



UPPSALA  
UNIVERSITET

# Calibration of probabilistic predictive models

Machine Learning Journal Club, Gatsby Unit

David Widmann

Department of Information Technology, Uppsala University, Sweden  
Centre for Interdisciplinary Mathematics, Uppsala University, Sweden

28 March 2022

# About me

## TL;DR 📖

- ▶ 31 year old PhD student at Uppsala University
- ▶ On parental leave since September 2021
- ▶ Research on uncertainty quantification of probabilistic models
- ▶ Active member in the Julia community



# About me

## Education

2017—now: PhD student (Uppsala University)

2016—2017: MSc Mathematics (TU Munich)

2013—2016: BSc Mathematics (TU Munich)

2007—2013: Human medicine (LMU and TU Munich)

# About me

## Education

2017—now: PhD student (Uppsala University)

2016—2017: MSc Mathematics (TU Munich)

2013—2016: BSc Mathematics (TU Munich)

2007—2013: Human medicine (LMU and TU Munich)

## Research interests

- ▶ Research topic: "Uncertainty-aware deep learning"
- ▶ Statistics, probability theory, scientific machine learning, and computer science
- ▶ Julia programming, e.g., SciML and Turing

## Papers

- ▶ J. Vaicenavicius et al. “Evaluating model calibration in classification.” In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019
  - ▶ Focus on multi-class classification, calibration lenses, calibration estimation and tests with ECE

# Papers

- ▶ J. Vaicenavicius et al. “Evaluating model calibration in classification.” In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019
  - ▶ Focus on multi-class classification, calibration lenses, calibration estimation and tests with ECE
- ▶ D. Widmann, F. Lindsten, and D. Zachariah. “Calibration tests in multi-class classification: A unifying framework.” In: *Advances in Neural Information Processing Systems 32*. 2019
  - ▶ Calibration errors and tests for multi-class classification based on matrix-valued kernels

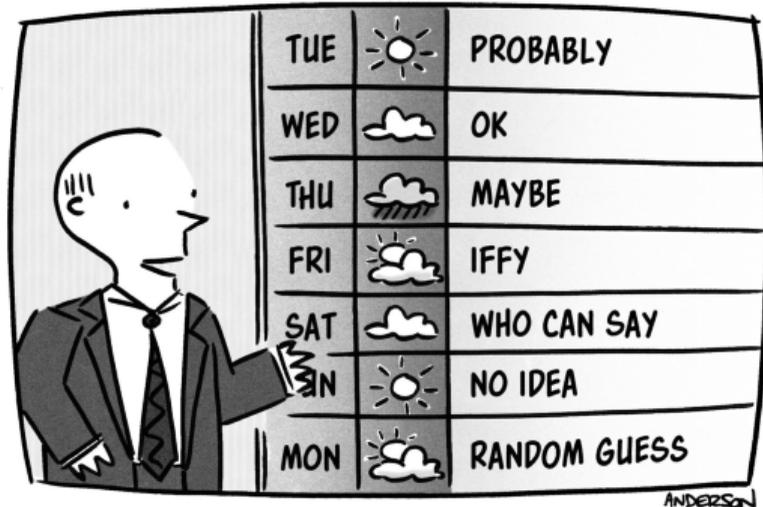
# Papers

- ▶ J. Vaicenavicius et al. “Evaluating model calibration in classification.” In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019
  - ▶ Focus on multi-class classification, calibration lenses, calibration estimation and tests with ECE
- ▶ D. Widmann, F. Lindsten, and D. Zachariah. “Calibration tests in multi-class classification: A unifying framework.” In: *Advances in Neural Information Processing Systems 32*. 2019
  - ▶ Calibration errors and tests for multi-class classification based on matrix-valued kernels
- ▶ D. Widmann, F. Lindsten, and D. Zachariah. “Calibration tests beyond classification.” In: *International Conference on Learning Representations*. 2021
  - ▶ Calibration errors and tests for probabilistic predictive models based on scalar-valued kernels

## Calibration: Motivation and definition

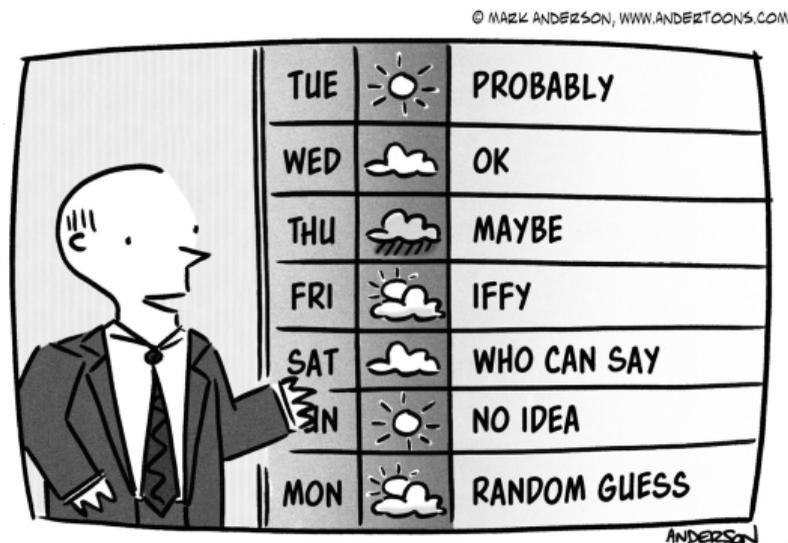
## Example: Weather forecasts

© MARK ANDERSON, WWW.ANDERTOONS.COM



"And now the 7-day forecast..."

## Example: Weather forecasts

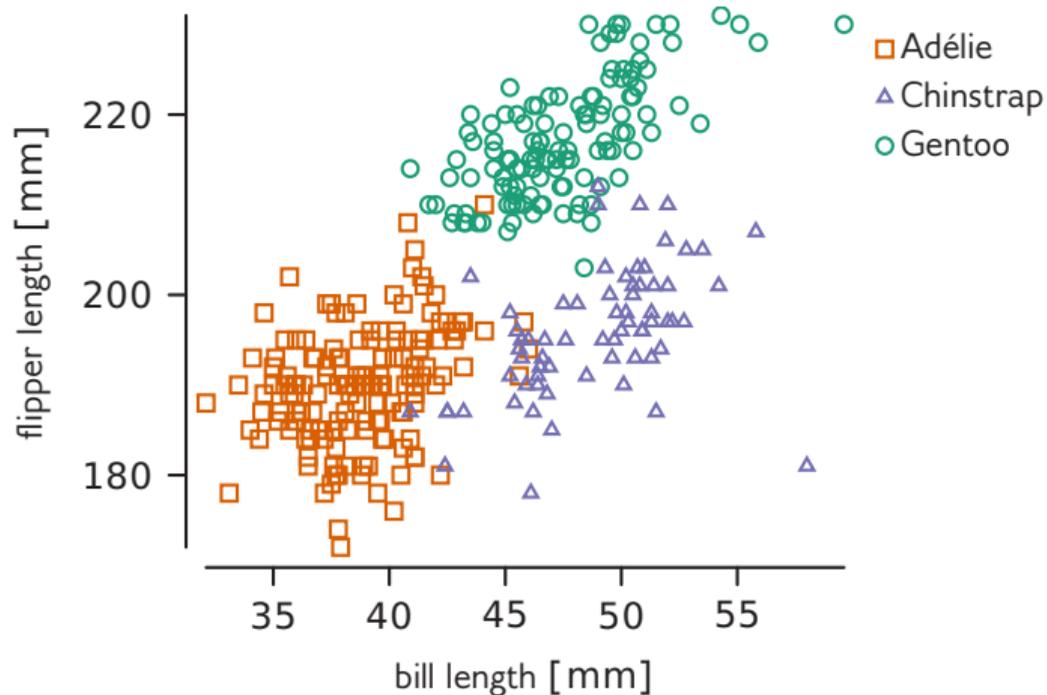


"And now the 7-day forecast..."

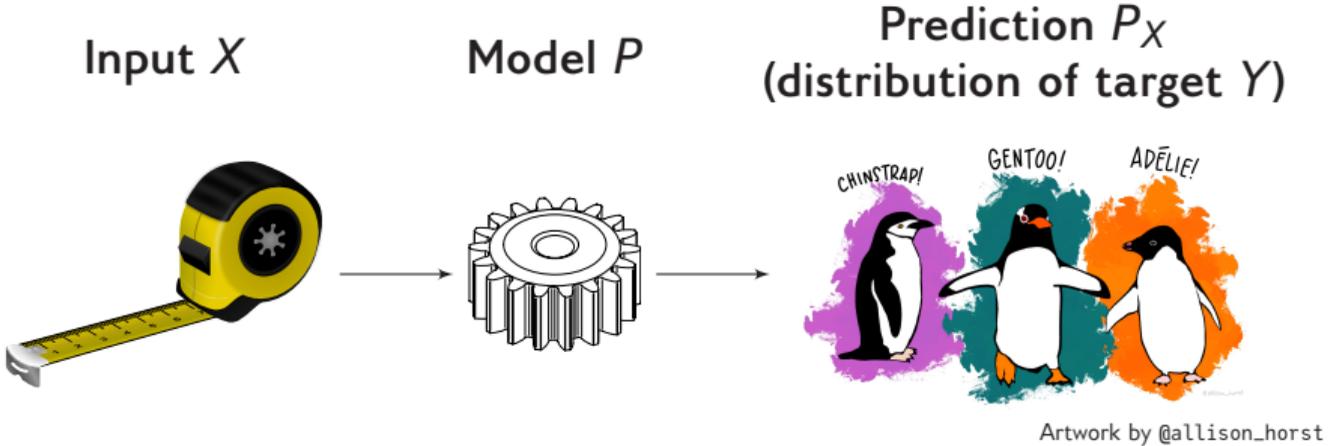
"Those forecasts which were marked 'doubtful' were the *best I could frame* under the circumstances. [...] If I make no distinction between these and others, I degrade the whole."

—E. Cooke

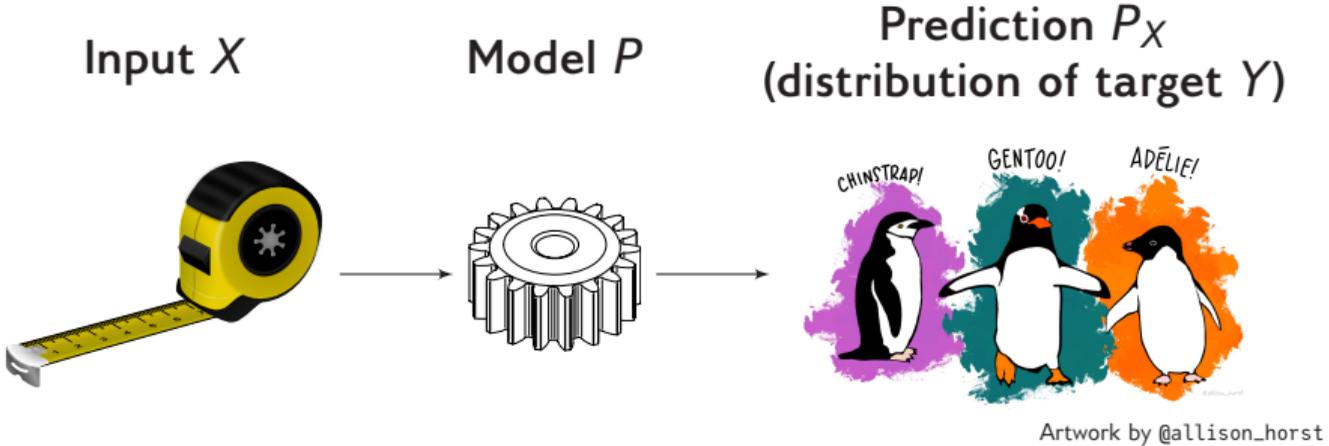
## Motivation: Classification example



# Motivation: Classification example



# Motivation: Classification example

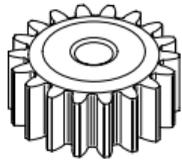


## Example: Prediction $P_X$

| Adélie | Chinstrap | Gentoo |
|--------|-----------|--------|
| 80%    | 10%       | 10%    |

