

Part III

Model Calibration

Adrian Galdran, MSC Research Fellow
Universitat Pompeu Fabra, Barcelona, Spain
University of Adelaide, Australia

Meritxell Riera i Marin, Researcher
Sycal Technologies, Barcelona, Spain
Universitat Pompeu Fabra, Barcelona, Spain



Contents

1. Calibration: what, when and why?
2. Visualizing & Measuring Calibration
3. Improving Calibration
4. Practical Hands-On Session



Palamós, ES >



NOW

HOURLY

10 DAY

MAPS

Tuesday, 6 June 2023



16

24°



60%

Precip Prob

17

23°



30%

Precip Prob

18

23°



10%

Precip Prob

19

22°



0%

Precip Prob

One of the most important tests of a forecast — I would argue that it is the single most important one — is called calibration. Out of all the times you said there was a 40 percent chance of rain, how often did rain actually occur? If, over the long run, it really did rain about 40 percent of the time, that means your forecasts were well calibrated. If it wound up raining just 20 percent of the time instead, or 60 percent of the time, they weren't.

Nate Silver, The Signal and the Noise: Why So Many Predictions Fail – but Some Don't

1. Calibration: what, when and why?

p	y
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	1

p	y
$\frac{1}{3}$	0
$\frac{1}{3}$	0
$\frac{1}{3}$	1
$\frac{1}{2}$	0
$\frac{1}{2}$	1

p	y
$\frac{3}{4}$	0
$\frac{3}{4}$	1
$\frac{3}{4}$	1
$\frac{3}{4}$	1
1	1

1. Calibration: what, when and why?

p	y
$\frac{1}{5}$	0
$\frac{1}{5}$	0
$\frac{1}{5}$	0
$\frac{1}{5}$	0
$\frac{1}{5}$	1

p	y
$\frac{1}{3}$	0
$\frac{1}{3}$	0
$\frac{1}{3}$	1
$\frac{1}{2}$	0
$\frac{1}{2}$	1

p	y
$\frac{3}{4}$	0
$\frac{3}{4}$	1
$\frac{3}{4}$	1
$\frac{3}{4}$	1
1	1

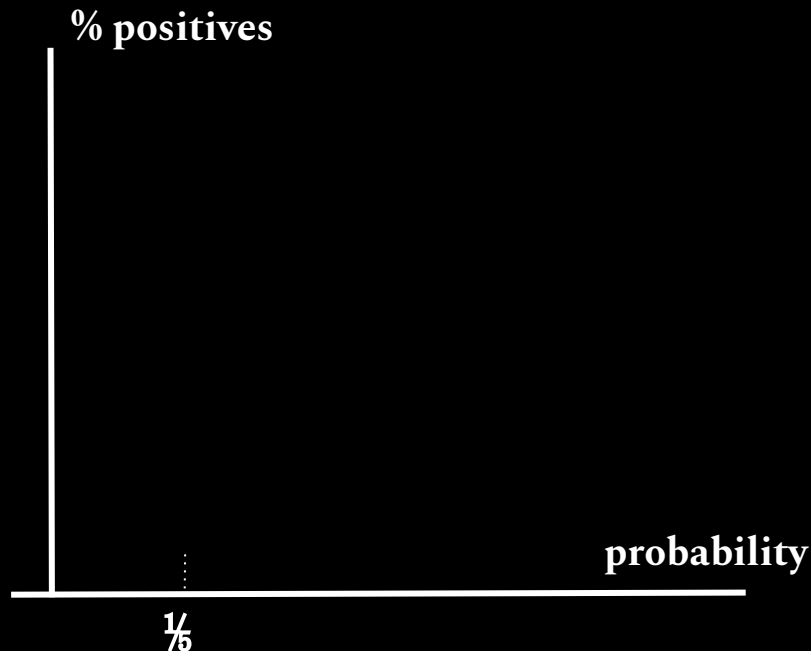
1. Calibration: what, when and why?

