

---

# Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with Dirichlet calibration

*Supplementary material*

---

**Meelis Kull**  
Department of Computer Science  
University of Tartu  
meelis.kull@ut.ee

**Miquel Perello-Nieto**  
Department of Computer Science  
University of Bristol  
miquel.perellonieto@bris.ac.uk

**Markus Kängsepp**  
Department of Computer Science  
University of Tartu  
markus.kangsepp@ut.ee

**Telmo Silva Filho**  
Department of Statistics  
Universidade Federal da Paraíba  
telmo@de.ufpb.br

**Hao Song**  
Department of Computer Science  
University of Bristol  
hao.song@bristol.ac.uk

**Peter Flach**  
Department of Computer Science  
University of Bristol and  
The Alan Turing Institute  
peter.flach@bristol.ac.uk

## Contents

<b>A</b>	<b>Source code</b>	<b>2</b>
<b>B</b>	<b>Proofs</b>	<b>2</b>
<b>C</b>	<b>Dirichlet calibration</b>	<b>4</b>
	C.1 Reliability diagram examples . . . . .	4
<b>D</b>	<b>Experimental setup</b>	<b>6</b>
	D.1 Datasets and performance estimation . . . . .	6
	D.2 Full example of statistical analysis . . . . .	7
<b>E</b>	<b>Results</b>	<b>9</b>
	E.1 Final ranking tables for all metrics . . . . .	9
	E.2 Final critical difference diagrams for every metric . . . . .	9
	E.3 Best calibrator hyperparameters . . . . .	10
	E.4 Comparison of classifiers . . . . .	10
	E.5 Deep neural networks . . . . .	10

## A Source code

The instructions and code for the experiments can be found on <https://dirichletcal.github.io/>.

## B Proofs

**Theorem 1** (Equivalence of generative, linear and canonical parametrisations). *The parametric families  $\hat{\mu}_{DirGen}(\mathbf{q}; \boldsymbol{\alpha}, \boldsymbol{\pi})$ ,  $\hat{\mu}_{DirLin}(\mathbf{q}; \mathbf{W}, \mathbf{b})$  and  $\hat{\mu}_{Dir}(\mathbf{q}; \mathbf{A}, \mathbf{c})$  are equal, i.e. they contain exactly the same calibration maps.*

*Proof.* We will prove that:

1. every function in  $\hat{\mu}_{DirGen}(\mathbf{q}; \boldsymbol{\alpha}, \boldsymbol{\pi})$  belongs to  $\hat{\mu}_{DirLin}(\mathbf{q}; \mathbf{W}, \mathbf{b})$ ;
2. every function in  $\hat{\mu}_{DirLin}(\mathbf{q}; \mathbf{W}, \mathbf{b})$  belongs to  $\hat{\mu}_{Dir}(\mathbf{q}; \mathbf{A}, \mathbf{c})$ ;
3. every function in  $\hat{\mu}_{Dir}(\mathbf{q}; \mathbf{A}, \mathbf{c})$  belongs to  $\hat{\mu}_{DirGen}(\mathbf{q}; \boldsymbol{\alpha}, \boldsymbol{\pi})$ .

**1.** Consider a function  $\hat{\mu}(\mathbf{q}) = \hat{\mu}_{DirGen}(\mathbf{q}; \boldsymbol{\alpha}, \boldsymbol{\pi})$ . Let us start with an observation that any vector  $\mathbf{x} = (x_1, \dots, x_k) \in (0, \infty)^k$  with only positive elements can be renormalised to add up to 1 using the expression  $\boldsymbol{\sigma}(\ln(\mathbf{x}))$ , since:

$$\boldsymbol{\sigma}(\ln(\mathbf{x})) = \mathbf{exp}(\ln(\mathbf{x})) / \left( \sum_i \exp(\ln(x_i)) \right) = \mathbf{x} / \left( \sum_i x_i \right)$$

where  $\mathbf{exp}$  is an operator applying exponentiation element-wise. Therefore,

$$\hat{\mu}(\mathbf{q}) = \boldsymbol{\sigma}(\ln(\pi_1 f_1(\mathbf{q}), \dots, \pi_k f_k(\mathbf{q})))$$

where  $f_i(\mathbf{q})$  is the probability density function of the distribution  $Dir(\boldsymbol{\alpha}^{(i)})$  where  $\boldsymbol{\alpha}^{(i)}$  is the  $i$ -th row of matrix  $\boldsymbol{\alpha}$ . Hence,  $f_i(\mathbf{q}) = \frac{1}{B(\boldsymbol{\alpha}^{(i)})} \prod_{j=1}^k q_j^{\alpha_{ij}-1}$ , where  $B(\cdot)$  denotes the multivariate beta function. Let us define a matrix  $\mathbf{W}$  and vector  $\mathbf{b}$  as follows:

$$w_{ij} = \alpha_{ij} - 1, \quad b_i = \ln(\pi_i) - \ln(B(\boldsymbol{\alpha}^{(i)}))$$

with  $w_{ij}$  and  $\alpha_{ij}$  denoting elements of matrices  $\mathbf{W}$  and  $\boldsymbol{\alpha}$ , respectively, and  $b_i, \pi_i$  denoting elements of vectors  $\mathbf{b}$  and  $\boldsymbol{\pi}$ . Now we can write

$$\begin{aligned} \ln(\pi_i f_i(\mathbf{q})) &= \ln(\pi_i) - \ln(B(\boldsymbol{\alpha}^{(i)})) + \ln \prod_{j=1}^k q_j^{\alpha_{ij}-1} \\ &= \ln(\pi_i) - \ln(B(\boldsymbol{\alpha}^{(i)})) + \sum_{j=1}^k (\alpha_{ij} - 1) \ln(q_j) \\ &= b_i + \sum_{j=1}^k w_{ij} \ln(q_j) \end{aligned}$$

and substituting this back into  $\hat{\mu}(\mathbf{q})$  we get:

$$\begin{aligned} \hat{\mu}(\mathbf{q}) &= \boldsymbol{\sigma}(\ln(\pi_1 f_1(\mathbf{q}), \dots, \pi_k f_k(\mathbf{q}))) \\ &= \boldsymbol{\sigma}(\mathbf{b} + \mathbf{W} \ln(\mathbf{q})) = \hat{\mu}_{DirLin}(\mathbf{q}; \mathbf{W}, \mathbf{b}) \end{aligned}$$

**2.** Consider a function  $\hat{\mu}(\mathbf{q}) = \hat{\mu}_{DirLin}(\mathbf{q}; \mathbf{W}, \mathbf{b})$ . Let us define a matrix  $\mathbf{A}$  and vector  $\mathbf{c}$  as follows:

$$a_{ij} = w_{ij} - \min_i w_{ij}, \quad \mathbf{c} = \boldsymbol{\sigma}(\mathbf{W} \ln \mathbf{u} + \mathbf{b})$$

with  $a_{ij}$  and  $w_{ij}$  denoting elements of matrices  $\mathbf{A}$  and  $\mathbf{W}$ , respectively, and  $\mathbf{u} = (1/k, \dots, 1/k)$  is a column vector of length  $k$ . Note that  $\mathbf{A} \mathbf{x} = \mathbf{W} \mathbf{x} + \mathbf{const}_1$  and  $\ln \boldsymbol{\sigma}(\mathbf{x}) = \mathbf{x} + \mathbf{const}_2$  for any  $\mathbf{x}$  where  $\mathbf{const}_1$  and  $\mathbf{const}_2$  are constant vectors (all elements are equal), but the constant depends on

x. Taking into account that  $\sigma(\mathbf{v} + \text{const}) = \sigma(\mathbf{v})$  for any vector  $\mathbf{v}$  and constant vector  $\text{const}$ , we obtain:

$$\begin{aligned}\hat{\mu}_{Dir}(\mathbf{q}; \mathbf{A}, \mathbf{c}) &= \sigma\left(\mathbf{A} \ln \frac{\mathbf{q}}{1/k} + \ln \mathbf{c}\right) = \sigma\left(\mathbf{W} \ln \frac{\mathbf{q}}{1/k} + \text{const}_1 + \ln \mathbf{c}\right) \\ &= \sigma\left(\mathbf{W} \ln \mathbf{q} - \mathbf{W} \ln \mathbf{u} + \text{const}_1 + \ln \sigma(\mathbf{W} \ln \mathbf{u} + \mathbf{b})\right) \\ &= \sigma\left(\mathbf{W} \ln \mathbf{q} - \mathbf{W} \ln \mathbf{u} + \text{const}_1 + \mathbf{W} \ln \mathbf{u} + \mathbf{b} + \text{const}_2\right) \\ &= \sigma\left(\mathbf{W} \ln \mathbf{q} + \mathbf{b} + \text{const}_1 + \text{const}_2\right) = \sigma\left(\mathbf{W} \ln \mathbf{q} + \mathbf{b}\right) = \hat{\mu}_{DirLin}(\mathbf{q}; \mathbf{W}, \mathbf{b}) \\ &= \hat{\mu}(\mathbf{q})\end{aligned}$$

3. Consider a function  $\hat{\mu}(\mathbf{q}) = \hat{\mu}_{Dir}(\mathbf{q}; \mathbf{A}, \mathbf{c})$ . Let us define a matrix  $\alpha$ , vector  $\mathbf{b}$  and vector  $\pi$  as follows:

$$\alpha_{ij} = a_{ij} + 1, \quad \mathbf{b} = \ln \mathbf{c} - \mathbf{A} \ln \mathbf{u}, \quad \pi_i = \exp(b_i) \cdot B(\alpha^{(i)})$$

with  $\alpha_{ij}$  and  $a_{ij}$  denoting elements of matrices  $\alpha$  and  $\mathbf{A}$ , respectively, and  $\mathbf{u} = (1/k, \dots, 1/k)$  is a column vector of length  $k$ . We can now write:

$$\begin{aligned}\hat{\mu}(\mathbf{q}) &= \hat{\mu}_{Dir}(\mathbf{q}; \mathbf{A}, \mathbf{c}) = \sigma\left(\mathbf{A} \ln \frac{\mathbf{q}}{1/k} + \ln \mathbf{c}\right) = \sigma\left(\mathbf{A} \ln \mathbf{q} - \mathbf{A} \ln \mathbf{u} + \ln \mathbf{c}\right) \\ &= \sigma\left((\alpha - \mathbf{1}) \ln \mathbf{q} + \mathbf{b}\right)\end{aligned}$$

Element  $i$  in the vector within the softmax is equal to:

$$\begin{aligned}\sum_{j=1}^k (\alpha_{ij} - 1) \ln(q_j) + b_j &= \sum_{j=1}^k (\alpha_{ij} - 1) \ln(q_j) + \ln\left(\pi_i \cdot \frac{1}{B(\alpha^{(i)})}\right) \\ &= \ln\left(\pi_i \cdot \frac{1}{B(\alpha^{(i)})} \prod_{j=1}^k q_j^{\alpha_{ij} - 1}\right) \\ &= \ln(\pi_i \cdot f_i(\alpha^{(i)}))\end{aligned}$$

and therefore:

$$\hat{\mu}(\mathbf{q}) = \sigma\left((\alpha - \mathbf{1}) \ln(\mathbf{q}) + \mathbf{b}\right) = \sigma\left(\ln(\pi_i \cdot f_i(\alpha^{(i)}))\right) = \hat{\mu}_{DirGen}(\mathbf{q}; \alpha, \pi)$$

□

The following proposition proves that temperature scaling can be viewed as a general-purpose calibration method, being a special case within the Dirichlet calibration map family.

**Proposition 1.** Let us denote the temperature scaling family by  $\hat{\mu}'_{TempS}(\mathbf{z}; t) = \sigma(\mathbf{z}/t)$  where  $\mathbf{z}$  are the logits. Then for any  $t$ , temperature scaling can be expressed as

$$\hat{\mu}'_{TempS}(\mathbf{z}; t) = \hat{\mu}_{DirLin}\left(\sigma(\mathbf{z}); \frac{1}{t} \mathbf{I}, \mathbf{0}\right)$$

where  $\mathbf{I}$  is the identity matrix and  $\mathbf{0}$  is the vector of zeros.

