

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

Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva  
Filho, Hao Song, Peter Flach

NeurIPS 2019



UNIVERSITY OF TARTU



University of  
BRISTOL



Universidade Federal da Paraíba  
ESTATÍSTICA



The  
Alan Turing  
Institute



# Notion of Calibration

A probabilistic classifier  $p^\wedge$  is:

- Multiclass-calibrated if for any prediction vector  $q = (q_1, \dots, q_k) \in \nabla_k$ , proportions of classes

$$P(Y = i \mid \hat{\mathbf{p}}(X) = \mathbf{q}) = q_i \quad \text{for } i = 1, \dots, k.$$

- Classwise-calibrated if

$$P(Y = i \mid \hat{p}_i(X) = q_i) = q_i.$$

- Confidence-calibrated if

$$P\left(Y = \operatorname{argmax}(\hat{\mathbf{p}}(X)) \mid \max(\hat{\mathbf{p}}(X)) = c\right) = c.$$

# Evaluation of Calibration

Some empirical measures of calibration are

confidence-ECE (place here equation)

$$\text{classwise-ECE} = \frac{1}{k} \sum_{j=1}^k \sum_{i=1}^m \frac{|B_{i,j}|}{n} |y_j(B_{i,j}) - \hat{p}_j(B_{i,j})|$$

Ever proper loss is minimised by the canonical calibration function

P-test versions for conf-ECE and cw-ECE

# Temperature Scaling

Temperature scaling proposed by [9] learns a temperature-parameter  $t > 0$  which decreases ( $t > 1$ ) or increases ( $t < 1$ ) the confidence of a model with softmax output.

$$\sigma_{\text{SM}}(\mathbf{z}_i)^{(k)} = \frac{\exp(z_i^{(k)})}{\sum_{j=1}^K \exp(z_i^{(j)})}, \quad \hat{p}_i = \max_k \sigma_{\text{SM}}(\mathbf{z}_i)^{(k)}$$

$$\hat{q}_i = \max_k \sigma_{\text{SM}}(\mathbf{z}_i / T)^{(k)}.$$

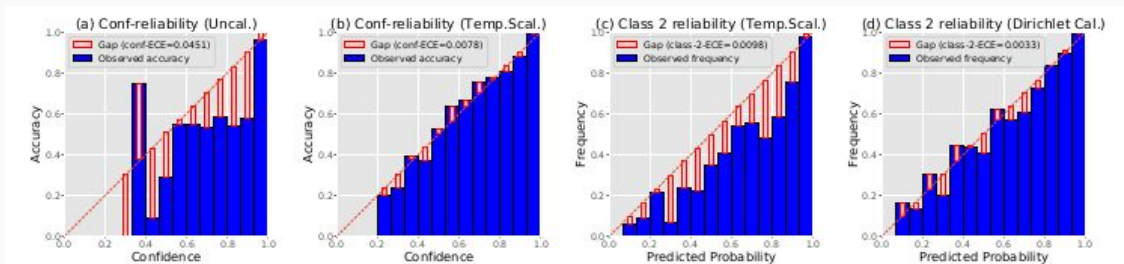


Figure 1: Reliability diagrams of c10\_resnet\_wide32 on CIFAR-10: (a) confidence-reliability before calibration; (b) confidence-reliability after temperature scaling; (c) classwise-reliability for class 2 after temperature scaling; (d) classwise-reliability for class 2 after Dirichlet calibration.