p	y
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	1



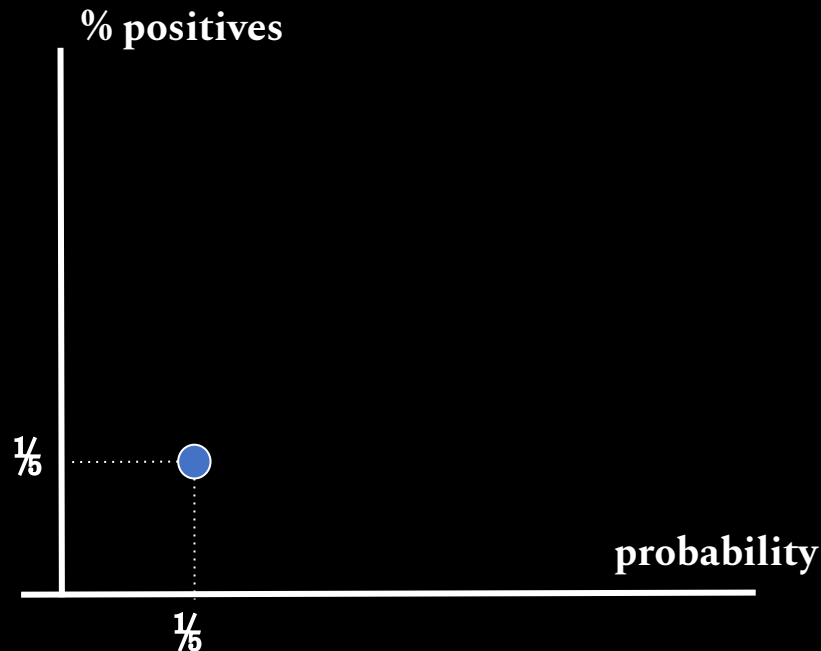
1. Calibration: what, when and why?

p	y
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	1



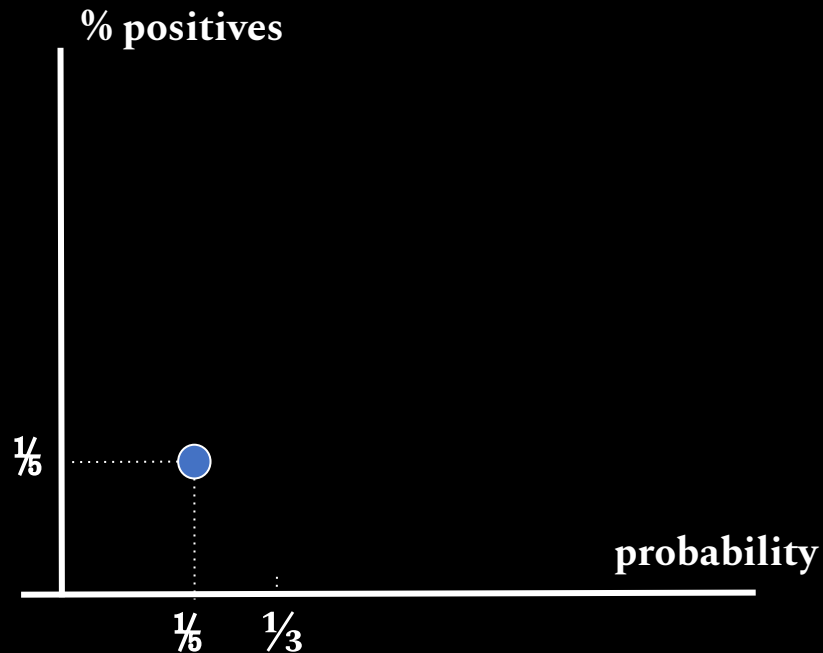
1. Calibration: what, when and why?