*Proof.* Let us first observe that for any  $\mathbf{x} \in \mathbb{R}^k$  there exists a constant vector  $\text{const}$  (all elements are equal) such that  $\ln \sigma(\mathbf{x}) = \mathbf{x} + \text{const}$ . Furthermore,  $\sigma(\mathbf{v} + \text{const}) = \sigma(\mathbf{v})$  for any vector  $\mathbf{v}$  and any constant vector  $\text{const}$ . Therefore,

$$\begin{aligned}\hat{\mu}_{DirLin}\left(\sigma(\mathbf{z}); \frac{1}{t} \mathbf{I}, \mathbf{0}\right) &= \sigma\left(\frac{1}{t} \mathbf{I} \ln \sigma(\mathbf{z})\right) \\ &= \sigma\left(\frac{1}{t} \mathbf{I} (\mathbf{z} + \text{const})\right) \\ &= \sigma\left(\frac{1}{t} \mathbf{I} \mathbf{z} + \frac{1}{t} \mathbf{I} \text{const}\right) \\ &= \sigma(\mathbf{z}/t + \text{const}') \\ &= \sigma(\mathbf{z}/t) \\ &= \hat{\mu}'_{TempS}(\mathbf{z}; t)\end{aligned}$$

where  $\text{const}' = \frac{1}{t} \mathbf{I} \text{const}$  is a constant vector as a product of a diagonal matrix with a constant vector. □

## C Dirichlet calibration

In this section we show some examples of reliability diagrams and other plots that can help to understand the representational power of Dirichlet calibration compared with other calibration methods.

### C.1 Reliability diagram examples

We will look at two examples of reliability diagrams on the original classifier and after applying 6 calibration methods. Figure 1 shows the first example for the 3 class classification dataset *balance-scale* and the classifier MLP. This figure shows the confidence-reliability diagram in the first column and the classwise-reliability diagrams in the other columns. Figure 1a shows how posterior probabilities from the MLP have small gaps between the true class proportions and the predicted means. This visualisation may indicate that the original classifier is already well calibrated. However, when we separate the reliability diagram per class, we notice that the predictions for the first class are underconfident, as indicated by high proportions of the true class at low mean predictions. On the other hand, the predictions on classes 2 and 3 are overconfident in the regions of posterior probabilities compressed between  $[0.2, 0.5]$  while being underconfident in higher regions.

Table 1: Averaged results for the confidence-ECE and classwise-ECE metrics of 6 calibrators applied on an MLP trained on the *balance-scale* dataset.

	DirL2	Beta	FreqB	Isot	WidB	TempS	Uncal
conf-ECE	<b>0.04</b> <sub>1</sub>	0.05 <sub>3</sub>	0.13 <sub>7</sub>	0.05 <sub>2</sub>	0.08 <sub>6</sub>	0.05 <sub>4</sub>	0.08 <sub>5</sub>
cw-ECE	0.12 <sub>2</sub>	0.13 <sub>3</sub>	0.29 <sub>7</sub>	<b>0.12</b> <sub>1</sub>	0.17 <sub>5</sub>	0.15 <sub>4</sub>	0.20 <sub>6</sub>

The following subfigures show how the different calibration methods try to reduce ECE, occasionally increasing the error. As can be seen in Table 1, Dirichlet L2 and One-vs.Rest isotonic regression obtain the lowest ECE while One-vs.Rest frequency binning makes the original calibration worse. Looking at Figure 1i it is possible to see how temperature scaling manages to reduce the overall overconfidence in the higher range of probabilities for classes 2 and 3, but makes the situation worse in the interval  $[0.2, 0.6]$ . However, it manages to reduce the overall ECE.

In the second example we show 3 calibration methods for a 4 class classification problem (car dataset) applied on the scores of an Adaboost SAMME classifier. Figure 2 shows one reliability diagram per class ( $C_1$  *acceptable*,  $C_2$  *good*,  $C_3$  *unacceptable*, and  $C_4$  *very good*).

From this Figure we can see that the uncalibrated model is underconfident for classes 1, 2 and 3, showing posterior probabilities never higher than 0.7, while having true class proportions higher than 0.7 in the mentioned interval. We can see that after applying some of the calibration models the posterior probabilities reach higher probability values. As can be seen in Table 2, Dirichlet L2

Table 2: Averaged results for the confidence-ECE and classwise-ECE metrics of 6 calibrators applied on an Adaboost SAMME trained on the *car* dataset.

	DirL2	Beta	FreqB	Isot	WidB	TempS	Uncal
conf-ECE	<b>0.07</b> <sub>1</sub>	0.10 <sub>4</sub>	0.12 <sub>5</sub>	0.07 <sub>2</sub>	0.09 <sub>3</sub>	0.14 <sub>7</sub>	0.14 <sub>6</sub>
cw-ECE	0.18 <sub>2</sub>	0.23 <sub>3</sub>	0.29 <sub>5</sub>	<b>0.18</b> <sub>1</sub>	0.25 <sub>4</sub>	0.32 <sub>7</sub>	0.29 <sub>6</sub>

and One-vs.Rest Isotonic Regression obtain the lowest ECE while Temperature Scaling makes the original calibration worse. Figure 2d shows how Dirichlet calibration with L2 regularisation achieved the largest spread of probabilities, also reducing the error mean gap with the predictions and the true class proportions. On the other hand, temperature scaling reduced ECE for class 1, but hurt the overall performance for the other classes.

A more detailed depiction of the previous reliability diagrams can be seen in Figure 3. In this case, the posterior probabilities are not introduced in bins, but a boxplot summarises their full distribution. The first observation here is, for the *good* and *very good* classes, the uncalibrated model tends to predict probability vectors with small variance, i.e. the outputs do not change much among different instances. Among the calibration approaches, temperature scaling still maintains this low level of

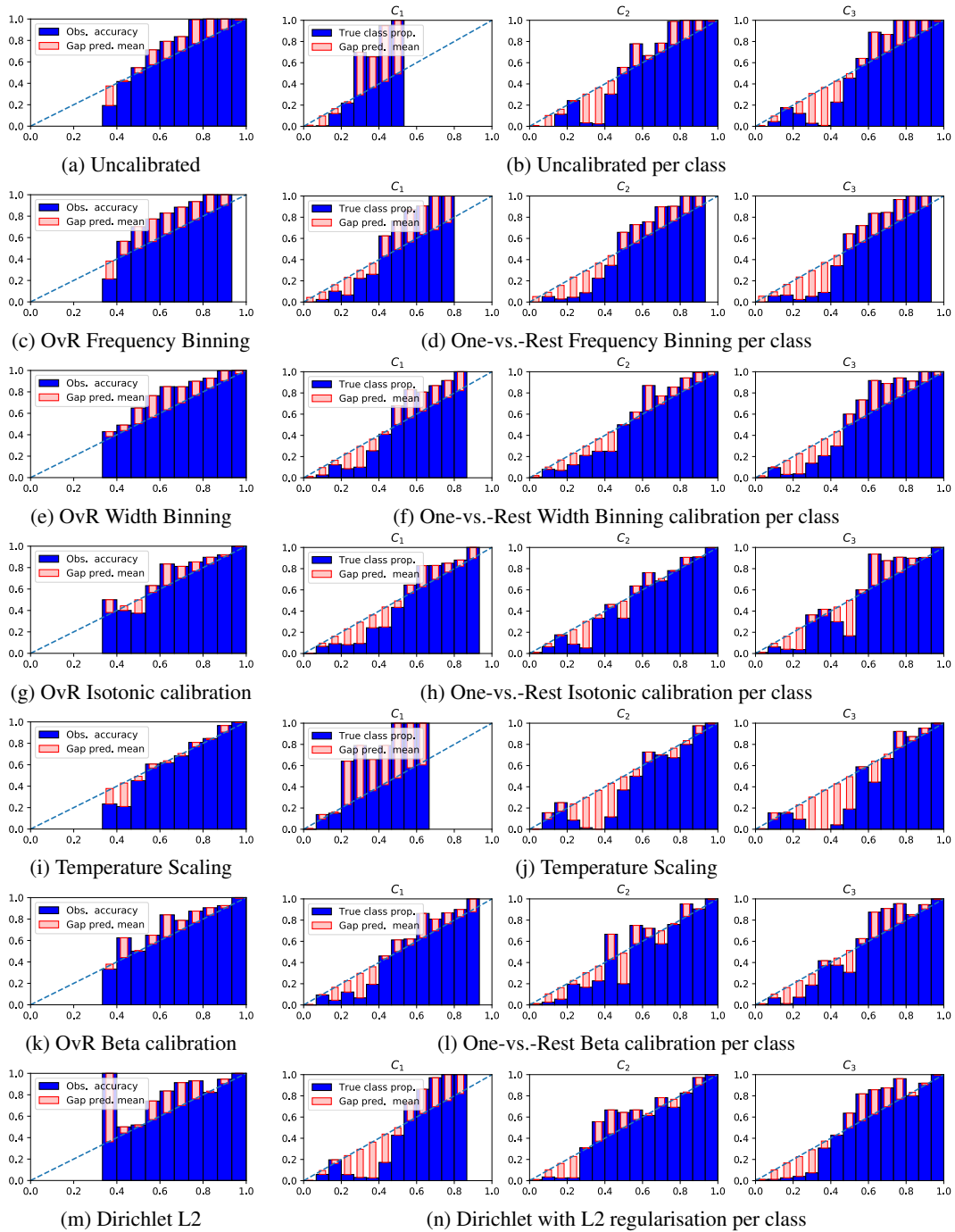


Figure 1: Confidence-reliability diagrams in the first column and classwise-reliability diagrams in the remaining columns, for a real experiment with the multilayer perceptron classifier on the balance-scale dataset and a subset of the calibrators. All the test partitions from the 5 times 5-fold-cross-validation have been aggregated to draw every plot.

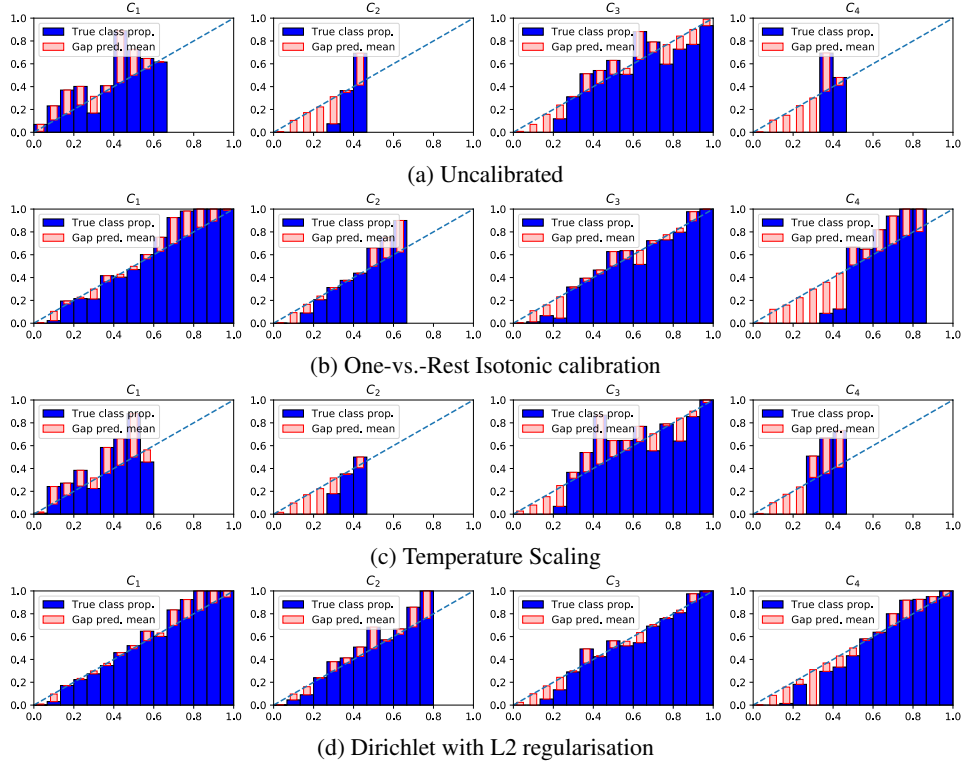


Figure 2: Reliability diagrams per class for a real experiment with the classifier Ada boost SAMME on the car dataset and 3 calibrators. The test partitions from the 5 times 5-fold-cross-validation have been aggregated to draw every plot.

variance, while both isotonic and Dirichlet L2 manage to show a higher variance on the outputs. While this observation cannot be justified here without quantitative analysis, another observation clearly shows an advantage of using Dirichlet L2. For the *acceptable* class, only Dirichlet L2 is capable of providing the highest mean probability for the correct class, while the other three methods tend to put higher probability mass on the *unacceptable* class on average.

## D Experimental setup

In this section we provide the detailed description of the experimental setup on a variety of non-neural classifiers and datasets. While our implementation of Dirichlet calibration is based on standard Newton-Raphson with multinomial logistic loss and L2 regularisation, as mentioned at the end of Section 3, existing implementations of logistic regression (e.g. scikit-learn) with the log transformed predicted probabilities can also be easily applied.