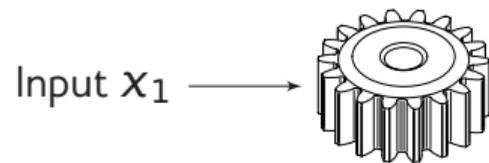
## Calibration: Intuition

**Model  $P$**

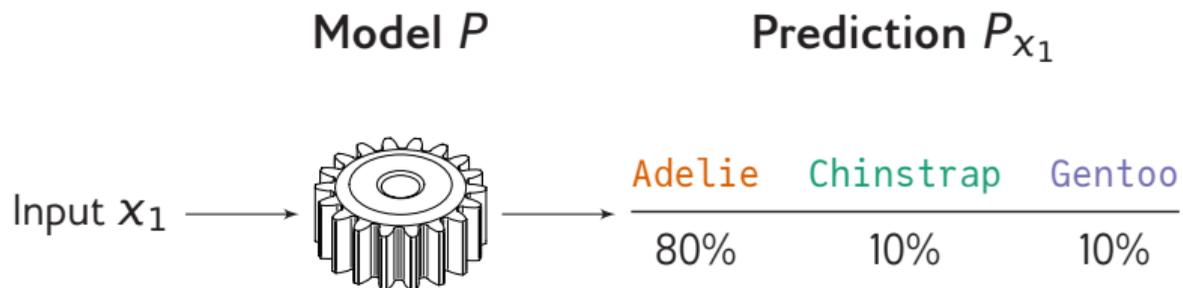


## Calibration: Intuition

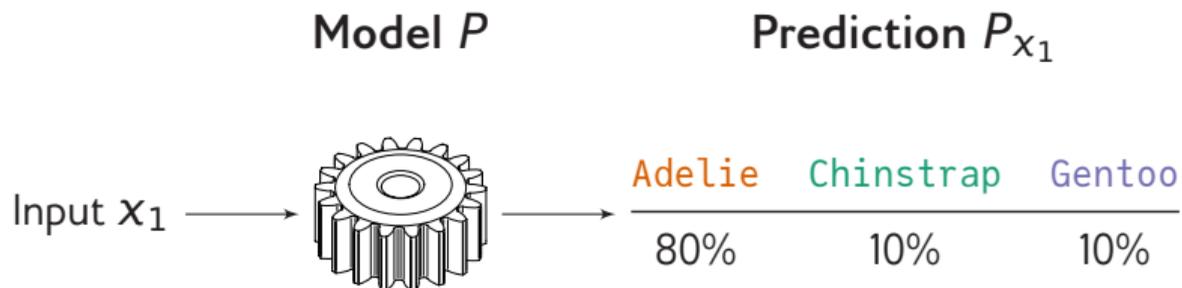
**Model  $P$**



## Calibration: Intuition



## Calibration: Intuition



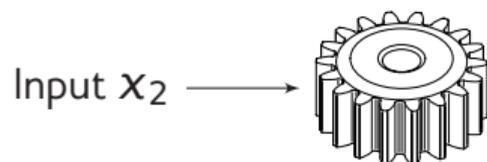
Empirical frequency

|        |           |        |
|--------|-----------|--------|
| Adelie | Chinstrap | Gentoo |
| <hr/>  |           |        |
|        |           |        |

The empirical frequency is shown as a single vertical bar below a horizontal line, indicating a uniform distribution across the three categories.

## Calibration: Intuition

Model  $P$



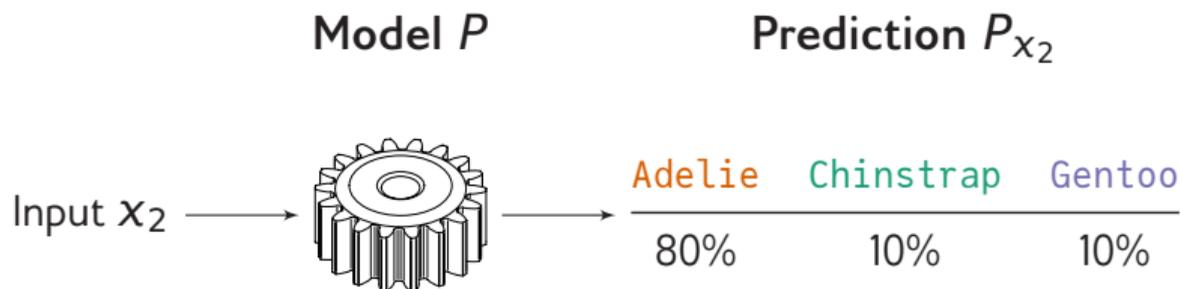
Empirical frequency

Adelie Chinstrap Gentoo

---

|

## Calibration: Intuition

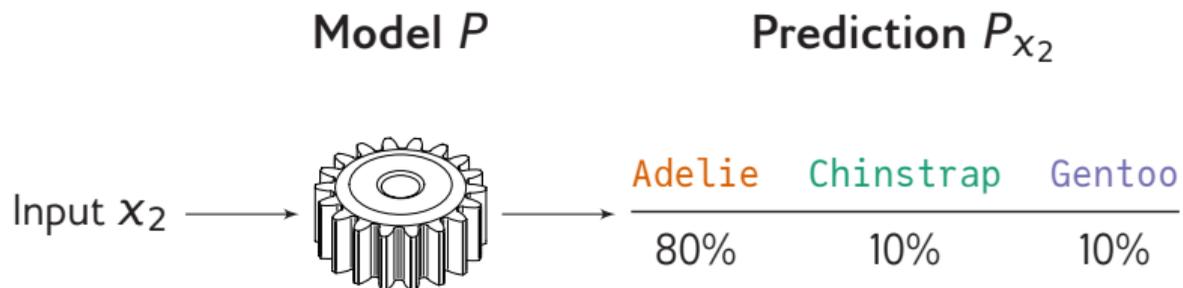


Empirical frequency

|        |           |        |
|--------|-----------|--------|
| Adelie | Chinstrap | Gentoo |
| <hr/>  |           |        |
|        |           |        |

The empirical frequency is shown as a single vertical bar below a horizontal line, indicating that the observed frequency for all categories is 1.

## Calibration: Intuition



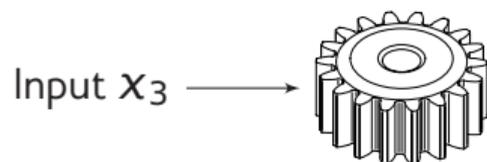
Empirical frequency

| Adelie | Chinstrap | Gentoo |
|--------|-----------|--------|
|        |           |        |

The empirical frequency table shows the observed counts for each category: Adelie (1), Chinstrap (1), and Gentoo (0).

## Calibration: Intuition

Model  $P$



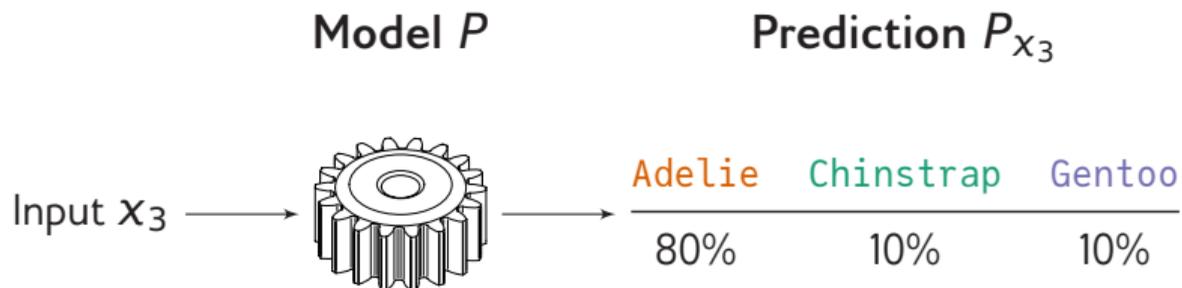
Empirical frequency

Adelie Chinstrap Gentoo

---

| |

## Calibration: Intuition

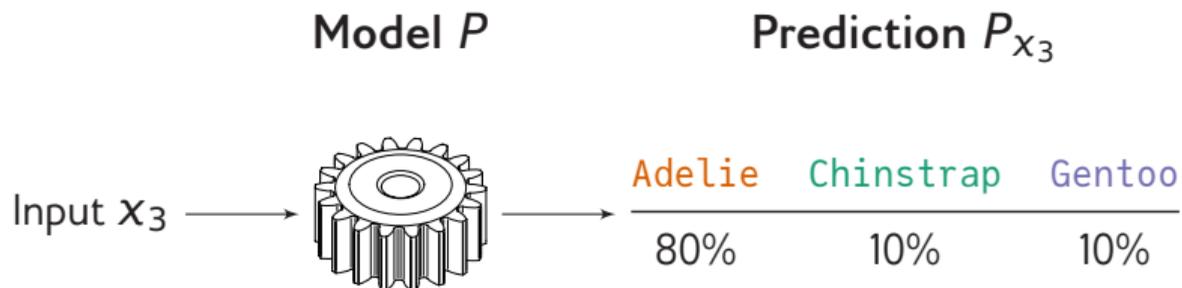


Empirical frequency

| Adelie | Chinstrap | Gentoo |
|--------|-----------|--------|
|        |           |        |

The empirical frequency table shows the observed counts for each category: Adelie (1), Chinstrap (1), and Gentoo (0).

## Calibration: Intuition

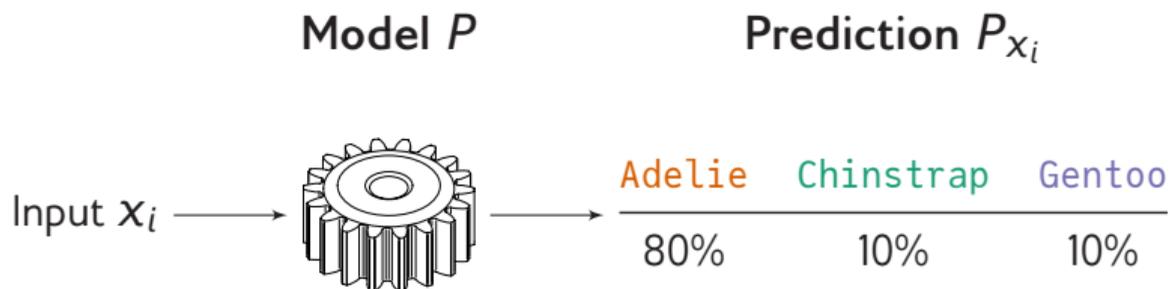


Empirical frequency

| Adelie | Chinstrap | Gentoo |
|--------|-----------|--------|
| //     | /         |        |

The empirical frequency is shown as a horizontal line with two vertical bars under 'Adelie' and one vertical bar under 'Chinstrap', representing the observed counts for each category.

## Calibration: Intuition



Empirical frequency

| Adelie | Chinstrap | Gentoo |
|--------|-----------|--------|
| ...    | ...       | ...    |

The empirical frequency is represented by a horizontal bar divided into three segments, with the number of ticks below each segment indicating the frequency of each class: Adelie (4 ticks), Chinstrap (2 ticks), and Gentoo (1 tick).

# Calibration

## Prediction $P_X$

| Adélie | Chinstrap | Gentoo |
|--------|-----------|--------|
| 80%    | 10%       | 10%    |

## Empirical frequency law( $Y | P_X$ )

| Adélie | Chinstrap | Gentoo |
|--------|-----------|--------|
| ...    | ...       | ...    |

---

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

# Calibration

Predictions consistent with empirically observed frequencies?

| Prediction $P_X$ |           |        |
|------------------|-----------|--------|
| Adélie           | Chinstrap | Gentoo |
| 80%              | 10%       | 10%    |

?