p	y
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	0
$\frac{1}{6}$	1



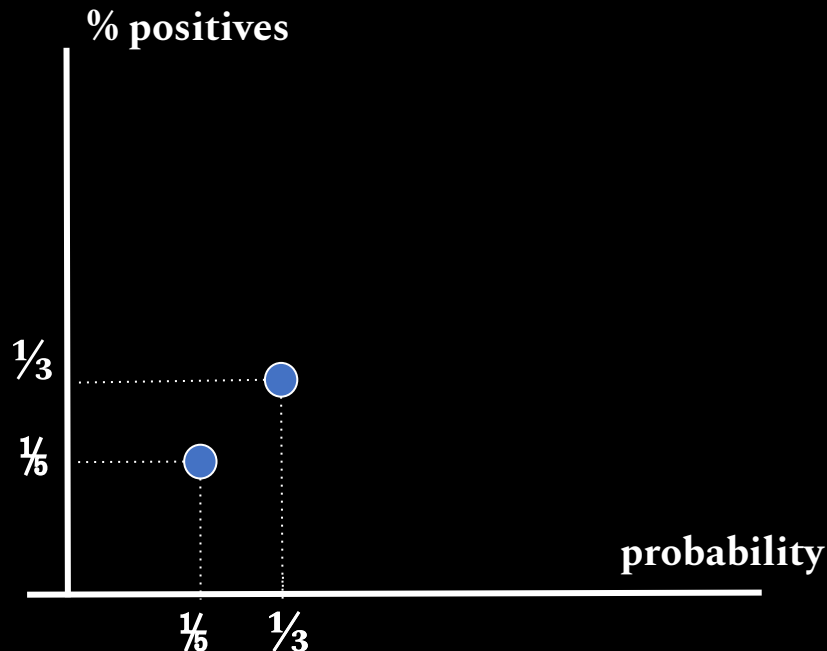
1. Calibration: what, when and why?

p	y
$\frac{1}{3}$	0
$\frac{1}{3}$	0
$\frac{1}{3}$	1



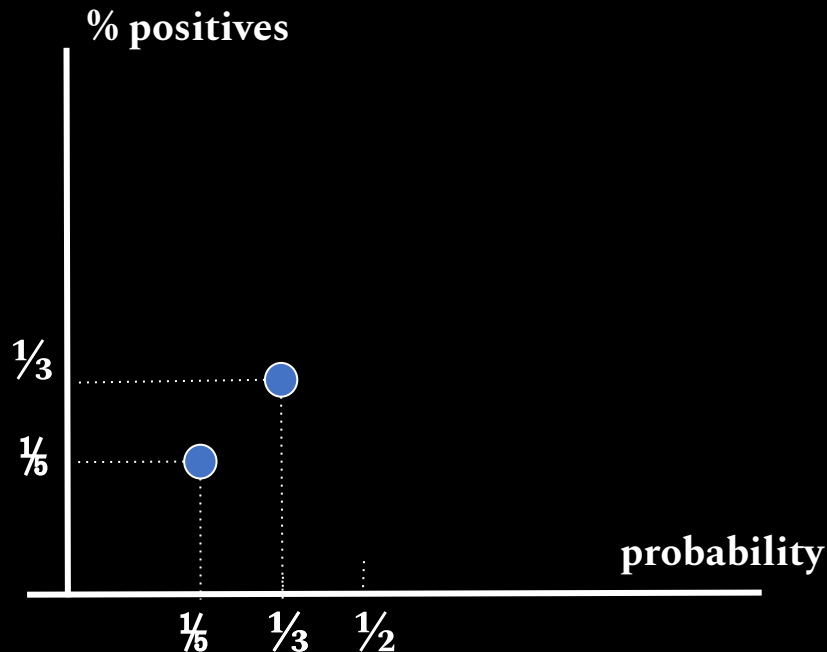
1. Calibration: what, when and why?

