

An Experimental Evaluation of Mixup Regression Forests

Rodriguez, Juan; Juez-Gil, Mario; Arnaiz-Gonzalez, Alvar; Kuncheva, Ludmila

Expert Systems with Applications

DOI:

<https://doi.org/10.1016/j.eswa.2020.113376>

Published: 01/08/2020

Peer reviewed version

[Cyswllt i'r cyhoeddiad / Link to publication](#)

Dyfyniad o'r fersiwn a gyhoeddwyd / Citation for published version (APA):

Rodriguez, J., Juez-Gil, M., Arnaiz-Gonzalez, A., & Kuncheva, L. (2020). An Experimental Evaluation of Mixup Regression Forests. *Expert Systems with Applications*, 151, [113376]. <https://doi.org/10.1016/j.eswa.2020.113376>

Hawliau Cyffredinol / General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

An Experimental Evaluation of Mixup Regression Forests

Juan J. Rodríguez^a, Mario Juez-Gil^a, Álvaro Arnaiz-González^a, Ludmila I. Kuncheva^b

^a*Escuela Politécnica Superior, Universidad de Burgos, 09006 Burgos, SPAIN*

^b*School of Computer Science and Electronic Engineering, Bangor University, Dean Street, Bangor LL57 1UT, UK*

Abstract

Over the past few decades, the remarkable prediction capabilities of ensemble methods have been used within a wide range of applications. Maximization of base-model ensemble accuracy and diversity are the keys to the heightened performance of these methods. One way to achieve diversity for training the base models is to generate artificial/synthetic instances for their incorporation with the original instances. Recently, the *mixup* method was proposed for improving the classification power of deep neural networks (Zhang et al., 2017). *Mixup* method generates artificial instances by combining pairs of instances and their labels, these new instances are used for training the neural networks promoting its regularization. In this paper, new regression tree ensembles trained with mixup, which we will refer to as Mixup Regression Forest, are presented and tested. The experimental study with 61 datasets showed that the mixup approach improved the results of both Random Forest and Rotation Forest.

Keywords: Mixup, Regression, Random Forest, Rotation Forest

1. Introduction

The idea that motivates this study, in relation to problems that ensemble techniques can solve, is that an increase in base-model diversity will improve ensemble performance, generalization, and robustness. Diversity is a key attribute of an ensemble, without which ensemble methods would not be as successful as they are (Kuncheva & Whitaker, 2003). It

Email addresses: jjrodriguez@ubu.es (Juan J. Rodríguez), mariojg@ubu.es (Mario Juez-Gil), alvarag@ubu.es (Álvar Arnaiz-González), l.i.kuncheva@bangor.ac.uk (Ludmila I. Kuncheva)

6 can be achieved in several ways: by using different methods for building the classifiers in
7 the ensemble (heterogeneous ensemble), by using methods that build classifiers with random
8 components, and by using different training sets. The focus of this paper rests on the last
9 strategy, in particular, in making new instances that not found in the original set for creating
10 different training sets.

Mixup has recently been proposed by Zhang et al. (2017) for training deep neural networks using combinations of pairs of examples and their labels. Given a training set where each example is (x, y) , with an input, x , and a corresponding output, y , then the combined examples (\tilde{x}, \tilde{y}) are generated as

$$\begin{aligned}\tilde{x} &= \lambda x_i + (1 - \lambda)x_j \\ \tilde{y} &= \lambda y_i + (1 - \lambda)y_j\end{aligned}$$

11 where (x_i, y_i) and (x_j, y_j) are two examples, drawn at random from the training data, and
12 $\lambda \in [0, 1]$. The values of λ were obtained using the Beta distribution: $\lambda \sim \text{Beta}(\alpha, \alpha)$, with
13 $\alpha \in (0, \infty)$.

14 Some example mixup data projections can be seen in figures 1 and 2. Figure 1 shows a
15 single input dataset where the input variable and the output variable are represented on the
16 x axis and the y axis, respectively, and the instances are generated with mixup. Figure 2
17 shows a couple of examples: two two-input datasets and the mixup-generated instances. The
18 output values of the original datasets are in $\{-1, 1\}$ and the output values of the datasets
19 that are generated are in $[-1, 1]$. Figure 3 shows the predictions of a single random tree for
20 the datasets shown in Figure 2.

21 Mixup differs from other data augmentation approaches, in so far as its outputs are also
22 combined. The combination of the outputs to address regression problems is a straightforward
23 procedure.

24 As shown in Figure 1, some of the examples generated with mixup are clearly noise.
25 Although it can be detrimental, noise injection has previously been used as a strategy
26 for building successful ensembles (Melville & Mooney, 2005; Frank & Pfahringer, 2006;
27 Martínez-Muñoz & Suárez, 2005; González et al., 2017). In mixup forests, the prevalence

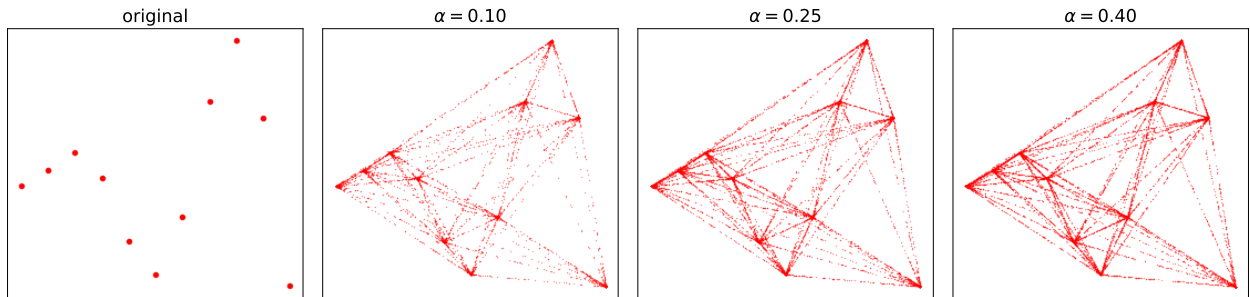


Figure 1: A regression problem dataset with a single input (x axis), and a single continuous output (y axis). Artificial instances are generated with mixup for $\alpha \in \{0.1, 0.25, 0.4\}$.

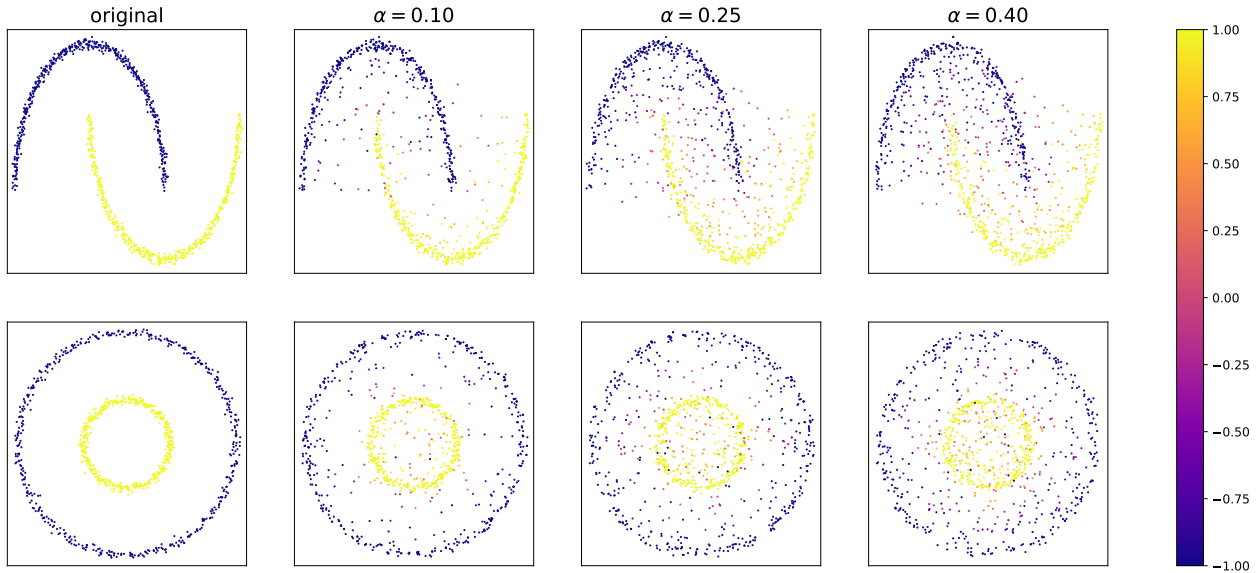


Figure 2: Two two-inputs datasets and the datasets generated with mixup for $\alpha \in \{0.1, 0.25, 0.4\}$. The output variables are shown in yellow and in blue.