---

| Empirical frequency law( $Y   P_X$ ) |           |        |
|--------------------------------------|-----------|--------|
| Adélie                               | Chinstrap | Gentoo |
| ...                                  | ...       | ...    |

# Calibration

Predictions consistent with empirically observed frequencies?

| Prediction $P_X$ |           |        | ?            | Empirical frequency law $\text{law}(Y   P_X)$ |           |        |
|------------------|-----------|--------|--------------|---|-----------|--------|
| Adélie           | Chinstrap | Gentoo |              | Adélie  | Chinstrap | Gentoo |
| 80%              | 10%       | 10%    | <u>=====</u> | ...   | ...       | ...    |

## Definition

A probabilistic predictive model  $P$  is calibrated if

$$\text{law}(Y | P_X) = P_X \quad \text{almost surely.}$$

D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

# Calibration

Predictions consistent with empirically observed frequencies?

| Prediction $P_X$ |           |        | ?             | Empirical frequency law $\text{law}(Y   P_X)$ |               |        |
|------------------|-----------|--------|---------------|---|---------------|--------|
| Adélie           | Chinstrap | Gentoo |               | Adélie  | Chinstrap     | Gentoo |
| 80%              | 10%       | 10%    | <u>      </u> | <u>      </u>                                 | <u>      </u> |        |
|                  |           |        | ...           | ...   | ...           |        |

## Definition

A probabilistic predictive model  $P$  is calibrated if

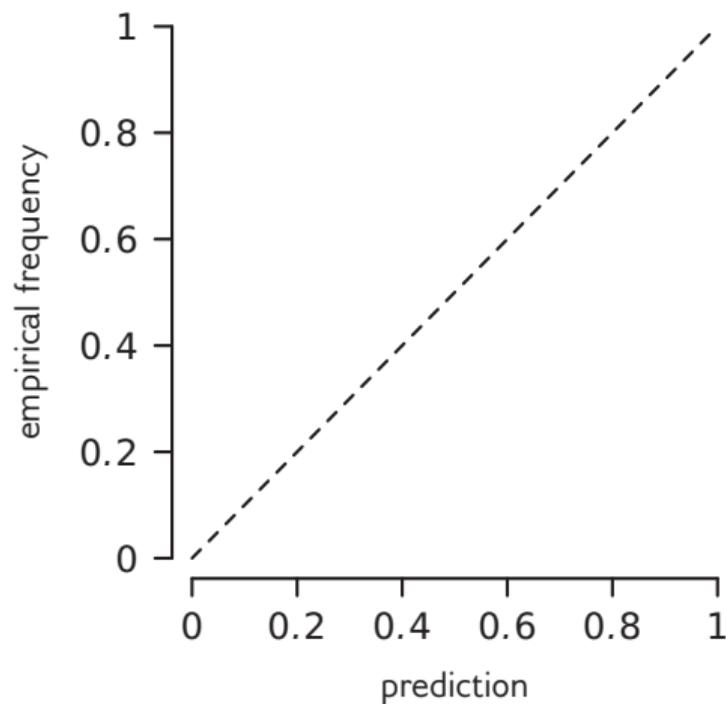
$$\text{law}(Y | P_X) = P_X \quad \text{almost surely.}$$

Notion captures also weaker confidence calibration

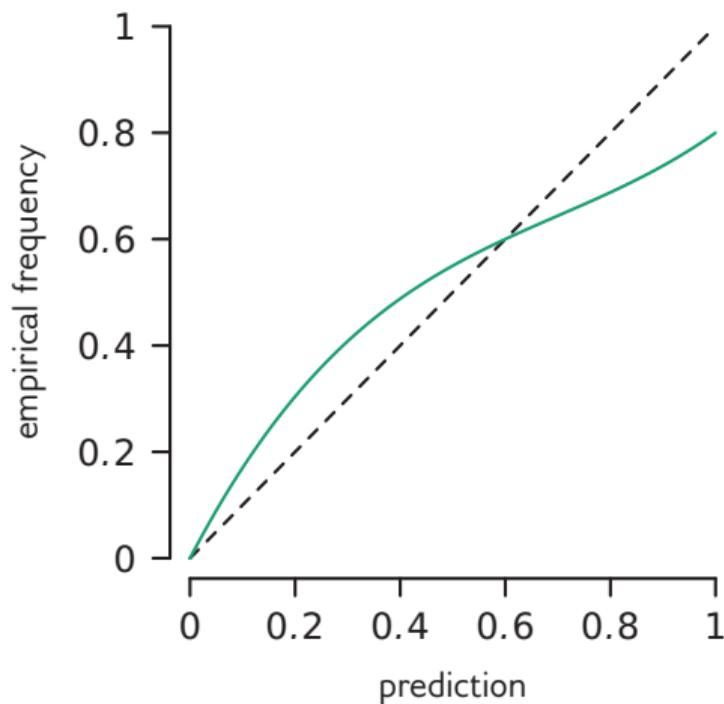
D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

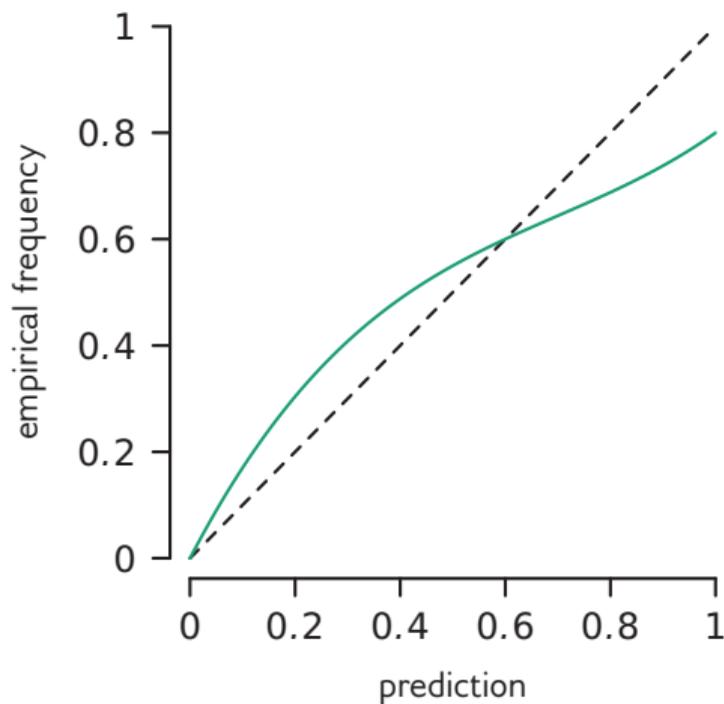
## Binary classification: Reliability diagram



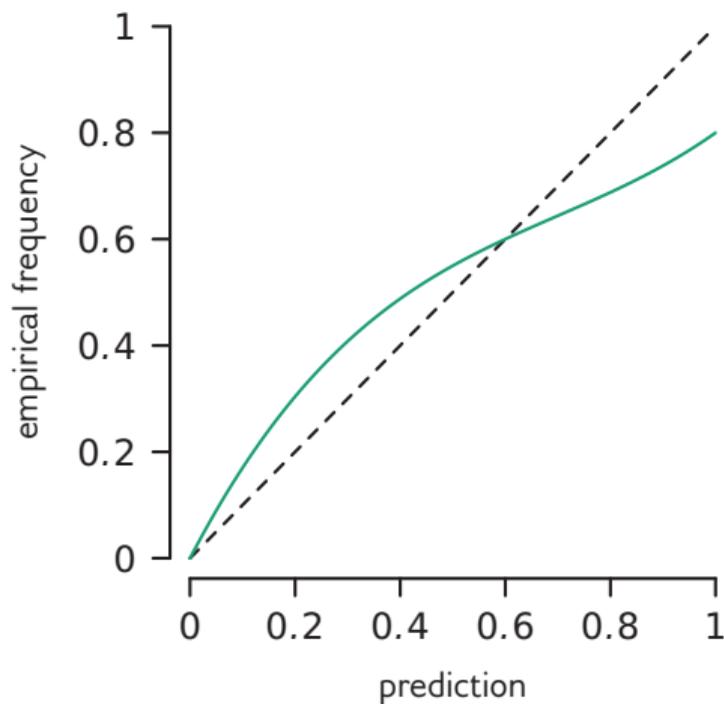
## Binary classification: Reliability diagram



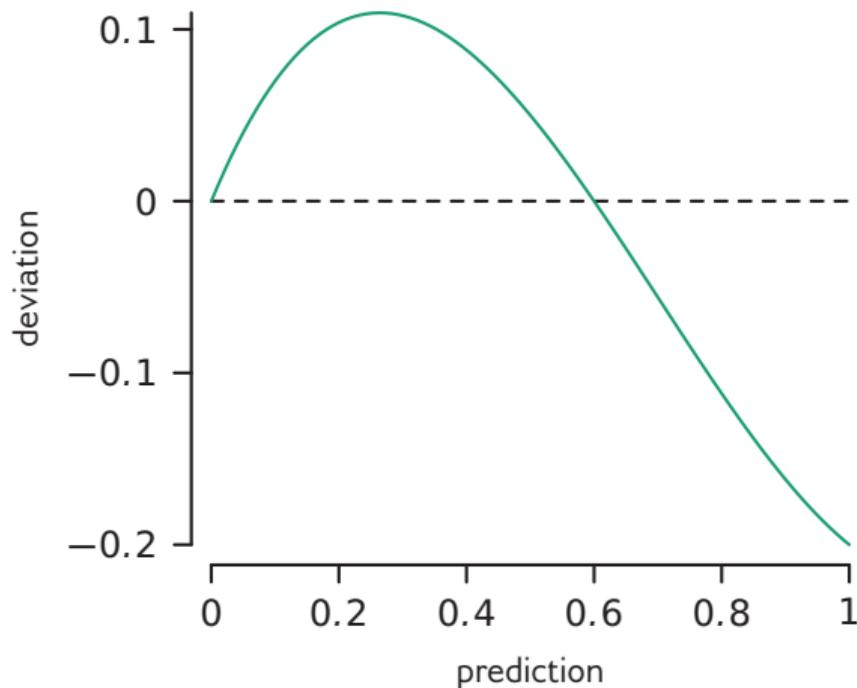
## Binary classification: Reliability diagram



## Binary classification: Reliability diagram



## Binary classification: Reliability diagram



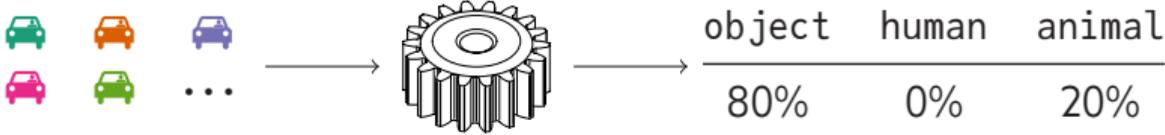
# Multi-class classification: All scores matter!



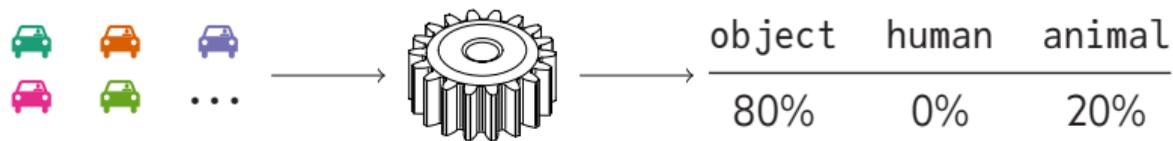
# Multi-class classification: All scores matter!



# Multi-class classification: All scores matter!

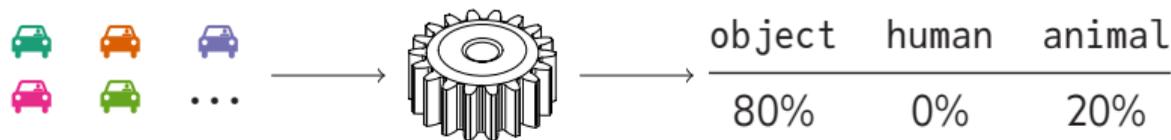


## Multi-class classification: All scores matter!



Common calibration evaluation techniques consider only the most-confident score

## Multi-class classification: All scores matter!



Common calibration evaluation techniques consider only the most-confident score

Common approaches do not distinguish between the two predictions even though the control actions based on these might be very different!

| object     | human | animal |
|------------|-------|--------|
| <b>80%</b> | 0%    | 20%    |
| <b>80%</b> | 20%   | 0%     |

# Weaker notions of calibration and calibration lenses

## Weaker notions

Weaker notions of calibration such as confidence calibration or calibration of marginal classifiers can be analyzed by considering calibration of induced predictive models.

# Weaker notions of calibration and calibration lenses

## Weaker notions

Weaker notions of calibration such as confidence calibration or calibration of marginal classifiers can be analyzed by considering calibration of induced predictive models.

## Definition (Calibration lenses)

Let  $\psi$  be a measurable function that defines targets  $Z := \psi(Y, P_X)$ . Then  $\psi$  induces a predictive model  $Q$  for targets  $Z$  with predictions

$$Q_X := \text{law}(\psi(\tilde{Y}, P_X))$$

where  $\tilde{Y} \sim P_X$ . Function  $\psi$  is called a *calibration lens*.