### D.1 Datasets and performance estimation

The full list of datasets, and a brief description of each one including the number of samples, features and classes is presented in Table 3.

Figure 4 shows how every dataset was divided in order to get an estimated performance for every combination of dataset, classifier and calibrator. Each dataset was divided using 5 times 5-fold-cross-validation to create 25 test partitions. For each of the 25 partitions the corresponding training set was divided further with a 3-fold-cross-validation for which the bigger portions were used to train the classifiers (and validate the calibrators if they had hyperparameters), and the small portion was used to train the calibrators. The 3 calibrators trained in the inner 3-folds were used to predict the corresponding test partition, and their predictions were averaged in order to obtain better estimates of their performance with the 7 different metrics (accuracy, Brier score, log-loss, maximum calibration

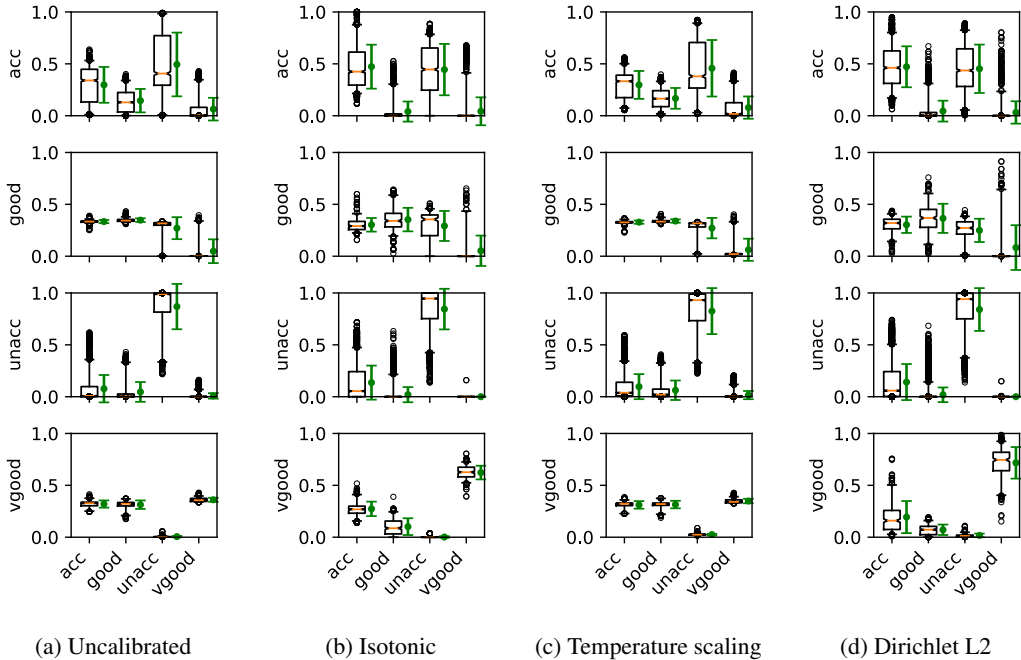


Figure 3: Effect of Dirichlet Calibration on the scores of Ada boost SAMME on the *car* dataset which is composed of 4 classes (*acceptable*, *good*, *unacceptable*, and *very good*). The whiskers of each box indicate the 5th and 95th percentile, the notch around the median indicates the confidence interval. The green error bar to the right of each box indicates one standard deviation on each side of the mean. In each subfigure, the first boxplot corresponds to the posterior probabilities for the samples of class 1, divided in 4 boxes representing the posterior probabilities for each class. A good classifier should have the highest posterior probabilities in the box corresponding to the true class. In Figure 3a it is possible to see that the first class (*acceptable*) is missclassified as belonging to the third class (*unacceptable*) with high probability values, while Dirichlet Calibration is able to alleviate that problem. Also, for the second and fourth true classes (*good*, and *very good*) the original classifier uses a reduced domain of probabilities (indicative of underconfidence), while Dirichlet calibration is able to spread these probabilities with more meaningful values (as indicated by a reduction of the calibration losses; See Figure 2).

error, confidence-ECE, classwise-ECE and the p test statistic of the ECE metrics). Finally, the 25 resulting measures were averaged.

## D.2 Full example of statistical analysis

The following is a full example of how the final rankings and statistical tests are computed. For this example, we will focus on the metric log-loss, and we will start with the naive Bayes classifier. Table 4 shows the estimated log-loss by averaging the 5-times 5-fold cross-validation log-losses of the inner 3-fold aggregated predictions. The sub-indices are the ranking of every calibrator for each dataset (ties in the ranking share the averaged rank). The resulting table of sub-indices is used to compute the Friedman test statistic, resulting in a value of 97.9 and a p-value of  $1.14e^{-17}$  indicating statistical difference between the different calibration methods. The last row contains the average ranks of the full table, which is shown in the corresponding critical difference diagram in Figure 5a. The critical difference uses the Bonferroni-Dunn one-tailed statistical test to compute the minimum ranking distance that is shown in the Figure, indicating that for this particular classifier and metric the Dirichlet calibrator with L2 regularisation is significantly better than the other methods, with the exception of Dirichlet with ODIR regularisation.

The same process is applied to each of the 11 classifiers for every metric. Table 6 shows the final average results of all classifiers. Notice that the row corresponding to naive Bayes has the rounded average rankings from Figure 5a.

dataset	n_samples	n_features	n_classes
abalone	4177	8	3
balance-scale	625	4	3
car	1728	6	4
cleveland	297	13	5
dermatology	358	34	6
glass	214	9	6
iris	150	4	3
landsat-satellite	6435	36	6
libras-movement	360	90	15
mfeat-karhunen	2000	64	10
mfeat-morphological	2000	6	10
mfeat-zernike	2000	47	10
optdigits	5620	64	10
page-blocks	5473	10	5
pendigits	10992	16	10
segment	2310	19	7
shuttle	101500	9	7
vehicle	846	18	4
vowel	990	10	11
waveform-5000	5000	40	3
yeast	1484	8	10

Table 3: Datasets used for the large-scale empirical comparison.

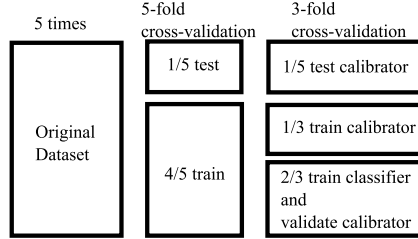
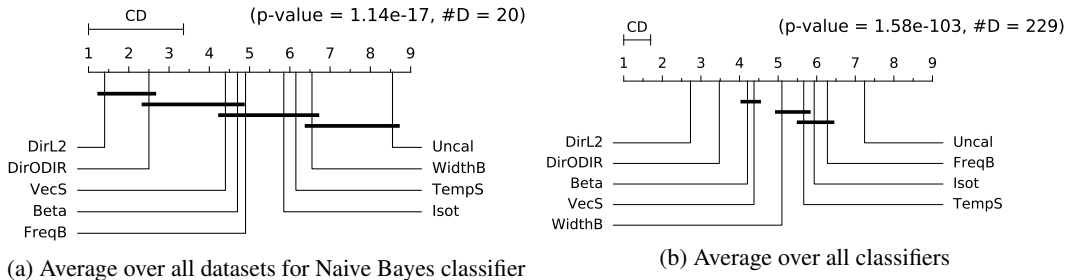


Figure 4: Partitions of each dataset in order to estimate out-of-sample measures.

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
abalone	1.95 <sub>9</sub>	<b>0.89<sub>1</sub></b>	0.89 <sub>2</sub>	0.89 <sub>5</sub>	0.89 <sub>6</sub>	0.89 <sub>4</sub>	0.90 <sub>7</sub>	0.89 <sub>3</sub>	0.92 <sub>8</sub>
balance-sc	0.48 <sub>9</sub>	0.22 <sub>2</sub>	<b>0.22<sub>1</sub></b>	0.31 <sub>4</sub>	0.41 <sub>8</sub>	0.23 <sub>3</sub>	0.38 <sub>7</sub>	0.36 <sub>6</sub>	0.35 <sub>5</sub>
car	1.57 <sub>9</sub>	0.38 <sub>2</sub>	<b>0.38<sub>1</sub></b>	0.59 <sub>6</sub>	0.92 <sub>8</sub>	0.43 <sub>3</sub>	0.56 <sub>5</sub>	0.55 <sub>4</sub>	0.67 <sub>7</sub>
cleveland	2.72 <sub>9</sub>	1.02 <sub>2</sub>	<b>1.00<sub>1</sub></b>	1.30 <sub>6</sub>	1.36 <sub>8</sub>	1.08 <sub>3</sub>	1.32 <sub>7</sub>	1.12 <sub>4</sub>	1.14 <sub>5</sub>
dermatolog	2.57 <sub>9</sub>	<b>0.20<sub>1</sub></b>	0.41 <sub>7</sub>	0.36 <sub>5</sub>	0.31 <sub>3</sub>	0.44 <sub>8</sub>	0.33 <sub>4</sub>	0.23 <sub>2</sub>	0.40 <sub>6</sub>
glass	2.93 <sub>9</sub>	<b>1.12<sub>1</sub></b>	1.57 <sub>6</sub>	1.62 <sub>8</sub>	1.35 <sub>4</sub>	1.48 <sub>5</sub>	1.58 <sub>7</sub>	1.13 <sub>3</sub>	1.12 <sub>2</sub>
iris	0.13 <sub>4</sub>	<b>0.11<sub>1</sub></b>	0.12 <sub>3</sub>	0.13 <sub>5</sub>	0.12 <sub>2</sub>	0.13 <sub>6</sub>	0.30 <sub>8</sub>	0.33 <sub>9</sub>	0.21 <sub>7</sub>
landsat-sa	3.87 <sub>9</sub>	<b>0.35<sub>1</sub></b>	0.36 <sub>2</sub>	0.55 <sub>4</sub>	0.61 <sub>7</sub>	0.55 <sub>3</sub>	0.58 <sub>5</sub>	0.58 <sub>6</sub>	0.74 <sub>8</sub>
libras-mov	4.99 <sub>9</sub>	<b>0.94<sub>1</sub></b>	1.73 <sub>7</sub>	1.32 <sub>3</sub>	1.16 <sub>2</sub>	1.33 <sub>4</sub>	2.04 <sub>8</sub>	1.66 <sub>6</sub>	1.43 <sub>5</sub>
mfeat-karh	0.45 <sub>9</sub>	<b>0.20<sub>1</sub></b>	0.20 <sub>2</sub>	0.23 <sub>3</sub>	0.24 <sub>5</sub>	0.23 <sub>4</sub>	0.39 <sub>8</sub>	0.38 <sub>7</sub>	0.29 <sub>6</sub>
mfeat-morp	1.78 <sub>9</sub>	<b>0.71<sub>1</sub></b>	0.79 <sub>2</sub>	0.91 <sub>7</sub>	1.01 <sub>8</sub>	0.86 <sub>5</sub>	0.85 <sub>4</sub>	0.82 <sub>3</sub>	0.88 <sub>6</sub>
mfeat-zern	1.73 <sub>9</sub>	0.60 <sub>2</sub>	<b>0.59<sub>1</sub></b>	0.71 <sub>3</sub>	0.76 <sub>5</sub>	0.71 <sub>4</sub>	0.85 <sub>8</sub>	0.82 <sub>6</sub>	0.84 <sub>7</sub>
optdigits	3.23 <sub>9</sub>	0.46 <sub>3</sub>	0.46 <sub>2</sub>	0.54 <sub>6</sub>	0.83 <sub>7</sub>	0.52 <sub>5</sub>	<b>0.43<sub>1</sub></b>	0.47 <sub>4</sub>	0.84 <sub>8</sub>
page-block	0.76 <sub>9</sub>	<b>0.17<sub>1</sub></b>	0.17 <sub>3</sub>	0.19 <sub>4</sub>	0.48 <sub>8</sub>	0.21 <sub>6</sub>	0.17 <sub>2</sub>	0.20 <sub>5</sub>	0.21 <sub>7</sub>
pendigits	1.29 <sub>9</sub>	0.19 <sub>2</sub>	<b>0.19<sub>1</sub></b>	0.46 <sub>3</sub>	0.53 <sub>7</sub>	0.46 <sub>4</sub>	0.47 <sub>5</sub>	0.48 <sub>6</sub>	0.58 <sub>8</sub>
segment	1.39 <sub>9</sub>	<b>0.28<sub>1</sub></b>	0.31 <sub>2</sub>	0.46 <sub>4</sub>	0.53 <sub>7</sub>	0.46 <sub>5</sub>	0.47 <sub>6</sub>	0.46 <sub>3</sub>	0.56 <sub>8</sub>
vehicle	2.34 <sub>9</sub>	<b>0.98<sub>1</sub></b>	0.99 <sub>2</sub>	1.08 <sub>5</sub>	1.16 <sub>7</sub>	1.07 <sub>4</sub>	1.17 <sub>8</sub>	1.04 <sub>3</sub>	1.10 <sub>6</sub>
vowel	0.83 <sub>5</sub>	<b>0.55<sub>1</sub></b>	0.62 <sub>2</sub>	0.80 <sub>3</sub>	0.85 <sub>6</sub>	0.80 <sub>4</sub>	1.04 <sub>8</sub>	1.06 <sub>9</sub>	0.89 <sub>7</sub>
waveform-5	0.80 <sub>9</sub>	<b>0.33<sub>1</sub></b>	0.33 <sub>2</sub>	0.37 <sub>4</sub>	0.43 <sub>7</sub>	0.35 <sub>3</sub>	0.39 <sub>6</sub>	0.38 <sub>5</sub>	0.46 <sub>8</sub>
yeast	5.12 <sub>9</sub>	1.26 <sub>2</sub>	<b>1.22<sub>1</sub></b>	1.33 <sub>6</sub>	2.06 <sub>8</sub>	1.32 <sub>5</sub>	1.29 <sub>3</sub>	1.31 <sub>4</sub>	1.43 <sub>7</sub>
avg rank	8.55	<b>1.40</b>	2.50	4.70	6.15	4.40	5.85	4.90	6.55