# Dirichlet Calibration

We consider the distribution of prediction vectors  $\hat{\mathbf{p}}(\mathbf{x})$  separately on instances of each class, and assume these are Dirichlet distributions with different parameters:

$$\hat{\mathbf{p}}(X) \mid Y = j \sim \text{Dir}(\boldsymbol{\alpha}^{(j)})$$

generative parametrisation:  $\hat{\boldsymbol{\mu}}_{\text{DirGen}}(\mathbf{q}; \boldsymbol{\alpha}, \boldsymbol{\pi}) = (\pi_1 f_1(\mathbf{q}), \dots, \pi_k f_k(\mathbf{q})) / z$

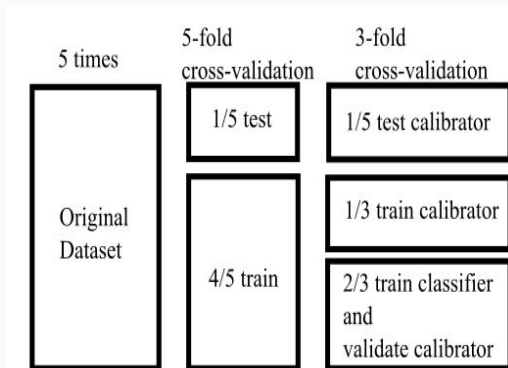
linear parametrisation:  $\hat{\boldsymbol{\mu}}_{\text{DirLin}}(\mathbf{q}; \mathbf{W}, \mathbf{b}) = \sigma(\mathbf{W} \ln \mathbf{q} + \mathbf{b})$

canonical parametrisation:  $\hat{\boldsymbol{\mu}}_{\text{Dir}}(\mathbf{q}; \mathbf{A}, \mathbf{c}) = \sigma(\mathbf{A} \ln \frac{\mathbf{q}}{1/k} + \ln \mathbf{c})$

All parameterizations are equal, i.e. they contain exactly the same calibration maps.

# Experiment non-neural: Setting

- 21 datasets and 11 classifiers = 231 settings
  - Logistic, nbayes, forest, adas, lda, qda, tree, knn, mlp, svc-linear, svc-rbf
- 6 calibration methods:
  - OvR\_Isotonic, OvR\_Width\_Bin, OvR\_Freq\_Bin, OvR\_Beta, Temp\_Scaling, Dirichlet\_L2
- 8 evaluation measures
- 5 times 5-fold-crossval.
  - Inner 3-fold-crossval.



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

# Experiment non-neural: Results

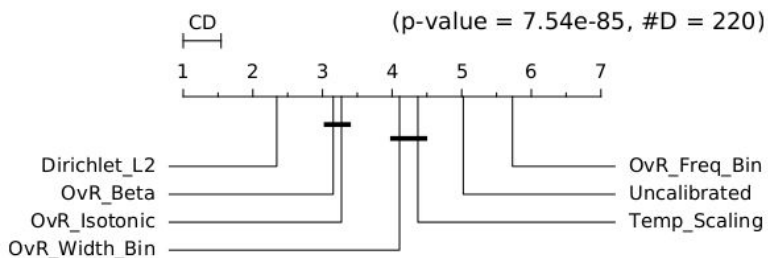
Table 1: Ranking of calibration methods for **p-cw-ECE** (Friedman's test significant with p-value  $7.54e^{-85}$ ).

	DirL2	Beta	FreqB	Isot	WidB	TempS	Uncal
adas	<b>2.4</b>	3.2	4.1	4.2	3.9	5.0	5.2
forest	3.5	<b>2.3</b>	5.7	3.0	3.6	5.0	5.0
knn	2.5	4.0	4.5	<b>2.1</b>	3.2	5.8	6.0
lda	<b>1.9</b>	3.1	5.8	3.0	3.5	5.0	5.8
logistic	<b>2.2</b>	2.8	6.4	3.0	4.2	3.9	5.5
mlp	<b>2.2</b>	2.9	6.7	4.0	5.2	3.0	4.1
nbayes	<b>1.4</b>	3.6	4.8	2.6	4.2	5.3	6.1
qda	<b>2.2</b>	2.8	6.3	2.5	3.8	4.8	5.6
svc-linear	<b>2.3</b>	2.7	6.7	3.8	4.0	3.7	4.8
svc-rbf	<b>2.9</b>	3.0	6.3	3.5	4.1	3.9	4.3
tree	<b>2.4</b>	4.3	5.9	4.2	5.2	3.0	3.0
avg rank	<b>2.34</b>	3.15	5.73	3.27	4.11	4.37	5.02