# Beyond classification

## Definition (reminder)

A probabilistic predictive model  $P$  is calibrated if

$$\text{law}(Y | P_X) = P_X \quad \text{almost surely.}$$

# Beyond classification

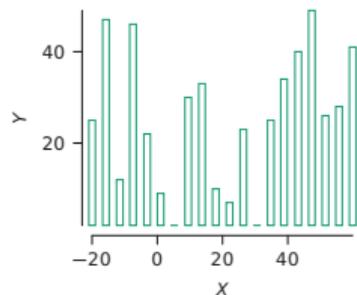
## Definition (reminder)

A probabilistic predictive model  $P$  is calibrated if

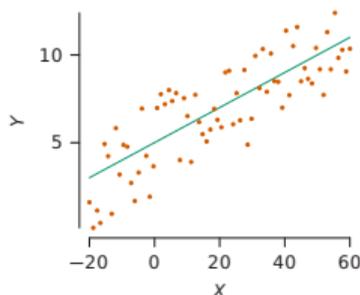
$$\text{law}(Y | P_X) = P_X \quad \text{almost surely.}$$

## Examples of other target spaces

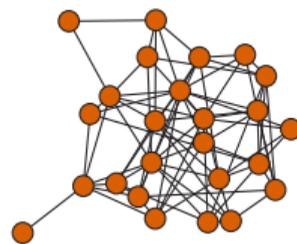
$\mathbb{N}_0$



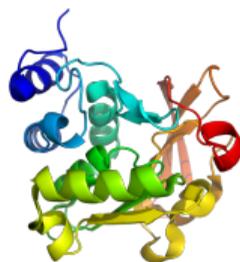
$\mathbb{R}^d$



graphs



protein structure



Calibration errors

## Expected calibration error (ECE)

### Definition

The expected calibration error (ECE) with respect to distance measure  $d$  is defined as

$$\text{ECE}_d := \mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X)).$$

## Expected calibration error (ECE)

### Definition

The expected calibration error (ECE) with respect to distance measure  $d$  is defined as

$$\text{ECE}_d := \mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X)).$$

### Choice of distance measure $d$

- ▶ For classification typically (semi-)metrics on the probability simplex (e.g., cityblock, Euclidean, or squared Euclidean distance)

# Expected calibration error (ECE)

## Definition

The expected calibration error (ECE) with respect to distance measure  $d$  is defined as

$$\text{ECE}_d := \mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X)).$$

## Choice of distance measure $d$

- ▶ For classification typically (semi-)metrics on the probability simplex (e.g., cityblock, Euclidean, or squared Euclidean distance)
- ▶ For general probabilistic predictive models **statistical divergences**

# Statistical divergences

## Definition

Let  $\mathcal{P}$  be a space of probability distributions. A function  $d: \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$  that satisfies

- ▶  $d(P, Q) \geq 0$  for all  $P, Q \in \mathcal{P}$ ,
- ▶  $d(P, Q) = 0$  if and only if  $P = Q$ ,

is a statistical divergence.

# Statistical divergences

## Definition

Let  $\mathcal{P}$  be a space of probability distributions. A function  $d: \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$  that satisfies

- ▶  $d(P, Q) \geq 0$  for all  $P, Q \in \mathcal{P}$ ,
- ▶  $d(P, Q) = 0$  if and only if  $P = Q$ ,

is a statistical divergence.

## Note

- ▶  $d$  does not need to be symmetric
- ▶  $d$  does not need to satisfy the triangle inequality

# Statistical divergences

## Definition

Let  $\mathcal{P}$  be a space of probability distributions. A function  $d: \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$  that satisfies

- ▶  $d(P, Q) \geq 0$  for all  $P, Q \in \mathcal{P}$ ,
- ▶  $d(P, Q) = 0$  if and only if  $P = Q$ ,

is a statistical divergence.

## Note

- ▶  $d$  does not need to be symmetric
- ▶  $d$  does not need to satisfy the triangle inequality

## Examples

- ▶  $f$ -divergences, e.g., Kullback-Leibler divergence or total variation distance
- ▶ Wasserstein distance

## Scoring rules: Definition

### Definition

The expected score of a probabilistic predictive model  $P$  is defined as

$$\mathbb{E}_{P_X, Y} S(P_X, Y)$$

where **scoring rule**  $s(p, y)$  is the reward of prediction  $p$  if the true outcome is  $y$ .

## Scoring rules: Definition

### Definition

The expected score of a probabilistic predictive model  $P$  is defined as

$$\mathbb{E}_{P_X, Y} s(P_X, Y)$$

where **scoring rule**  $s(p, y)$  is the reward of prediction  $p$  if the true outcome is  $y$ .

### Examples for classification

- ▶ Brier score:  $s(p, y) = - \int_{\Omega} ((\delta_y - p)^2)(d\omega)$
- ▶ Logarithmic score:  $s(p, y) = \log p(\{y\})$

## Scoring rules: Decomposition

For proper scoring rules

$$\begin{aligned}\mathbb{E}_{P_X, Y} s(P_X, Y) &= \mathbb{E}_{P_X} d(\text{law}(Y), \text{law}(Y | P_X)) \\ &\quad - \mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X)) - S(\text{law}(Y), \text{law}(Y))\end{aligned}$$

Expected score of  $P$  under  $Q$

$$S(P, Q) := \int_{\Omega} s(P, \omega) Q(d\omega)$$

Score divergence

$$d(P, Q) = S(Q, Q) - S(P, Q)$$

## Scoring rules: Decomposition

For proper scoring rules

$$\mathbb{E}_{P_X, Y} s(P_X, Y) = \underbrace{\mathbb{E}_{P_X} d(\text{law}(Y), \text{law}(Y | P_X))}_{\text{resolution}} - \mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X)) - S(\text{law}(Y), \text{law}(Y))$$

Expected score of  $P$  under  $Q$

$$S(P, Q) := \int_{\Omega} s(P, \omega) Q(d\omega)$$

Score divergence

$$d(P, Q) = S(Q, Q) - S(P, Q)$$

## Scoring rules: Decomposition

For proper scoring rules

$$\mathbb{E}_{P_X, Y} s(P_X, Y) = \underbrace{\mathbb{E}_{P_X} d(\text{law}(Y), \text{law}(Y | P_X))}_{\text{resolution}} - \underbrace{\mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X))}_{\text{calibration}} - S(\text{law}(Y), \text{law}(Y))$$

Expected score of  $P$  under  $Q$

$$S(P, Q) := \int_{\Omega} s(P, \omega) Q(d\omega)$$

Score divergence

$$d(P, Q) = S(Q, Q) - S(P, Q)$$

# Scoring rules: Decomposition

For proper scoring rules

$$\mathbb{E}_{P_X, Y} s(P_X, Y) = \underbrace{\mathbb{E}_{P_X} d(\text{law}(Y), \text{law}(Y | P_X))}_{\text{resolution}} - \underbrace{\mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X))}_{\text{calibration}} - \underbrace{S(\text{law}(Y), \text{law}(Y))}_{\text{uncertainty of } Y}$$

Expected score of  $P$  under  $Q$

$$S(P, Q) := \int_{\Omega} s(P, \omega) Q(d\omega)$$

Score divergence

$$d(P, Q) = S(Q, Q) - S(P, Q)$$

# Scoring rules: Decomposition

For proper scoring rules

$$\mathbb{E}_{P_X, Y} s(P_X, Y) = \underbrace{\mathbb{E}_{P_X} d(\text{law}(Y), \text{law}(Y | P_X))}_{\text{resolution}} - \underbrace{\mathbb{E}_{P_X} d(P_X, \text{law}(Y | P_X))}_{\text{calibration}} - \underbrace{S(\text{law}(Y), \text{law}(Y))}_{\text{uncertainty of } Y}$$

Expected score of  $P$  under  $Q$

$$S(P, Q) := \int_{\Omega} s(P, \omega) Q(d\omega)$$

Score divergence

$$d(P, Q) = S(Q, Q) - S(P, Q)$$

**Models can trade off calibration for resolution!**

# An alternative definition of calibration

## Theorem

A probabilistic predictive model  $P$  is calibrated if

$$(P_X, Y) \stackrel{d}{=} (P_X, Z_X),$$

where  $Z_X \sim P_X$ .

# An alternative definition of calibration

## Theorem

A probabilistic predictive model  $P$  is calibrated if

$$(P_X, Y) \stackrel{d}{=} (P_X, Z_X),$$

where  $Z_X \sim P_X$ .

Calibration error as distance between  $\text{law}((P_X, Y))$  and  $\text{law}((P_X, Z_X))$

## Calibration error: Integral probability metric

$$\text{CE}_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{P_{X,Y}} f(P_X, Y) - \mathbb{E}_{P_{X,Z_X}} f(P_X, Z_X) \right|$$

---

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

# Calibration error: Integral probability metric

$$\text{CE}_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{P_X, Y} f(P_X, Y) - \mathbb{E}_{P_X, Z_X} f(P_X, Z_X) \right|$$

## Examples

- ▶ 1-Wasserstein distance:  $\mathcal{F} = \{f : \|f\|_{\text{Lip}} \leq 1\}$
- ▶ Total variation distance:  $\mathcal{F} = \{f : \|f\|_{\infty} \leq 1\}$

## Calibration error: Integral probability metric

$$\text{CE}_{\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{P_{X,Y}} f(P_X, Y) - \mathbb{E}_{P_{X,Z_X}} f(P_X, Z_X) \right|$$

### Examples

- ▶ 1-Wasserstein distance:  $\mathcal{F} = \{f : \|f\|_{\text{Lip}} \leq 1\}$
- ▶ Total variation distance:  $\mathcal{F} = \{f : \|f\|_{\infty} \leq 1\}$

Common choices of  $\text{ECE}_d$  in classification can be formulated in this way

## Kernel calibration error: Maximum mean discrepancy (MMD)

Choose  $\mathcal{F} = \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq 1\}$  for some reproducing kernel Hilbert space  $\mathcal{H}$

---

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

## Kernel calibration error: Maximum mean discrepancy (MMD)

Choose  $\mathcal{F} = \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq 1\}$  for some reproducing kernel Hilbert space  $\mathcal{H}$

### Reproducing kernel Hilbert space (RKHS)

- ▶ Hilbert space of functions that satisfy  $f$  close to  $g \Rightarrow f(\mathbf{x})$  close to  $g(\mathbf{x})$

---

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

## Kernel calibration error: Maximum mean discrepancy (MMD)

Choose  $\mathcal{F} = \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq 1\}$  for some reproducing kernel Hilbert space  $\mathcal{H}$

### Reproducing kernel Hilbert space (RKHS)

- ▶ Hilbert space of functions that satisfy  $f$  close to  $g \Rightarrow f(\mathbf{x})$  close to  $g(\mathbf{x})$
- ▶ Possesses a positive-definite function  $k$  as reproducing kernel

---

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

## Kernel calibration error: Maximum mean discrepancy (MMD)

Choose  $\mathcal{F} = \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq 1\}$  for some reproducing kernel Hilbert space  $\mathcal{H}$

### Reproducing kernel Hilbert space (RKHS)

- ▶ Hilbert space of functions that satisfy  $f$  close to  $g \Rightarrow f(\mathbf{x})$  close to  $g(\mathbf{x})$
- ▶ Possesses a positive-definite function  $k$  as reproducing kernel

### Definition

The kernel calibration error (KCE) of a model  $P$  with respect to kernel  $k$  is defined as

$$\text{KCE}_k^2 := \text{CE}_{\mathcal{F}}^2 = \int k((p, y), (\tilde{p}, \tilde{y})) \mu(d(p, y)) \mu(d(\tilde{p}, \tilde{y})),$$

where  $\mu = \text{law}((P_X, Y)) - \text{law}((P_X, Z_X))$ .