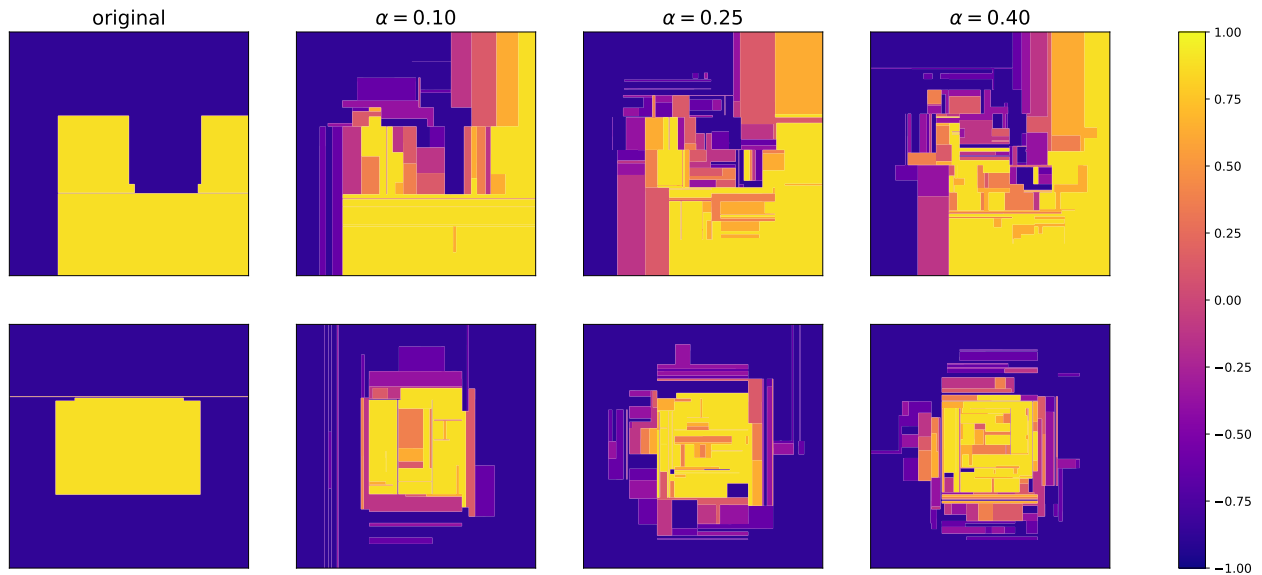


Figure 3: Predictions given by a single random tree trained with the corresponding datasets from Figure 2.

28 of these noisy examples can be controlled with the α value and the number of artificial
 29 examples that are generated.

30 Ensemble techniques have successfully been applied in various domains over the past
 31 few decades. Many works and several literature reviews have been published on both clas-
 32 sification (Kuncheva, 2014) and regression (Mendes-Moreira et al., 2012) ensembles. Some
 33 illustrative examples of ensemble applications are detailed below.

34 In industrial environments, ensembles can be used as predictive models with adaptive ca-
 35 pabilities, for example, to respond to incidences at processing plants (Soares & Araújo, 2015).
 36 Financial forecasting with ensembles has also been a very frequent research topic, among
 37 other examples, for the prediction of trading in stocks (Weng et al., 2018) and bankruptcy
 38 trends (Chen et al., 2020). It is also of great industrial interest, for example, in the construc-
 39 tion industry, where ensembles have been used for the prediction of financial distress (Choi
 40 et al., 2018). Many techniques for credit risk assessment have been proposed, based on both
 41 statistics and Artificial Intelligence (AI) models; a task in which ensembles have demon-
 42 strated good performance (Marqués et al., 2012). In biometrics, improved recognition rates
 43 can be achieved using multimodal biometric systems that capture multiple biometric traits,

44 e.g. fingerprint, iris and facial features; multimodal data learning in those fields can be
45 addressed by using ensembles (Ross & Jain, 2003). The advantages and the convenience of
46 ensemble learning to learn from multimodal features have likewise benefited several clinical
47 practices (Tay et al., 2013). The sort of highly robust system required for image recognition
48 tasks, such as facial recognition, can be provided by ensembles, to address the diversity of
49 facial expressions and aging effects (Sirlantzis et al., 2008). Real-life problems, such as spam
50 detection (Geng et al., 2007), translation of DNA sequences (García-Pedrajas et al., 2012),
51 and the detection of credit-card fraud (Panigrahi et al., 2009), are known as imbalanced
52 learning problems that can also be solved using ensemble techniques (Galar et al., 2012).
53 The mixup data augmentation strategy proposed in this paper, might therefore lead to even
54 better ensemble models for the aforementioned applications, as the artificial generation of
55 instances has the potential to improve the performance of almost any ensemble method.

56 The contribution of this study relates to the novel use of the mixup approach. It demon-
57 strates that artificial examples generated by mixup contribute to improved ensemble perfor-
58 mance in regression tasks. Mixup is therefore considered for regression, mainly because of
59 its simplicity: it can be used with all data types and needs no adjustments to the model.

60 The rest of the paper will be organized as follows. In section 2, a brief literature review of
61 the most relevant works in this field will be presented. In section 3, the experimental setup
62 will be described. Then the results will be presented and analyzed in Section 4. Finally, some
63 concluding remarks and suggestions for future research work will be outlined in section 5.

64 **2. Related works**

65 Diversity between the members of an ensemble means that those ensembles are capable
66 of better predictions than the individual ensemble members. One way to achieve diversity
67 is by introducing artificial examples for training, for example through the mixup approach.
68 Data augmentation with artificial examples has previously been used in many ensemble
69 algorithms, some of which are detailed below.

70 In DECORATE (Diverse Ensemble Creation by Oppositional Relabeling of Artificial
71 Training Examples) (Melville & Mooney, 2003, 2005), instances are generated based on the

72 distribution of the data. The labels of the new instances are assigned with a probability
73 that is proportional to the inverse of the probability assigned by the current ensemble,
74 because the purpose of the artificial instances is to increase diversity. In Bagging with Input
75 Smearing (Frank & Pfahringer, 2006), the generation of artificial instances add noise to
76 actual instances.

77 In imbalanced classification problems¹, artificial examples are commonly used for increas-
78 ing the number of instances of the minority class/es. As with mixup, in SMOTE (Chawla
79 et al., 2002), artificial instances are also obtained by combining pairs of instances. In this
80 case, as both instances in a pair are of the same class, the label of the artificial instances
81 is the same as the instances used to generate them. SMOTE was not originally proposed
82 as an ensemble method and can in fact be used as a pre-processing step before the con-
83 struction of a model. Nevertheless, it can also be directly used in ensembles, by training
84 each base classifier with a different set of original and artificial instances. SMOTE has been
85 combined with generic ensemble methods giving rise to SMOTEBoost (Chawla et al., 2003)
86 and SMOTEBagging (Wang & Yao, 2009), among others.

87 There are many other methods for balancing datasets by augmenting the minority classes
88 with artificial instances (Han et al., 2005; He et al., 2008; Menardi & Torelli, 2014; Zhu
89 et al., 2017). Some of these methods, such as SMOTE, have also been adapted to regression
90 problems (Torgo et al., 2013).

91 Likewise, highly sophisticated approaches exist for augmenting datasets. Most of those
92 have been specifically designed for a given data type, for example, images (Tokozume et al.,
93 2017, 2018; Inoue, 2018; Summers & Dinneen, 2019). Such approaches require training and
94 adjusting a model, in order to generate the artificial instances (Mayo & Frank, 2017; Verma
95 et al., 2018; Guo et al., 2018; Lindenbaum et al., 2018; Beckham et al., 2019).

96 Here, mixup was chosen as the simplest augmentation method and the significant advan-
97 tage of its use with regression ensembles of random trees (Mixup Regression Forests) will
98 be demonstrated in the following section.

¹Imbalanced classification problems are those related to datasets and domains where one class has a much greater number of examples than another (Haixiang et al., 2017).

99 **3. Experimental setting**

100 The purpose of this experiment is to demonstrate the advantage of the mixup augmen-
101 tation step. Two of the best state-of-the-art ensemble methods (singled out by extensive
102 experimental studies (Random Forest (Breiman, 2001; Fernández-Delgado et al., 2014) and
103 Rotation Forest (Rodríguez et al., 2006; Pardo et al., 2013; Bagnall et al., 2018)) are tested
104 with and without the mixup step over a large collection of datasets. The experimental setup
105 is presented below.

106 *3.1. Datasets*

107 Table 1 shows the main characteristics of the 61 regression datasets used in the experi-
108 ments. All of them are available in the format used by Weka² (Hall et al., 2009). Thirty of
109 the 61 datasets were collected by Luís Torgo³.

110 *3.2. Methods*

111 The mixup method is used in combination with Random Forest (Breiman, 2001) and
112 Rotation Forest (Rodríguez et al., 2006; Pardo et al., 2013). Both Random and Rotation
113 Forest are used to transform the training dataset. In Random Forest, the dataset is sampled,
114 whereas in Rotation Forest, it is rotated and then sampled. The mixup transformation can
115 be done before or after the two above-mentioned ensemble transformations. Four methods
116 are therefore available:

- 117 • MixRandFor: The dataset is augmented with mixup and then sampled.
- 118 • RandMixFor: The dataset is sampled and then the sample is augmented with mixup.
- 119 • MixRotFor: The dataset is augmented with mixup and then rotated.
- 120 • RotMixFor: The dataset is rotated and then augmented with mixup.