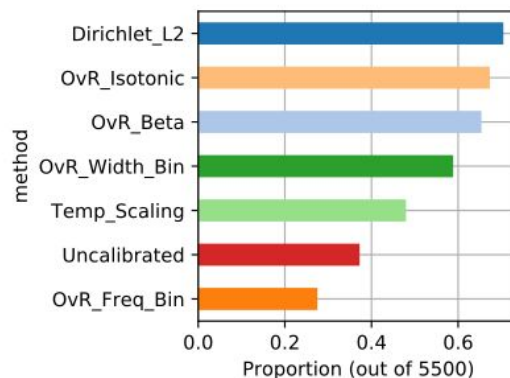
Table 2: Ranking of calibration methods for **log-loss** (p-value  $4.39e^{-77}$ ).

	DirL2	Beta	FreqB	Isot	WidB	TempS	Uncal
	<b>1.4</b>	3.1	3.2	4.3	3.5	5.9	6.6
	4.2	<b>1.9</b>	4.7	4.1	2.9	5.2	5.2
	3.8	4.8	3.0	<b>1.6</b>	2.0	6.5	6.5
	<b>1.6</b>	2.2	5.2	5.2	3.5	4.6	5.7
	<b>1.3</b>	2.1	5.8	6.1	3.5	3.6	5.6
	<b>2.2</b>	2.3	6.5	6.2	4.7	2.9	3.4
	<b>1.1</b>	3.4	3.4	4.0	4.4	5.5	6.3
	<b>1.7</b>	2.7	5.6	4.6	3.4	4.2	5.8
	<b>1.3</b>	2.3	6.1	6.1	4.3	3.0	4.8
	2.6	<b>2.2</b>	4.3	4.8	4.5	4.0	5.6
	3.9	5.1	3.4	<b>2.1</b>	2.4	5.6	5.6
	<b>2.25</b>	2.92	4.66	4.48	3.54	4.61	5.54

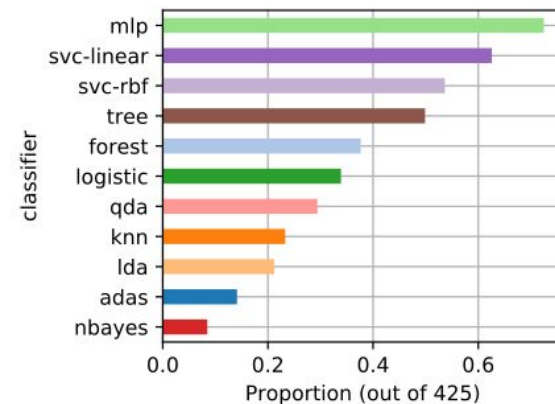
# Experiment non-neural: Results



(a) p-cw-ECE critical difference



(b) p-cw-ECE for calibrators



(c) p-cw-ECE for classifiers



# Experiment Deep Neural Networks: Setting

- 3 datasets
  - CIFAR-10, CIFAR-100, SVHN
- 11 convolutional NNs
  - ResNet 110, ResNet 110 SD, ResNet 152 SD, WideNet 32, LeNet 5
- 3 pretrained CNNs
- 5 calibration methods:
  - TempS, Dir-L2, Dir-ODIR, VecS, MS-ODIR
- 8 evaluation measures
- 5-fold-crossval.
  - Optimal regularisation parameters

# Experiment Deep Neural Networks: Results

Table 3: Scores and ranking of calibration methods for **cw-ECE**.

