

# Model calibration

Alexander Lyzhov

# Why calibration?



# Overview

Classifier calibration definitions

Histogram regression estimators

Unbiased calibration estimation

Empirical results

Bonus: calibration methods

# How to define calibration?

$$\mathbb{P}[Y = \arg \max_y g_y(X) \mid \max_y g_y(X)] = \max_y g_y(X)^1 \quad \text{conf}$$

$$\mathbb{P}[Y = y \mid g_y(X)] = g_y(X) \quad \text{marginal}$$

$$\mathbb{P}[Y = y \mid g(X)] = g_y(X) \quad \text{joint}$$

$$\text{conf} \not\Rightarrow \text{marginal}$$

$$\text{conf} + \text{marginal} \not\Rightarrow \text{joint}$$

---

<sup>1</sup>Nixon, Jeremy, et al. "Measuring calibration in deep learning." arXiv preprint arXiv:1904.01685 (2019).

$$\mathbb{P}[Y = \arg \max_y g_y(X) \mid \max_y g_y(X)] = \max_y g_y(X) \quad \text{conf}$$

$$\mathbb{P}[Y = y \mid g_y(X)] = g_y(X) \quad \text{marginal}$$

conf  $\not\Rightarrow$  marginal

$g(X)$	$\mathbb{P}[Y \in \cdot \mid g(X)]$
( <u>0.6</u> , 0.1, 0.3)	( <u>0.7</u> , 0.2, 0.1)
(0.4, <u>0.6</u> , 0.0)	(0.3, <u>0.5</u> , 0.2)

$$\mathbb{P}[Y = \arg \max_y g_y(X) \mid \max_y g_y(X)] = \max_y g_y(X) \quad \text{conf}$$

$$\mathbb{P}[Y = y \mid g_y(X)] = g_y(X) \quad \text{marginal}$$

$$\mathbb{P}[Y = y \mid g(X)] = g_y(X) \quad \text{joint}$$

$$\text{conf} + \text{marginal} \not\Rightarrow \text{joint}$$

$g(X)$	$\mathbb{P}[Y \in \cdot \mid g(X)]$
( <u>0.1</u> , 0.3, 0.6)	( <u>0.2</u> , 0.2, 0.6)
( <u>0.1</u> , 0.6, 0.3)	( <u>0.0</u> , 0.7, 0.3)
(0.3, 0.1, 0.6)	(0.2, 0.2, 0.6)
(0.3, 0.6, 0.1)	(0.4, 0.5, 0.1)
(0.6, 0.1, 0.3)	(0.7, 0.0, 0.3)
(0.6, 0.3, 0.1)	(0.5, 0.4, 0.1)

$$\mathbb{P}[Y = y \mid g(X)] = g_y(X)$$

joint

joint  $\Rightarrow$  high accuracy

$g(X)$	$\mathbb{P}[Y \in \cdot \mid g(X)]$
(0.5, 0.5)	(0.5, 0.5)

	$g(x)$	$y$
$x_1$	(0.5, 0.5)	(1, 0)
$x_2$	(0.5, 0.5)	(0, 1)

# Calibration errors

$$r(\xi) := (\mathbb{P}[Y = 1 \mid g(X) = \xi], \dots, \mathbb{P}[Y = m \mid g(X) = \xi])$$

**Joint** calibration  $\iff r(\xi) = \xi$

$$\text{ECE} = \mathbb{E}[d(r(g(X)), g(X))]$$

$$\text{MCE} = \sup[d(r(g(X)), g(X))]$$

**Joint** calibration  $\iff \text{ECE}=0, \text{MCE}=0$

$r(\xi)$  is unknown. *How to estimate?*



# Histogram regression

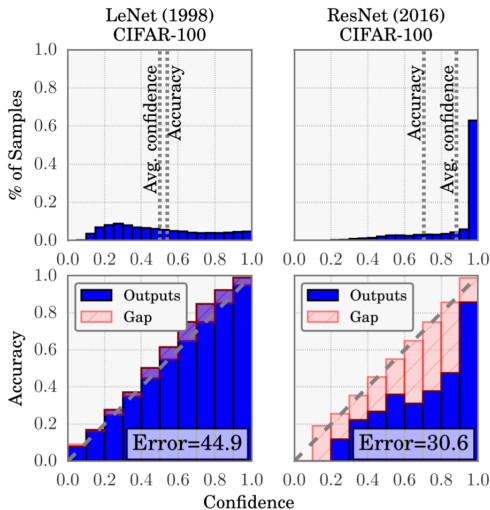
Partition simplex and let  $b[\xi]$  be a partition corresponding to  $\xi$ .