p	y
$\frac{1}{3}$	0
$\frac{1}{3}$	0
$\frac{1}{3}$	1



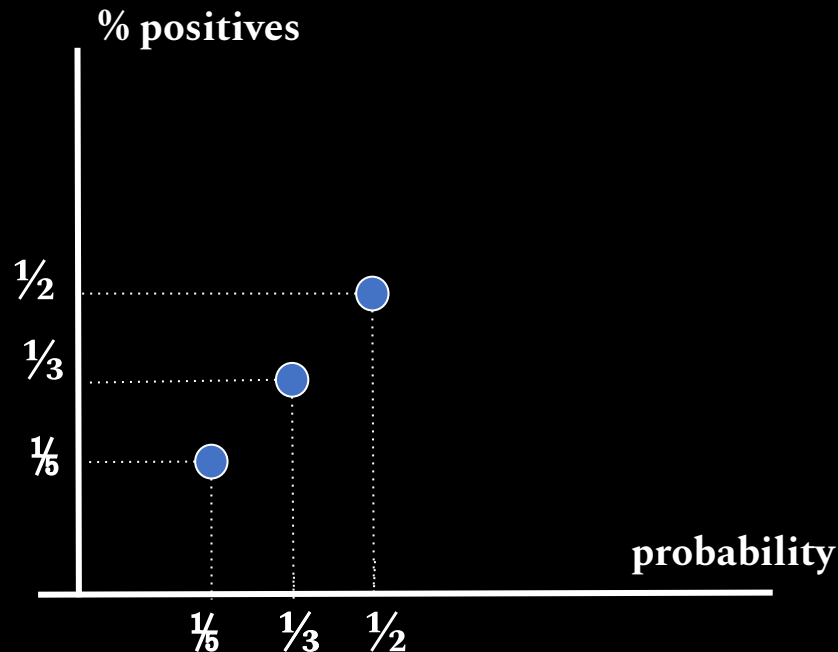
1. Calibration: what, when and why?

p	y
$\frac{1}{3}$	0
$\frac{1}{3}$	0
$\frac{1}{3}$	1
$\frac{1}{2}$	0
$\frac{1}{2}$	1



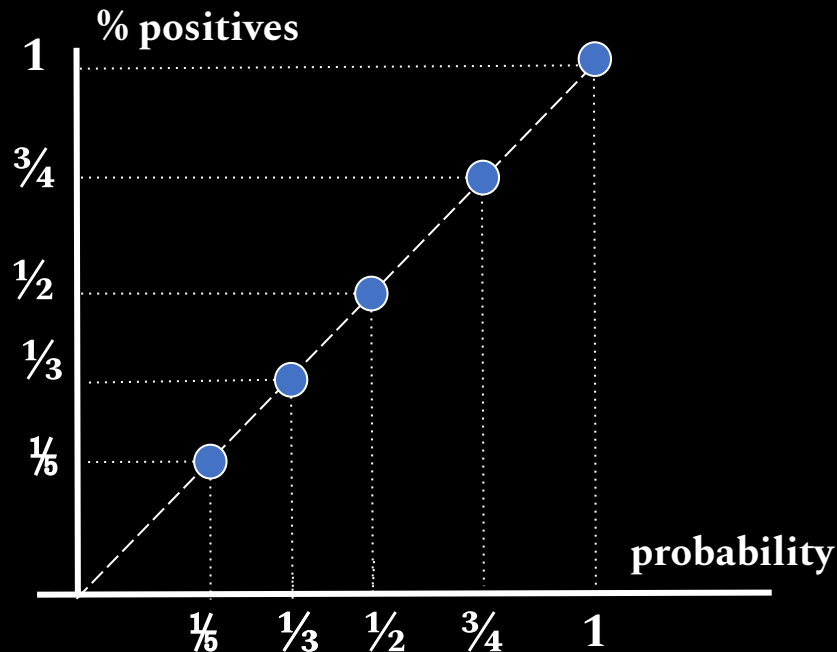
1. Calibration: what, when and why?

p	y
$\frac{1}{3}$	0
$\frac{1}{3}$	0
$\frac{1}{3}$	1
$\frac{1}{2}$	0
$\frac{1}{2}$	1