	Uncal	general-purpose calibrators			calibrators using logits	
		TempS	Dir-L2	Dir-ODIR	VecS	MS-ODIR
c10_convnet	0.104 <sub>6</sub>	0.044 <sub>4</sub>	0.043 <sub>2</sub>	0.045 <sub>5</sub>	<b>0.043<sub>1</sub></b>	0.044 <sub>3</sub>
c10_densenet40	0.114 <sub>6</sub>	0.040 <sub>5</sub>	<b>0.034<sub>1</sub></b>	0.037 <sub>4</sub>	0.036 <sub>2</sub>	0.037 <sub>3</sub>
c10_lenet5	0.198 <sub>6</sub>	0.171 <sub>5</sub>	<b>0.052<sub>1</sub></b>	0.059 <sub>4</sub>	0.057 <sub>2</sub>	0.059 <sub>3</sub>
c10_resnet110	0.098 <sub>6</sub>	0.043 <sub>5</sub>	<b>0.032<sub>1</sub></b>	0.039 <sub>4</sub>	0.037 <sub>3</sub>	0.036 <sub>2</sub>
c10_resnet110_SD	0.086 <sub>6</sub>	0.031 <sub>4</sub>	0.031 <sub>5</sub>	0.029 <sub>3</sub>	0.027 <sub>2</sub>	<b>0.027<sub>1</sub></b>
c10_resnet_wide32	0.095 <sub>6</sub>	0.048 <sub>5</sub>	0.032 <sub>3</sub>	0.029 <sub>2</sub>	0.032 <sub>4</sub>	<b>0.029<sub>1</sub></b>
c100_convnet	0.424 <sub>6</sub>	<b>0.227<sub>1</sub></b>	0.402 <sub>5</sub>	0.240 <sub>3</sub>	0.241 <sub>4</sub>	0.240 <sub>2</sub>
c100_densenet40	0.470 <sub>6</sub>	0.187 <sub>2</sub>	0.330 <sub>5</sub>	<b>0.186<sub>1</sub></b>	0.189 <sub>3</sub>	0.191 <sub>4</sub>
c100_lenet5	0.473 <sub>6</sub>	0.385 <sub>5</sub>	0.219 <sub>4</sub>	0.213 <sub>2</sub>	<b>0.203<sub>1</sub></b>	0.214 <sub>3</sub>
c100_resnet110	0.416 <sub>6</sub>	0.201 <sub>3</sub>	0.359 <sub>5</sub>	<b>0.186<sub>1</sub></b>	0.194 <sub>2</sub>	0.203 <sub>4</sub>
c100_resnet110_SD	0.375 <sub>6</sub>	0.203 <sub>4</sub>	0.373 <sub>5</sub>	0.189 <sub>3</sub>	<b>0.170<sub>1</sub></b>	0.186 <sub>2</sub>
c100_resnet_wide32	0.420 <sub>6</sub>	0.186 <sub>4</sub>	0.333 <sub>5</sub>	0.180 <sub>2</sub>	<b>0.171<sub>1</sub></b>	0.180 <sub>3</sub>
SVHN_convnet	0.159 <sub>6</sub>	0.038 <sub>4</sub>	0.043 <sub>5</sub>	0.026 <sub>2</sub>	<b>0.025<sub>1</sub></b>	0.027 <sub>3</sub>
SVHN_resnet152_SD	0.019 <sub>2</sub>	<b>0.018<sub>1</sub></b>	0.022 <sub>6</sub>	0.020 <sub>3</sub>	0.021 <sub>5</sub>	0.021 <sub>4</sub>
Average rank	5.71	3.71	3.79	2.79	2.29	2.71