²http://www.cs.waikato.ac.nz/ml/weka/index_datasets.html

³<http://www.dcc.fc.up.pt/~ltorgo/Regression/DataSets.html>

Dataset	Examples	Numeric	Nominal	Dataset	Examples	Numeric	Nominal
2d-planes	40768	10	0	house-16H	22784	16	0
abalone	4177	7	1	house-8L	22784	8	0
aileron	13750	40	0	housing	506	12	1
auto-horse	205	17	8	hungarian	294	6	7
auto-mpg	398	4	3	kin8nm	8192	8	0
auto-price	159	15	0	longley	16	6	0
auto93	93	16	6	lowbwt	189	2	7
bank-32nh	8192	32	0	machine-cpu	209	6	0
bank-8FM	8192	8	0	mbagrade	61	1	1
basketball	96	4	0	meta	528	19	2
bodyfat	252	14	0	mv	40768	7	3
bolts	40	7	0	pbc	418	10	8
breast-tumor	286	1	8	pharynx	195	1	10
cal-housing	20640	8	0	pole	15000	48	0
cholesterol	303	6	7	pollution	60	15	0
cleveland	303	6	7	puma32H	8192	32	0
cloud	108	4	2	puma8NH	8192	8	0
cpu	209	6	1	pw-linear	200	10	0
cpu-act	8192	21	0	pyrimidines	74	27	0
cpu-small	8192	12	0	quake	2178	3	0
delta-aileron	7129	5	0	schlvote	38	4	1
delta-elevators	9517	6	0	sensory	576	0	11
detroit	13	13	0	servo	167	0	4
diabetes-numeric	43	2	0	sleep	62	7	0
echo-months	130	6	3	stock	950	9	0
elevators	16599	18	0	strike	625	5	1
elusage	55	1	1	triazines	186	60	0
fishcatch	158	5	2	veteran	137	3	4
friedman	40768	10	0	vineyard	52	3	0
fruitfly	125	2	2	wisconsin	194	32	0
gascons	27	4	0				

Table 1: Experimental dataset characteristics.

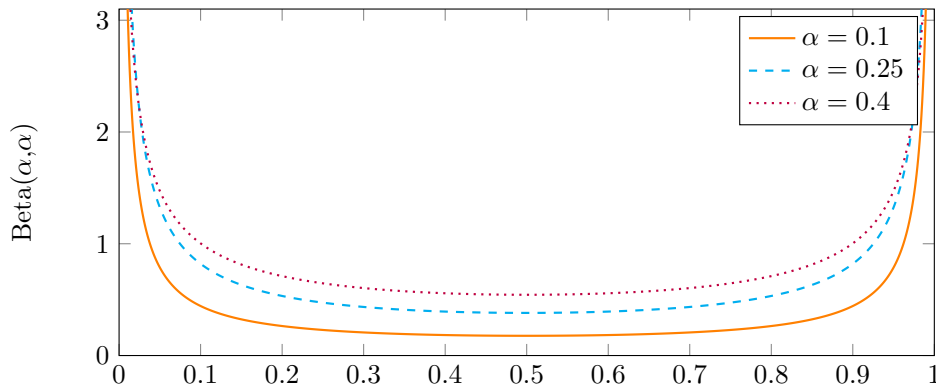


Figure 4: Beta distribution of the α values under consideration.

121 3.3. Settings

122 The experiments were performed using Weka (Hall et al., 2009). The default parameter’s
 123 values of Random Forest and Rotation Forest were used, unless otherwise specified. For
 124 Random Forest, the default number of random attributes is $\log_2(m) + 1$ where m is the
 125 number of attributes. For Rotation Forest, the default size for each group of attributes
 126 is 3. The default method for constructing the trees in Rotation Forest, which only works
 127 for classification, is J48. Hence, REPTree, a tree method for regression, was used with no
 128 pruning, as ensembles generally work better with unstable models and pruning increases
 129 stability.

130 The results were generated using a 5×2 -fold cross validation. The reported values are
 131 therefore averaged values from the 10 experiments. Three performance measures were cal-
 132 culated: RMSE (Root Mean Squared Error), MAE (Mean Absolute Error), and correlation.

133 The size of each ensemble was set at 100. The number of artificial examples to be
 134 generated was set at 50% of the training data size. Three values were applied (0.10, 0.25,
 135 and 0.40) for the α values (in the Beta distribution), from the recommended range of $[0.1, 0, 4]$
 136 in (Zhang et al., 2017). Figure 4 plots the Beta distribution for these α values.

137 One option for using mixup with nominal attributes is to transform them into numeric
 138 attributes. For example, one approach is to turn them into numerical values (that introduces
 139 an artificial order), and another is to turn them into binary attributes (greatly multiplying
 140 the attributes when there are many nominal values per attribute). Nevertheless, the mixing

141 of two nominal value attributes was done in the experiments, by randomly selecting a single
142 one. The probability of selecting the first nominal value is λ .

143 The number of artificial examples and the α value are hyper-parameters that can poten-
144 tially improve the results when adjusted for each dataset.

145 4. Results and discussion

146 Tables 2, 3, and 4 show the results for RMSE, for MAE, and for correlation, respectively.

147 *Pairwise comparisons.* Tables 5 and 6 show the number of datasets for which the column
148 method achieved better results than the row method. As 61 datasets were used in the
149 experiments, a value greater than or equal to 31 will indicate that the column method has
150 better results than the row method. It can be seen that the results are favorable for variants
151 with mixup, especially for RMSE and correlation.

152 *Relative scores.* Figure 5 shows the boxplots of the relative scores, comparing the original
153 method (Random or Rotation Forest) with the variants with mixup. The relative score for a
154 given measure is defined as $(b - a)/a$ where a and b represent the performance of the original
155 method and the performance of the variant method, respectively. When the measure is an
156 error (RMSE or MAE), negative values of the score indicate that the variant is better. In
157 contrast, positive values for correlation indicate that the variant is better. Each boxplot was
158 obtained from the relative scores of the 61 datasets. The outliers were not included in the
159 boxplots for the relative scores, as their inclusion would leave the boxes very small, because
160 the relative scores of these few datasets (outliers) are much larger.

161 The boxplots and the signs of the median values are generally favorable for the variants
162 with mixup. The only exceptions are RandMixFor and RotMixFor with $\alpha \in \{0.25, 0.40\}$
163 for MAE.

164 *Influence of α .* The following approach shows how the α values can affect the performance
165 measures. For a given dataset, method and performance measure, the values of the measure
166 were calculated for $\alpha = 0.1, 0.25, 0.4$ and then scaled to the interval $[0, 1]$. Then, a parabola

Table 5: Comparisons of Random Forest variants. Each cell shows the number of datasets where the column method is better than the row method.

(a) RMSE								
	<i>Rand For</i>	<i>MixRand For-0.10</i>	<i>MixRand For-0.25</i>	<i>MixRand For-0.40</i>	<i>RandMix For-0.10</i>	<i>RandMix For-0.25</i>	<i>RandMix For-0.40</i>	<i>Total</i>
RandFor		36	39	45	39	39	32	230
MixRandFor-0.10	24		29	30	27	28	28	166
MixRandFor-0.25	21	30		31	29	25	25	161
MixRandFor-0.40	15	29	28		25	25	20	142
RandMixFor-0.10	22	34	31	36		25	27	175
RandMixFor-0.25	21	31	34	34	36		33	189
RandMixFor-0.40	28	31	34	39	34	26		192
Total	131	191	195	215	190	168	165	

(b) MAE								
	<i>Rand For</i>	<i>MixRand For-0.10</i>	<i>MixRand For-0.25</i>	<i>MixRand For-0.40</i>	<i>RandMix For-0.10</i>	<i>RandMix For-0.25</i>	<i>RandMix For-0.40</i>	<i>Total</i>
RandFor		35	36	33	33	27	28	192
MixRandFor-0.10	25		25	28	26	22	22	148
MixRandFor-0.25	24	34		25	29	17	19	148
MixRandFor-0.40	27	31	35		28	16	18	155
RandMixFor-0.10	27	34	30	31		19	18	159
RandMixFor-0.25	33	38	42	44	40		27	224
RandMixFor-0.40	33	38	42	42	43	33		231
Total	169	210	210	203	199	134	132	

(c) Correlation								
	<i>Rand For</i>	<i>MixRand For-0.10</i>	<i>MixRand For-0.25</i>	<i>MixRand For-0.40</i>	<i>RandMix For-0.10</i>	<i>RandMix For-0.25</i>	<i>RandMix For-0.40</i>	<i>Total</i>
RandFor		40	40	46	44	44	43	257
MixRandFor-0.10	21		33	39	31	33	33	190
MixRandFor-0.25	21	28		40	34	31	35	189
MixRandFor-0.40	15	22	21		24	30	30	142
RandMixFor-0.10	17	30	27	36		30	31	171
RandMixFor-0.25	17	28	30	31	31		37	174
RandMixFor-0.40	18	28	26	31	30	24		157
Total	109	176	177	223	194	192	209	

Table 6: Comparisons of Rotation Forest variants. Each cell shows the number of datasets where the column method is better than the row method.

