

# Background Check: A general technique to build more reliable and versatile classifiers.

## Supplementary material

Miquel Perello Nieto<sup>\*†</sup>, Telmo M. Silva Filho<sup>\*‡</sup>, Meelis Kull<sup>†</sup> and Peter Flach<sup>†</sup>

<sup>†</sup>Intelligent Systems Laboratory, University of Bristol, UK

<sup>‡</sup>Centro de Informatica, Universidade Federal de Pernambuco, Brazil

Email: <sup>†</sup>firstname.lastname@bristol.ac.uk, <sup>‡</sup>tmsf@cin.ufpe.br

### REFERENCES

- [1] M. Lichman, "UCI machine learning repository," 2013.
- [2] K. Hempstalk, E. Frank, and I. H. Witten, *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2008, Antwerp, Belgium, September 15-19, 2008, Proceedings, Part I*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, ch. One-Class Classification by Combining Density and Class Probability Estimation, pp. 505–519.
- [3] C. Ferri and J. Hernández-Orallo, "Cautious classifiers," *Proceedings of ROC Analysis in Artificial Intelligence, 1st International Workshop (ROCAI- 2004)*, vol. 4, pp. 27–36, 2004.
- [4] L. Li, Q. Hu, X. Wu, and D. Yu, "Exploration of classification confidence in ensemble learning," *Pattern Recognition*, vol. 47, no. 9, pp. 3120 – 3131, 2014.
- [5] D. Tax and R. Duin, "Growing a multi-class classifier with a reject option," *Pattern Recognition Letters*, vol. 29, no. 10, pp. 1565–1570, jul 2008.
- [6] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

### I. PROOFS

#### Proposition 1.

$$p(b|x) = \frac{1}{1+r(x)},$$

$$p(f_c|x) = \frac{p(f_c|f, x)r(x)}{1+r(x)} \quad \text{for } c = 1, \dots, k$$

*Proof.*  $\frac{1}{1+r(x)} = \frac{1}{(p(b|x)+p(f|x))/p(b|x)} = \frac{1}{p(b|x)} = p(b|x)$ ;  
 $\frac{p(f_c|f, x)r(x)}{1+r(x)} = p(f_c|f, x) \frac{p(f|x)}{p(b|x)} p(b|x) = p(f_c, f|x) = p(f_c|x)$ .  $\square$

#### Proposition 2.

$$q_f(x) = \frac{p(x|f)}{\max_x p(x|f)}, p(x|f) = \frac{q_f(x)}{\int_x q_f(x) dx}.$$

*Proof.*  $q_f(x) = \frac{p(x|f)}{\max_x p(x|f)} = \frac{p(x|f)p(f)}{\max_x p(x|f)p(f)} = \frac{p(x|f)}{\max_x p(x|f)}$ ,  
 $\int_x q_f(x) dx = \int_x \frac{p(x|f)}{\max_x p(x|f)} dx = \frac{\int_x p(x|f) dx}{\max_x p(x|f)} = \frac{1}{\max_x p(x|f)}$ ,  
 $p(x|f) = q_f(x) \max_x p(x|f) = q_f(x) \frac{1}{\int_x q_f(x) dx}$ .  $\square$

**Proposition 3.** If  $\mu$  is the affine background bias with  $\mu(0) = 0$ , then  $p(f|x)$  is a monotonically decreasing function of  $\mu(1)$  of the form  $p(f|x) = 1/(\mu(1) + 1)$ .

*Proof.* We have that:

$$p(f|x) = \frac{q_f}{q_b + q_f}.$$

Applying the affine background bias we get:

$$p(f|x) = \frac{q_f}{(1 - q_f)\mu(0) + q_f\mu(1) + q_f}.$$

Finally, with  $\mu(0) = 0$  and eliminating  $q_f$ , we arrive at:

$$p(f|x) = \frac{1}{\mu(1) + 1}.$$

$\square$

**Proposition 4.** Let  $\mu$  be the affine background bias with  $\mu(0) = 0$ , then for a given rejection threshold  $\theta$ ,  $\mu(1) = \theta$ .