---

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests beyond classification." In: *International Conference on Learning Representations*. 2021

# Choice of kernel

## Observations

- ▶ Kernel  $k$  defined on the product space of predictions and targets

# Choice of kernel

## Observations

- ▶ Kernel  $k$  defined on the product space of predictions and targets
- ▶ In multi-class classification,  $k$  can be identified with a matrix-valued kernel on the space of predictions

# Choice of kernel

## Observations

- ▶ Kernel  $k$  defined on the product space of predictions and targets
- ▶ In multi-class classification,  $k$  can be identified with a matrix-valued kernel on the space of predictions
- ▶ For specific kernel choices,  $Z_X$  can be integrated out analytically

# Choice of kernel

## Observations

- ▶ Kernel  $k$  defined on the product space of predictions and targets
- ▶ In multi-class classification,  $k$  can be identified with a matrix-valued kernel on the space of predictions
- ▶ For specific kernel choices,  $Z_X$  can be integrated out analytically
- ▶ Otherwise numerical integration methods (e.g., Monte Carlo integration) can be used to integrate out  $Z_X$

# Choice of kernel

## Observations

- ▶ Kernel  $k$  defined on the product space of predictions and targets
- ▶ In multi-class classification,  $k$  can be identified with a matrix-valued kernel on the space of predictions
- ▶ For specific kernel choices,  $Z_X$  can be integrated out analytically
- ▶ Otherwise numerical integration methods (e.g., Monte Carlo integration) can be used to integrate out  $Z_X$
- ▶ Suggestive to use tensor product kernels  $k = k_{\mathcal{P}} \otimes k_{\mathcal{Y}}$ , where  $k_{\mathcal{P}}$  and  $k_{\mathcal{Y}}$  are kernels on the space of predictions and targets, respectively

# Tensor product kernel

## Construction of $k_{\mathcal{P}}$ with Hilbertian metrics

- ▶ For Hilbertian metrics of form  $d_{\mathcal{P}}(p, \tilde{p}) = \|\phi(p) - \phi(\tilde{p})\|_2$  for some  $\phi: \mathcal{P} \rightarrow \mathbb{R}^d$ ,

$$k_{\mathcal{P}}(p, \tilde{p}) = \exp(-\lambda d_{\mathcal{P}}^{\nu}(p, \tilde{p})), \quad (1)$$

is valid kernel on the space of predictions for  $\lambda > 0$  and  $\nu \in (0, 2]$

# Tensor product kernel

## Construction of $k_{\mathcal{P}}$ with Hilbertian metrics

- ▶ For Hilbertian metrics of form  $d_{\mathcal{P}}(p, \tilde{p}) = \|\phi(p) - \phi(\tilde{p})\|_2$  for some  $\phi: \mathcal{P} \rightarrow \mathbb{R}^d$ ,

$$k_{\mathcal{P}}(p, \tilde{p}) = \exp(-\lambda d_{\mathcal{P}}^{\nu}(p, \tilde{p})), \quad (1)$$

is valid kernel on the space of predictions for  $\lambda > 0$  and  $\nu \in (0, 2]$

- ▶ Parameterization of predictions gives rise to  $\phi$  naturally

# Tensor product kernel

## Construction of $k_{\mathcal{P}}$ with Hilbertian metrics

- ▶ For Hilbertian metrics of form  $d_{\mathcal{P}}(p, \tilde{p}) = \|\phi(p) - \phi(\tilde{p})\|_2$  for some  $\phi: \mathcal{P} \rightarrow \mathbb{R}^d$ ,

$$k_{\mathcal{P}}(p, \tilde{p}) = \exp(-\lambda d_{\mathcal{P}}^{\nu}(p, \tilde{p})), \quad (1)$$

is valid kernel on the space of predictions for  $\lambda > 0$  and  $\nu \in (0, 2]$

- ▶ Parameterization of predictions gives rise to  $\phi$  naturally
- ▶ For many mixture models, Hilbertian metrics of model components can be lifted to Hilbertian metric of mixture models

Estimation of calibration errors

## Estimation of calibration errors

### Task

Estimate the calibration error of a model  $P$  from a validation dataset  $(X_i, Y_i)_{i=1, \dots, n}$  of features and corresponding targets.

# Estimation of calibration errors

## Task

Estimate the calibration error of a model  $P$  from a validation dataset  $(X_i, Y_i)_{i=1, \dots, n}$  of features and corresponding targets.

## Dataset of predictions and targets sufficient

- ▶ Calibration (errors) defined based only on predictions and targets
- ▶ Estimation can be performed with dataset  $(P_{X_i}, Y_i)$  of predictions and corresponding targets instead
- ▶ Highlights that structure of features and model is not relevant for calibration estimation

# ECE: Estimation

## Problem

The estimation of  $\text{law}(Y | P_X)$  is challenging.

---

 J. Vaicenavicius et al. "Evaluating model calibration in classification." In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

# ECE: Estimation

## Problem

The estimation of  $\text{law}(Y | P_X)$  is challenging.

## Binning predictions

- ▶ Common approach in classification
- ▶ Often leads to **biased and inconsistent** estimators

---

 J. Vaicenavicius et al. "Evaluating model calibration in classification." In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019

 D. Widmann, F. Lindsten, and D. Zachariah. "Calibration tests in multi-class classification: A unifying framework." In: *Advances in Neural Information Processing Systems* 32. 2019

# ECE: Experiments

## 10-class classification

For three models **M1**, **M2** and **M3**,  $10^4$  synthetic datasets  $(P_{X_i}, Y_i)_{i=1, \dots, 250}$  are sampled according to

- ▶  $P_{X_i} = \text{Cat}(p_i)$  with  $p_i \sim \text{Dir}(0.1, \dots, 0.1)$ ,

# ECE: Experiments

## 10-class classification

For three models **M1**, **M2** and **M3**,  $10^4$  synthetic datasets  $(P_{X_i}, Y_i)_{i=1, \dots, 250}$  are sampled according to

▶  $P_{X_i} = \text{Cat}(p_i)$  with  $p_i \sim \text{Dir}(0.1, \dots, 0.1)$ ,

▶  $Y_i$  conditionally on  $P_{X_i}$  from

**M1**:  $P_{X_i}$ ,      **M2**:  $0.5P_{X_i} + 0.5\delta_1$ ,      **M3**:  $U(\{1, \dots, 10\})$ .

# ECE: Experiments

## 10-class classification

For three models **M1**, **M2** and **M3**,  $10^4$  synthetic datasets  $(P_{X_i}, Y_i)_{i=1, \dots, 250}$  are sampled according to

▶  $P_{X_i} = \text{Cat}(p_i)$  with  $p_i \sim \text{Dir}(0.1, \dots, 0.1)$ ,

▶  $Y_i$  conditionally on  $P_{X_i}$  from

**M1**:  $P_{X_i}$ ,      **M2**:  $0.5P_{X_i} + 0.5\delta_1$ ,      **M3**:  $U(\{1, \dots, 10\})$ .

Model **M1** is calibrated, and models **M2** and **M3** are uncalibrated.

# ECE: Experiments

## 10-class classification

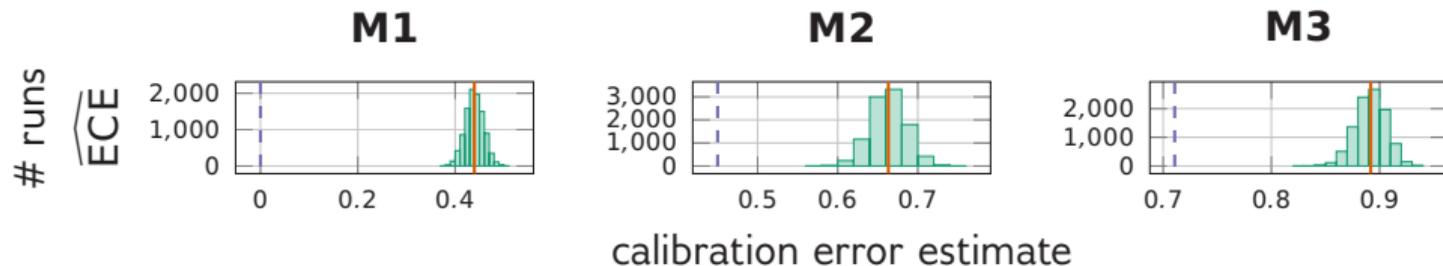
For three models **M1**, **M2** and **M3**,  $10^4$  synthetic datasets  $(P_{X_i}, Y_i)_{i=1, \dots, 250}$  are sampled according to

▶  $P_{X_i} = \text{Cat}(p_i)$  with  $p_i \sim \text{Dir}(0.1, \dots, 0.1)$ ,

▶  $Y_i$  conditionally on  $P_{X_i}$  from

**M1**:  $P_{X_i}$ ,      **M2**:  $0.5P_{X_i} + 0.5\delta_1$ ,      **M3**:  $U(\{1, \dots, 10\})$ .

Model **M1** is calibrated, and models **M2** and **M3** are uncalibrated.



## Kernel calibration error: Estimation

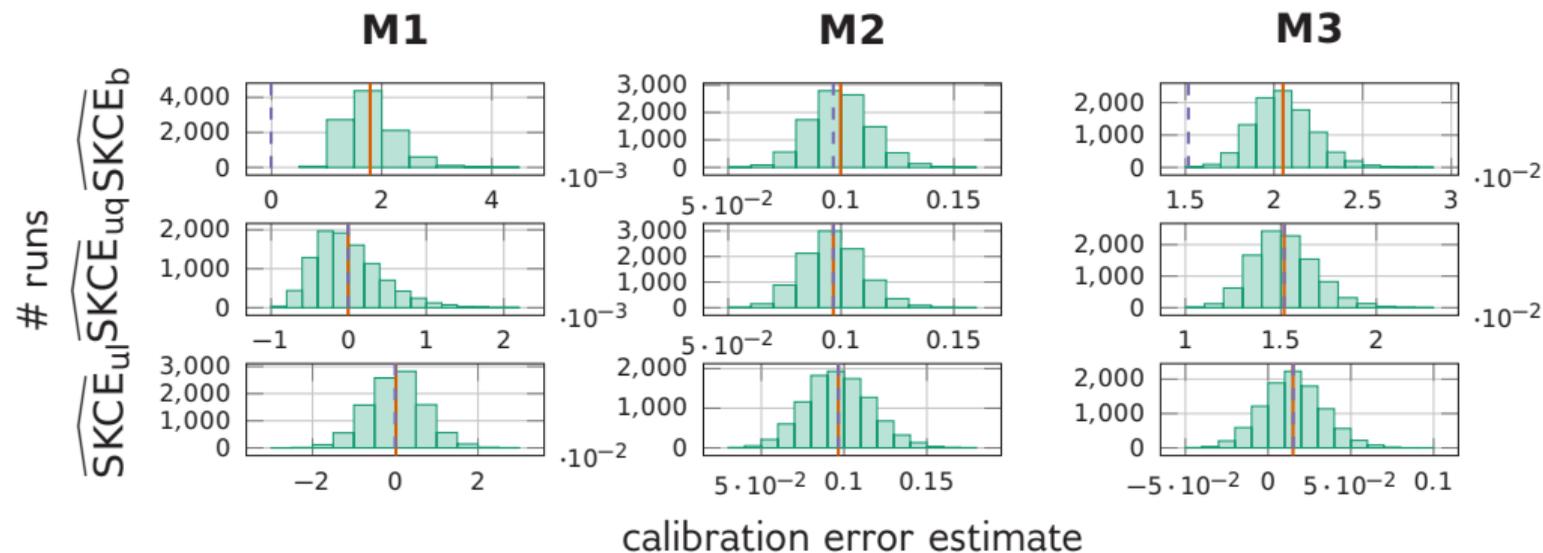
- ▶ For the MMD unbiased and consistent estimators are available

## Kernel calibration error: Estimation

- ▶ For the MMD unbiased and consistent estimators are available
- ▶ Variance can be reduced by marginalizing out  $Z_X$

# Kernel calibration error: Estimation

- ▶ For the MMD unbiased and consistent estimators are available
- ▶ Variance can be reduced by marginalizing out  $Z_X$