(a) RMSE

	<i>Rot For</i>	<i>MixRot For-0.10</i>	<i>MixRot For-0.25</i>	<i>MixRot For-0.40</i>	<i>RotMix For-0.10</i>	<i>RotMix For-0.25</i>	<i>RotMix For-0.40</i>	Total
RotFor		38	35	36	37	34	36	216
MixRotFor-0.10	23		24	24	28	22	22	143
MixRotFor-0.25	26	36		28	36	26	26	178
MixRotFor-0.40	25	36	30		37	30	26	184
RotMixFor-0.10	24	30	24	23		20	24	145
RotMixFor-0.25	27	37	33	29	39		24	189
RotMixFor-0.40	25	38	33	33	36	35		200
Total	150	215	179	173	213	167	158	

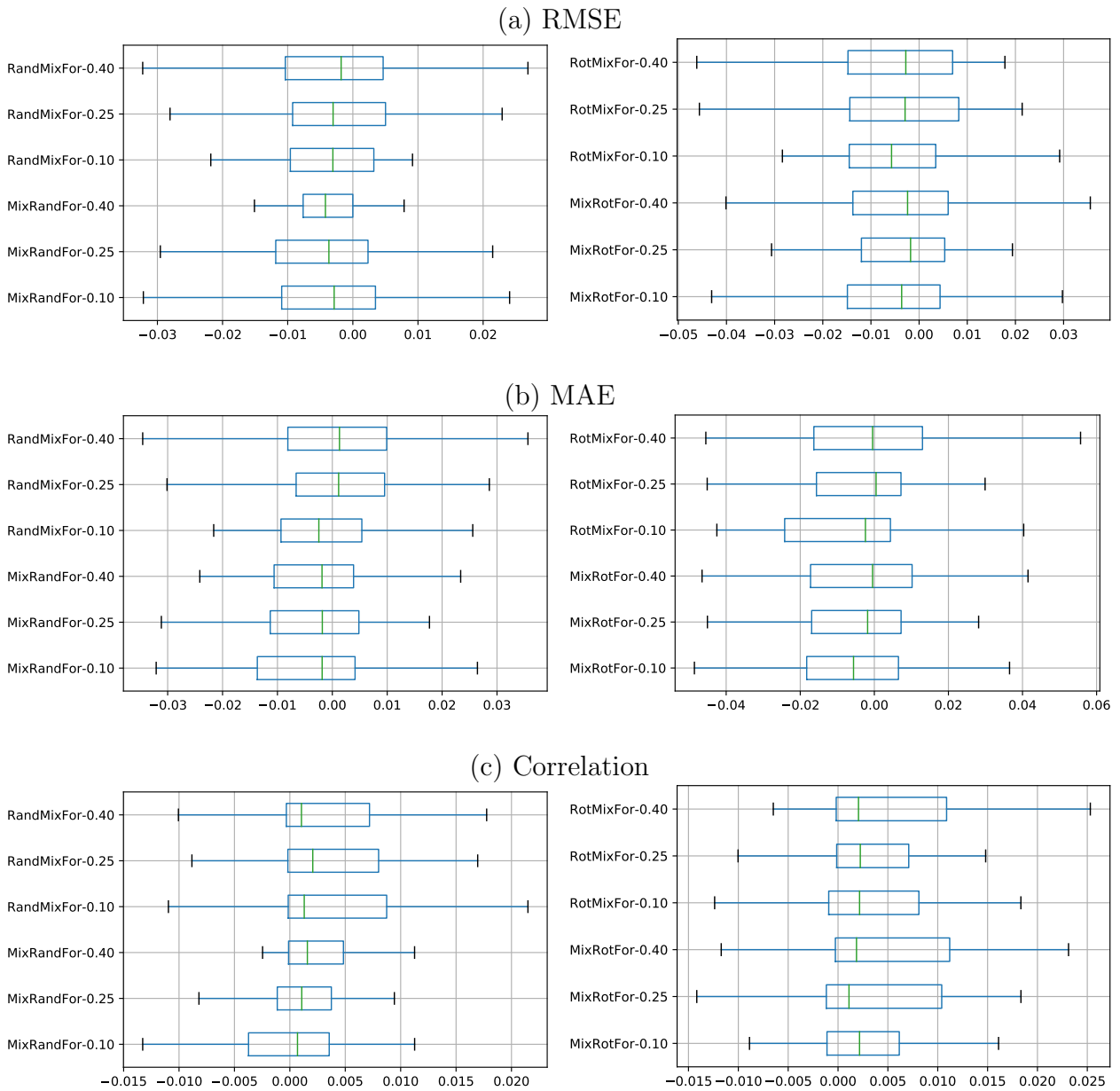
(b) MAE

	<i>Rot For</i>	<i>MixRot For-0.10</i>	<i>MixRot For-0.25</i>	<i>MixRot For-0.40</i>	<i>RotMix For-0.10</i>	<i>RotMix For-0.25</i>	<i>RotMix For-0.40</i>	Total
RotFor		39	31	31	35	29	32	197
MixRotFor-0.10	22		20	22	34	20	21	139
MixRotFor-0.25	30	40		27	38	29	22	186
MixRotFor-0.40	30	38	32		39	32	23	194
RotMixFor-0.10	26	25	21	21		19	20	132
RotMixFor-0.25	32	40	30	28	40		24	194
RotMixFor-0.40	29	39	38	36	40	35		217
Total	169	221	172	165	226	164	142	

(c) Correlation

	<i>Rot For</i>	<i>MixRot For-0.10</i>	<i>MixRot For-0.25</i>	<i>MixRot For-0.40</i>	<i>RotMix For-0.10</i>	<i>RotMix For-0.25</i>	<i>RotMix For-0.40</i>	Total
RotFor		40	41	41	40	42	43	247
MixRotFor-0.10	21		30	31	33	30	29	174
MixRotFor-0.25	20	31		31	35	25	34	176
MixRotFor-0.40	20	30	30		31	21	29	161
RotMixFor-0.10	21	27	26	29		23	32	158
RotMixFor-0.25	19	31	35	40	38		32	195
RotMixFor-0.40	18	32	27	32	29	29		167
Total	119	191	189	204	206	170	199	

Figure 5: Boxplots of relative performances. The start and end of the box are the first and third quartiles, the band inside the box is the median. Outliers are not shown.



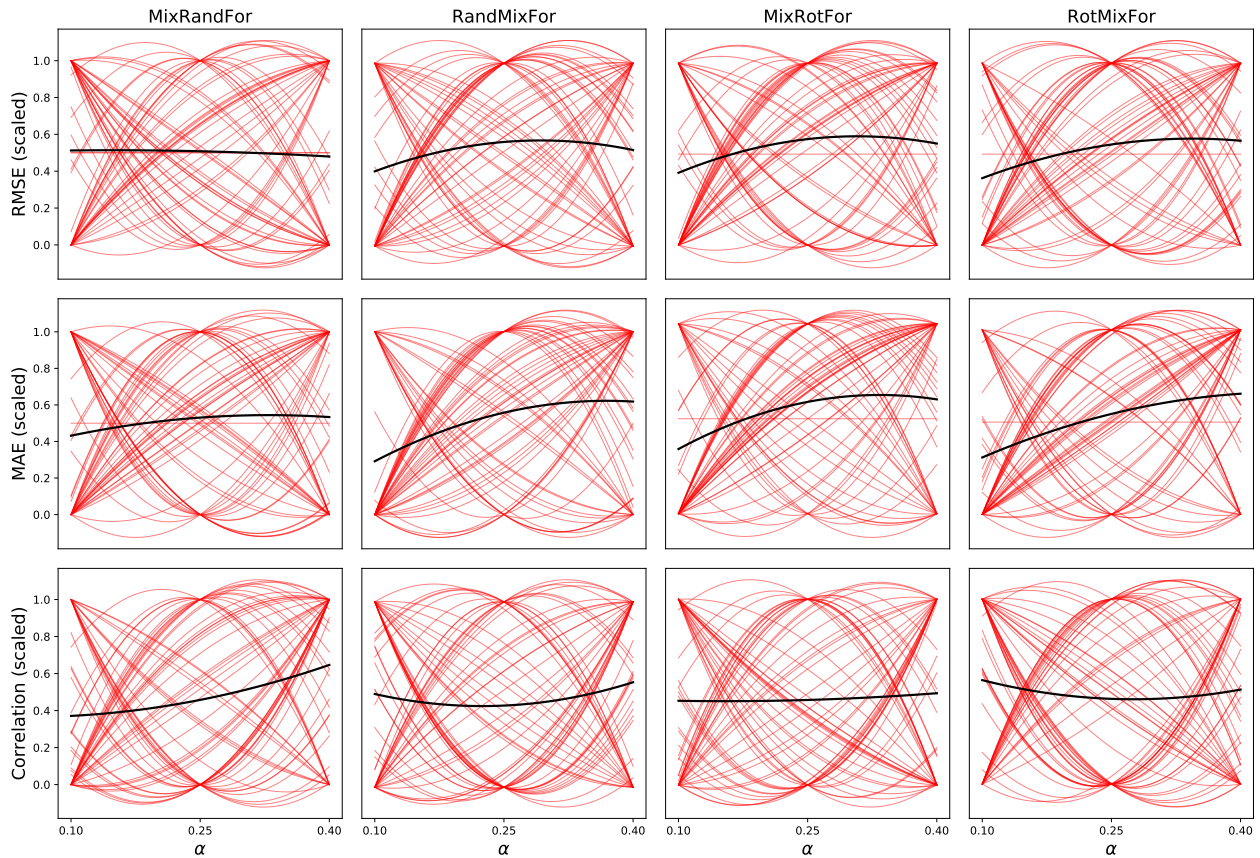


Figure 6: Scaled measures as a function of α . Each red parabola corresponds to a single dataset; the black parabola plots the average values.