Approximate  $r_y(\xi) = \mathbb{P}(Y = y \mid g(X) = \xi)$  with:

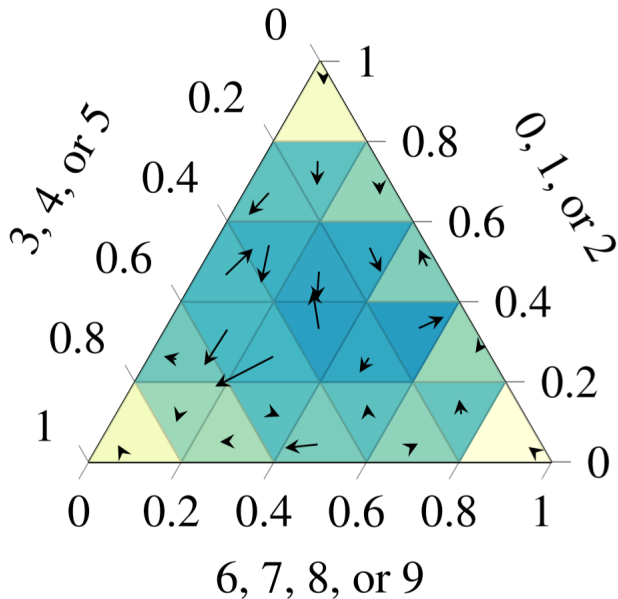
$$\hat{r}_y(\xi) := |\{i : g(x_i) \in b[\xi] \wedge y_i = y\}| / |\{i : g(x_i) \in b[\xi]\}|$$

Approximate  $\text{ECE} = \mathbb{E}[d(r(g(X)), g(X))]$  with:

$$\widehat{\text{ECE}} = \sum_b \frac{n_b}{N} d\left(\frac{1}{n_b} \sum_b \hat{r}_y(g(x)), \frac{1}{n_b} \sum_b g(x)\right)$$



<sup>1</sup>Guo, Chuan, et al. "On calibration of modern neural networks." Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017.



---

<sup>1</sup>Nixon, Jeremy, et al. "Measuring calibration in deep learning." arXiv preprint arXiv:1904.01685 (2019).

# Histogram regression pathologies

$$\widehat{\text{ECE}} = \sum_b \frac{n_b}{N} d \left( \frac{1}{n_b} \sum_b \hat{r}_y(g(x)), \frac{1}{n_b} \sum_b g(x) \right)$$

$$\lim_{\max_b \sup_{p,q \in b} \|p-q\|_2 \rightarrow 0} \lim_{n \rightarrow \infty} \widehat{\text{ECE}} = \text{ECE}$$

$$\lim_{n \rightarrow \infty} \widehat{\text{ECE}} \leq \text{ECE}^1$$

1. Biases for different models can be different
2.  $\widehat{\text{ECE}}$  is random, distribution is not taken into account
3. No common unit or scale
4. Depends a lot on binning scheme
5. Number of bins scales exponentially with number of classes

---

<sup>1</sup>Nixon, Jeremy, et al. "Measuring calibration in deep learning." arXiv preprint arXiv:1904.01685 (2019).

# Histogram regression pathologies

Example of  $\widehat{\text{ECE}}$  bias <sup>1</sup>:

$$\mathbb{P}(Y = 1) = \mathbb{P}(Y = 2) = 1/2$$

$$\mathbb{P}(Y = 0 \mid X = x) \in \{0, 1\} \quad \forall x$$

$$g_{\text{opt}}(x) = \mathbb{P}(Y), g_{\text{const}}(x) \equiv (1/2, 1/2)$$

$$\text{ECE}(g_{\text{opt}}) = 0, \text{ECE}(g_{\text{const}}) = 0$$

One bin, single data point.

$$\mathbb{E}(\widehat{\text{ECE}}(g_{\text{opt}})) = 0, \mathbb{E}(\widehat{\text{ECE}}(g_{\text{const}})) = 1/2$$

---

<sup>1</sup>Nixon, Jeremy, et al. "Measuring calibration in deep learning." arXiv preprint arXiv:1904.01685 (2019).