*Proof.* Following Chow's rule, for a  $k$ -class cautious classification problem, the minimum condition such that an instance  $x$  can be accepted and classified by the model is:

$$p(f_c|x) = \theta,$$

where  $f_c$  represents the foreground class with the highest class conditional probability for instance  $x$ . In our  $(k + 1)$ -class setting, with the extra class being background, this condition is rewritten as:

$$p(f_c|f, x) = p(b|x), \text{ and } p(f_c|x)p(f|x) = p(b|x)$$

Substituting  $p(f_c|x) = \theta$  and  $p(b|x) = 1 - p(f|x)$  and isolating  $p(f|x)$ , we get:

$$p(f|x) = \frac{1}{\theta + 1}$$

Then, from Proposition 3, we arrive at  $\mu(1) = \theta$   $\square$

### II. ALGORITHMS

#### III. DATA PREPROCESSING

In order to demonstrate the versatility of BC we selected 41 datasets from UCI [1]. Half of them have been used previously in publications [2], [3], [4], [5] which we cite and/or compare against. Because of our interest in multiclass classification problems we selected 20 additional datasets with more than 3 classes. When the datasets were available from the dataset repository mldata.org we used the Python library scikit-learn

---

**Algorithm 1** Training BCD

---

**Require:**

Number of *foreground* classes  $k$ ;  
If  $k > 1$ , the  $k$ -class *foreground* classifier;

**Algorithm:**

- 1: Uniformly generate artificial *background* data around *foreground* data;
- 2: Train a binary discriminative classifier of *foreground* vs *background*;
- 3: **if**  $k > 1$  **then**
- 4:   Combine classifiers into a  $(k + 1)$ -class posterior probability estimator;
- 5: **end if**

**return**  $(k + 1)$ -class posterior probability estimator.

---

---

**Algorithm 2** Testing BCR

---

**Require:**

Number of *foreground* classes  $k$ ;  
If  $k > 1$ , the  $k$ -class *foreground* classifier;  
background bias  $\mu$ ;  
One-class model trained on *foreground* data;

**Algorithm:**

- 1: Obtain  $q_f$  from the one-class model;
- 2: Estimate  $q_b$  as  $\mu(q_f)$ ;
- 3: Estimate posterior probabilities  $p(b|x)$  and  $p(f|x)$ ;
- 4: **if**  $k > 1$  **then**
- 5:   Obtain  $k$ -class probability vector from *foreground* classifier;
- 6:   Combine calibrated probabilities into a  $(k + 1)$ -vector;
- 7: **end if**

**return**  $(k + 1)$ -class posterior probability estimates.

---

---

**Algorithm 3** Cautious classification with BC

---

**Require:**

$k$ -class *foreground* classifier;  
 $k$ -class rejection threshold  $\theta$ ;

**Algorithm:**

- 1: Set  $\mu(1) = \theta$ ;
- 2: Estimate posterior probabilities  $p(b|x)$  and  $p(f|x)$ ;
- 3: Obtain  $k$ -class probability vector from *foreground* classifier;
- 4: Combine probabilities into a  $(k + 1)$ -vector;
- 5: For every instance  $x$  predict  $\hat{y} = \text{argmax}_i(p(y = i|x))$ ;
- 6: Reject  $x$  if  $\hat{y} = (k + 1)$ ;

**return** Predictions.

---

---

**Algorithm 4** Outlier detection with BC–training phase

---

**Require:**

Number of *foreground* classes  $k$ ;  
 $k$ -class *foreground* classifier;  
background bias  $\mu$ ;

**Algorithm:**

- 1: **if**  $\mu(0) = \mu(1) = 0.5$  **then**
- 2:   Obtain  $(k + 1)$ -class posterior probability estimator with BCD;
- 3: **else**
- 4:   Obtain  $(k + 1)$ -class posterior probability estimator with BCR;
- 5: **end if**

**return**  $(k + 1)$ -class posterior probability estimator.

---

---

**Algorithm 5** Outlier detection with BC–test phase

---

**Require:**

Number of *foreground* classes  $k$ ;  
 $(k + 1)$ -class posterior probability estimator BC;

**Algorithm:**

- 1: Obtain  $(k + 1)$ -class posterior probability estimates from BC;
- 2: For every instance  $x$  predict  $\hat{y} = \text{argmax}_i(p(y = i|x))$ ;
- 3: Mark  $x$  as outlier if  $\hat{y} = (k + 1)$ ;

**return** Predictions.

---

nominal features were transformed into numerical values. We chose this option instead of transforming them into sparse binary representations in order to reduce the computational cost of the experiments. For each nominal feature all its values were sorted alphabetically, next they were substituted by their corresponding index in the sorted list, starting from zero. In case of missing values a special number was assigned to them and they were preprocessed as missing values in the next step.