167 was fitted to the three points. Figure 6 shows these parabolas, and a final parabola (shown in
 168 black) obtained by averaging the scaled values across all the datasets. There is no consistent
 169 pattern of the parabolas for the individual datasets, indicating that the optimal value of α
 170 depends on the dataset.

171 *Average ranks.* Figure 7 shows the average ranks for Random Forest and its variants with
 172 mixup. The best method is assigned rank 1, the second is assigned rank 2, and so on.
 173 The worst method is assigned rank 7, as we are comparing 7 alternatives for each ensemble
 174 method (the original ensemble, MixXXX for three values of α , and XXXMix for three
 175 values of α .) With the aim of evaluating whether some variants are significantly better than
 176 the starting method (without mixup), the Bonferroni-Dunn test was performed over the

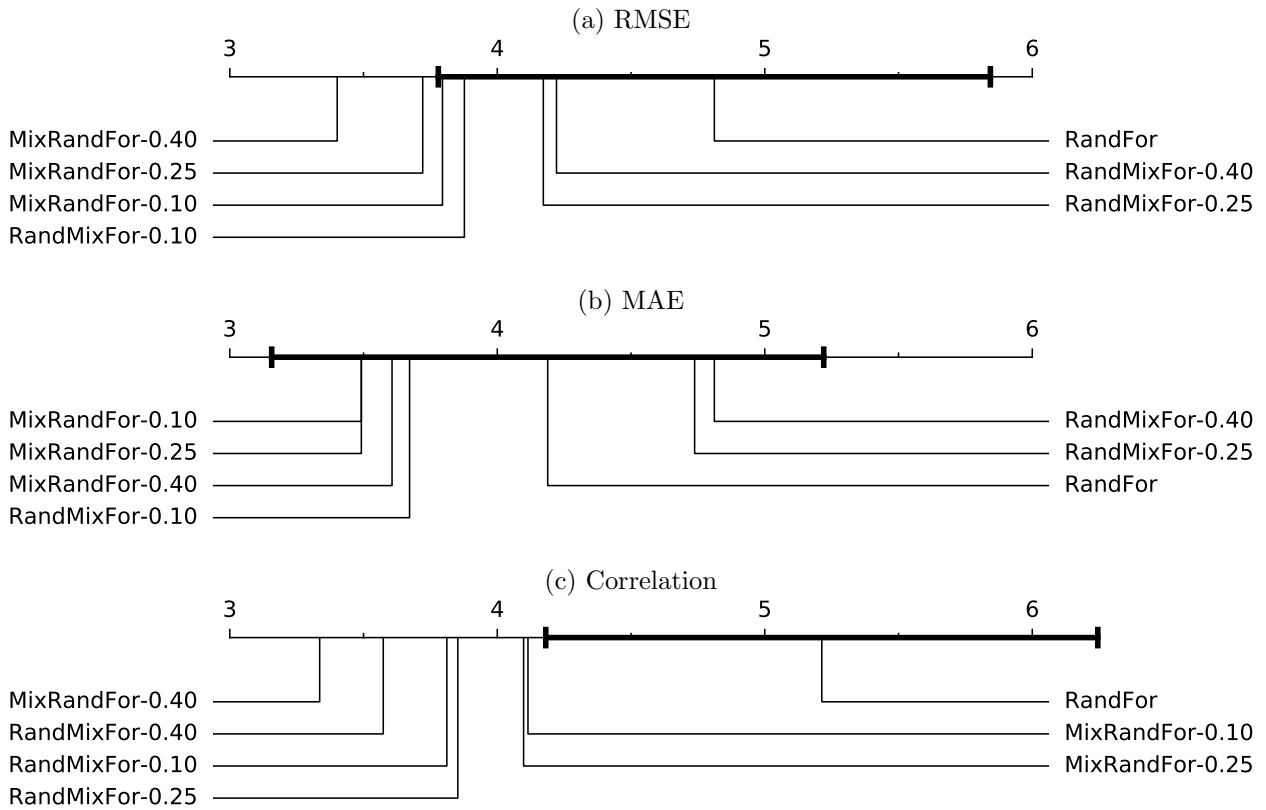


Figure 7: Comparison of Random Forest against variants with Mixup, with the Bonferroni-Dunn test. The marked interval spans the critical value and is centered at the mean rank for Random Forest. Variants with ranks outside the marked interval are significantly different ($p < 0.05$) than Random Forest.

177 ranks (Demšar, 2006) using Random or Rotation Forest as the control classifier. Random
 178 Forest without mixup had the worst average rank for RMSE and correlation. The advantage
 179 of mixup for MAE was less clear, as two variants with mixup were worse.

180 Figure 8 shows the average ranks for Rotation Forest and its mixup variants. In the
 181 same way as Random Forest, Rotation Forest without mixup shows the worst average rank
 182 for RMSE and correlation. The three variants with mixup were worse for MAE, while the
 183 other three were better.

184 Table 7 shows the average ranks for Random Forest, Rotation Forest, and their variants
 185 with mixup. Instead of having two independent ranks, one for Random Forest and the other
 186 for Rotation Forest, as with the two previous Figures (7 and 8), these tables show the ranks

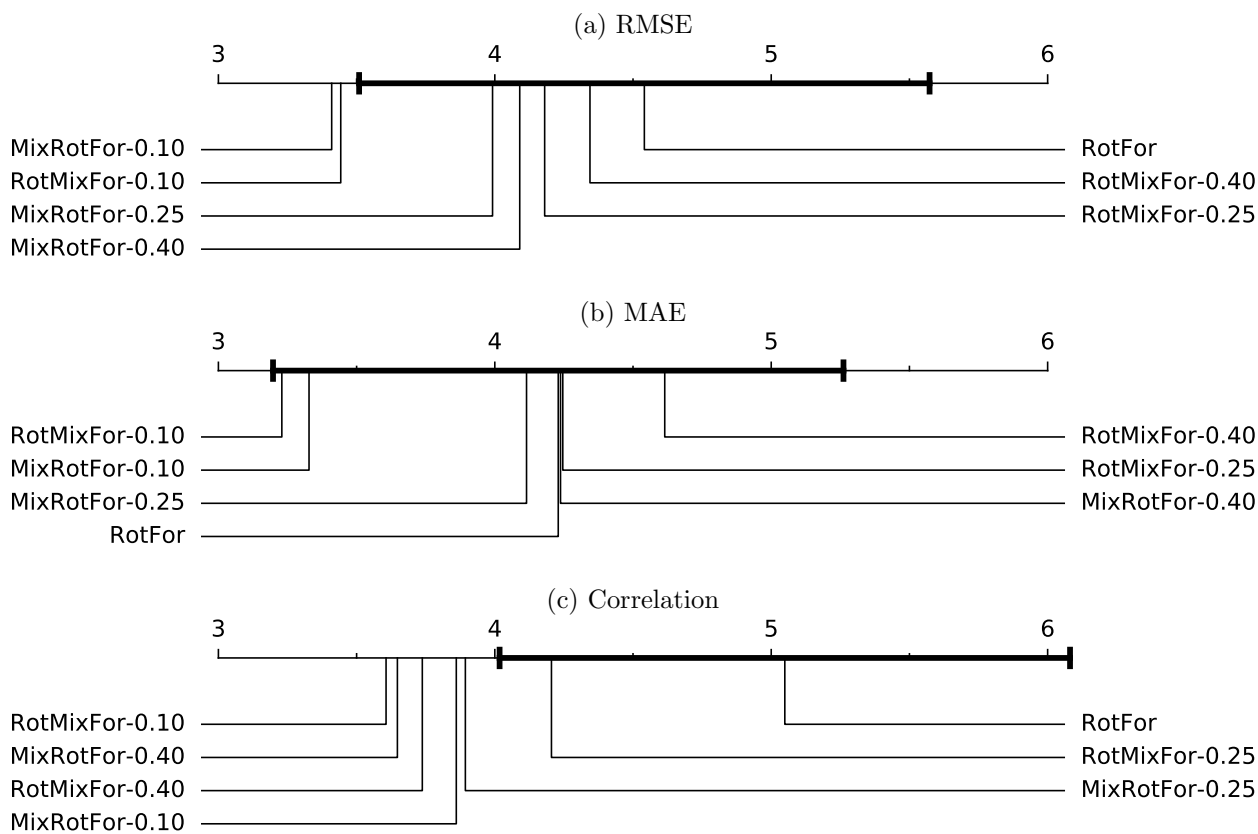


Figure 8: Comparison of Rotation Forest against variants with Mixup, with the Bonferroni-Dunn test. The marked interval spans the critical value and is centered at the mean rank for Rotation Forest. Variants with ranks outside the marked interval are significantly different ($p < 0.05$) from Rotation Forest.