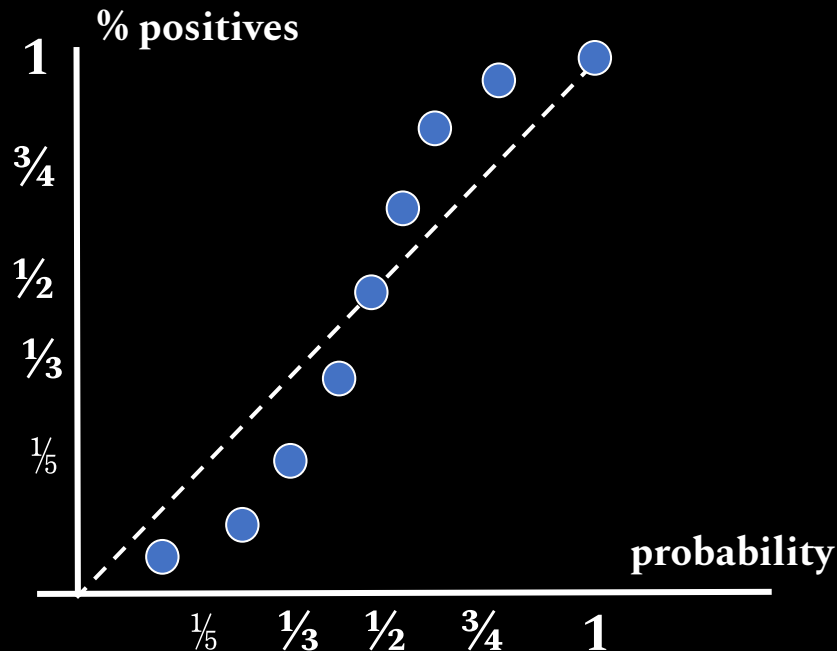
1. Calibration: what, when and why?

p	y
$\frac{3}{4}$	0
$\frac{3}{4}$	1
$\frac{3}{4}$	1
$\frac{3}{4}$	1
1	1



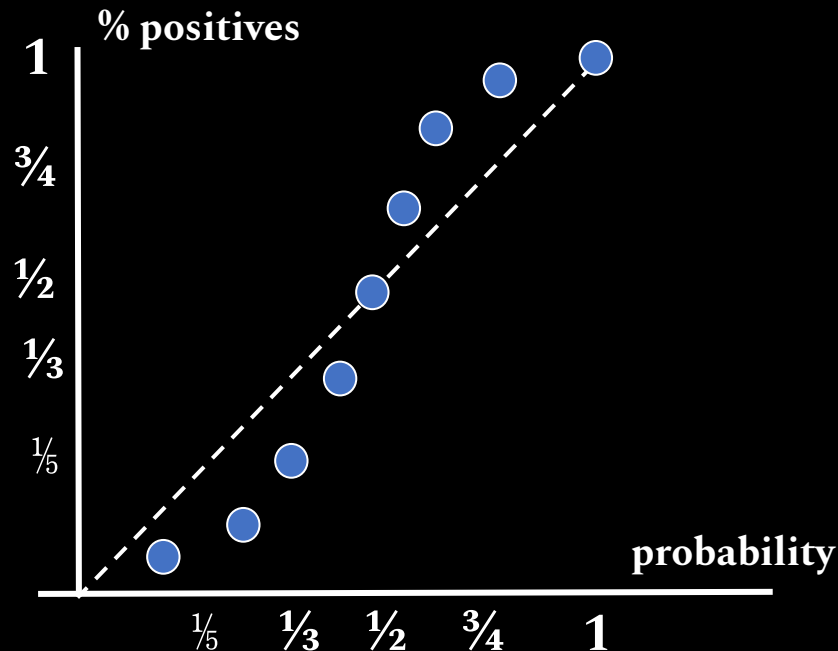
1. Calibration: what, when and why?

QUESTION:
Are these predictions
under-confident
or
over-confident?



1. Calibration: what, when and why?

QUESTION:
Are these predictions
under-confident
or
over-confident?



2. Visualizing & Measuring Calibration

- **Reliability Plots**

Not enough items with a given confidence to estimate population statistics decently:

model predicts with $p=0.2 \rightarrow 20\%$ positives”

What if you only have 2 items predicted with $p=0.2$?

We can group predictions in bins, and **plot them against $y=x$** .