Secondly, all samples that contained missing values were preprocessed in two different ways, depending on the proportion of instances with missing values. In datasets where less than 25% of the instances had missing values these samples were removed (46 samples from autos, 6 from cleveland, 31 from credit-approval, 8 from dermatology and 4 from wpbc). In the other cases the missing values were substituted by the mean of their corresponding feature (167 values from hepatitis, 1605 from horse and 2480 from mushroom).

Thirdly, datasets with more than 30 000 instances were reduced to 10% of their original size (letter and shuttle) in order to reduce the computational cost. Finally, all features were standardised with mean zero and variance one. Table I summarises the datasets in terms of number of samples, features and classes after preprocessing.

#### IV. TABLES

[6] to download them directly from that repository. Otherwise, we downloaded the data from the UCI webpage.

Because of the large variety of datasets, we had to preprocess and standardise them to run all our experiments. First,

Name	Samples	Features	Classes
abalone	4177	8	3
autos	159	25	6
balance-scale	625	4	3
car	1728	6	4
cleveland	297	13	5
credit-approval	653	15	2
dermatology	358	34	6
diabetes	768	8	2
ecoli	336	7	8
flare	1389	10	6
german	1000	20	2
glass	214	9	6
heart-statlog	270	13	2
hepatitis	155	19	2
horse	300	27	2
ionosphere	351	34	2
iris	150	4	3
landsat-satellite	6435	36	6
letter	3511	16	26
libras-movement	360	90	15
lung-cancer	96	7129	2
mfeat-karhunen	2000	64	10
mfeat-morphological	2000	6	10
mfeat-zernike	2000	47	10
mushroom	8124	22	2
optdigits	5620	64	10
page-blocks	5473	10	5
pendigits	10992	16	10
scene-classification	2407	294	2
segment	2310	19	7
shuttle	10154	9	7
sonar	208	60	2
spambase	4601	57	2
tic-tac	958	9	2
vehicle	846	18	4
vowel	990	10	11
waveform-5000	5000	40	3
wdbc	569	30	2
wdbc	194	33	2
yeast	1484	8	10
zoo	101	16	7

TABLE I: Description of the 41 classification datasets from UCI used for the experiments

	BC	O-norm	T-norm
abalone	48.90(3)	49.08(2)	<b>49.94(1)</b>
autos	71.75(3)	<b>74.96(1)</b>	74.75(2)
balance-car	62.36(3)	93.68(2)	<b>93.86(1)</b>
cleveland	88.54(2)	<b>91.53(1)</b>	79.96(3)
credit-a	<b>67.54(1)</b>	44.76(3)	63.63(2)
dermatol	64.27(3)	79.90(2)	<b>81.01(1)</b>
diabetes	<b>82.51(1)</b>	82.33(2)	82.26(3)
ecoli	<b>78.92(1)</b>	75.10(3)	78.44(2)
flare	83.83(2)	82.53(3)	<b>84.47(1)</b>
german	<b>59.21(1)</b>	57.36(3)	58.30(2)
glass	78.61(2)	77.71(3)	<b>79.44(1)</b>
heart-st	65.08(2)	64.73(3)	<b>71.07(1)</b>
hepatiti	<b>80.66(1)</b>	80.03(2)	78.93(3)
horse	<b>84.99(1)</b>	66.16(3)	84.21(2)
ionosphe	78.63(2)	69.48(3)	<b>82.07(1)</b>
iris	<b>87.64(1)</b>	82.15(3)	83.15(2)
landsat-	<b>80.08(1)</b>	79.66(3)	79.8(2)
letter	66.25(3)	<b>84.13(1)</b>	83.13(2)
libras-m	72.01(3)	<b>79.52(1)</b>	77.12(2)
lung-can	<b>46.01(1)</b>	43.38(2.5)	43.38(2.5)
mfeat-ka	<b>34.58(1)</b>	34.20(2.5)	34.20(2.5)
mfeat-mo	<b>84.11(1)</b>	33.41(2)	33.39(3)
mfeat-ze	71.42(3)	76.25(2)	<b>77.45(1)</b>
mushroom	<b>75.88(1)</b>	60.30(2)	59.98(3)
optdigit	88.05(3)	<b>99.77(1)</b>	99.61(2)
page-blo	87.25(3)	<b>90.88(1)</b>	87.82(2)
pendigit	<b>94.13(1)</b>	73.70(3)	90.85(2)
scene-cl	78.29(3)	<b>91.99(1)</b>	86.58(2)
segment	<b>84.81(1)</b>	33.37(2.5)	33.37(2.5)
shuttle	82.80(3)	<b>91.80(1)</b>	90.63(2)
sonar	78.66(3)	82.43(2)	<b>83.93(1)</b>
spambase	<b>65.00(1)</b>	36.07(2.5)	36.07(2.5)
tic-tac	78.36(3)	<b>85.88(1)</b>	82.55(2)
vehicle	75.25(2)	72.81(3)	<b>77.49(1)</b>
vowel	63.89(3)	<b>72.73(1)</b>	69.18(2)
waveform	71.58(3)	<b>74.80(1)</b>	72.91(2)
wdbc	<b>86.44(1)</b>	53.54(3)	53.66(2)
wdbc	<b>88.57(1)</b>	84.72(2)	82.81(3)
wdbc	<b>64.29(1)</b>	61.60(3)	61.92(2)
yeast	59.03(2)	53.74(3)	<b>67.47(1)</b>
zoo	<b>86.75(1)</b>	86.70(2)	85.34(3)
Average	<b>74.32(1.90)</b>	70.95(2.14)	72.59(1.95)