Table 4: Ranking of calibration methods applied on the classifier nbayes with the measure=loss(Friedman statistic test = 9.79E+01, p-value = 1.14E-17)



(a) Average over all datasets for Naive Bayes classifier

(b) Average over all classifiers

Figure 5: Critical Difference diagrams for the averaged ranking results of the metric Log-loss.



## E Results

In this Section we present all the final results, including ranking tables for every metric, critical difference diagrams, the best hyperparameters selected for Dirichlet calibration with L2 regularisation, Frequency binning and Width binning; a comparison of how calibrated the 11 classifiers are, and additional results on deep neural networks.

### E.1 Final ranking tables for all metrics

We present here all the final ranking tables for all metrics (Tables 5, 6, 7, 8, 9, 10, 11, and 12). For each ranking, a lower value is indicative of a better metric value (eg. a higher accuracy corresponds to a lower ranking, while a lower log-loss corresponds to a lower ranking as well). Additional details on how to interpret the tables can be found in Section D.2.

Table 5: Rankings for Accuracy

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	7.5	3.1	<b>2.5</b>	5.7	7.5	5.5	4.5	4.2	4.4
forest	5.5	5.4	5.7	<b>3.6</b>	4.8	3.9	4.5	6.9	4.7
knn	5.8	6.4	4.4	4.3	6.4	4.1	3.7	6.4	<b>3.5</b>
lda	6.6	<b>3.5</b>	<b>2.9</b>	5.1	6.1	4.7	4.4	6.5	5.2
logistic	6.5	3.1	<b>2.7</b>	4.4	6.4	5.7	5.1	6.2	4.9
mlp	4.9	4.9	4.6	<b>3.7</b>	5.0	4.6	4.8	7.0	5.5
nbayes	8.2	2.2	<b>1.6</b>	5.5	7.8	3.9	4.3	5.0	6.5
qda	7.3	3.5	<b>2.6</b>	4.6	7.2	4.6	3.4	5.7	6.0
svc-linear	6.1	<b>3.1</b>	3.3	4.6	6.6	4.8	4.2	6.7	5.5
svc-rbf	7.4	4.9	4.5	4.1	6.7	3.7	<b>3.1</b>	4.9	5.7
tree	3.8	6.9	5.5	4.6	<b>3.6</b>	5.5	4.5	5.3	5.3
avg rank	6.31	4.30	<b>3.67</b>	4.56	6.18	4.64	4.24	5.90	5.20

Table 6: Rankings for log-loss

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	8.4	2.1	<b>2.0</b>	4.7	7.8	5.1	4.9	4.8	5.1
forest	7.0	5.7	4.4	<b>3.1</b>	3.9	3.4	5.8	6.8	4.8
knn	8.2	2.8	4.9	6.4	7.2	5.4	<b>2.7</b>	4.4	3.0
lda	7.4	<b>1.9</b>	3.0	3.5	5.6	4.0	7.2	7.1	5.2
logistic	7.6	<b>1.7</b>	2.5	3.2	5.3	4.0	8.1	7.6	5.0
mlp	4.5	<b>2.7</b>	4.5	3.1	3.5	4.0	8.0	8.6	6.2
nbayes	8.6	<b>1.4</b>	2.5	4.7	6.2	4.4	5.8	4.9	6.5
qda	7.5	<b>2.2</b>	3.0	3.8	5.4	4.5	6.7	7.3	4.6
svc-linear	6.7	<b>1.9</b>	2.3	3.7	4.4	3.7	8.1	8.1	6.0
svc-rbf	7.6	4.0	3.4	3.7	5.7	<b>2.9</b>	6.0	5.4	6.3
tree	6.2	3.6	5.7	6.2	7.4	6.7	<b>2.0</b>	4.0	3.2
avg rank	7.24	<b>2.73</b>	3.48	4.21	5.66	4.38	5.93	6.28	5.10

Table 7: Rankings for Brier score

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	8.5	2.8	<b>2.0</b>	4.5	8.3	4.8	4.0	4.7	5.5
forest	7.3	6.1	4.9	2.7	4.4	<b>2.2</b>	3.5	7.8	6.2
knn	6.4	5.5	4.2	3.8	5.9	<b>2.8</b>	3.6	7.3	5.5
lda	6.9	3.0	<b>2.3</b>	3.6	6.8	4.1	4.0	8.0	6.2
logistic	8.3	<b>2.3</b>	2.4	3.9	6.2	4.3	3.6	8.0	6.0
mlp	5.5	3.8	3.9	<b>3.5</b>	5.4	3.8	3.6	8.7	7.0
nbayes	8.5	2.0	<b>1.4</b>	5.0	7.2	3.4	4.0	6.3	7.1
qda	7.5	3.0	<b>2.0</b>	4.2	6.8	4.1	3.4	7.8	6.2
svc-linear	7.7	2.4	<b>2.3</b>	4.3	5.5	4.0	4.0	8.6	6.3
svc-rbf	7.8	4.4	3.8	3.7	6.4	<b>3.1</b>	3.5	5.8	6.5
tree	<b>1.7</b>	7.1	5.1	3.9	2.8	3.6	5.1	8.4	7.1
avg rank	6.90	3.88	<b>3.13</b>	3.90	5.97	3.65	3.85	7.40	6.32

Table 8: Rankings for MCE

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	8.2	3.8	<b>3.7</b>	4.8	6.2	4.3	5.2	4.0	4.7
forest	6.3	5.9	4.8	4.0	4.9	5.4	4.7	4.9	<b>4.0</b>
knn	3.5	5.4	6.4	4.7	6.7	6.5	3.9	4.6	<b>3.4</b>
lda	6.6	5.1	4.2	5.1	4.6	<b>4.1</b>	5.6	5.0	4.7
logistic	5.2	4.3	4.5	5.8	4.7	5.6	6.1	4.7	<b>4.0</b>
mlp	5.6	5.0	6.4	4.0	5.2	4.5	5.6	5.2	<b>3.4</b>
nbayes	7.8	5.5	5.5	4.6	4.5	4.0	5.0	<b>4.0</b>	4.1
qda	6.8	4.9	5.2	5.3	4.0	5.1	4.8	5.6	<b>3.4</b>
svc-linear	5.3	5.0	5.0	5.0	4.8	4.9	5.9	4.9	<b>4.2</b>
svc-rbf	4.4	5.9	5.7	5.4	4.4	6.0	<b>3.9</b>	4.1	5.1
tree	5.4	6.3	5.0	5.1	5.8	5.6	4.0	<b>3.8</b>	
avg rank	5.91	5.18	5.14	4.90	5.09	5.11	4.97	4.64	<b>4.07</b>

Table 9: Rankings for confidence-ECE

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	8.2	3.1	<b>2.5</b>	4.0	7.3	2.9	4.2	6.4	6.3
forest	7.3	5.9	3.0	4.0	4.4	<b>2.5</b>	2.8	8.3	6.7
knn	4.0	4.7	<b>3.2</b>	4.5	4.1	3.4	5.5	8.1	7.5
lda	6.7	4.4	<b>3.0</b>	3.4	6.0	4.0	3.3	8.0	6.2
logistic	8.0	3.6	3.4	4.0	5.0	3.5	<b>3.2</b>	8.0	6.3
mlp	5.7	4.2	3.2	3.8	5.6	3.8	<b>2.3</b>	8.7	7.7
nbayes	8.0	3.8	<b>2.3</b>	4.3	5.6	3.5	3.6	7.5	6.3
qda	6.5	4.3	2.9	4.0	6.2	4.0	<b>2.2</b>	8.4	6.5
svc-linear	7.0	4.0	3.2	4.2	4.6	3.2	3.2	8.5	6.9
svc-rbf	5.7	4.5	4.4	4.6	4.0	4.0	<b>3.2</b>	8.6	6.0
tree	<b>2.0</b>	6.2	4.1	3.0	3.7	3.7	5.9	8.8	7.6
avg rank	6.27	4.44	<b>3.20</b>	4.00	5.12	3.50	3.59	8.12	6.75

Table 10: Rankings for classwise-ECE

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	6.7	<b>3.3</b>	<b>2.9</b>	5.4	6.3	3.4	4.0	6.4	6.5
forest	7.6	5.6	3.4	3.6	5.0	<b>2.1</b>	2.7	8.3	6.7
knn	4.1	4.6	3.1	4.2	5.5	<b>2.5</b>	5.4	7.9	7.6
lda	6.3	4.1	3.4	3.4	6.8	3.8	<b>3.0</b>	7.9	6.3
logistic	8.0	3.1	3.6	3.8	6.1	3.4	3.1	7.8	5.9
mlp	5.5	4.1	3.2	3.2	6.3	3.8	<b>2.7</b>	8.6	7.6
nbayes	7.4	2.8	<b>1.9</b>	5.0	7.8	3.4	3.8	7.0	6.0
qda	6.5	3.6	2.6	4.0	7.0	4.0	<b>2.5</b>	8.3	6.3
svc-linear	7.1	3.7	3.8	4.0	4.8	3.2	3.2	8.5	6.7
svc-rbf	6.0	4.0	4.0	4.2	5.0	3.7	<b>3.3</b>	8.6	6.2
tree	<b>2.2</b>	6.3	4.0	3.7	4.2	3.4	5.0	8.5	7.7
avg rank	6.12	4.12	<b>3.28</b>	4.05	5.89	3.34	3.54	7.99	6.67

### E.2 Final critical difference diagrams for every metric

In order to perform a final comparison between calibration methods, we considered every combination of dataset and classifier as a group  $n = \#datasets \times \#classifiers$ , and ranked the results of the  $k$  calibration methods. With this setting, we have performed the Friedman statistical test followed by the one-tailed Bonferroni-Dunn test to obtain critical differences (CDs) for every metric (See Figure

Table 11: Rankings for p-confidence-ECE

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	7.9	2.8	<b>2.5</b>	4.3	6.9	3.8	4.5	6.3	6.1
forest	7.0	6.0	<b>2.9</b>	3.6	3.6	3.2	3.2	8.4	7.1
knn	6.2	3.0	<b>2.9</b>	4.2	4.8	3.9	4.6	7.9	7.5
lda	7.3	3.7	<b>3.1</b>	3.9	5.3	4.1	3.5	8.0	6.2
logistic	7.8	3.5	4.0	3.7	4.7	3.6	<b>3.4</b>	8.0	6.4
mlp	5.1	4.3	3.7	3.6	4.7	4.1	<b>3.2</b>	8.7	7.6
nbayes	8.1	3.2	<b>2.5</b>	4.4	5.3	3.6	3.3	7.6	6.9
qda	6.8	3.9	3.4	3.6	5.5	4.3	<b>2.6</b>	8.3	6.6
svc-linear	6.8	4.0	3.5	3.8	4.3	<b>3.3</b>	4.1	8.5	6.8
svc-rbf	5.6	4.5	4.2	4.6	4.2	4.0	<b>3.2</b>	8.4	6.2
tree	<b>2.9</b>	5.3	4.0	3.3	4.0	4.5	5.2	8.2	7.5
avg rank	6.50	4.01	<b>3.33</b>	3.90	4.84	3.87	3.71	8.03	6.81

