
Classifier Calibration, Why and How?

No Trustworthy Machine Learning Without It

Peter Flach

Peter.Flach@bristol.ac.uk

[classifier-calibration.github.io/](https://github.com/peterflach/classifier-calibration)



What's in this module?

For a machine learning classifier to be trustworthy, its outputs need to be meaningful.

In particular, if a classifier outputs probabilities they need to correspond to relative frequencies of events in the world.

We will take a closer look at what this means, and how this can be achieved.



Table of Contents

Let's talk about the weather...

Why are we interested in calibration?

Common sources of miscalibration

A first look at some calibration techniques

Calibrating multi-class classifiers



Let's talk about the weather...

Part of: Classifier Calibration, Why and How?

Peter Flach

Peter.Flach@bristol.ac.uk

classifier-calibration.github.io/



Taking inspiration from forecasting

Weather forecasters started thinking about calibration a long time ago (Brier, 1950).

- ✍ A forecast '70% chance of rain' should be followed by rain 70% of the time.

This is immediately applicable to binary classification:

- ✍ A prediction '70% chance of spam' should be spam 70% of the time.

and to multi-class classification:

- ✍ A prediction '70% chance of setosa, 10% chance of versicolor and 20% chance of virginica' should be setosa/versicolor/virginica 70/10/20% of the time.

In general:

- ✍ A predicted probability (vector) should match empirical (observed) probabilities.

Q: What does 'x% of the time' mean?



Forecasting example

Let's consider a small toy example:

- ✂ Two predictions of '10% chance of rain' were both followed by 'no rain'.
- ✂ Two predictions of '40% chance of rain' were once followed by 'no rain', and once by 'rain'.
- ✂ Three predictions of '70% chance of rain' were once followed by 'no rain', and twice by 'rain'.
- ✂ One prediction of '90% chance of rain' was followed by 'rain'.

Q: Is this forecaster well-calibrated?



Over- and under-estimates

	\hat{p}	y
0	0.1	0
1	0.1	0
2	0.4	0
3	0.4	1
4	0.7	0
5	0.7	1
6	0.7	1
7	0.9	1

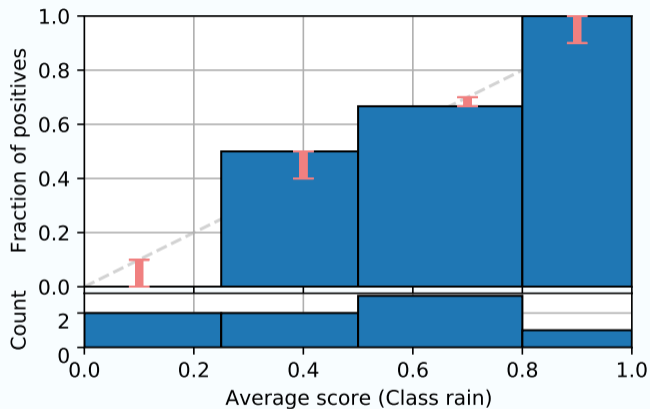
This forecaster is doing a pretty decent job:

- '10% chance of rain' was a slight over-estimate ($\bar{y} = 0/2 = 0\%$);
- '40% chance of rain' was a slight under-estimate ($\bar{y} = 1/2 = 50\%$);
- '70% chance of rain' was a slight over-estimate ($\bar{y} = 2/3 = 67\%$);
- '90% chance of rain' was a slight under-estimate ($\bar{y} = 1/1 = 100\%$).