## Calibration tests

## Problems with calibration errors

- ▶ Calibration errors have no meaningful unit or scale

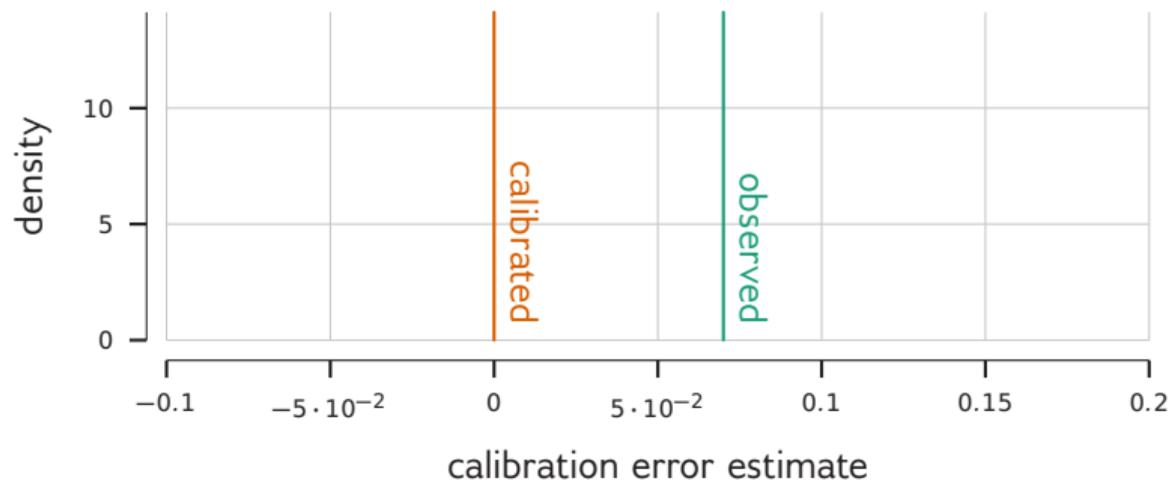
## Problems with calibration errors

- ▶ Calibration errors have no meaningful unit or scale
- ▶ Different calibration errors rank models differently

## Problems with calibration errors

- ▶ Calibration errors have no meaningful unit or scale
- ▶ Different calibration errors rank models differently
- ▶ Calibration error estimators are random variables

# Calibration tests



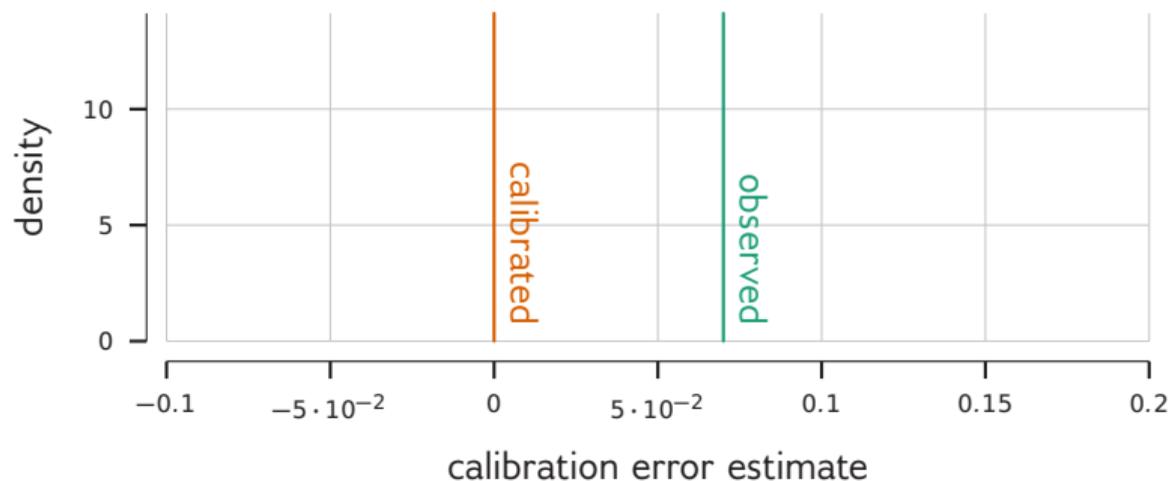
---

 J. Bröcker and L. A. Smith. "Increasing the reliability of reliability diagrams." In: *Weather and Forecasting* (2007)

 J. Vaicenavicius et al. "Evaluating model calibration in classification." In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019

# Calibration tests

Null hypothesis  $H_0 :=$  “model is calibrated”

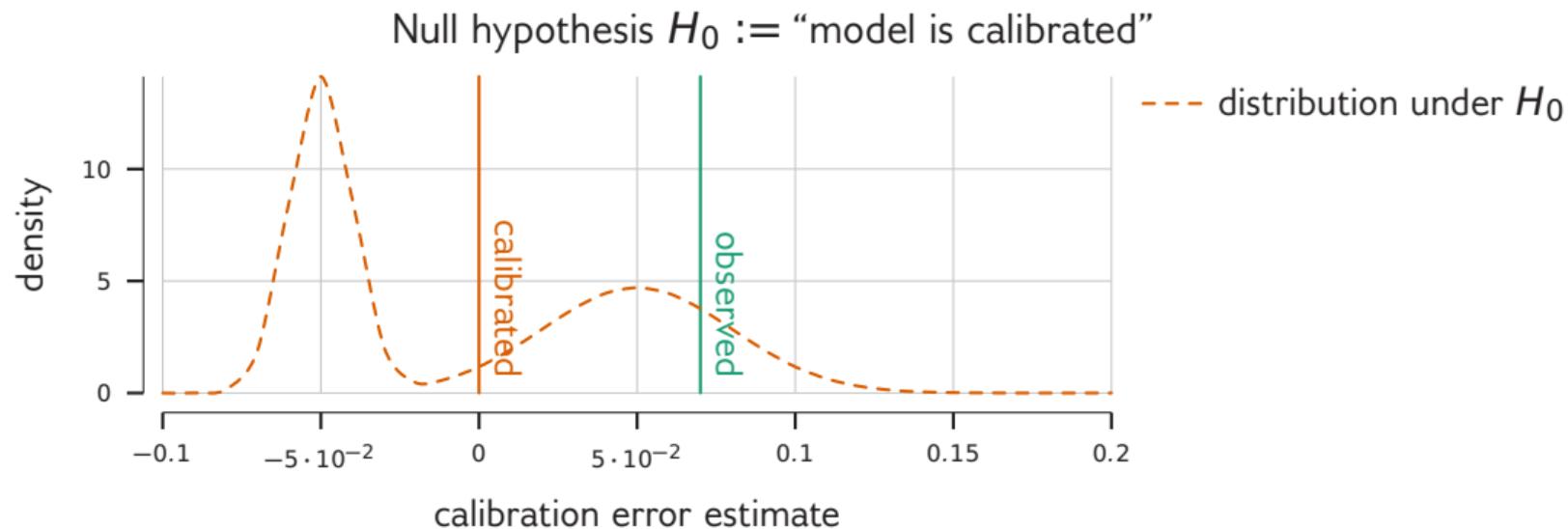


---

 J. Bröcker and L. A. Smith. “Increasing the reliability of reliability diagrams.” In: *Weather and Forecasting* (2007)

 J. Vaicenavicius et al. “Evaluating model calibration in classification.” In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019

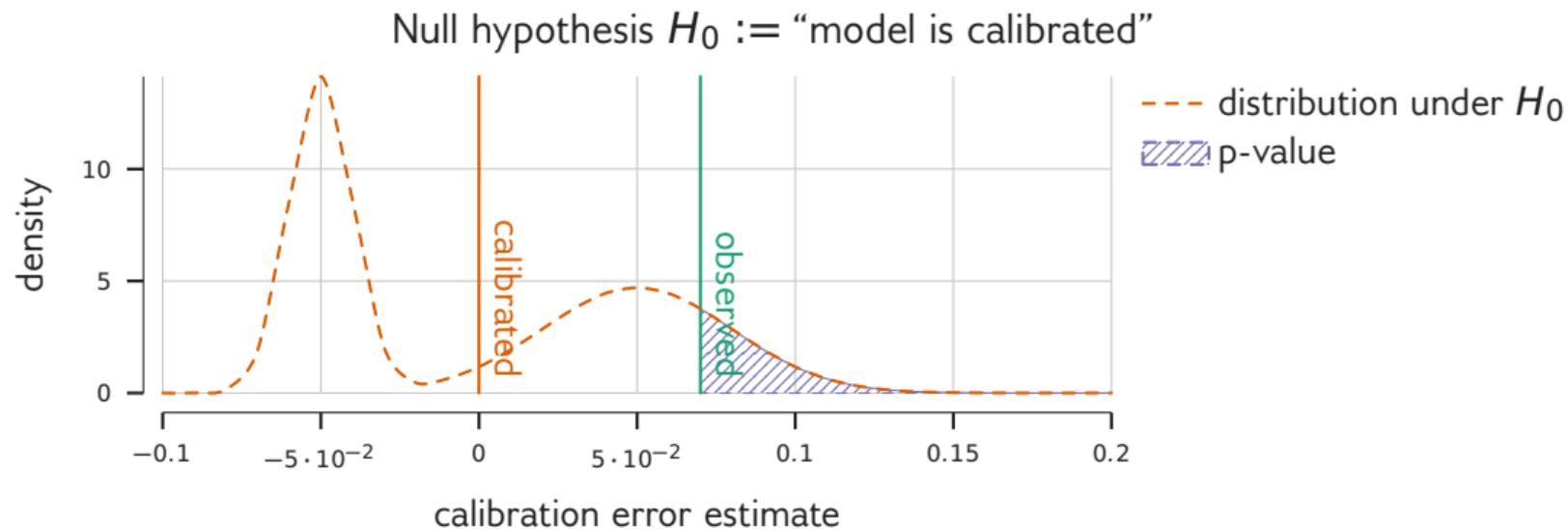
# Calibration tests



J. Bröcker and L. A. Smith. “Increasing the reliability of reliability diagrams.” In: *Weather and Forecasting* (2007)

J. Vaicenavicius et al. “Evaluating model calibration in classification.” In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019

# Calibration tests

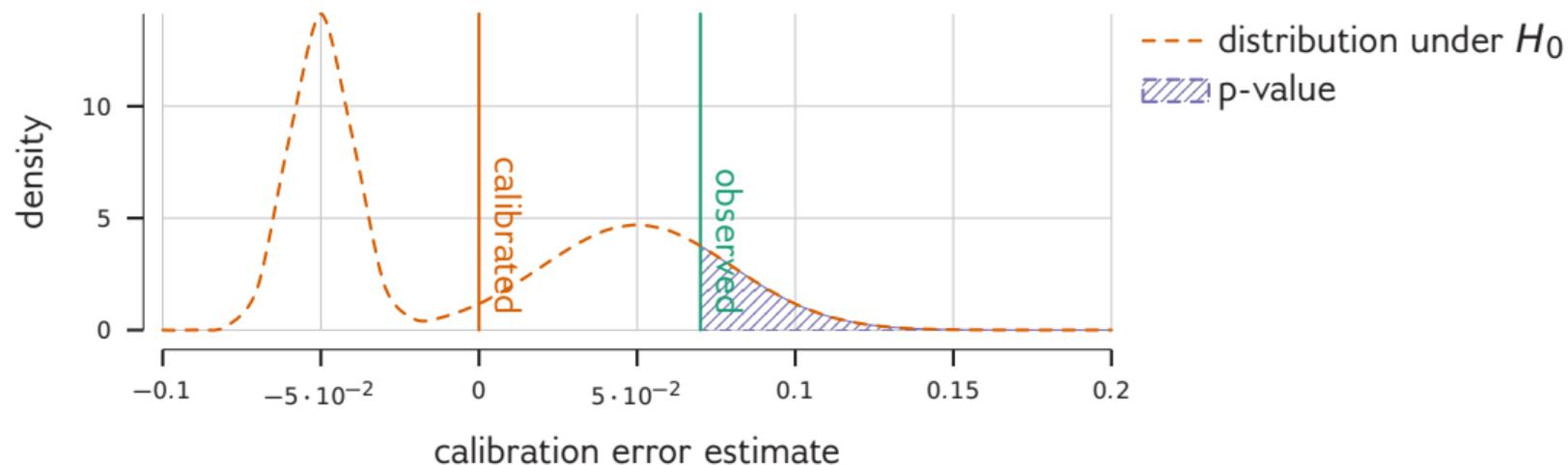


▣ J. Bröcker and L. A. Smith. “Increasing the reliability of reliability diagrams.” In: *Weather and Forecasting* (2007)

▣ J. Vaicenavicius et al. “Evaluating model calibration in classification.” In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019

# Calibration tests

Null hypothesis  $H_0 :=$  “model is calibrated”