- **Expected Calibration Error**

The average of gaps across bins, weighted by bin population:

$$\text{ECE} = \frac{1}{M} \sum_{i=1}^M \frac{1}{|B_i|} |\text{prob}(B_i) - \text{pos}(B_i)|$$

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	2

p	y
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	3

p	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	1

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	2

p	y
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	3

p	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	1

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	2

p	y
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	3

p	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	2

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	2

p	y
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	3

p	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	1

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	2

p	y
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	3

p	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	1

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	2

p	y
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	3

p	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	2

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	2

p	y
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	3

p	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	1

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	(\hat{y} , c)	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	(1, $\frac{2}{3}$)	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	(1, $\frac{2}{3}$)	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	(1, $\frac{2}{3}$)	2

p	(\hat{y} , c)	y
$[0, \frac{2}{3}, \frac{1}{3}]$	(2, $\frac{2}{3}$)	2
$[0, \frac{2}{3}, \frac{1}{3}]$	(2, $\frac{2}{3}$)	2
$[0, \frac{2}{3}, \frac{1}{3}]$	(2, $\frac{2}{3}$)	3

p	(\hat{y} , c)	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	(3, $\frac{2}{3}$)	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	(3, $\frac{2}{3}$)	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	(3, $\frac{2}{3}$)	1

2. Visualizing & Measuring Calibration

- Generalizing from Binary to Multi-Class classifiers

Full-calibration: consider the whole probability vector.

Class-wise calibration: only consider marginal probabilities.

Confidence calibration: only consider highest probability.

p	(\hat{y} , c)	y
$[\frac{2}{3}, \frac{1}{3}, 0]$	$(1, \frac{2}{3})$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	$(1, \frac{2}{3})$	1
$[\frac{2}{3}, \frac{1}{3}, 0]$	$(1, \frac{2}{3})$	2

p	(\hat{y} , c)	y
$[0, \frac{2}{3}, \frac{1}{3}]$	$(2, \frac{2}{3})$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	$(2, \frac{2}{3})$	2
$[0, \frac{2}{3}, \frac{1}{3}]$	$(2, \frac{2}{3})$	3

p	(\hat{y} , c)	y
$[\frac{1}{3}, 0, \frac{2}{3}]$	$(3, \frac{2}{3})$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	$(3, \frac{2}{3})$	3
$[\frac{1}{3}, 0, \frac{2}{3}]$	$(3, \frac{2}{3})$	2

2. Visualizing & Measuring Calibration

- Expected Calibration Error (binary)

$$\text{bin-ECE} = \frac{1}{M} \sum_{i=1}^M \frac{1}{|B_i|} |\text{prob}(B_i) - \text{pos}(B_i)|$$

- Expected Class-Wise Calibration Error (multi-class)

$$\text{cw-ECE} = \frac{1}{K} \sum_{k=1}^K \text{bin-ECE}_k \quad [\text{one-vs-rest}]$$

- Expected Confidence-Calibration Error (multi-class)

$$\text{conf-ECE} = \frac{1}{M} \sum_{i=1}^M \frac{1}{|B_i|} |\text{conf}(B_i) - \text{acc}(B_i)|$$

2. Visualizing & Measuring Calibration

- Alternative Calibration Measures: **Proper Scoring Rules**

p	y
$(\frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon, \frac{1}{3}-\epsilon)$	1
$(\frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon)$	1
$(\frac{1}{3}-\epsilon, \frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon)$	2

p	y
$(\frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon)$	2
$(\frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon)$	3
$(\frac{1}{3}-\epsilon, \frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon)$	3

2. Visualizing & Measuring Calibration

- Alternative Calibration Measures: **Proper Scoring Rules**

p	y
$(\frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon, \frac{1}{3}-\epsilon)$	1
$(\frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon)$	1
$(\frac{1}{3}-\epsilon, \frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon)$	2

p	y
$(\frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon)$	2
$(\frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon, \frac{1}{3}-\epsilon)$	3
$(\frac{1}{3}-\epsilon, \frac{1}{3}-\epsilon, \frac{1}{3}+2\epsilon)$	3

This classifier predicts a random class with full uncertainty. It always has a confidence of $\frac{1}{3}$, and it has an accuracy of $\frac{1}{3}$. Therefore it is **perfectly conf-calibrated**, but **useless**.