Visualising forecasts: the reliability diagram

	\hat{p}	y
0	0.1	0
1	0.1	0
2	0.4	0
3	0.4	1
4	0.7	0
5	0.7	1
6	0.7	1
7	0.9	1

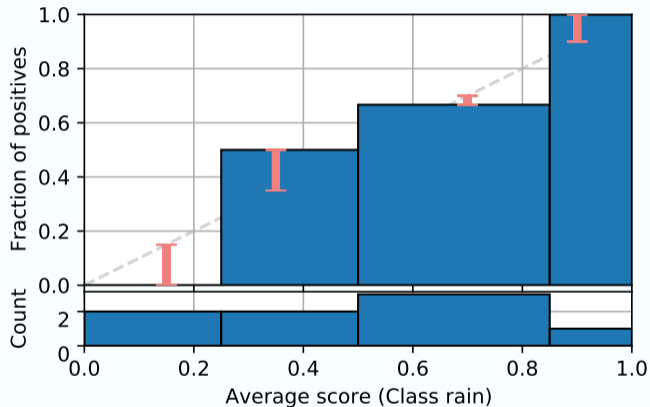


Figures generated with PyCalib (Perello-Nieto et al., 2021)



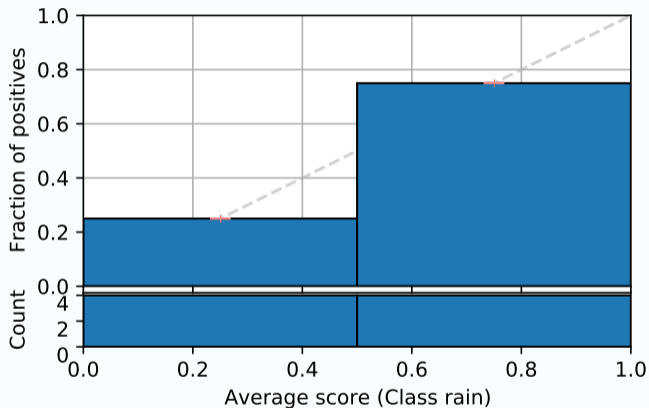
Changing the numbers slightly

	\hat{p}	y
0	0.1	0
1	0.2	0
2	0.3	0
3	0.4	1
4	0.6	0
5	0.7	1
6	0.8	1
7	0.9	1



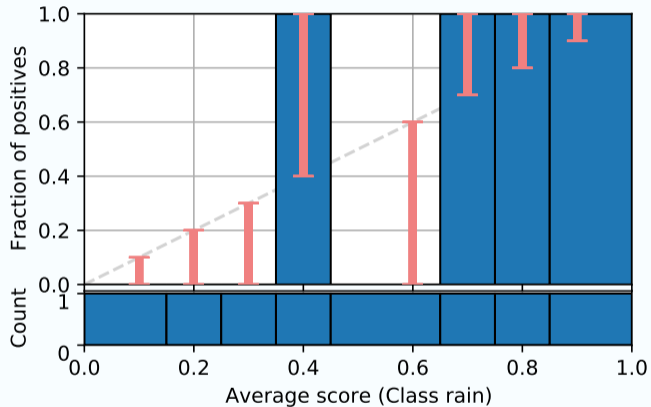
Or should we group forecasts differently?

	\hat{p}	y
0	0.1	0
1	0.2	0
2	0.3	0
3	0.4	1
4	0.6	0
5	0.7	1
6	0.8	1
7	0.9	1



Or not at all?

	\hat{p}	y
0	0.1	0
1	0.2	0
2	0.3	0
3	0.4	1
4	0.6	0
5	0.7	1
6	0.8	1
7	0.9	1



Binning or pooling predictions is a fundamental notion

We need bins to **evaluate** the degree of calibration:

- In order to decide whether a weather forecaster is well-calibrated, we need to look at a good number of forecasts, say over one year.
- We also need to make sure that there are a reasonable number of forecasts for separate probability values, so we can obtain reliable empirical estimates.
 - ✍ Trade-off: large bins give better empirical estimates, small bins allows a more fine-grained assessment of calibration.

But adjusting forecasts in groups also gives rise to practical calibration **methods**:

- empirical binning
- isotonic regression (aka ROC convex hull)



References



Brier, G. W. (1950). Verification of forecasts expressed in terms of probabilities.
Monthly Weather Review, 78(1), 1–3.