# Unbiased calibration estimation

Instead of  $\text{ECE} = \mathbb{E}[d(r(g(X)), g(X))]$ , define error as:  
 $\text{CE}[\mathcal{F}, g] := \sup_{f \in \mathcal{F}} \mathbb{E}[(r(g(X)) - g(X))^\top f(g(X))]^1$

**joint**  $\implies \text{CE} = 0$

$\text{CE} = 0 \not\Rightarrow$  **joint**

$\text{CE} = 0, \mathcal{F} = \mathcal{C}(\Delta^m \rightarrow \mathbb{R}^m) \implies$  **joint**

---

<sup>1</sup>Widmann, David, Fredrik Lindsten, and Dave Zachariah.

"Calibration tests in multi-class classification: A unifying framework."  
arXiv preprint arXiv:1910.11385 (2019).

# Unbiased calibration estimation

RKHS  $\mathcal{H} \iff$  kernel

Kernel  $k : X \times X \rightarrow \mathbb{R}$ ,

$$k(s, t) = k(t, s), \sum_{i,j=1}^n u_i u_j k(t_j, t_j) \geq 0$$

Hilbert space  $\{f : X \rightarrow \mathbb{R}\}$  is RKHS if  $E_t f = f(t)$  is continuous.

Reproducing property:  $f(t) = \langle f, k(\cdot, t) \rangle_{\mathcal{H}}$

Matrix-valued kernel  $k : \Delta^m \times \Delta^m \rightarrow \mathbb{R}^{m \times m}$ ,

$$k(s, t) = k(t, s)^T, \sum_{i,j=1}^n u_i^T k(t_i, t_j) u_j \geq 0$$

Hilbert space  $\{f : \Delta^m \rightarrow \mathbb{R}^m\}$  is RKHS if  $E_{t,u} f = \langle u, f(t) \rangle_{\mathbb{R}^m}$  is continuous.

Reproducing property:  $\langle u, f(t) \rangle_{\mathbb{R}^m} = \langle k(\cdot, t)u, f \rangle_{\mathcal{H}}$

# Unbiased calibration estimation

Let  $\mathcal{F}$  be unit ball in RKHS of  $k$ .

Kernel calibration error:

$$\text{KCE}[k, g] := \text{CE}[\mathcal{F}, g] = \sup_{f \in \mathcal{F}} \mathbb{E} \left[ (r(g(X)) - g(X))^{\top} f(g(X)) \right]$$

$$\text{KCE} = 0 \iff \text{joint}$$

$$\text{KCE} = \left( \mathbb{E} \left[ (e_Y - g(X))^{\top} k(g(X), g(X')) (e_{Y'} - g(X')) \right] \right)^{\frac{1}{2}}$$

No  $r(g(X))$  now!



# Unbiased calibration estimation

We had  $\text{CE}[\mathcal{F}, g] := \sup_{f \in \mathcal{F}} \mathbb{E} \left[ (r(g(X)) - g(X))^{\top} f(g(X)) \right]$

$\text{KCE}[k, g] := \text{CE}[\mathcal{F}, g]$

$\text{KCE} = \left( \mathbb{E} \left[ (e_Y - g(X))^{\top} k(g(X), g(X')) (e_{Y'} - g(X')) \right] \right)^{\frac{1}{2}}$

$h_{i,j} := (e_{Y_i} - g(X_i))^{\top} k(g(X_i), g(X_j)) (e_{Y_j} - g(X_j))$

Notation	Definition	Properties	Complexity
$\widehat{\text{SKCE}}_{\text{b}}$	$n^{-2} \sum_{i,j=1}^n h_{i,j}$	biased	$O(n^2)$
$\widehat{\text{SKCE}}_{\text{uq}}$	$\binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} h_{i,j}$	unbiased	$O(n^2)$
$\widehat{\text{SKCE}}_{\text{ul}}$	$\lfloor n/2 \rfloor^{-1} \sum_{i=1}^{\lfloor n/2 \rfloor} h_{2i-1, 2i}$	unbiased	$O(n)$

Potentially trainable.