Table 12: Rankings for p-classwise-ECE

	Uncal	DirL2	DirODIR	Beta	TempS	VecS	Isot	FreqB	WidthB
adas	7.4	3.1	<b>3.0</b>	5.1	6.8	3.7	4.1	6.0	5.7
forest	6.1	4.9	4.7	<b>3.2</b>	4.4	4.3	3.6	8.0	5.8
knn	7.7	<b>2.0</b>	4.2	5.0	6.1	4.6	3.0	7.0	5.5
lda	7.3	<b>2.9</b>	4.1	3.5	6.0	3.8	3.7	7.9	5.8
logistic	6.9	<b>3.0</b>	4.3	3.6	5.1	4.3	4.1	8.0	5.6
mlp	5.0	<b>2.5</b>	5.2	2.9	4.3	4.5	5.0	8.6	6.9
nbayes	8.2	<b>1.6</b>	2.3	4.9	7.0	3.2	4.5	6.8	6.5
qda	7.5	<b>2.2</b>	3.6	3.8	6.0	4.6	3.7	8.2	5.5
svc-linear	6.2	<b>3.0</b>	4.4	4.0	4.5	4.0	4.8	8.5	5.5
svc-rbf	5.5	3.8	4.1	4.4	5.0	4.8	<b>3.4</b>	8.3	5.8
tree	3.7	4.6	<b>3.6</b>	5.0	4.2	4.4	4.2	8.0	7.3
avg rank	6.50	<b>3.07</b>	3.97	4.12	5.40	4.21	4.01	7.75	5.98

6). The results showed Dirichlet L2 as the best calibration method for the measures accuracy, log-loss and p-cw-ece with statistical significance (See Figures 6a 6c, and 6h), and in the group of the best calibration methods in the rest of the metrics with statistical significance, but no difference within the group. It is worth mentioning that Figure 6c showed statistical difference between Dirichlet L2, OvR Beta, OvR width binning, and the rest of the calibrators in one group; in the mentioned order.

### E.3 Best calibrator hyperparameters

Figure 8 shows the best hyperparameters for every inner 3-fold-cross-validation. Dirichlet L2 (Figure 8a) shows a preference for regularisation hyperparameter  $\lambda = 1e^{-3}$  and lower values. Our current minimum regularisation value of  $1e^{-7}$  is also being selected multiple times, indicating that lower values may be optimal in several occasions. However, this fact did not seem to hurt the overall good results in our experiments. One-vs.-Rest frequency binning tends to prefer 10 bins of equal number of samples, while One-vs.-Rest width binning prefers 5 equal sized bins (See Figures 8b and 8c respectively).

### E.4 Comparison of classifiers

In this Section we compare all the classifiers without post-hoc calibration on 17 of the datasets; from the total of 21 datasets *shuttle*, *yeast*, *mfeat-karhunen* and *libras-movement* were removed from this analysis as at least one classifier was not able to complete the experiment.

Figure 9 shows the Critical Difference diagram for all the 8 metrics. In particular, the MLP and the SVC with linear kernel are always in the group with the higher rankings and never in the lowest. Similarly, random forest is consistently in the best group, but in the worst group as well in 4 of the measures. SVC with radial basis kernel is in the best group 6 times, but 3 times in the worst. On the other hand, naive Bayes and Adaboost SAMME are consistently in the worst group and never in the best one. The rest of the classifiers did not show a clear ranking position.

Figures 10b and 10a show the proportion of times each classifier passed the p-conf-ECE and p-cw-ECE statistical test for all datasets and cross-validation folds.

### E.5 Deep neural networks

In this section, we provide further discussion about results from the deep networks experiments. These are given in the form of critical difference diagrams (Figure 11) and tables (Tables 13-20) both including the following measures: error rate, log-loss, Brier score, maximum calibration error (MCE), confidence-ECE (conf-ECE), classwise-ECE (cw-ECE), as well as significance measures p-conf-ECE and p-cw-ECE.

In addition, Table 21 compares MS-ODIR and vector scaling on log-loss. On the table, we also added MS-ODIR-zero which was obtained from the respective MS-ODIR model by replacing the off-diagonal entries with zeroes. Each experiment is replicated three times with different splits on datasets. This is done to compare the stability of the methods. In each replication, the best scoring model is written in bold.

Finally, Figure 12 shows that temperature scaling systematically under-estimates class 4 probabilities on the model `c10_resnet_wide32` on CIFAR-10.

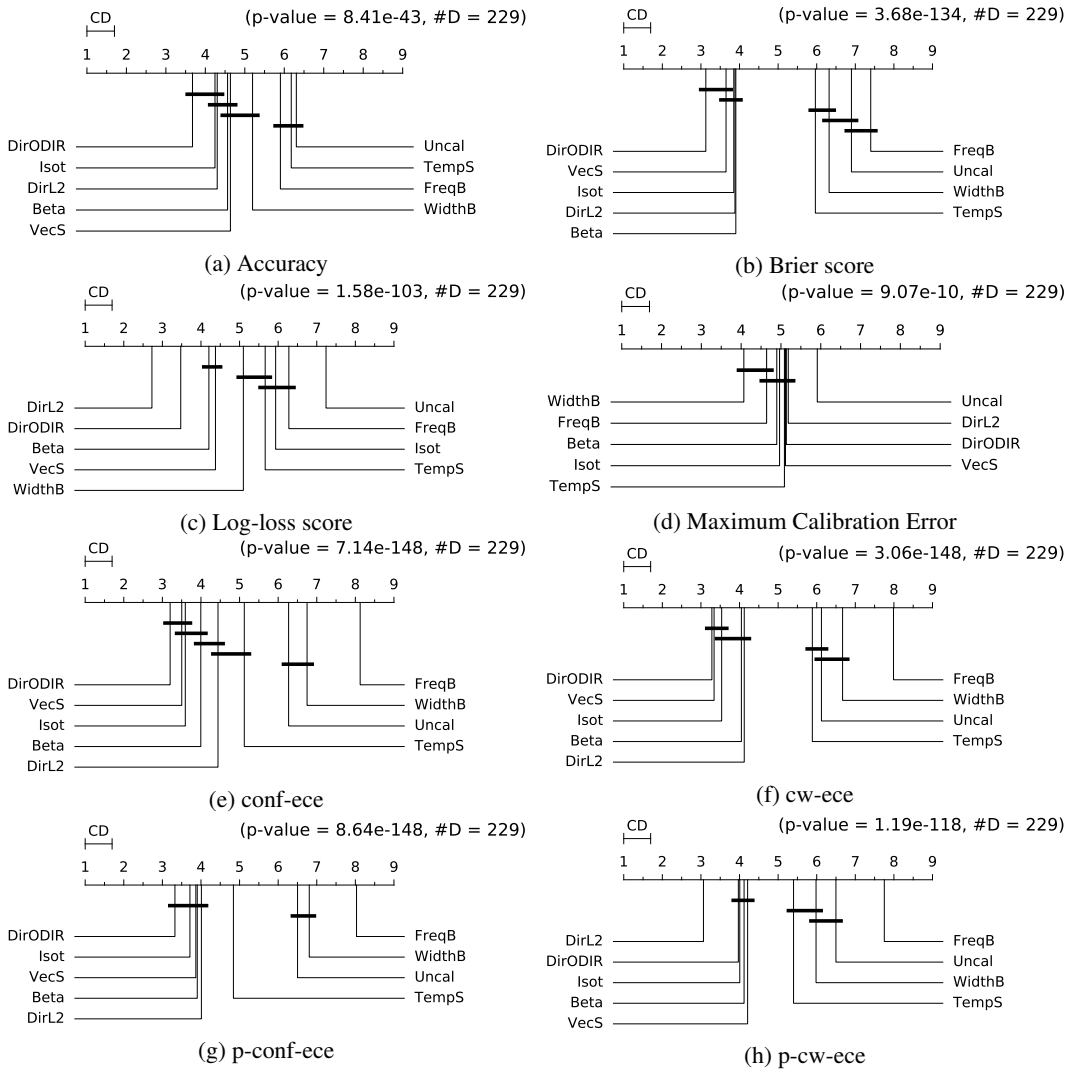


Figure 6: Critical difference of the average of multiclass classifiers.

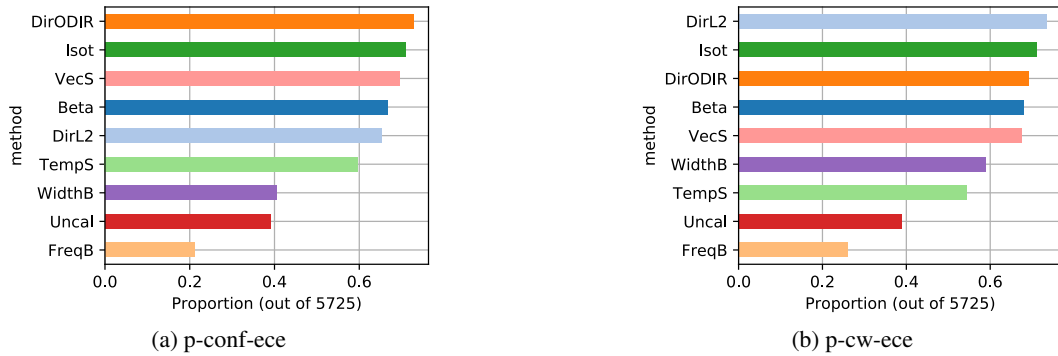


Figure 7: Proportion of times each calibrator passes a calibration p-test with a p-value higher than 0.05.

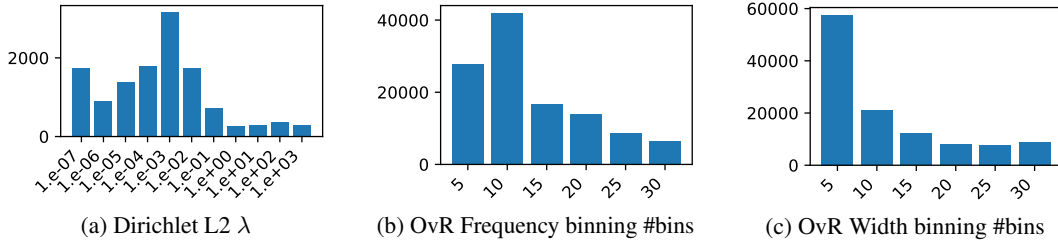


Figure 8: Histogram of the selected hyperparameters during the inner 3-fold-cross-validation

