48)
- Tao, D., Cheng, J., Yu, Z., Yue, K., and Wang, L. (2018). Domain-weighted majority voting for crowdsourcing. *IEEE transactions on neural networks and learning systems*, 30(1):163–174. (page 2)
- Teh, Y., Seeger, M., and Jordan, M. (2005). Semiparametric latent factor models. In *AISTATS 2005-Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics*. (pages 25, 89, and 92)

- Theodoridis, S. (2015). *Machine learning: a Bayesian and optimization perspective*. Academic press. (page 5)
- Tian, T. and Zhu, J. (2015). Max-margin majority voting for learning from crowds. In *Advances in neural information processing systems*, pages 1621–1629. (pages 11 and 12)
- Tikhonov, A. N. and Arsenin, V. Y. (1977). Solutions of ill-posed problems. *New York*, 1:30. (page 11)
- Titsias, M. (2009). Variational learning of inducing variables in sparse gaussian processes. In *Artificial intelligence and statistics*, pages 567–574. PMLR. (page 91)
- Vargas-Cardona, H. D., Orozco, Á. A., Álvarez, A. M., and Álvarez, M. A. (2019). Tensor decomposition processes for interpolation of diffusion magnetic resonance imaging. *Expert Systems with Applications*, 118:92–108. (page 1)
- Varma, M. and Babu, B. R. (2009). More generality in efficient multiple kernel learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 1065–1072. (page 29)
- Venanzi, M., Guiver, J., Kazai, G., Kohli, P., and Shokouhi, M. (2014). Community-based bayesian aggregation models for crowdsourcing. In *Proceedings of the 23rd international conference on World wide web*, pages 155–164. ACM. (pages 4 and 48)
- Wang, X. and Bi, J. (2016). Bi-convex optimization to learn classifiers from multiple biomedical annotations. *IEEE/ACM transactions on computational biology and bioinformatics*, 14(3):564–575. (page 4)
- Welinder, P., Branson, S., Perona, P., and Belongie, S. J. (2010). The multidimensional wisdom of crowds. In *Advances in neural information processing systems*, pages 2424–2432. (page 1)
- Xiao, H., Xiao, H., and Eckert, C. (2013). Learning from multiple observers with unknown expertise. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 595–606. Springer. (pages 17, 49, 62, 96, and 133)
- Xie, Y., Li, J., Zhu, T., and Liu, C. (2019). Continuous-valued annotations aggregation for heart rate detection. *IEEE Access*, 7:37664–37671. (page 4)
- Yan, Y., Rosales, R., Fung, G., Schmidt, M. W., Valadez, G. H., Bogoni, L., Moy, L., and Dy, J. G. (2010). Modeling annotator expertise: Learning when everybody knows a bit of something. In *AISTATS*, pages 932–939. (pages 15, 34, 35, 48, 53, 55, 62, and 77)
- Yan, Y., Rosales, R., Fung, G., Subramanian, R., and Dy, J. (2014). Learning from multiple annotators with varying expertise. *Machine learning*, 95(3):291–327. (pages 4, 11, 15, 28, 46, 47, 48, and 70)
- Zhai, X., Oliver, A., Kolesnikov, A., and Beyer, L. (2019). S4l: Self-supervised semi-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1476–1485. (page 122)

- Zhang, A. and Gao, X. (2018). Data-dependent kernel sparsity preserving projection and its application for semi-supervised classification. *Multimedia Tools and Applications*, pages 1–17. (pages 21 and 49)
- Zhang, J., Wu, X., and Sheng, V. S. (2014). Imbalanced multiple noisy labeling. *IEEE Transactions on Knowledge and Data Engineering*, 27(2):489–503. (page 3)
- Zhang, J., Wu, X., and Sheng, V. S. (2015). Imbalanced multiple noisy labeling. *IEEE Transactions on Knowledge and Data Engineering*, 27(2):489–503. (page 46)
- Zhang, J., Wu, X., and Sheng, V. S. (2016). Learning from crowdsourced labeled data: a survey. *Artificial Intelligence Review*, 46(4):543–576. (page 1)
- Zhang, P. and Obradovic, Z. (2011). Learning from inconsistent and unreliable annotators by a gaussian mixture model and bayesian information criterion. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 553–568. Springer. (pages 16 and 48)
- Zhu, T., Pimentel, M. A., Clifford, G. D., and Clifton, D. A. (2019). Unsupervised bayesian inference to fuse biosignal sensory estimates for personalising care. *IEEE journal of biomedical and health informatics*, 23(1):47. (pages 2, 4, 17, 79, 105, and 110)