TABLE II: Mean accuracies for each dataset and 20 iterations of 5-fold cross-validation for Background Check, O-norm and T-norm methods [5]. The number in brackets represent the rankings of the three methods per dataset.

method	Accuracy		Log-loss	
	EP-CC	BC	EP-CC	BC
abalone	55.06 ± 1.5	<b>55.36 ± 1.4*</b>	3.94 ± 0.7	<b>3.32 ± 0.7***</b>
autos	67.54 ± 9.2	<b>69.49 ± 7.2*</b>	1.03 ± 0.5	<b>0.46 ± 0.3***</b>
balance-sc	<b>91.21 ± 2.3***</b>	90.54 ± 2.1	0.96 ± 0.3	<b>0.54 ± 0.4***</b>
car	71.61 ± 1.7	<b>71.63 ± 0.9</b>	2.60 ± 0.2	<b>2.26 ± 0.3***</b>
cleveland	55.95 ± 5.0	<b>58.44 ± 3.1***</b>	1.97 ± 0.5	<b>1.36 ± 0.6***</b>
credit-app	85.85 ± 2.9	<b>86.14 ± 2.8*</b>	<b>9.40 ± 0.5</b>	9.45 ± 0.6
dermatolog	96.40 ± 2.2	<b>96.45 ± 2.2</b>	0.21 ± 0.1	<b>0.05 ± 0.1***</b>
diabetes	76.66 ± 2.6	<b>77.13 ± 2.9**</b>	10.30 ± 0.6	<b>10.14 ± 0.9</b>
ecoli	<b>85.23 ± 3.6</b>	84.20 ± 5.5	0.58 ± 0.2	<b>0.37 ± 0.2***</b>
flare	39.74 ± 2.8	<b>42.82 ± 2.3***</b>	2.96 ± 0.5	<b>2.19 ± 0.7***</b>
german	<b>75.12 ± 2.5</b>	74.90 ± 2.2	<b>3.18 ± 0.7</b>	3.30 ± 1.0
glass	<b>64.50 ± 6.8***</b>	62.02 ± 6.4	1.62 ± 0.5	<b>0.95 ± 0.5***</b>
heart-stat	81.85 ± 5.1	<b>83.13 ± 5.2***</b>	6.48 ± 1.0	<b>5.21 ± 1.3***</b>
hepatitis	82.21 ± 5.9	<b>83.65 ± 5.0*</b>	11.77 ± 1.8	<b>10.96 ± 2.2*</b>
horse	78.69 ± 5.5	<b>80.94 ± 4.1***</b>	4.80 ± 1.0	<b>2.47 ± 1.1***</b>
ionosphere	86.43 ± 3.4	<b>88.82 ± 3.2***</b>	11.81 ± 0.1	<b>9.48 ± 0.8***</b>
iris	96.43 ± 3.2	<b>96.73 ± 3.1</b>	0.41 ± 0.4	<b>0.20 ± 0.3***</b>
landsat-sa	86.45 ± 0.8	<b>86.79 ± 0.8***</b>	0.53 ± 0.1	<b>0.36 ± 0.1***</b>
letter	80.96 ± 1.5	<b>81.40 ± 1.3**</b>	0.15 ± 0.0	<b>0.09 ± 0.0***</b>
libras-mov	<b>79.31 ± 4.0***</b>	76.54 ± 4.5	0.27 ± 0.1	<b>0.16 ± 0.1***</b>
lung-cance	98.70 ± 2.8	<b>99.42 ± 1.6*</b>	<b>16.50 ± 0.0</b>	<b>16.50 ± 0.0</b>