Perello-Nieto, M., Song, H., Silva Filho, T. & Kängsepp, M. (2021). *Pycalib a library for classifier calibration* (Version 0.1.0.dev3). Zenodo.
<https://doi.org/10.5281/zenodo.5518877>



Why are we interested in calibration?

Part of: Classifier Calibration, Why and How?

Peter Flach

Peter.Flach@bristol.ac.uk

classifier-calibration.github.io/



Why are we interested in calibration?

To calibrate means **to employ a known scale with known properties**.

✈ E.g., additive scale with a well-defined zero, so that ratios are meaningful.

For classifiers we want to use the probability scale, so that we can

- justifiably use default decision rules (e.g., maximum posterior probability);
- adjust these decision rules in a straightforward way to account for different class priors or misclassification costs;
- combine probability estimates in a well-founded way.

Q: Is the probability scale additive?

Q: How would you combine probability estimates from several well-calibrated models?



Optimal decisions I

Denote the cost of predicting class j for an instance of true class i as $C(\hat{Y} = j|Y = i)$.
The expected cost of predicting class j for instance x is

$$C(\hat{Y} = j|X = x) = \sum_i P(Y = i|X = x)C(\hat{Y} = j|Y = i)$$

where $P(Y = i|X = x)$ is the probability of instance x having true class i (as would be given by the Bayes-optimal classifier).

The optimal decision is then to predict the class with lowest expected cost:

$$\hat{Y}^* = \operatorname{argmin}_j C(\hat{Y} = j|X = x) = \operatorname{argmin}_j \sum_i P(Y = i|X = x)C(\hat{Y} = j|Y = i)$$



Optimal decisions II

In binary classification we have:

$$C(\hat{Y} = +|X = x) = P(+|x)C(+|+) + (1 - P(+|x))C(+|-)$$

$$C(\hat{Y} = -|X = x) = P(+|x)C(-|+) + (1 - P(+|x))C(-|-)$$

On the optimal decision boundary these two expected costs are equal, which gives

$$P(+|x) = \frac{C(+|-) - C(-|-)}{C(+|-) - C(-|-) + C(-|+) - C(+|+)} \triangleq c$$

This gives the optimal threshold on the hypothetical Bayes-optimal probabilities. It is also the best thing to do in practice – as long as the probabilities are well-calibrated!

Optimal decisions III

Without loss of generality we can set the cost of true positives and true negatives to zero; $c = \frac{c_{FP}}{c_{FP} + c_{FN}}$ is then the cost of a false positive in proportion to the combined cost of one false positive and one false negative.

- ✎ E.g., if false positives are 4 times as costly as false negatives then we set the decision threshold to $4/(4 + 1) = 0.8$ in order to only make positive predictions if we're pretty certain.

Similar reasoning applies to changes in class priors:

- if we trained on balanced classes but want to deploy with 4 times as many positives compared to negatives, we lower the decision threshold to 0.2;
- more generally, if we trained for class ratio r and deploy for class ratio r' we set the decision threshold to $r/(r + r')$.

Cost and class prior changes can be combined in the obvious way.

Common sources of miscalibration

Part of: Classifier Calibration, Why and How?

Peter Flach

Peter.Flach@bristol.ac.uk

classifier-calibration.github.io/



Common sources of miscalibration

Underconfidence: a classifier thinks it's **worse** at separating classes than it actually is.

- Hence we need to *pull predicted probabilities away from the centre*.

Overconfidence: a classifier thinks it's **better** at separating classes than it actually is.

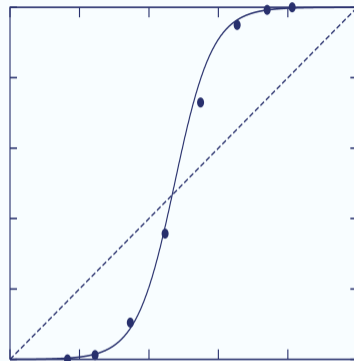
- Hence we need to *push predicted probabilities toward the centre*.