2. Visualizing & Measuring Calibration

- Alternative Calibration Measures: **Proper Scoring Rules**

p	y
$(\frac{2}{3}, 0, \frac{1}{3})$	1
$(0, \frac{1}{3}, \frac{2}{3})$	1
$(\frac{1}{3}, \frac{2}{3}, 0)$	2

p	y
$(0, \frac{1}{3}, \frac{2}{3})$	2
$(\frac{1}{3}, 0, \frac{2}{3})$	3
$(0, \frac{1}{3}, \frac{2}{3})$	3

2. Visualizing & Measuring Calibration

- Alternative Calibration Measures: **Proper Scoring Rules**

p	y
$(\frac{2}{3}, 0, \frac{1}{3})$	1
$(0, \frac{1}{3}, \frac{2}{3})$	1
$(\frac{1}{3}, \frac{2}{3}, 0)$	2

p	y
$(0, \frac{1}{3}, \frac{2}{3})$	2
$(\frac{1}{3}, 0, \frac{2}{3})$	3
$(0, \frac{1}{3}, \frac{2}{3})$	3

This classifier always predicts with $\sim \frac{2}{3}$ confidence. Also, it has an accuracy of $\frac{2}{3}$. It is **perfectly conf-calibrated**, but it has more **discrimination ability** than random guessing.

2. Visualizing & Measuring Calibration

- Alternative Calibration Measures: **Proper Scoring Rules**

p	y
(1, 0, 0)	1
(1, 0, 0)	1
(0, 1, 0)	2

p	y
(0, 1, 0)	2
(0, 0, 1)	3
(0, 0, 1)	3

2. Visualizing & Measuring Calibration

- Alternative Calibration Measures: **Proper Scoring Rules**

p	y
(1, 0, 0)	1
(1, 0, 0)	1
(0, 1, 0)	2

p	y
(0, 1, 0)	2
(0, 0, 1)	3
(0, 0, 1)	3

This is a god-like classifier. It is always 100% confident, and always right. It is totally **calibrated** and **perfectly discriminative**.

**PSRs are a tool for measuring
calibration & discrimination jointly.**

2. Visualizing & Measuring Calibration

- **Proper Scoring Rules**

Measure discrimination+calibration at individual item level

Most popular: Brier Score, Negative Log-Likelihood

$$\text{Brier}(\mathbf{p}, \mathbf{y}) = \|\mathbf{p} - \mathbf{y}\|_2^2$$

$$\text{NLL}(\mathbf{p}, \mathbf{y}) = -\log(\mathbf{p}_y)$$

Example: $\mathbf{y} = \mathbf{3}$, $\mathbf{y} = (0, 0, 1)$, $\mathbf{p} = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)$, $\mathbf{q} = \left(0, \frac{1}{3}, \frac{2}{3}\right)$

$$\text{Brier}(\mathbf{p}, \mathbf{y}) = 2/3 \quad \text{Brier}(\mathbf{q}, \mathbf{y}) = 2/9 \quad \text{Brier}(\mathbf{y}, \mathbf{y}) = 0$$

$$\text{NLL}(\mathbf{p}, \mathbf{y}) \approx 0.477 \quad \text{NLL}(\mathbf{q}, \mathbf{y}) \approx 0.176 \quad \text{NLL}(\mathbf{y}, \mathbf{y}) = 0$$

Note that a fully uncertain prediction (\mathbf{p}) does not score well.

3. Improving Calibration

• Post-Training Calibration

Classic methods: **Platt Scaling & Isotonic Regression**:

- Platt: Fits a logistic regression model using validation set.
- Isotonic: Fits a monotonic piecewise constant mapping, optimizing bins to maximize calibration.

Both designed for binary models

Temperature Scaling: Uses a validation set to learn a scalar T dividing logits before applying softmax and tempers their value:

$$p_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} \mapsto p_j = \frac{e^{(z_j/T)}}{\sum_{k=1}^N e^{(z_k/T)}}$$

We will see examples in the hands-on

3. Improving Calibration

- **Model Ensembling**

Ensembling several diverse models can improve calibration. Of course it comes with a computational overhead.

- **Training Time Calibration**

Over-parametrized NNs can keep on learning the training set until they are fully confident, minimizing NLL indefinitely. We can avoid this by **regularizing** so as to **disencourage confidence**.

Label Smoothing, MixUp, Focal Loss... Careful of **underfitting!**
Always report also a PSR, not only ECE.

We will see examples in the hands-on

4. Hands-On

<https://shorturl.at/hyFO3>