Table 7: Average ranks.

RMSE		MAE		Correlation	
Method	Rank	Method	Rank	Method	Rank
RotMixFor-0.10	5.655738	RotMixFor-0.10	5.581967	RotMixFor-0.10	6.016393
MixRotFor-0.10	5.786885	MixRotFor-0.10	5.827869	RotMixFor-0.40	6.114754
MixRotFor-0.25	6.467213	RotMixFor-0.25	6.762295	MixRotFor-0.40	6.188525
RotMixFor-0.25	6.508197	MixRotFor-0.25	6.811475	MixRotFor-0.10	6.418033
MixRotFor-0.40	6.598361	MixRotFor-0.40	6.950820	MixRotFor-0.25	6.549180
RotMixFor-0.40	6.770492	RotFor	7.106557	RotMixFor-0.25	6.795082
RotFor	7.016393	RotMixFor-0.40	7.163934	MixRandFor-0.40	7.778689
MixRandFor-0.40	7.893443	MixRandFor-0.10	7.754098	RotFor	7.852459
MixRandFor-0.25	8.196721	MixRandFor-0.25	7.852459	RandMixFor-0.40	7.983607
MixRandFor-0.10	8.221311	RandMixFor-0.10	7.885246	RandMixFor-0.10	8.155738
RandMixFor-0.10	8.418033	MixRandFor-0.40	8.098361	RandMixFor-0.25	8.213115
RandMixFor-0.25	8.893443	RandFor	8.540984	MixRandFor-0.10	8.491803
RandMixFor-0.40	9.040984	RandMixFor-0.25	9.311475	MixRandFor-0.25	8.491803
RandFor	9.532787	RandMixFor-0.40	9.352459	RandFor	9.950820

187 for all the methods together. With regard to RMSE, all the Rotation Forest variants are
188 above all the Random Forest variants. Moreover, the two original methods (without mixup)
189 are the last methods in their respective sets. Likewise, with regard to MAE, the Rotation
190 Forest variants are above all the Random Forest variants, although there are a few variants
191 with mixup below the method without mixup. The methods without mixup for correlation
192 are below all the other methods in their set, although there is some overlap between the two
193 sets, because RandMixFor-0.40 is above RotFor.

194 Figures 9 and 10 show boxplots for the ranks of the different datasets. Both the Random
195 Forest and the Rotation Forest variants are independently depicted in Figure 9, so the rank
196 values range from 1 to 7. The Random Forest and the Rotation Forest variants are jointly
197 depicted in Figure 10, so the rank values range from 1 to 14. These figures support the idea
198 that the use of mixup variants is advisable.

199 Overall, Rotation Forest shows better performance compared to Random Forest, and
200 mixup offers an advantage for both ensemble methods, which has been empirically demon-
201 strated in our experiment.

Figure 9: Boxplots for the ranks. The boxplots to the left refer to the Random Forest variants and those to the right refer to the Rotation Forest variants.

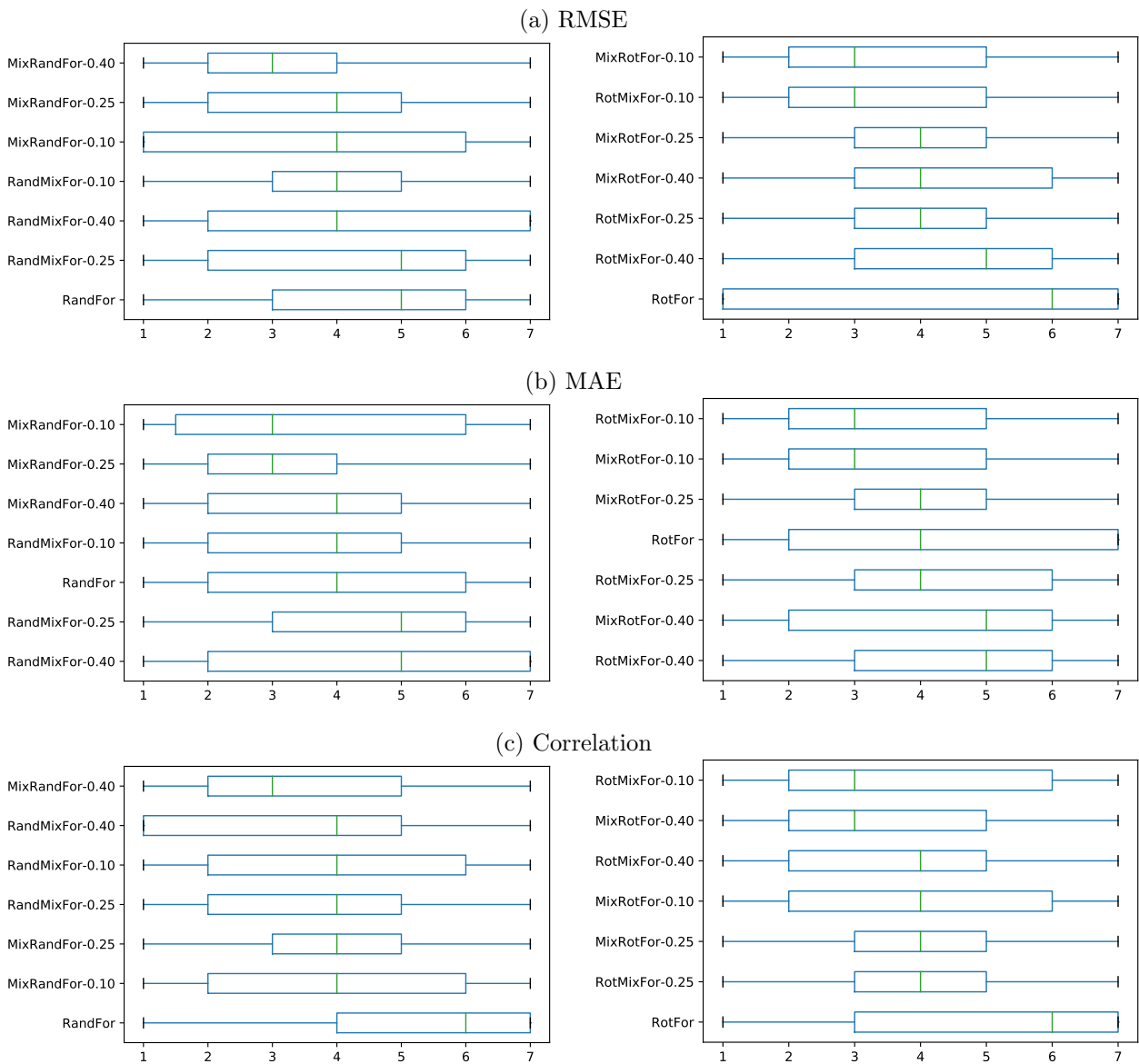
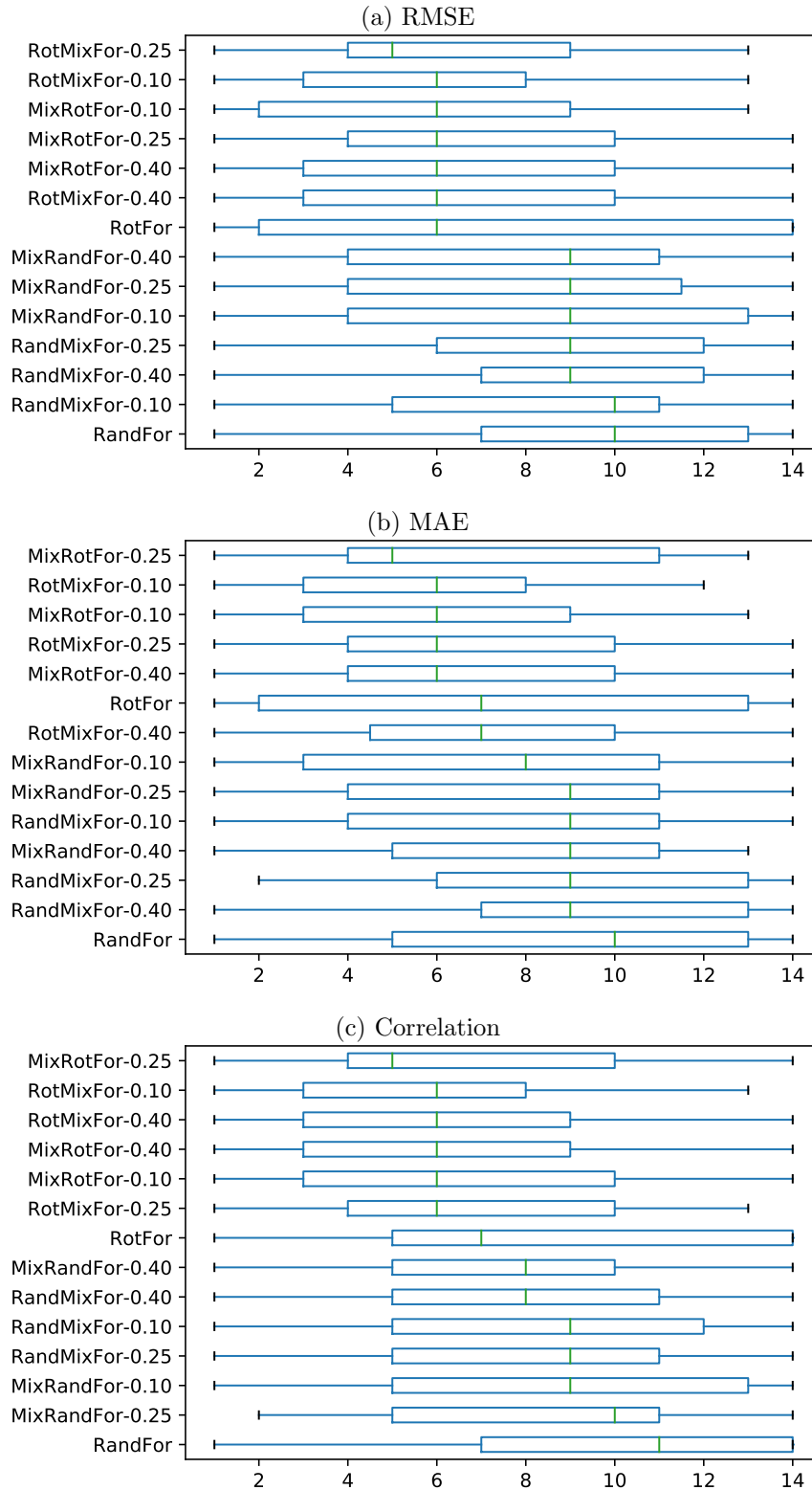