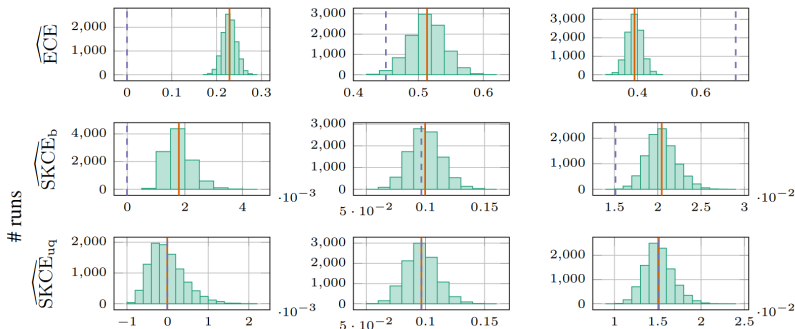
# Results

$$k(x, y) = \exp(-\|x - y\|/\nu)$$

M1

M2

M3



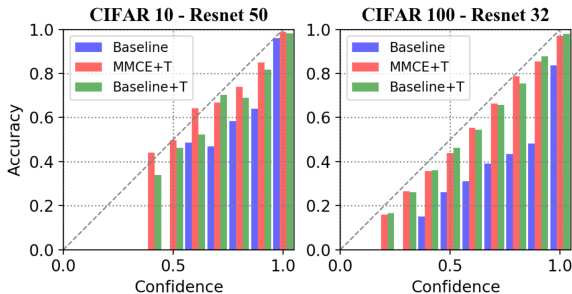
# Results

$$\text{conf}_i = \max g(y_i | x_i)$$

$$\text{correct}_i = \mathbb{I} [\hat{y}_i = \arg \max g(y_i | x_i)]$$

$$\text{SKCE}_\theta(B) = \sum_{i,j \in B} \frac{(\text{conf}_i - \text{correct}_i)(\text{conf}_j - \text{correct}_j)k(r_i, r_j)}{|B|^2} \quad 1$$

$$\min_\theta \sum_{(x_i, y_i) \in B} \log g_\theta(y_i | x_i) + \lambda (\text{SKCE}_\theta(B))^{\frac{1}{2}}$$



---

<sup>1</sup>Kumar, Aviral, Sunita Sarawagi, and Ujjwal Jain. "Trainable calibration measures for neural networks from kernel mean embeddings." International Conference on Machine Learning. 2018.

# Bonus: calibration methods

## 1. Non-parametric methods

Histogram binning<sup>1</sup>:

$$\min_{\theta_1, \dots, \theta_M} \sum_{m=1}^M \sum_{i=1}^n \mathbf{1}(a_m \leq \hat{p}_i < a_{m+1}) (\theta_m - y_i)^2$$

Isotonic regression<sup>2</sup>:

$$\min_{\theta, a} \sum_{m=1}^M \sum_{i=1}^n \mathbf{1}(a_m \leq \hat{p}_i < a_{m+1}) (\theta_m - y_i)^2$$
$$0 = a_1 \leq a_2 \leq \dots \leq a_{M+1} = 1, \theta_1 \leq \theta_2 \leq \dots \leq \theta_M$$

---

<sup>1</sup>Zadrozny, Bianca, and Charles Elkan. "Obtaining calibrated probability estimates from decision trees and naive Bayesian classifiers." *ICML*. Vol. 1. 2001.

<sup>2</sup>Zadrozny, Bianca, and Charles Elkan. "Transforming classifier scores into accurate multiclass probability estimates." *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*. 2002.

# Bonus: calibration methods

## 2. Parametric methods

Platt scaling:  $\text{Softmax}(Wz + b)^1$

Temperature scaling:  $\text{Softmax}(z/T)$

+ of temperature scaling:

Easy and very effective for **conf**

Doesn't affect accuracy

- of temperature scaling:

Requires validation set

Changes confidence uniformly

Doesn't work well for **marginal**

---

<sup>1</sup>Platt, John. "Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods." Advances in large margin classifiers 10.3 (1999): 61-74.

## Bonus: calibration methods

### 3. Calibration during training

Optimize kernel calibration error

Regularize entropy of predictive distribution<sup>1</sup>

Use focal loss instead of NLL:  $FL(p) = -(1-g_{y_i}(x))^{\gamma} \log g_{y_i}(x)$ <sup>2</sup>

---

<sup>1</sup>Pereyra, Gabriel, et al. "Regularizing neural networks by penalizing confident output distributions." arXiv preprint arXiv:1701.06548 (2017).

<sup>2</sup>Mukhoti, Jishnu, et al. "Calibrating Deep Neural Networks using Focal Loss." arXiv preprint arXiv:2002.09437 (2020).

# Open questions

- ▶ How to optimize estimates of marginal/joint calibration while training?
- ▶ How to pick the kernel? Can we learn it?
- ▶ Does this help in safety-critical applications?
- ▶ Do existing recalibration techniques for NNs result in better marginal/joint calibration?
- ▶ How do NN hyperparameters influence kernel calibration error?