A classifier can be overconfident for one class and underconfident for the other, in which case all predicted probabilities need to be increased or decreased.



Underconfidence example

- Underconfidence typically gives **sigmoidal** distortions.
- To calibrate these means to *pull predicted probabilities away from the centre.*

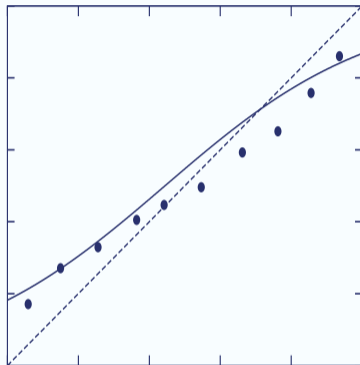


Source: Niculescu-Mizil and Caruana (2005)



Overconfidence example

- Overconfidence is very common, and usually a consequence of over-counting evidence.
- Here, distortions are **inverse-sigmoidal**
- Calibrating these means to *push predicted probabilities toward the centre*.



Source: Niculescu-Mizil and Caruana (2005)

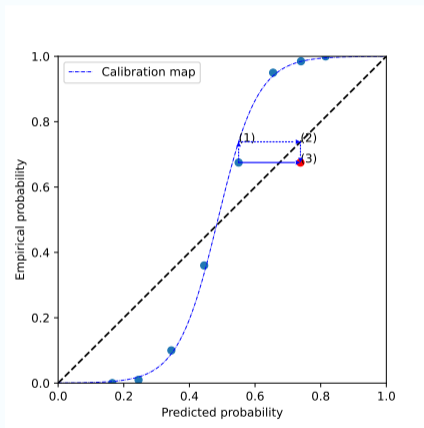


Why fitting the distortions helps with calibration


In clockwise direction, the dotted arrows show:

- (1) using a point's uncalibrated score on the x-axis as input to the calibration map,
- (2) mapping the resulting output back to the diagonal, and
- (3) combine with the empirical probability of the point we started from.

The closer the original point is to the fitted calibration map, the closer the calibrated point (in red) will be to the diagonal.



References

-  Niculescu-Mizil, A. & Caruana, R. (2005). Predicting good probabilities with supervised learning. In *22nd international conference on machine learning (icml'05)*, ACM Press.



A first look at some calibration techniques

Part of: Classifier Calibration, Why and How?

Peter Flach

Peter.Flach@bristol.ac.uk

classifier-calibration.github.io/



A first look at some calibration techniques

Parametric calibration involves modelling the score distributions within each class.

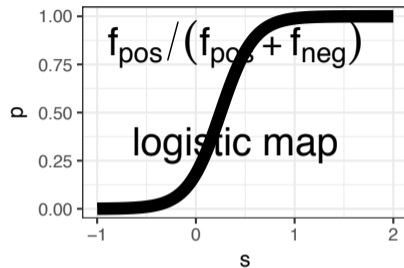
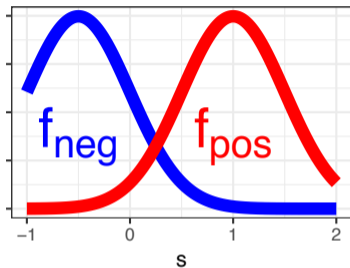
- ✍ **Platt scaling** = Logistic calibration can be derived by assuming that the scores within both classes are normally distributed with the same variance (Platt, 2000).
- ✍ **Beta calibration** employs Beta distributions instead, to deal with scores already on a $[0, 1]$ scale (Kull et al., 2017).
- ✍ **Dirichlet calibration** for more than two classes (Kull et al., 2019).

Non-parametric calibration often ignores scores and employs ranks instead.

- ✍ E.g., **isotonic regression** = pool adjacent violators = ROC convex hull (Fawcett & Niculescu-Mizil, 2007; Zadrozny & Elkan, 2001).