Table 4: 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	0.391 <sub>6</sub>	<b>0.195<sub>1</sub></b>	0.197 <sub>4</sub>	0.195 <sub>2</sub>	0.197 <sub>5</sub>	0.196 <sub>3</sub>
c10_densenet40	0.428 <sub>6</sub>	0.225 <sub>5</sub>	<b>0.220<sub>1</sub></b>	0.224 <sub>4</sub>	0.223 <sub>3</sub>	0.222 <sub>2</sub>
c10_lenet5	0.823 <sub>6</sub>	0.800 <sub>5</sub>	0.744 <sub>2</sub>	0.744 <sub>3</sub>	0.747 <sub>4</sub>	<b>0.743<sub>1</sub></b>
c10_resnet110	0.358 <sub>6</sub>	0.209 <sub>5</sub>	<b>0.203<sub>1</sub></b>	0.205 <sub>3</sub>	0.206 <sub>4</sub>	0.204 <sub>2</sub>
c10_resnet110_SD	0.303 <sub>6</sub>	0.178 <sub>5</sub>	0.177 <sub>4</sub>	0.176 <sub>3</sub>	0.175 <sub>2</sub>	<b>0.175<sub>1</sub></b>
c10_resnet_wide32	0.382 <sub>6</sub>	0.191 <sub>5</sub>	0.185 <sub>4</sub>	0.182 <sub>2</sub>	0.183 <sub>3</sub>	<b>0.182<sub>1</sub></b>
c100_convnet	1.641 <sub>6</sub>	<b>0.942<sub>1</sub></b>	1.189 <sub>5</sub>	0.961 <sub>2</sub>	0.964 <sub>4</sub>	0.961 <sub>3</sub>
c100_densenet40	2.017 <sub>6</sub>	1.057 <sub>2</sub>	1.253 <sub>5</sub>	1.059 <sub>4</sub>	1.058 <sub>3</sub>	<b>1.051<sub>1</sub></b>
c100_lenet5	2.784 <sub>6</sub>	2.650 <sub>5</sub>	2.595 <sub>4</sub>	2.490 <sub>2</sub>	2.516 <sub>3</sub>	<b>2.487<sub>1</sub></b>
c100_resnet110	1.694 <sub>6</sub>	1.092 <sub>3</sub>	1.212 <sub>5</sub>	1.096 <sub>4</sub>	1.089 <sub>2</sub>	<b>1.074<sub>1</sub></b>
c100_resnet110_SD	1.353 <sub>6</sub>	0.942 <sub>3</sub>	1.198 <sub>5</sub>	0.945 <sub>4</sub>	<b>0.923<sub>1</sub></b>	0.927 <sub>2</sub>
c100_resnet_wide32	1.802 <sub>6</sub>	0.945 <sub>3</sub>	1.087 <sub>5</sub>	0.953 <sub>4</sub>	0.937 <sub>2</sub>	<b>0.933<sub>1</sub></b>
SVHN_convnet	0.205 <sub>6</sub>	0.151 <sub>5</sub>	0.142 <sub>3</sub>	0.138 <sub>2</sub>	0.144 <sub>4</sub>	<b>0.138<sub>1</sub></b>
SVHN_resnet152_SD	0.085 <sub>6</sub>	<b>0.079<sub>1</sub></b>	0.085 <sub>5</sub>	0.080 <sub>2</sub>	0.081 <sub>4</sub>	0.081 <sub>3</sub>
Average rank	6.0	3.5	3.79	2.93	3.14	1.64

# Conclusion

1. Dirichlet calibration: New parametric general-purpose multiclass calibration method
  - a. Natural extension of two-class Beta calibration
  - b. Easy to implement with multinomial logistic regression on log-transformed class probabilities
2. Best or tied best performance with 21 datasets x 11 classifiers
3. Advances state-of-the-art on Neural Networks by introducing ODIR regularisation



UNIVERSITY OF TARTU



# Future work

1. Which neural architectures and training methods have temperature scaling as a canonical calibration function
2. Use other distributions of the exponential family
3. Investigate scores coming from mixtures of distributions per class



UNIVERSITY OF TARTU



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

Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva  
Filho, Hao Song, Peter Flach

NeurIPS 2019



UNIVERSITY OF TARTU



University of  
BRISTOL



Universidade Federal da Paraíba  
ESTATÍSTICA



The  
Alan Turing  
Institute