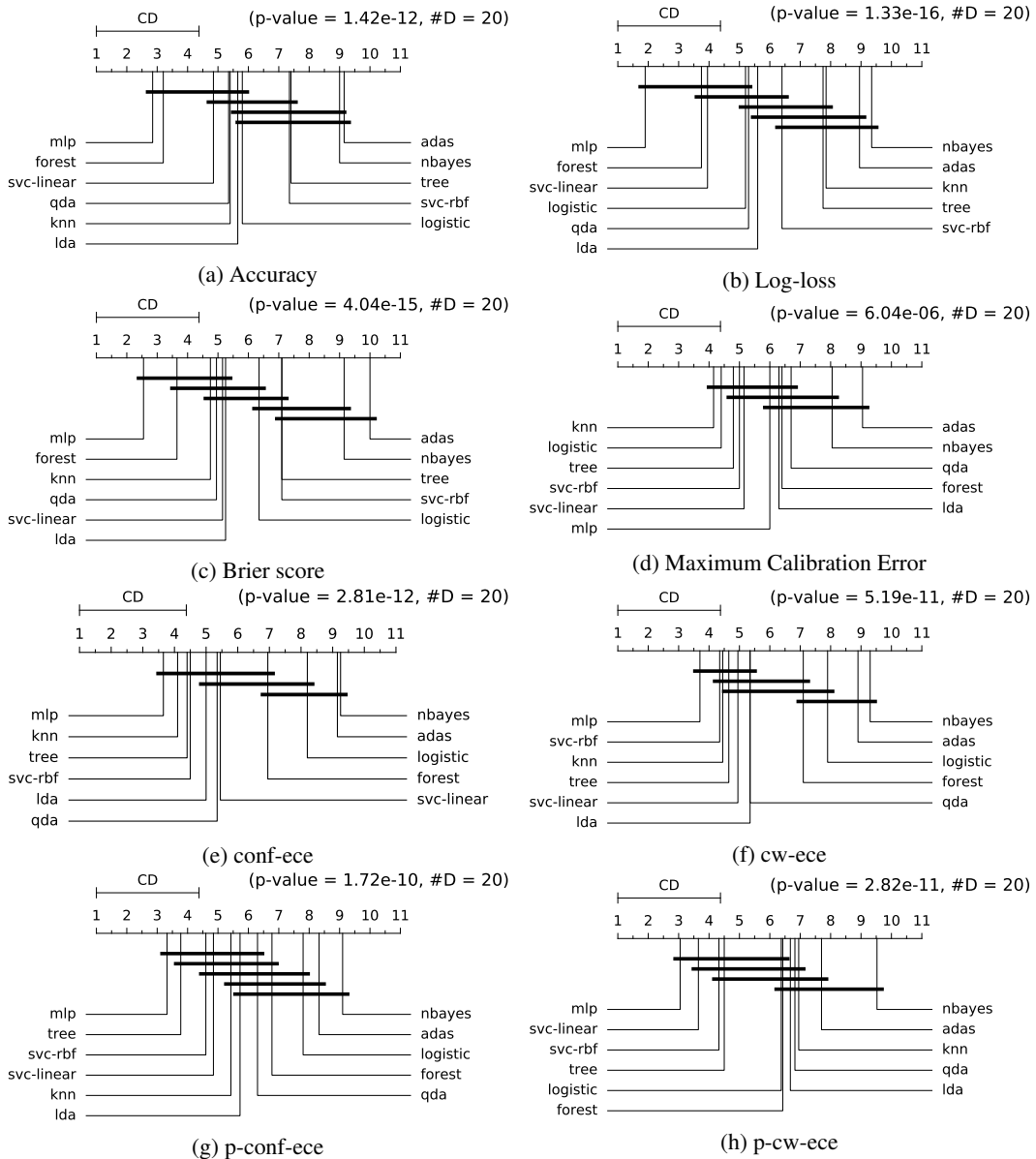


Figure 9: Critical difference of uncalibrated classifiers.

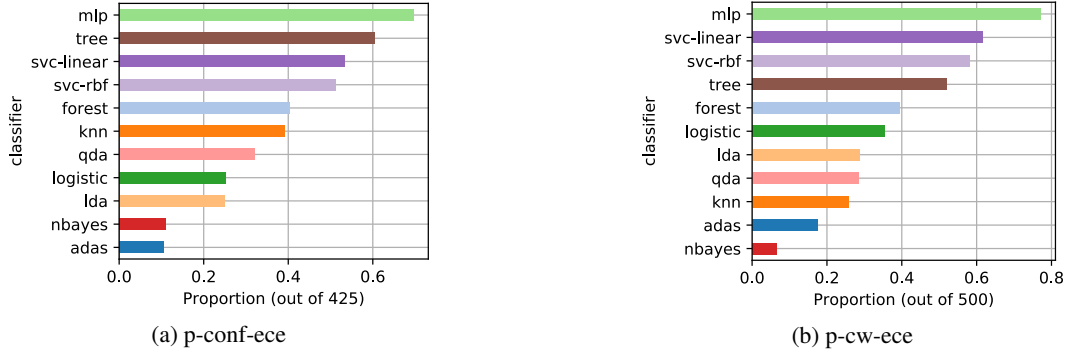


Figure 10: Proportion of times each classifier is already calibrated with different p-tests.

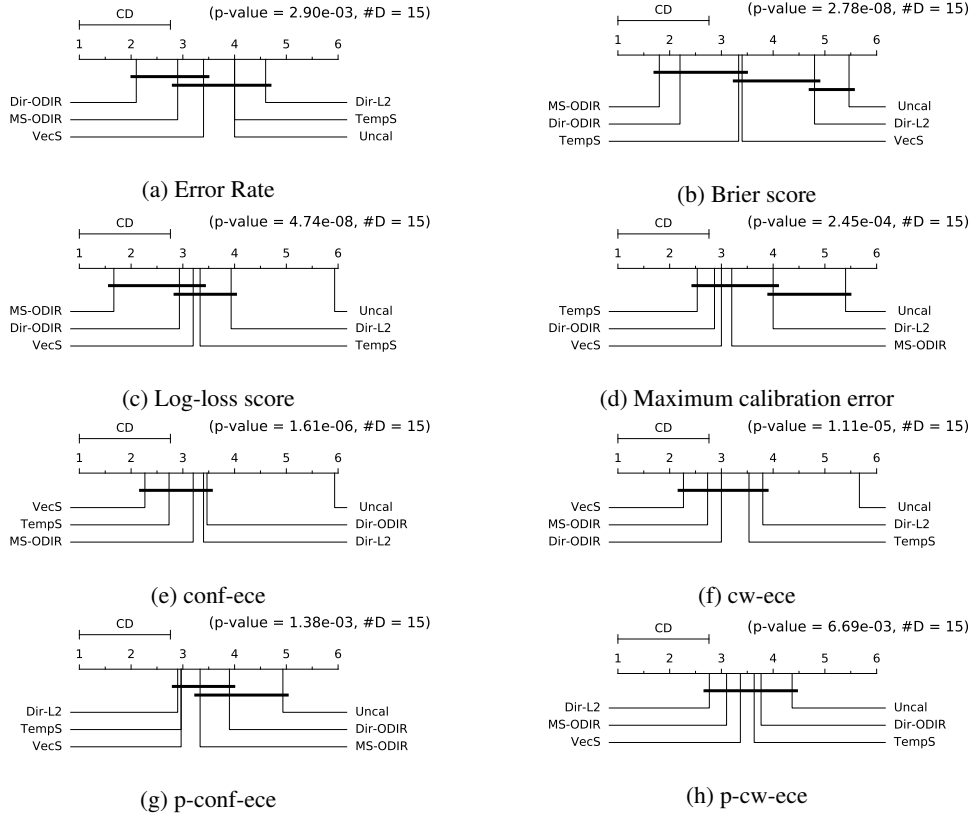


Figure 11: Critical difference of the deep neural networks.

Table 13: Scores and ranking of calibration methods for **log-loss**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.39098 <sub>6</sub>	<b>0.19497</b> <sub>1</sub>	0.19692 <sub>4</sub>	0.19536 <sub>2</sub>	0.19743 <sub>5</sub>	0.19634 <sub>3</sub>
c10_densenet40	0.42821 <sub>6</sub>	0.22509 <sub>5</sub>	<b>0.22048</b> <sub>1</sub>	0.22371 <sub>4</sub>	0.22270 <sub>3</sub>	0.22240 <sub>2</sub>
c10_lenet5	0.82326 <sub>6</sub>	0.80031 <sub>5</sub>	0.74418 <sub>2</sub>	0.74441 <sub>3</sub>	0.74704 <sub>4</sub>	<b>0.74262</b> <sub>1</sub>
c10_resnet110	0.35827 <sub>6</sub>	0.20926 <sub>5</sub>	<b>0.20303</b> <sub>1</sub>	0.20511 <sub>3</sub>	0.20595 <sub>4</sub>	0.20375 <sub>2</sub>
c10_resnet110_SD	0.30325 <sub>6</sub>	0.17760 <sub>5</sub>	0.17694 <sub>4</sub>	0.17608 <sub>3</sub>	0.17549 <sub>2</sub>	<b>0.17537</b> <sub>1</sub>
c10_resnet_wide32	0.38170 <sub>6</sub>	0.19148 <sub>5</sub>	0.18464 <sub>4</sub>	0.18203 <sub>2</sub>	0.18276 <sub>3</sub>	<b>0.18165</b> <sub>1</sub>
c100_convnet	1.64120 <sub>6</sub>	<b>0.94162</b> <sub>1</sub>	1.18945 <sub>5</sub>	0.96121 <sub>2</sub>	0.96369 <sub>4</sub>	0.96141 <sub>3</sub>
c100_densenet40	2.01740 <sub>6</sub>	1.05713 <sub>2</sub>	1.25293 <sub>5</sub>	1.05909 <sub>4</sub>	1.05831 <sub>3</sub>	<b>1.05084</b> <sub>1</sub>
c100_lenet5	2.78365 <sub>6</sub>	2.64979 <sub>5</sub>	2.59482 <sub>4</sub>	2.48951 <sub>2</sub>	2.51590 <sub>3</sub>	<b>2.48670</b> <sub>1</sub>
c100_resnet110	1.69371 <sub>6</sub>	1.09169 <sub>3</sub>	1.21239 <sub>5</sub>	1.09607 <sub>4</sub>	1.08916 <sub>2</sub>	<b>1.07370</b> <sub>1</sub>
c100_resnet110_SD	1.35250 <sub>6</sub>	0.94214 <sub>3</sub>	1.19837 <sub>5</sub>	0.94477 <sub>4</sub>	<b>0.92341</b> <sub>1</sub>	0.92731 <sub>2</sub>
c100_resnet_wide32	1.80215 <sub>6</sub>	0.94453 <sub>3</sub>	1.08711 <sub>5</sub>	0.95288 <sub>4</sub>	0.93650 <sub>2</sub>	<b>0.93273</b> <sub>1</sub>
SVHN_convnet	0.20460 <sub>6</sub>	0.15142 <sub>5</sub>	0.14246 <sub>3</sub>	0.13791 <sub>2</sub>	0.14388 <sub>4</sub>	<b>0.13760</b> <sub>1</sub>
SVHN_resnet152_SD	0.08542 <sub>6</sub>	<b>0.07861</b> <sub>1</sub>	0.08463 <sub>5</sub>	0.08038 <sub>2</sub>	0.08124 <sub>4</sub>	0.08100 <sub>3</sub>
avg rank	6.0	3.5	3.79	2.93	3.14	1.64

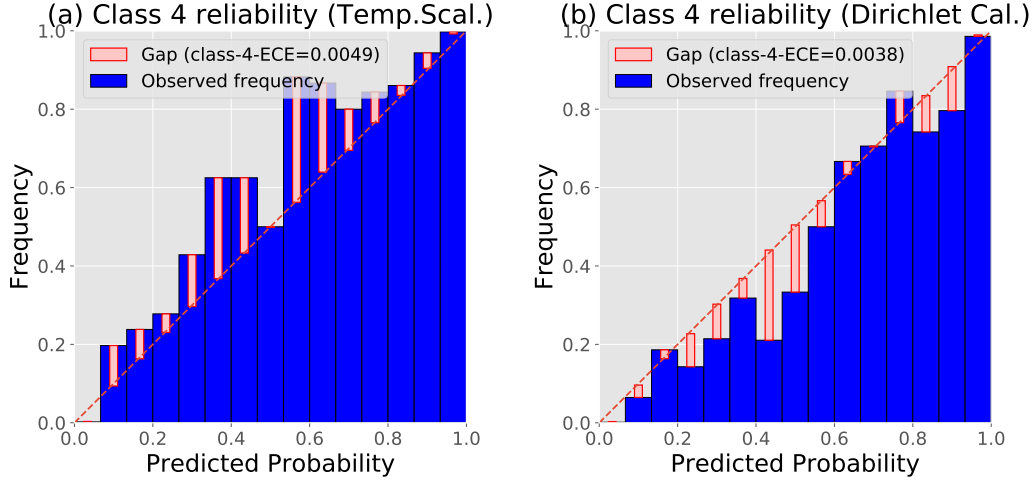


Figure 12: Reliability diagrams of c10\_resnet\_wide32 on CIFAR-10: (a) classwise-reliability for class 4 after temperature scaling; (b) classwise-reliability for class 4 after Dirichlet calibration.