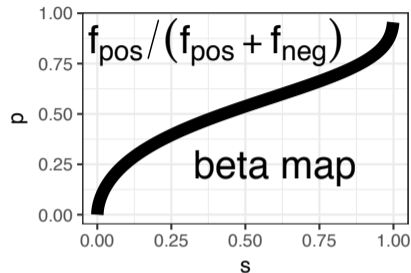
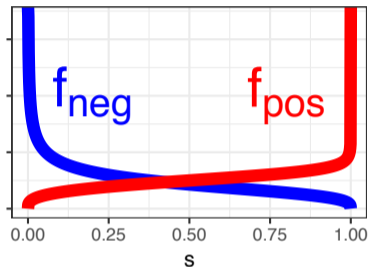
Platt scaling



$$p(s; w, m) = \frac{1}{1 + \exp(-w(s - m))}$$
$$w = (\mu_{pos} - \mu_{neg}) / \sigma^2, m = (\mu_{pos} + \mu_{neg}) / 2$$



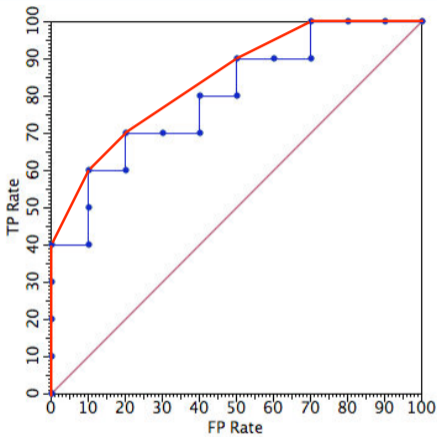
Beta calibration



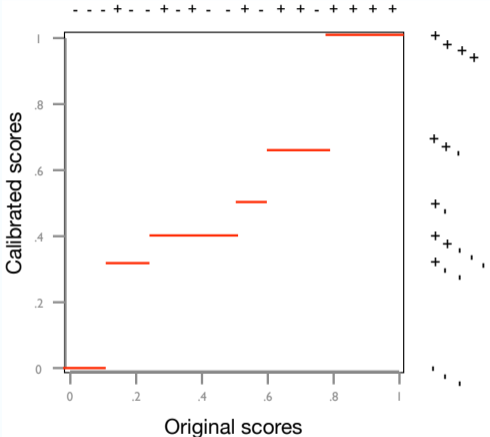
$$p(s; a, b, c) = \frac{1}{1 + \exp(-a \ln s - b \ln(1 - s) - c)}$$
$$a = \alpha_{pos} - \alpha_{neg}, b = \beta_{neg} - \beta_{pos}$$








Isotonic regression



Source: Flach (2016)



References I

-  Fawcett, T. & Niculescu-Mizil, A. (2007). PAV and the ROC convex hull. *Machine Learning*, 68(1), 97–106.
-  Flach, P. A. (2016). ROC analysis. In *Encyclopedia of machine learning and data mining*. Springer.
-  Kull, M., Perello-Nieto, M., Kängsepp, M., Filho, T. S., Song, H. & Flach, P. (2019). Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with Dirichlet calibration. In *Advances in neural information processing systems (nips'19)*.
-  Kull, M., Silva Filho, T. M. & Flach, P. (2017). Beyond Sigmoids: How to obtain well-calibrated probabilities from binary classifiers with beta calibration. *Electronic Journal of Statistics*, 11(2), 5052–5080.
-  Platt, J. (2000). Probabilities for SV Machines. In A. J. Smola, P. Bartlett, B. Schölkopf & D. Schuurmans (Eds.), *Advances in large-margin classifiers* (pp. 61–74). MIT Press.



References II



Zadrozny, B. & Elkan, C. (2001). Obtaining calibrated probability estimates from decision trees and naive Bayesian classifiers. In *18th international conference on machine learning (icml'01)*.

Calibrating multi-class classifiers

Part of: Classifier Calibration, Why and How?