mfeat-karh	95.50 ± 1.0	<b>96.64 ± 0.9***</b>	0.08 ± 0.0	<b>0.04 ± 0.0***</b>
mfeat-morp	<b>73.63 ± 1.8</b>	73.46 ± 1.8	0.59 ± 0.1	<b>0.49 ± 0.1***</b>
mfeat-zern	81.38 ± 1.2	<b>82.96 ± 1.3***</b>	0.37 ± 0.1	<b>0.09 ± 0.0***</b>
mushroom	<b>98.86 ± 0.3***</b>	98.44 ± 0.5	8.83 ± 0.1	<b>8.65 ± 0.2***</b>
optdigits	97.72 ± 0.5	<b>98.51 ± 0.3***</b>	0.04 ± 0.0	<b>0.02 ± 0.0***</b>
page-block	96.00 ± 0.5	<b>96.14 ± 0.4**</b>	0.28 ± 0.0	<b>0.22 ± 0.0***</b>
pendigits	98.08 ± 0.3	<b>98.25 ± 0.2***</b>	0.06 ± 0.0	<b>0.03 ± 0.0***</b>
scene-clas	76.33 ± 2.0	<b>80.10 ± 1.3***</b>	4.04 ± 0.2	<b>1.72 ± 0.8***</b>
segment	94.80 ± 1.0	<b>94.85 ± 0.9</b>	0.24 ± 0.1	<b>0.15 ± 0.1***</b>
shuttle	<b>97.68 ± 0.3***</b>	97.03 ± 0.4	<b>0.10 ± 0.0***</b>	0.13 ± 0.0
sonar	73.57 ± 6.4	<b>77.76 ± 5.7***</b>	4.62 ± 1.2	<b>2.19 ± 0.8***</b>
spambase	92.71 ± 0.8	<b>92.83 ± 0.8</b>	6.82 ± 0.3	<b>6.12 ± 0.4***</b>
tic-tac	<b>67.44 ± 3.0***</b>	65.42 ± 0.9	<b>9.44 ± 1.6***</b>	11.70 ± 0.8
vehicle	78.70 ± 2.7	<b>79.64 ± 2.6***</b>	1.21 ± 0.4	<b>0.60 ± 0.3***</b>
vowel	77.58 ± 3.3	<b>78.62 ± 2.6***</b>	0.38 ± 0.2	<b>0.14 ± 0.1***</b>
waveform-5	84.95 ± 1.0	<b>86.17 ± 0.9***</b>	1.40 ± 0.4	<b>0.59 ± 0.2***</b>
wdbc	96.07 ± 1.7	<b>97.29 ± 1.5***</b>	6.61 ± 0.3	<b>6.05 ± 0.4***</b>
wdbc	74.68 ± 7.1	<b>78.52 ± 4.2***</b>	<b>2.07 ± 1.0</b>	2.37 ± 1.3
yeast	58.89 ± 2.4	<b>59.15 ± 2.4</b>	1.29 ± 0.1	<b>1.02 ± 0.2***</b>
zoo	92.71 ± 3.1	<b>95.14 ± 3.6***</b>	0.38 ± 0.2	<b>0.08 ± 0.1***</b>
Average	81.54 ± 14.19	<b>82.27 ± 13.89***</b>	3.42 ± 4.12	<b>2.98 ± 4.12***</b>

TABLE III: Mean and standard deviation of accuracy and log-loss on 41 datasets. Obtained from 20 iterations of 5-fold cross-validation. A Wilcoxon signed rank-sum test was performed for each metric and dataset; \* significant at  $p < 0.05$ ; \*\* significant at  $p < 0.005$ ; \*\*\* significant at  $p < 0.001$ .