Figure 10: Boxplots for the ranks. The ranks are obtained using both Random and Rotation Forests variants.



202 *Limitations.* The scope of this study is nevertheless limited. The two parameters of the
203 method, the α value for the Beta distribution, and the number of synthetic examples that
204 are generated were not adjusted for each dataset. Only three values of α were considered
205 and the number of synthetic examples was arbitrarily fixed at 50%. Ensemble size is another
206 parameter that can affect the results and that can interact with the previous parameters.
207 Moreover, the default parameter’s values for Random Forest and Rotation Forest were used
208 with no previous adjustment for the study.

209 The mixup approach has been applied to only two ensemble methods, Random Forest
210 and Rotation Forest, although it could be applied to other methods. For instance, another
211 very successful ensemble method, although not commonly used for regression, is boost-
212 ing (Solomatine & Shrestha, 2004). The mixup approach can also be used with ensembles
213 by combining other regression methods rather than classification trees. The usefulness of
214 the mixup approach for regression ensembles with other ensembles and base methods is as
215 yet unproven.

216 The mixup method was the only method considered for generating artificial instances.
217 Other methods for generating artificial instances might be better suited for a given dataset.

218 **5. Conclusions and future research**

219 The mixup strategy has been previously used for regularizing deep neural networks,
220 although this method can also be used for increasing diversity in ensembles. In this paper,
221 we have shown that the performance of regression forest methods can be improved by using
222 the mixup strategy, which introduces artificial instances in the datasets used for training each
223 regression tree. The advantages of the mixup method have been experimentally shown for
224 both Random Forest and Rotation Forest over a broad set of 61 datasets. Our experimental
225 results favored the Rotation Forest and its improved variants.

226 Some limitations of the study can be approached in future works. The mixup method has
227 one parameter, α . We found no clear pattern of influence for the three experimental values
228 (0.1, 0.25, and 0.4). Adjusting α for each dataset and varying the number of generated
229 artificial instances can both potentially improve the results.

230 Mixup forest can be applied to other ensemble methods, such as boosting variants. It
231 can also be used with ensembles formed by other regression models instead of trees.

232 A future research line is the adaptation of the mixup method for classification datasets.
233 As mentioned earlier, the use of mixup for regression is straightforward, because the output
234 value is continuous. Nevertheless, the application of this method to classification requires
235 a previous decision on the best way of combining different nominal classes. The method
236 could also be useful in problems with several outputs, such as multi-label classification and
237 multi-target regression.

238 The distribution of the instances can make the mixup strategy counterproductive, be-
239 cause it may add noise in a localized region of the space. With this in mind, further research
240 on the convexity of the space could help clarify the advisability of applying mixup. More-
241 over, more advanced data augmentation techniques that take into account the manifold of
242 the actual instances would be interesting to explore (Guo et al., 2018; Verma et al., 2018).

243 Recently, imbalance for regression has been studied (Torgo et al., 2013). The evaluation
244 of whether mixup can be used to work with imbalanced datasets is also a promising area for
245 future research.

246 **Acknowledgments**

247 This work was supported through project TIN2015-67534-P (MINECO/FEDER, UE) of
248 the *Ministerio de Economía y Competitividad* of the Spanish Government, project BU085P17
249 (JCyL/FEDER, UE) of the *Junta de Castilla y León* (both projects co-financed through Eu-
250 ropean Union FEDER funds), and by the *Consejería de Educación* of the *Junta de Castilla*
251 *y León* and the European Social Fund through a pre-doctoral grant (EDU/1100/2017). The
252 second author is grateful for a Mobility Grant from the *Universidad de Burgos*. The third
253 author is grateful for a Mobility Grant (CAS19/00100) from the *Ministerio de Ciencia,*
254 *Innovación y Universidades* of the Spanish Government. The authors gratefully acknowl-
255 edge the support of NVIDIA Corporation and its donation of the TITAN Xp GPUs that
256 facilitated this research.

257 **References**

- 258 Bagnall, A., Bostrom, A., Cawley, G., Flynn, M., Large, J., & Lines, J. (2018). Is rotation forest the best
259 classifier for problems with continuous features? *arXiv preprint arXiv:1809.06705*, .
- 260 Beckham, C., Honari, S., Lamb, A., Verma, V., Ghadiri, F., Devon Hjelm, R., & Pal, C. (2019). Adversarial
261 mixup resynthesizers. *arXiv e-prints*, .
- 262 Breiman, L. (2001). Random forests. *Machine learning*, *45*, 5–32.
- 263 Chawla, N., Bowyer, K., Hall, L., & Kegelmeyer, W. (2002). Smote: Synthetic minority over-sampling
264 technique. *Journal of Artificial Intelligence Research*, *16*, 321–357.
- 265 Chawla, N., Lazarevic, A., Hall, L., & Bowyer, K. (2003). SMOTEBoost: Improving prediction of the
266 minority class in boosting. In *7th European Conference on Principles and Practice of Knowledge Discovery
267 in Databases(PKDD 2003)* (pp. 107–119).
- 268 Chen, Z., Chen, W., & Shi, Y. (2020). Ensemble learning with label proportions for bankruptcy prediction.
269 *Expert Systems with Applications*, *146*, 113155. doi:10.1016/j.eswa.2019.113155.
- 270 Choi, H., Son, H., & Kim, C. (2018). Predicting financial distress of contractors in the construction industry
271 using ensemble learning. *Expert Systems with Applications*, *110*, 1 – 10. doi:10.1016/j.eswa.2018.05.
272 026.
- 273 Demšar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning
274 research*, *7*, 1–30.
- 275 Fernández-Delgado, M., Cernadas, E., Barro, S., & Amorim, D. (2014). Do we need hundreds of classifiers
276 to solve real world classification problems? *Journal of Machine Learning Research*, *15*, 3133–3181. URL:
277 <http://jmlr.org/papers/v15/delgado14a.html>.
- 278 Frank, E., & Pfahringer, B. (2006). Improving on bagging with input smearing. In W.-K. Ng, M. Kit-
279 suregawa, J. Li, & K. Chang (Eds.), *Advances in Knowledge Discovery and Data Mining* (pp. 97–106).
280 Berlin, Heidelberg: Springer Berlin Heidelberg.
- 281 Galar, M., Fernandez, A., Barrenechea, E., Bustince, H., & Herrera, F. (2012). A review on ensembles
282 for the class imbalance problem: Bagging-, boosting-, and hybrid-based approaches. *IEEE Transactions
283 on Systems, Man and Cybernetics Part C: Applications and Reviews*, *42*, 463–484. doi:10.1109/TSMCC.
284 2011.2161285.
- 285 García-Pedrajas, N., Pérez-Rodríguez, J., García-Pedrajas, M., Ortiz-Boyer, D., & Fyfe, C. (2012). Class im-
286 balance methods for translation initiation site recognition in DNA sequences. *Knowledge-Based Systems*,
287 *25*, 22–34. doi:10.1016/j.knosys.2011.05.002.
- 288 Geng, G. G., Wang, C. H., Li, Q. D., Xu, L., & Jin, X. B. (2007). Boosting the performance of web spam
289 detection with ensemble under-sampling classification. In *Proceedings - Fourth International Conference
290 on Fuzzy Systems and Knowledge Discovery, FSKD 2007* (pp. 583–587). volume 4. doi:10.1109/FSKD.