Peter Flach

Peter.Flach@bristol.ac.uk

classifier-calibration.github.io/



What's so special about multi-class calibration?

Similar to classification, some methods are inherently multi-class but most are not.

This leads to (at least) three different ways of **defining** what it means to be fully multiclass-calibrated.

✍ Many recent papers use the (weak) notion of confidence calibration.

Evaluating multi-class calibration is in its full generality still an open problem.



Definitions of calibration for more than two classes

The following definitions of calibration are equivalent for binary classification but increasingly stronger for more than two classes:

Confidence calibration: only consider the highest predicted probability.

Class-wise calibration: only consider marginal probabilities.

Multi-class calibration: consider the entire vector of predicted probabilities.



Confidence calibration

This was proposed by Guo et al. (2017), requiring that among all instances where the probability of **the most likely class** is predicted to be c , the expected accuracy is c . (We call this ‘confidence calibration’ to distinguish it from the stronger notions of calibration.)

Formally, a probabilistic classifier $\hat{\mathbf{p}} : \mathcal{X} \rightarrow \Delta_k$ is **confidence-calibrated**, if for any confidence level $c \in [0, 1]$, the actual proportion of the predicted class, among all possible instances \mathbf{x} being predicted this class with confidence c , is equal to c :

$$P(Y = i \mid \hat{p}_i(\mathbf{x}) = c) = c \quad \text{where } i = \underset{j}{\operatorname{argmax}} \hat{p}_j(\mathbf{x}).$$



Class-wise calibration

Originally proposed by Zadrozny and Elkan (2002), this requires that all **one-vs-rest** probability estimators obtained from the original multiclass model are calibrated.

Formally, a probabilistic classifier $\hat{\mathbf{p}} : \mathcal{X} \rightarrow \Delta_k$ is **classwise-calibrated**, if for any class i and any predicted probability q_i for this class, the actual proportion of class i , among all possible instances \mathbf{x} getting the same prediction $\hat{p}_i(\mathbf{x}) = q_i$, is equal to q_i :

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



Multi-class calibration

This is the **strongest form of calibration** for multiple classes, subsuming the previous two definitions.

A probabilistic classifier $\hat{\mathbf{p}} : \mathcal{X} \rightarrow \Delta_k$ is **multiclass-calibrated** if for any prediction vector $\mathbf{q} = (q_1, \dots, q_k) \in \Delta_k$, the proportions of classes among all possible instances \mathbf{x} getting the same prediction $\hat{\mathbf{p}}(\mathbf{x}) = \mathbf{q}$ are equal to the prediction vector \mathbf{q} :

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



Reminder: binning needed

For practical purposes, the conditions in these definitions need to be relaxed. This is where **binning** comes in.

Once we have the bins, we can draw a **reliability diagram** as in the two-class case. For class-wise calibration, we can show per-class reliability diagrams or a single averaged one.

The degree of calibration is assessed using the **gaps** in the reliability diagram. All of this will be elaborated in the next part of the tutorial.



Important points to remember

Only well-calibrated probability estimates are worthy to be called probabilities:
otherwise they are just scores that happen to be in the $[0, 1]$ range.

Binning will be required in some form:

instance-based probability evaluation metrics such as Brier score or log-loss
always measure calibration **plus something else.**

In multi-class settings, think carefully about which form of calibration you need:
e.g., confidence-calibration is too weak in a cost-sensitive setting.



References

-  Guo, C., Pleiss, G., Sun, Y. & Weinberger, K. Q. (2017). On Calibration of Modern Neural Networks. In *34th international conference on machine learning*.
-  Zadrozny, B. & Elkan, C. (2002). Transforming Classifier Scores into Accurate Multiclass Probability Estimates. In *8th acm sigkdd international conference on knowledge discovery and data mining - kdd '02*, ACM Press.



Classifier Calibration, Why and How?

No Trustworthy Machine Learning Without It

Peter Flach

Peter.Flach@bristol.ac.uk

classifier-calibration.github.io/