Reject  $H_0$  if p-value is small

 J. Bröcker and L. A. Smith. “Increasing the reliability of reliability diagrams.” In: *Weather and Forecasting* (2007)

 J. Vaicenavicius et al. “Evaluating model calibration in classification.” In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. Vol. 89. Apr. 2019

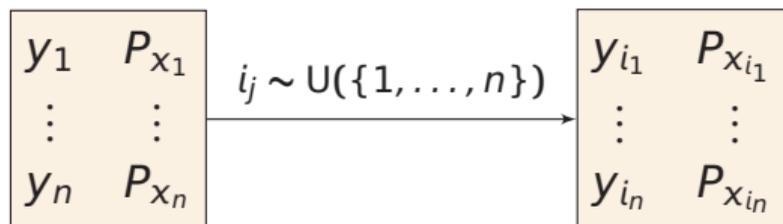
# Consistency resampling

Original dataset

|          |           |
|----------|-----------|
| $y_1$    | $P_{x_1}$ |
| $\vdots$ | $\vdots$  |
| $y_n$    | $P_{x_n}$ |

# Consistency resampling

Original dataset



# Consistency resampling

Original dataset

|          |           |
|----------|-----------|
| $y_1$    | $P_{x_1}$ |
| $\vdots$ | $\vdots$  |
| $y_n$    | $P_{x_n}$ |

$$i_j \sim U(\{1, \dots, n\})$$

|           |               |
|-----------|---------------|
| $y_{i_1}$ | $P_{x_{i_1}}$ |
| $\vdots$  | $\vdots$      |
| $y_{i_n}$ | $P_{x_{i_n}}$ |

$$\tilde{y}_j \sim P_{x_j}$$

Resampled dataset under  $H_0$

|                   |               |
|-------------------|---------------|
| $\tilde{y}_{i_1}$ | $P_{x_{i_1}}$ |
| $\vdots$          | $\vdots$      |
| $\tilde{y}_{i_n}$ | $P_{x_{i_n}}$ |

# Consistency resampling

Original dataset

|          |           |
|----------|-----------|
| $y_1$    | $P_{X_1}$ |
| $\vdots$ | $\vdots$  |
| $y_n$    | $P_{X_n}$ |

$i_j \sim U(\{1, \dots, n\})$

|           |               |
|-----------|---------------|
| $y_{i_1}$ | $P_{X_{i_1}}$ |
| $\vdots$  | $\vdots$      |
| $y_{i_n}$ | $P_{X_{i_n}}$ |

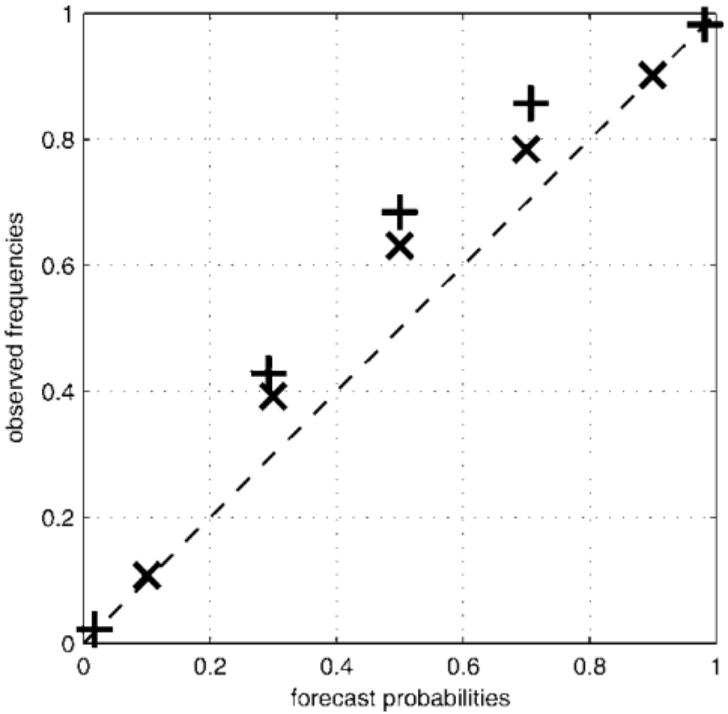
$\tilde{y}_j \sim P_{X_j}$

Resampled dataset under  $H_0$

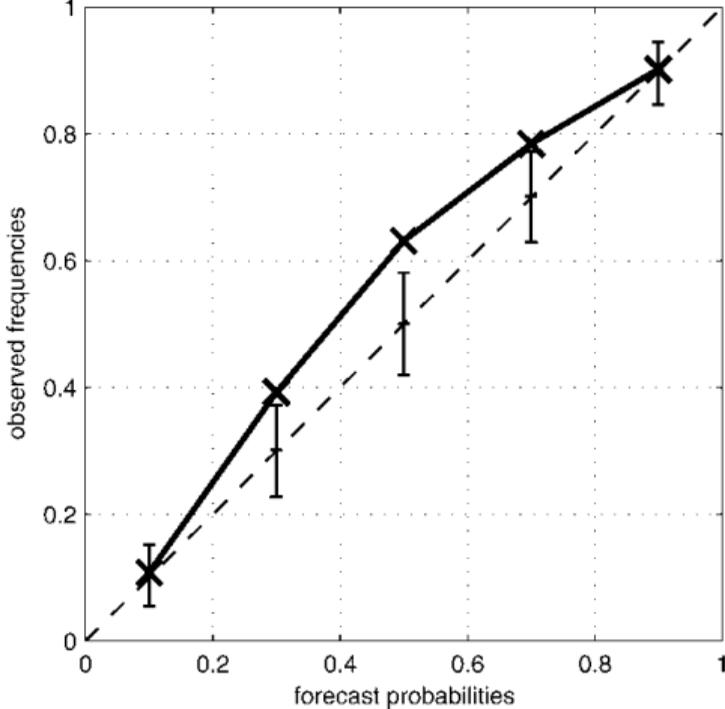
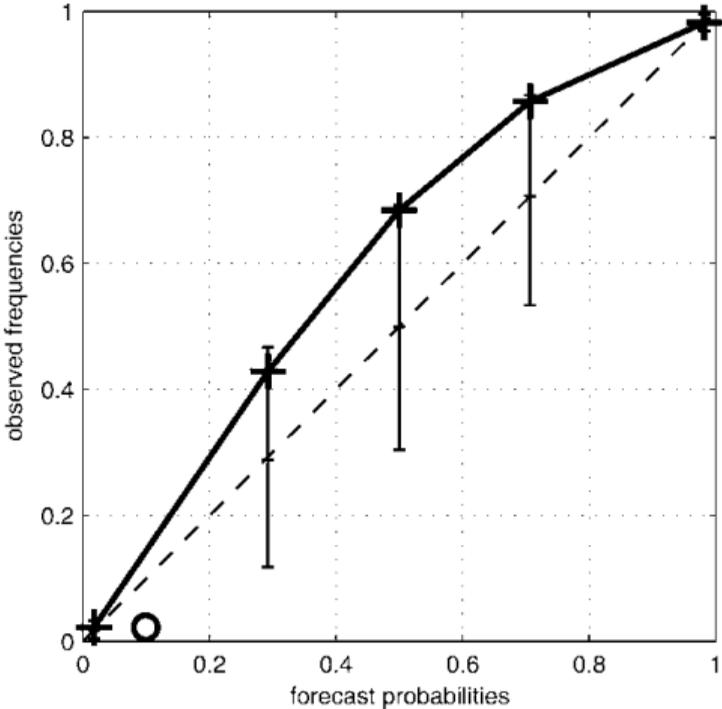
|                   |               |
|-------------------|---------------|
| $\tilde{y}_{i_1}$ | $P_{X_{i_1}}$ |
| $\vdots$          | $\vdots$      |
| $\tilde{y}_{i_n}$ | $P_{X_{i_n}}$ |