291 2007.207.

292 González, S., García, S., Lázaro, M., Figueiras-Vidal, A. R., & Herrera, F. (2017). Class switching according
 293 to nearest enemy distance for learning from highly imbalanced data-sets. *Pattern Recognition*, *70*, 12–24.

294 Guo, H., Mao, Y., & Zhang, R. (2018). Mixup as locally linear out-of-manifold regularization. *CoRR*,
 295 *abs/1809.02499*. URL: <http://arxiv.org/abs/1809.02499>.

296 Haixiang, G., Yijing, L., Shang, J., Mingyun, G., Yuanyue, H., & Bing, G. (2017). Learning from class-
 297 imbalanced data: Review of methods and applications. *Expert Systems with Applications*, *73*, 220 – 239.
 298 doi:10.1016/j.eswa.2016.12.035.

299 Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. H. (2009). The weka data
 300 mining software: an update. *SIGKDD Explor. Newsl.*, *11*, 10–18. doi:10.1145/1656274.1656278.

301 Han, H., Wang, W.-Y., & Mao, B.-H. (2005). Borderline-smote: A new over-sampling method in im-
 302 balanced data sets learning. In D.-S. Huang, X.-P. Zhang, & G.-B. Huang (Eds.), *Advances in Intel-*
 303 *ligent Computing: International Conference on Intelligent Computing, ICIC 2005, Hefei, China, Au-*
 304 *gust 23-26, 2005, Proceedings, Part I* (pp. 878–887). Berlin, Heidelberg: Springer Berlin Heidelberg.
 305 doi:10.1007/11538059_91.

306 He, H., Bai, Y., Garcia, E. A., & Li, S. (2008). ADASYN: Adaptive synthetic sampling approach for
 307 imbalanced learning. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World*
 308 *Congress on Computational Intelligence)* (pp. 1322–1328). doi:10.1109/IJCNN.2008.4633969.

309 Inoue, H. (2018). Data augmentation by pairing samples for images classification. *CoRR*, *abs/1801.02929*.
 310 URL: <http://arxiv.org/abs/1801.02929>.

311 Kuncheva, L. I. (2014). *Combining pattern classifiers: methods and algorithms*. John Wiley & Sons.

312 Kuncheva, L. I., & Whitaker, C. J. (2003). Measures of diversity in classifier ensembles and their relationship
 313 with the ensemble accuracy. *Machine Learning*, *51*, 181–207. doi:10.1023/A:1022859003006.

314 Lindenbaum, O., Stanley, J., Wolf, G., & Krishnaswamy, S. (2018). Geometry based data generation. In
 315 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances*
 316 *in Neural Information Processing Systems 31* (pp. 1400–1411). Curran Associates, Inc. URL: <http://papers.nips.cc/paper/7414-geometry-based-data-generation.pdf>.

317 //papers.nips.cc/paper/7414-geometry-based-data-generation.pdf.

318 Marqués, A., García, V., & Sánchez, J. (2012). Two-level classifier ensembles for credit risk assessment.
 319 *Expert Systems with Applications*, *39*, 10916 – 10922. doi:10.1016/j.eswa.2012.03.033.

320 Martínez-Muñoz, G., & Suárez, A. (2005). Switching class labels to generate classification ensembles. *Pattern*
 321 *Recognition*, *38*, 1483–1494.

322 Mayo, M., & Frank, E. (2017). Improving naive bayes for regression with optimised artificial surrogate data.
 323 *CoRR*, *abs/1707.04943*. URL: <http://arxiv.org/abs/1707.04943>.

324 Melville, P., & Mooney, R. J. (2003). Constructing diverse classifier ensembles using artificial training

325 examples. In *IJCAI* (pp. 505–510). volume 3.

326 Melville, P., & Mooney, R. J. (2005). Creating diversity in ensembles using artificial data. *Information*
327 *Fusion*, *6*, 99–111.

328 Menardi, G., & Torelli, N. (2014). Training and assessing classification rules with imbalanced data. *Data*
329 *Mining and Knowledge Discovery*, *28*, 92–122. doi:10.1007/s10618-012-0295-5.

330 Mendes-Moreira, J., Soares, C., Jorge, A. M., & Sousa, J. F. D. (2012). Ensemble approaches for regression:
331 A survey. *ACM computing surveys*, *45*, 10.

332 Panigrahi, S., Kundu, A., Sural, S., & Majumdar, A. K. (2009). Credit card fraud detection: A fusion
333 approach using Dempster-Shafer theory and Bayesian learning. *Information Fusion*, *10*, 354–363. doi:10.
334 1016/j.inffus.2008.04.001.

335 Pardo, C., Diez-Pastor, J. F., García-Osorio, C., & Rodríguez, J. J. (2013). Rotation forests for regression.
336 *Applied Mathematics and Computation*, *219*, 9914–9924.

337 Rodríguez, J. J., Kuncheva, L. I., & Alonso, C. J. (2006). Rotation forest: A new classifier ensemble
338 method. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *28*, 1619–1630. URL: [http:](http://doi.ieeecomputersociety.org/10.1109/TPAMI.2006.211)
339 [//doi.ieeecomputersociety.org/10.1109/TPAMI.2006.211](http://doi.ieeecomputersociety.org/10.1109/TPAMI.2006.211).

340 Ross, A., & Jain, A. (2003). Information fusion in biometrics. *Pattern Recognition Letters*, *24*, 2115 – 2125.
341 doi:10.1016/S0167-8655(03)00079-5.

342 Sirlantzis, K., Hoque, S., & Fairhurst, M. (2008). Diversity in multiple classifier ensembles based on binary
343 feature quantisation with application to face recognition. *Applied Soft Computing*, *8*, 437 – 445. doi:10.
344 1016/j.asoc.2005.08.002.

345 Soares, S. G., & Araújo, R. (2015). A dynamic and on-line ensemble regression for changing environments.
346 *Expert Systems with Applications*, *42*, 2935 – 2948. doi:10.1016/j.eswa.2014.11.053.

347 Solomatine, D. P., & Shrestha, D. L. (2004). Adaboost.RT: a boosting algorithm for regression problems. In
348 *2004 IEEE International Joint Conference on Neural Networks (IEEE Cat. No.04CH37541)* (pp. 1163–
349 1168 vol.2). volume 2. doi:10.1109/IJCNN.2004.1380102.

350 Summers, C., & Dinneen, M. J. (2019). Improved mixed-example data augmentation. In *2019 IEEE Winter*
351 *Conference on Applications of Computer Vision (WACV)* (pp. 1262–1270). IEEE.

352 Tay, W.-L., Chui, C.-K., Ong, S.-H., & Ng, A. C.-M. (2013). Ensemble-based regression analysis of
353 multimodal medical data for osteopenia diagnosis. *Expert Systems with Applications*, *40*, 811 – 819.
354 doi:10.1016/j.eswa.2012.08.031.

355 Tokozume, Y., Ushiku, Y., & Harada, T. (2017). Learning from between-class examples for deep sound
356 recognition. *arXiv preprint arXiv:1711.10282*, .

357 Tokozume, Y., Ushiku, Y., & Harada, T. (2018). Between-class learning for image classification. In *Pro-*
358 *ceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 5486–5494).

- 359 Torgo, L., Ribeiro, R. P., Pfahringer, B., & Branco, P. (2013). Smote for regression. In L. Correia, L. P.
360 Reis, & J. Cascalho (Eds.), *Progress in Artificial Intelligence* (pp. 378–389). Berlin, Heidelberg: Springer
361 Berlin Heidelberg.
- 362 Verma, V., Lamb, A., Beckham, C., Najafi, A., Mitliagkas, I., Courville, A., Lopez-Paz, D., & Bengio,
363 Y. (2018). Manifold mixup: Better representations by interpolating hidden states. *arXiv preprint*
364 *arXiv:1806.05236*, .
- 365 Wang, S., & Yao, X. (2009). Diversity analysis on imbalanced data sets by using ensemble models. In
366 *Computational Intelligence and Data Mining, 2009. CIDM'09. IEEE Symposium on* (pp. 324–331). IEEE.
- 367 Weng, B., Lu, L., Wang, X., Megahed, F. M., & Martinez, W. (2018). Predicting short-term stock prices
368 using ensemble methods and online data sources. *Expert Systems with Applications*, 112, 258 – 273.
369 doi:10.1016/j.eswa.2018.06.016.
- 370 Zhang, H., Cissé, M., Dauphin, Y. N., & Lopez-Paz, D. (2017). mixup: Beyond empirical risk minimization.
371 *CoRR*, abs/1710.09412. URL: <http://arxiv.org/abs/1710.09412v2>.
- 372 Zhu, T., Lin, Y., & Liu, Y. (2017). Synthetic minority oversampling technique for multiclass imbalance prob-
373 lems. *Pattern Recognition*, 72, 327 – 340. URL: [http://www.sciencedirect.com/science/article/
374 pii/S0031320317302947](http://www.sciencedirect.com/science/article/pii/S0031320317302947). doi:10.1016/j.patcog.2017.07.024.