Table 14: Scores and ranking of calibration methods for **Brier score**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.01090 <sub>6</sub>	<b>0.00952</b> <sub>1</sub>	0.00969 <sub>5</sub>	0.00955 <sub>3</sub>	0.00958 <sub>4</sub>	0.00953 <sub>2</sub>
c10_densenet40	0.01274 <sub>6</sub>	0.01100 <sub>4</sub>	0.01102 <sub>5</sub>	0.01097 <sub>2</sub>	0.01097 <sub>3</sub>	<b>0.01097</b> <sub>1</sub>
c10_lenet5	0.03788 <sub>6</sub>	0.03748 <sub>5</sub>	0.03510 <sub>2</sub>	0.03511 <sub>3</sub>	0.03523 <sub>4</sub>	<b>0.03502</b> <sub>1</sub>
c10_resnet110	0.01102 <sub>6</sub>	0.00979 <sub>4</sub>	0.00979 <sub>5</sub>	0.00977 <sub>2</sub>	0.00978 <sub>3</sub>	<b>0.00976</b> <sub>1</sub>
c10_resnet110_SD	0.00981 <sub>6</sub>	0.00874 <sub>4</sub>	0.00877 <sub>5</sub>	0.00867 <sub>3</sub>	0.00867 <sub>2</sub>	<b>0.00866</b> <sub>1</sub>
c10_resnet_wide32	0.01047 <sub>6</sub>	0.00924 <sub>5</sub>	0.00909 <sub>4</sub>	<b>0.00888</b> <sub>1</sub>	0.00891 <sub>3</sub>	0.00889 <sub>2</sub>
c100_convnet	0.00425 <sub>5</sub>	<b>0.00358</b> <sub>1</sub>	0.00441 <sub>6</sub>	0.00358 <sub>2</sub>	0.00362 <sub>4</sub>	0.00361 <sub>3</sub>
c100_densenet40	0.00491 <sub>6</sub>	0.00401 <sub>3</sub>	0.00468 <sub>5</sub>	0.00400 <sub>2</sub>	0.00403 <sub>4</sub>	<b>0.00400</b> <sub>1</sub>
c100_lenet5	0.00813 <sub>6</sub>	0.00792 <sub>5</sub>	0.00786 <sub>4</sub>	0.00760 <sub>2</sub>	0.00767 <sub>3</sub>	<b>0.00760</b> <sub>1</sub>
c100_resnet110	0.00453 <sub>6</sub>	0.00392 <sub>3</sub>	0.00438 <sub>5</sub>	0.00391 <sub>2</sub>	0.00393 <sub>4</sub>	<b>0.00391</b> <sub>1</sub>
c100_resnet110_SD	0.00418 <sub>5</sub>	0.00367 <sub>4</sub>	0.00456 <sub>6</sub>	0.00364 <sub>3</sub>	<b>0.00360</b> <sub>1</sub>	0.00361 <sub>2</sub>
c100_resnet_wide32	0.00432 <sub>6</sub>	0.00355 <sub>4</sub>	0.00401 <sub>5</sub>	0.00354 <sub>3</sub>	0.00352 <sub>2</sub>	<b>0.00351</b> <sub>1</sub>
SVHN_convnet	0.00776 <sub>6</sub>	0.00598 <sub>5</sub>	0.00555 <sub>3</sub>	<b>0.00530</b> <sub>1</sub>	0.00561 <sub>4</sub>	0.00532 <sub>2</sub>
SVHN_resnet152_SD	0.00297 <sub>3</sub>	<b>0.00291</b> <sub>1</sub>	0.00305 <sub>6</sub>	0.00293 <sub>2</sub>	0.00299 <sub>5</sub>	0.00298 <sub>4</sub>
avg rank	5.64	3.5	4.71	2.21	3.29	1.64

Table 15: Scores and ranking of calibration methods for **confidence-ECE**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.04760 <sub>6</sub>	0.01065 <sub>5</sub>	0.00769 <sub>2</sub>	0.00960 <sub>4</sub>	<b>0.00740</b> <sub>1</sub>	0.00782 <sub>3</sub>
c10_densenet40	0.05500 <sub>6</sub>	0.00946 <sub>2</sub>	<b>0.00568</b> <sub>1</sub>	0.01097 <sub>5</sub>	0.01018 <sub>4</sub>	0.00988 <sub>3</sub>
c10_lenet5	0.05180 <sub>6</sub>	0.01665 <sub>5</sub>	0.01383 <sub>3</sub>	0.01367 <sub>2</sub>	<b>0.01310</b> <sub>1</sub>	0.01468 <sub>4</sub>
c10_resnet110	0.04750 <sub>6</sub>	0.01132 <sub>5</sub>	<b>0.00680</b> <sub>1</sub>	0.01086 <sub>3</sub>	0.01130 <sub>4</sub>	0.01059 <sub>2</sub>
c10_resnet110_SD	0.04113 <sub>6</sub>	<b>0.00555</b> <sub>1</sub>	0.00646 <sub>4</sub>	0.00815 <sub>5</sub>	0.00579 <sub>3</sub>	0.00566 <sub>2</sub>
c10_resnet_wide32	0.04505 <sub>6</sub>	0.00784 <sub>4</sub>	<b>0.00524</b> <sub>1</sub>	0.00837 <sub>5</sub>	0.00769 <sub>3</sub>	0.00727 <sub>2</sub>
c100_convnet	0.17614 <sub>6</sub>	<b>0.01367</b> <sub>1</sub>	0.14347 <sub>5</sub>	0.02069 <sub>3</sub>	0.01965 <sub>2</sub>	0.02660 <sub>4</sub>
c100_densenet40	0.21156 <sub>6</sub>	<b>0.00902</b> <sub>1</sub>	0.12380 <sub>5</sub>	0.01138 <sub>2</sub>	0.01224 <sub>3</sub>	0.02197 <sub>4</sub>
c100_lenet5	0.12125 <sub>6</sub>	0.01499 <sub>4</sub>	0.01369 <sub>2</sub>	0.02003 <sub>5</sub>	<b>0.01294</b> <sub>1</sub>	0.01407 <sub>3</sub>
c100_resnet110	0.18480 <sub>6</sub>	<b>0.02380</b> <sub>1</sub>	0.14535 <sub>5</sub>	0.02822 <sub>4</sub>	0.02693 <sub>2</sub>	0.02735 <sub>3</sub>
c100_resnet110_SD	0.15861 <sub>5</sub>	<b>0.01214</b> <sub>1</sub>	0.15920 <sub>6</sub>	0.02283 <sub>4</sub>	0.01296 <sub>2</sub>	0.02246 <sub>3</sub>
c100_resnet_wide32	0.18784 <sub>6</sub>	<b>0.01472</b> <sub>1</sub>	0.13509 <sub>5</sub>	0.01891 <sub>3</sub>	0.01718 <sub>2</sub>	0.02581 <sub>4</sub>
SVHN_convnet	0.07755 <sub>6</sub>	0.01179 <sub>4</sub>	0.01910 <sub>5</sub>	0.00997 <sub>2</sub>	<b>0.00934</b> <sub>1</sub>	0.01037 <sub>3</sub>
SVHN_resnet152_SD	0.00862 <sub>6</sub>	0.00607 <sub>4</sub>	0.00691 <sub>5</sub>	<b>0.00582</b> <sub>1</sub>	0.00595 <sub>2</sub>	0.00604 <sub>3</sub>
avg rank	5.93	2.79	3.57	3.43	2.21	3.07

Table 16: Scores and ranking of calibration methods for **classwise-ECE**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.10375 <sub>6</sub>	0.04423 <sub>4</sub>	0.04262 <sub>2</sub>	0.04507 <sub>5</sub>	<b>0.04259</b> <sub>1</sub>	0.04352 <sub>3</sub>
c10_densenet40	0.11430 <sub>6</sub>	0.03977 <sub>5</sub>	<b>0.03412</b> <sub>1</sub>	0.03687 <sub>4</sub>	0.03609 <sub>2</sub>	0.03678 <sub>3</sub>
c10_lenet5	0.19849 <sub>6</sub>	0.17141 <sub>5</sub>	<b>0.05185</b> <sub>1</sub>	0.05891 <sub>4</sub>	0.05705 <sub>2</sub>	0.05862 <sub>3</sub>
c10_resnet110	0.09846 <sub>6</sub>	0.04344 <sub>5</sub>	<b>0.03206</b> <sub>1</sub>	0.03950 <sub>4</sub>	0.03653 <sub>3</sub>	0.03615 <sub>2</sub>
c10_resnet110_SD	0.08647 <sub>6</sub>	0.03071 <sub>4</sub>	0.03148 <sub>5</sub>	0.02937 <sub>3</sub>	0.02713 <sub>2</sub>	<b>0.02681</b> <sub>1</sub>
c10_resnet_wide32	0.09530 <sub>6</sub>	0.04775 <sub>5</sub>	0.03153 <sub>3</sub>	0.02947 <sub>2</sub>	0.03164 <sub>4</sub>	<b>0.02921</b> <sub>1</sub>
c100_convnet	0.42414 <sub>6</sub>	<b>0.22683</b> <sub>1</sub>	0.40185 <sub>5</sub>	0.24041 <sub>3</sub>	0.24063 <sub>4</sub>	0.23958 <sub>2</sub>
c100_densenet40	0.47026 <sub>6</sub>	0.18664 <sub>2</sub>	0.32985 <sub>5</sub>	<b>0.18630</b> <sub>1</sub>	0.18879 <sub>3</sub>	0.19112 <sub>4</sub>
c100_lenet5	0.47264 <sub>6</sub>	0.38481 <sub>5</sub>	0.21865 <sub>4</sub>	0.21348 <sub>2</sub>	<b>0.20293</b> <sub>1</sub>	0.21379 <sub>3</sub>
c100_resnet110	0.41644 <sub>6</sub>	0.20095 <sub>3</sub>	0.35885 <sub>5</sub>	<b>0.18639</b> <sub>1</sub>	0.19442 <sub>2</sub>	0.20270 <sub>4</sub>
c100_resnet110_SD	0.37518 <sub>6</sub>	0.20310 <sub>4</sub>	0.37346 <sub>5</sub>	0.18895 <sub>3</sub>	<b>0.17015</b> <sub>1</sub>	0.18552 <sub>2</sub>
c100_resnet_wide32	0.42027 <sub>6</sub>	0.18573 <sub>4</sub>	0.33258 <sub>5</sub>	0.17951 <sub>2</sub>	<b>0.17082</b> <sub>1</sub>	0.17966 <sub>3</sub>
SVHN_convnet	0.15935 <sub>6</sub>	0.03830 <sub>4</sub>	0.04276 <sub>5</sub>	0.02638 <sub>2</sub>	<b>0.02480</b> <sub>1</sub>	0.02750 <sub>3</sub>
SVHN_resnet152_SD	0.01940 <sub>2</sub>	<b>0.01849</b> <sub>1</sub>	0.02184 <sub>6</sub>	0.01988 <sub>3</sub>	0.02120 <sub>5</sub>	0.02088 <sub>4</sub>
avg rank	5.71	3.71	3.79	2.79	2.29	2.71

Table 17: Scores and ranking of calibration methods for **MCE**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.59173 <sub>6</sub>	0.23150 <sub>4</sub>	0.12432 <sub>2</sub>	0.24830 <sub>5</sub>	0.12831 <sub>3</sub>	<b>0.07621</b> <sub>1</sub>
c10_densenet40	0.33396 <sub>6</sub>	0.09929 <sub>2</sub>	0.11679 <sub>4</sub>	<b>0.07858</b> <sub>1</sub>	0.12046 <sub>5</sub>	0.11297 <sub>3</sub>
c10_lenet5	0.11281 <sub>6</sub>	0.09158 <sub>3</sub>	<b>0.05112</b> <sub>1</sub>	0.09009 <sub>2</sub>	0.09996 <sub>4</sub>	0.10061 <sub>5</sub>
c10_resnet110	0.29580 <sub>6</sub>	0.23639 <sub>4</sub>	0.24405 <sub>5</sub>	<b>0.08331</b> <sub>1</sub>	0.13130 <sub>2</sub>	0.22678 <sub>3</sub>
c10_resnet110_SD	0.32484 <sub>6</sub>	<b>0.07823</b> <sub>1</sub>	0.23064 <sub>5</sub>	0.13309 <sub>3</sub>	0.14276 <sub>4</sub>	0.08422 <sub>2</sub>
c10_resnet_wide32	0.37215 <sub>4</sub>	<b>0.07060</b> <sub>1</sub>	0.49283 <sub>6</sub>	0.41567 <sub>5</sub>	0.26539 <sub>3</sub>	0.26372 <sub>2</sub>
c100_convnet	0.36391 <sub>6</sub>	0.13689 <sub>4</sub>	0.23333 <sub>5</sub>	0.07235 <sub>2</sub>	<b>0.07043</b> <sub>1</sub>	0.08171 <sub>3</sub>
c100_densenet40	0.45400 <sub>6</sub>	<b>0.02213</b> <sub>1</sub>	0.19748 <sub>5</sub>	0.04074 <sub>2</sub>	0.04293 <sub>3</sub>	0.05004 <sub>4</sub>
c100_lenet5	0.20097 <sub>6</sub>	0.05836 <sub>3</sub>	<b>0.05678</b> <sub>1</sub>	0.06774 <sub>4</sub>	0.05749 <sub>2</sub>	0.08939 <sub>5</sub>
c100_resnet110	0.39882 <sub>6</sub>	0.07099 <sub>2</sub>	0.20732 <sub>5</sub>	0.08026 <sub>4</sub>	0.07354 <sub>3</sub>	<b>0.06678</b> <sub>1</sub>
c100_resnet110_SD	0.48291 <sub>6</sub>	0.04099 <sub>2</sub>	0.24578 <sub>5</sub>	0.05979 <sub>3</sub>	<b>0.04038</b> <sub>1</sub>	0.06612 <sub>4</sub>
c100_resnet_wide32	0.45639 <sub>6</sub>	<b>0.03606</b> <sub>1</sub>	0.19370 <sub>5</sub>	0.05521 <sub>2</sub>	0.06605 <sub>4</sub>	0.06468 <sub>3</sub>
SVHN_convnet	0.30011 <sub>5</sub>	0.40691 <sub>6</sub>	<b>0.16154</b> <sub>1</sub>	0.18458 <sub>3</sub>	0.16312 <sub>2</sub>	0.18588 <sub>4</sub>
SVHN_resnet152_SD	0.25032 <sub>5</sub>	<b>0.18244</b> <sub>1</sub>	0.23895 <sub>4</sub>	0.19649 <sub>2</sub>	0.23092 <sub>3</sub>	0.80082 <sub>6</sub>
avg rank	5.71	2.5	3.86	2.79	2.86	3.29

Table 18: Scores and ranking of calibration methods for **error rate (%)**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	6.18000 <sub>2</sub>	6.18000 <sub>2</sub>	6.38000 <sub>6</sub>	<b>6.12000</b> <sub>1</sub>	6.36000 <sub>5</sub>	6.32000 <sub>4</sub>
c10_densenet40	7.58000 <sub>5</sub>	7.58000 <sub>5</sub>	<b>7.49000</b> <sub>1</sub>	7.53000 <sub>4</sub>	7.52000 <sub>3</sub>	7.50000 <sub>2</sub>
c10_lenet5	27.26000 <sub>5</sub>	27.26000 <sub>5</sub>	<b>25.25000</b> <sub>1</sub>	25.44000 <sub>2</sub>	25.49000 <sub>3</sub>	25.50000 <sub>4</sub>
c10_resnet110	6.44000 <sub>1</sub>	6.44000 <sub>1</sub>	6.54000 <sub>6</sub>	6.49000 <sub>4</sub>	6.47000 <sub>3</sub>	6.49000 <sub>4</sub>
c10_resnet110_SD	5.96000 <sub>5</sub>	5.96000 <sub>5</sub>	5.90000 <sub>4</sub>	<b>5.77000</b> <sub>1</sub>	5.83000 <sub>3</sub>	5.81000 <sub>2</sub>
c10_resnet_wide32	6.07000 <sub>5</sub>	6.07000 <sub>5</sub>	5.94000 <sub>4</sub>	5.76000 <sub>2</sub>	<b>5.74000</b> <sub>1</sub>	5.81000 <sub>3</sub>
c100_convnet	26.12000 <sub>1</sub>	26.12000 <sub>1</sub>	30.96000 <sub>6</sub>	26.22000 <sub>3</sub>	26.56000 <sub>4</sub>	26.60000 <sub>5</sub>
c100_densenet40	30.00000 <sub>3</sub>	30.00000 <sub>3</sub>	33.48000 <sub>6</sub>	29.87000 <sub>2</sub>	30.16000 <sub>5</sub>	<b>29.61000</b> <sub>1</sub>
c100_lenet5	66.41000 <sub>5</sub>	66.41000 <sub>5</sub>	65.97000 <sub>4</sub>	62.53000 <sub>2</sub>	63.59000 <sub>3</sub>	<b>62.44000</b> <sub>1</sub>
c100_resnet110	28.52000 <sub>4</sub>	28.52000 <sub>4</sub>	30.04000 <sub>6</sub>	<b>28.36000</b> <sub>1</sub>	28.40000 <sub>2</sub>	28.45000 <sub>3</sub>
c100_resnet110_SD	27.17000 <sub>4</sub>	27.17000 <sub>4</sub>	31.43000 <sub>6</sub>	26.96000 <sub>3</sub>	26.50000 <sub>2</sub>	<b>26.42000</b> <sub>1</sub>
c100_resnet_wide32	26.18000 <sub>4</sub>	26.18000 <sub>4</sub>	27.69000 <sub>6</sub>	26.07000 <sub>2</sub>	26.08000 <sub>3</sub>	<b>26.06000</b> <sub>1</sub>
SVHN_convnet	3.82750 <sub>5</sub>	3.82750 <sub>5</sub>	3.42811 <sub>3</sub>	<b>3.34728</b> <sub>1</sub>	3.51845 <sub>4</sub>	3.37105 <sub>2</sub>
SVHN_resnet152_SD	1.84773 <sub>2</sub>	1.84773 <sub>2</sub>	1.90535 <sub>6</sub>	<b>1.80547</b> <sub>1</sub>	1.87462 <sub>4</sub>	1.87462 <sub>4</sub>
avg rank	4.14	4.14	4.64	2.11	3.25	2.71

Table 19: Scores and ranking of calibration methods for **p-confidence-ECE**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.06	0.032 <sub>4</sub>	0.363 <sub>2</sub>	0.019 <sub>5</sub>	<b>0.461</b> <sub>1</sub>	0.052 <sub>3</sub>
c10_densenet40	0.04	0.002 <sub>2</sub>	<b>0.525</b> <sub>1</sub>	0.000 <sub>4</sub>	0.000 <sub>4</sub>	0.000 <sub>4</sub>
c10_lenet5	0.06	0.008 <sub>5</sub>	0.027 <sub>4</sub>	0.084 <sub>3</sub>	<b>0.155</b> <sub>1</sub>	0.144 <sub>2</sub>
c10_resnet110	0.04	0.000 <sub>4</sub>	<b>0.246</b> <sub>1</sub>	0.000 <sub>4</sub>	0.000 <sub>4</sub>	0.000 <sub>4</sub>
c10_resnet110_SD	0.06	0.105 <sub>4</sub>	<b>0.179</b> <sub>1</sub>	0.003 <sub>5</sub>	0.114 <sub>3</sub>	0.124 <sub>2</sub>
c10_resnet_wide32	0.06	0.017 <sub>3</sub>	<b>0.281</b> <sub>1</sub>	0.005 <sub>4</sub>	0.005 <sub>4</sub>	0.076 <sub>2</sub>
c100_convnet	0.05	<b>0.174</b> <sub>1</sub>	0.000 <sub>5</sub>	0.049 <sub>2</sub>	0.021 <sub>3</sub>	0.000 <sub>5</sub>
c100_densenet40	0.05	<b>0.817</b> <sub>1</sub>	0.000 <sub>5</sub>	0.617 <sub>2</sub>	0.238 <sub>3</sub>	0.000 <sub>5</sub>
c100_lenet5	0.06	0.153 <sub>4</sub>	0.217 <sub>3</sub>	0.001 <sub>5</sub>	0.395 <sub>2</sub>	<b>0.422</b> <sub>1</sub>
c100_resnet110	0.03	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>
c100_resnet110_SD	0.04	0.009 <sub>2</sub>	0.000 <sub>4</sub>	0.000 <sub>4</sub>	<b>0.060</b> <sub>1</sub>	0.000 <sub>4</sub>
c100_resnet_wide32	0.04	<b>0.022</b> <sub>1</sub>	0.000 <sub>4</sub>	0.000 <sub>4</sub>	0.001 <sub>2</sub>	0.000 <sub>4</sub>
mmist_mlp	0.06	0.616 <sub>3</sub>	<b>0.948</b> <sub>1</sub>	0.486 <sub>4</sub>	0.455 <sub>5</sub>	0.677 <sub>2</sub>
SVHN_convnet	0.03	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>
SVHN_resnet152_SD	0.03	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>	0.000 <sub>3</sub>
avg rank	4.93	2.97	2.9	3.9	2.97	3.33

Table 20: Scores and ranking of calibration methods for **p-classwise-ECE**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.06	0.0104 <sub>4</sub>	<b>0.1276</b> <sub>1</sub>	0.0038 <sub>5</sub>	0.0340 <sub>2</sub>	0.0114 <sub>3</sub>
c10_densenet40	0.04	0.0000 <sub>4</sub>	<b>0.0093</b> <sub>1</sub>	0.0000 <sub>4</sub>	0.0000 <sub>4</sub>	0.0000 <sub>4</sub>
c10_lenet5	0.05	0.0000 <sub>5</sub>	<b>0.6014</b> <sub>1</sub>	0.0390 <sub>4</sub>	0.1230 <sub>2</sub>	0.0501 <sub>3</sub>
c10_resnet110	0.04	0.0000 <sub>4</sub>	<b>0.0088</b> <sub>1</sub>	0.0000 <sub>4</sub>	0.0000 <sub>4</sub>	0.0000 <sub>4</sub>
c10_resnet110_SD	0.06	0.0058 <sub>5</sub>	0.0105 <sub>3</sub>	0.0077 <sub>4</sub>	0.1816 <sub>2</sub>	<b>0.2196</b> <sub>1</sub>
c10_resnet_wide32	0.05	0.0000 <sub>5</sub>	0.0096 <sub>3</sub>	0.0158 <sub>2</sub>	0.0006 <sub>4</sub>	<b>0.0249</b> <sub>1</sub>
c100_convnet	0.04	<b>0.0770</b> <sub>1</sub>	0.0000 <sub>4</sub>	0.0000 <sub>4</sub>	0.0000 <sub>4</sub>	0.0000 <sub>4</sub>
c100_densenet40	0.03	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>
c100_lenet5	0.03	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>
c100_resnet110	0.03	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>
c100_resnet110_SD	0.03	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>
c100_resnet_wide32	0.03	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>
mnist_mlp	0.06	<b>0.5669</b> <sub>1</sub>	0.0842 <sub>3</sub>	0.0022 <sub>5</sub>	0.0280 <sub>4</sub>	0.1178 <sub>2</sub>
SVHN_convnet	0.03	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>
SVHN_resnet152_SD	0.03	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>	0.0000 <sub>3</sub>
avg rank	4.37	3.63	2.77	3.77	3.37	3.1

Table 21: Comparison of MS-ODIR and vector scaling for **log-loss**.

	VecS	Replication 1		VecS	Replication 2		VecS	Replication 3	
		MS-ODIR	MS-ODIR-zero		MS-ODIR	MS-ODIR-zero		MS-ODIR	MS-ODIR-zero
c10_convnet	0.19774	<b>0.19632</b>	0.19632	—	—	—	—	—	—
c10_densenet40	<b>0.22240</b>	0.22240	0.22240	0.21316	<b>0.21186</b>	0.21366	0.21350	<b>0.21325</b>	0.21327
c10_lenet5	0.74688	<b>0.74262</b>	0.74830	0.69392	<b>0.69287</b>	0.69335	<b>0.67955</b>	0.67974	0.68127
c10_resnet110	0.20624	<b>0.20375</b>	0.20537	0.20064	<b>0.19803</b>	0.20040	0.19655	<b>0.19536</b>	0.19739
c10_resnet110_SD	0.17545	<b>0.17537</b>	0.17539	0.18123	<b>0.18094</b>	0.18097	<b>0.17799</b>	0.17829	0.17829
c10_resnet_wide32	0.18274	<b>0.18165</b>	0.18302	0.18522	<b>0.18364</b>	0.18546	0.17431	<b>0.17274</b>	0.17448
c100_convnet	0.96311	<b>0.96141</b>	0.96149	—	—	—	—	—	—
c100_densenet40	1.05714	<b>1.05084</b>	1.06804	1.06366	<b>1.05456</b>	1.07107	1.07704	<b>1.06918</b>	1.08559
c100_lenet5	2.51695	<b>2.48670</b>	2.57932	2.21546	<b>2.20054</b>	2.22360	2.28054	<b>2.27887</b>	2.29485
c100_resnet110	1.08824	<b>1.07370</b>	1.10137	1.09066	<b>1.08267</b>	1.11116	1.11977	<b>1.10672</b>	1.13900
c100_resnet110_SD	<b>0.92275</b>	0.92731	0.92730	0.87758	<b>0.87698</b>	0.87701	<b>0.88523</b>	0.88731	0.88727
c100_resnet_wide32	0.93724	<b>0.93273</b>	0.94060	0.93291	<b>0.92531</b>	0.94854	0.93183	<b>0.92439</b>	0.94568
SVHN_convnet	0.14392	<b>0.13760</b>	0.14507	—	—	—	—	—	—
SVHN_resnet152_SD	0.08131	<b>0.08100</b>	0.08100	0.12728	<b>0.12723</b>	0.12639	0.12559	<b>0.12453</b>	0.12381