estimate p-value

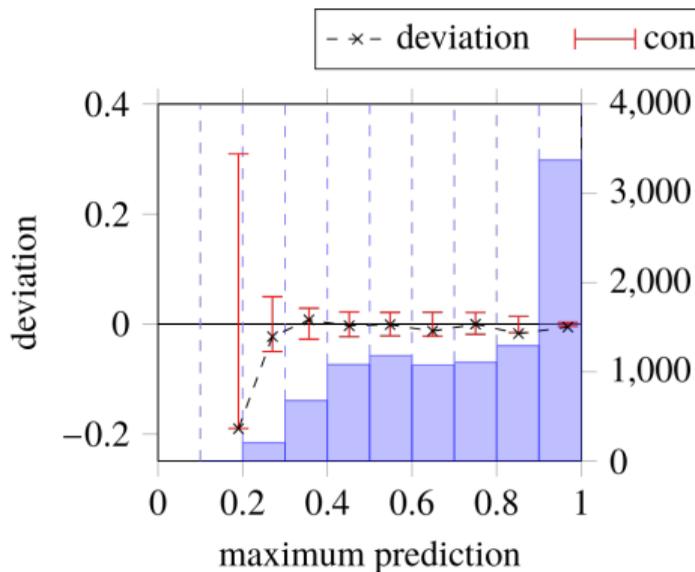
# Consistency bars



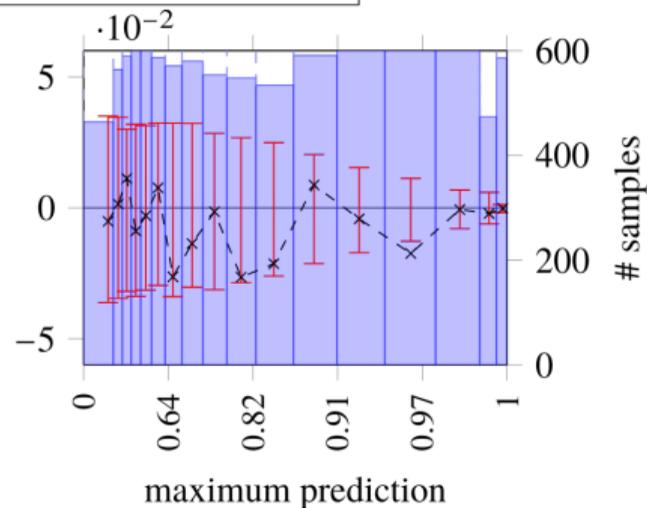
# Consistency bars



# Variants



(a) Equally-sized bins



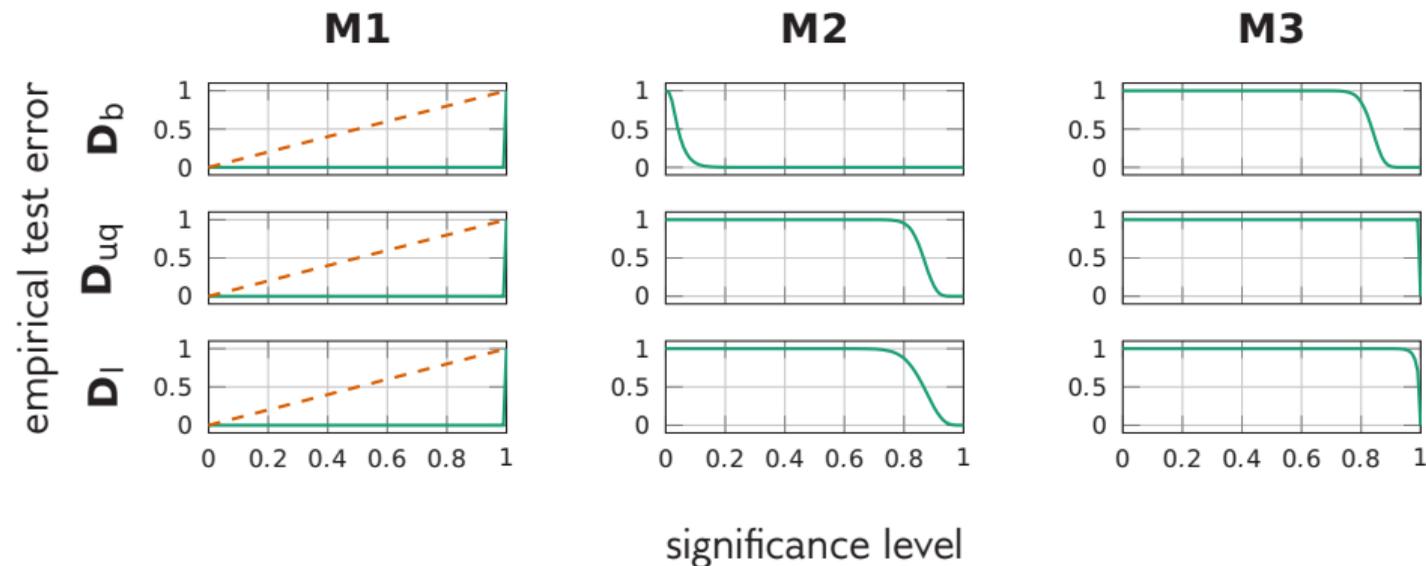
(b) Data-dependent bins

# Kernel calibration error: Distribution-free tests

**Upper bound** the p-value

# Kernel calibration error: Distribution-free tests

Upper bound the p-value

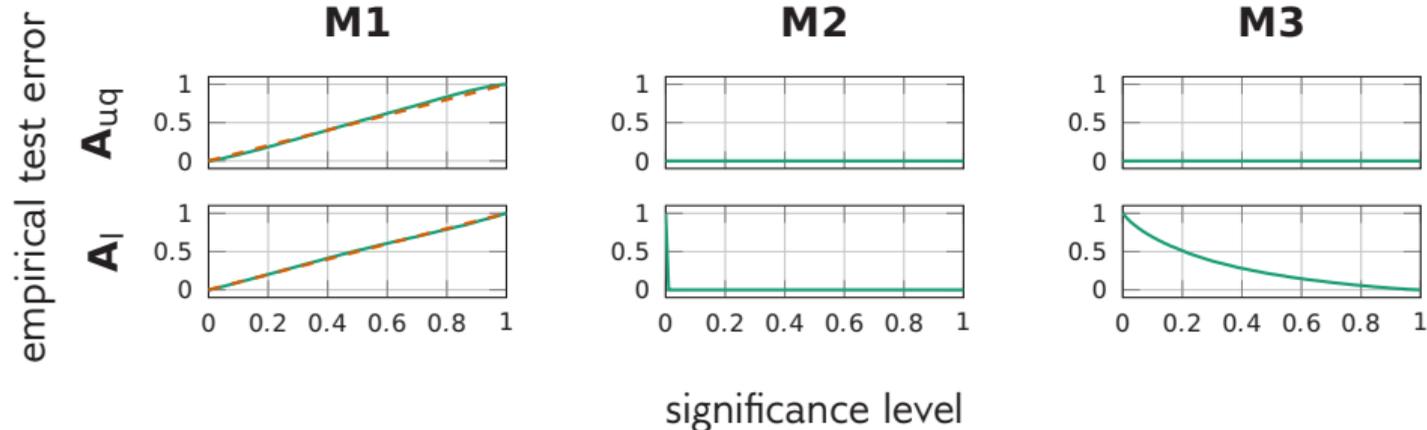


## Kernel calibration error: Asymptotic tests

**Approximate** the p-value based on the **asymptotic** distribution

# Kernel calibration error: Asymptotic tests

**Approximate** the p-value based on the **asymptotic** distribution



Calibration: Software packages

# CalibrationAnalysis.jl

## Summary

- ▶ Suite for analyzing calibration of probabilistic predictive models
- ▶ Written in Julia, with interfaces in Python (`pycalibration`) and R (`rcalibration`)

# CalibrationAnalysis.jl

## Summary

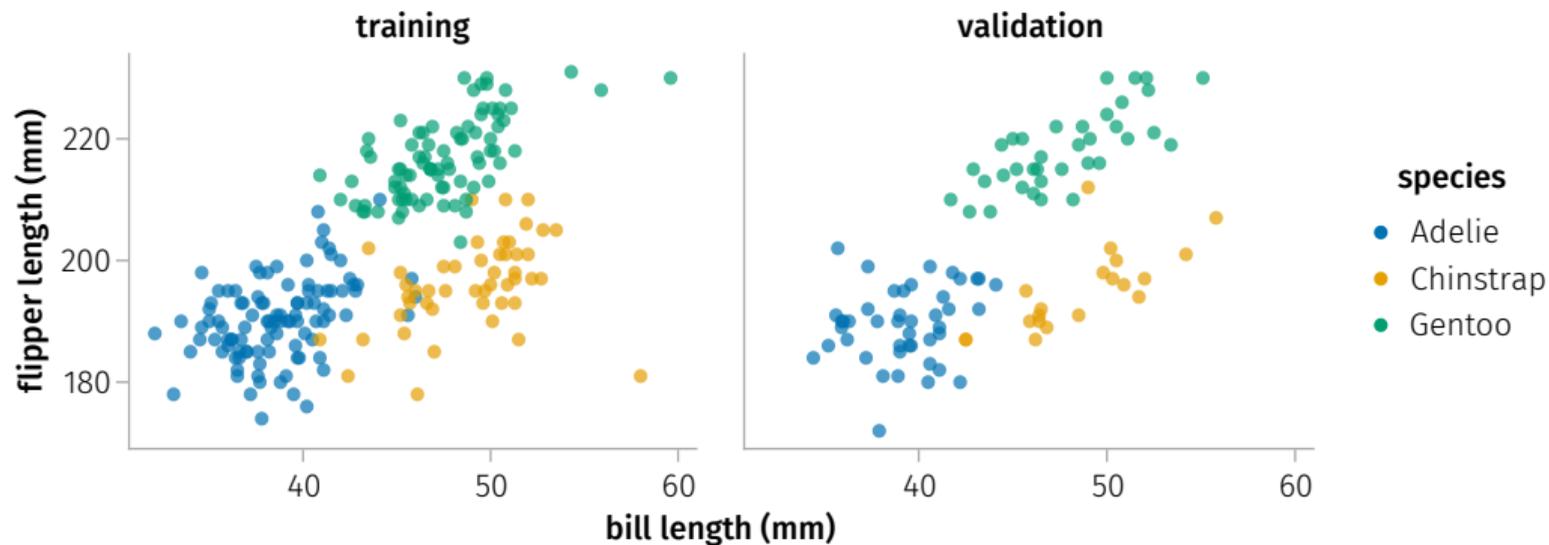
- ▶ Suite for analyzing calibration of probabilistic predictive models
- ▶ Written in Julia, with interfaces in Python (`pycalibration`) and R (`rcalibration`)

## Features

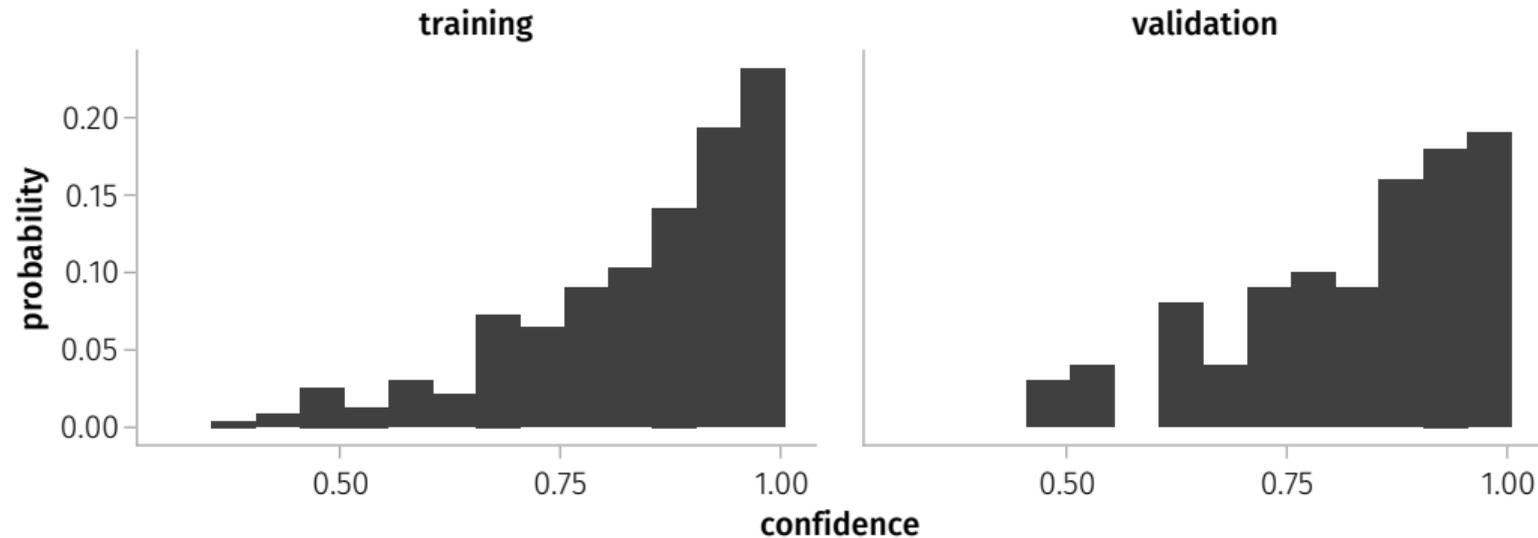
- ▶ Supports classification and regression models
- ▶ Reliability diagrams (`ReliabilityDiagrams.jl`)
- ▶ Estimation of calibration errors such as ECE and KCE (`CalibrationErrors.jl`)
- ▶ Calibration tests (`CalibrationTests.jl`)
- ▶ Integration with Julia ecosystem: Supports `Plots.jl` and `Makie.jl`, `KernelFunctions.jl`, and `HypothesisTests.jl`

## Calibration analysis: Penguins example

We train a naive Bayes classifier of penguin species based on bill depth, bill length, flipper length, and body mass.



# Binary predictions



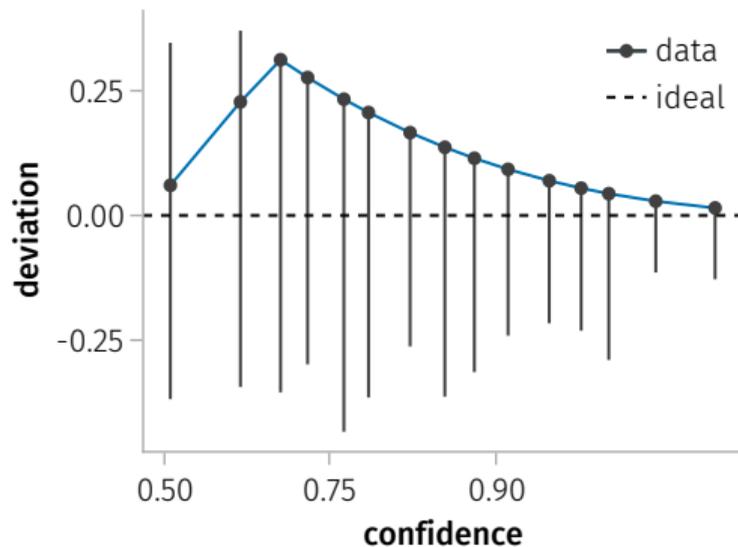
# Reliability diagram

## Code

```
julia> using CalibrationAnalysis, CairoMakie
```

```
julia> reliability(  
    confidence,  
    outcome;  
    binning=EqualMass(; n=15),  
    deviation=true,  
    ↪ consistencybars=ConsistencyBars(),  
)
```

## Polished result



## Expected calibration error: Code

```
julia> ece = ECE(UniformBinning(5), TotalVariation());
```

```
julia> ece(confidence, outcome)
```

```
0.06594437403598197
```

```
julia> ece(predictions, observations)
```

```
0.15789651955832515
```

## Kernel calibration error: Code

```
julia> kernel = GaussianKernel() ⊗ WhiteKernel();
```

```
julia> skce = SKCE(kernel);
```

```
julia> skce(predictions, observations)
```

```
0.0032631144705774404
```

```
julia> skce = SKCE(kernel; unbiased=false);
```

```
julia> skce(predictions, observations)
```

```
0.004202113116841622
```

```
julia> skce = SKCE(kernel; blocksize=5);
```

```
julia> skce(predictions, observations)
```

```
-0.005037270862051889
```

## Calibration test: Code

```
julia> AsymptoticSKCETest(kernel, predictions, observations)
```

```
Asymptotic SKCE test
```

```
-----  
Population details:
```

```
parameter of interest:  SKCE  
value under h_0:       0.0  
point estimate:        0.00326311
```

```
Test summary:
```

```
outcome with 95% confidence: reject h_0  
one-sided p-value:          0.0150
```

```
Details:
```

```
test statistic: -0.0009060378940361157
```

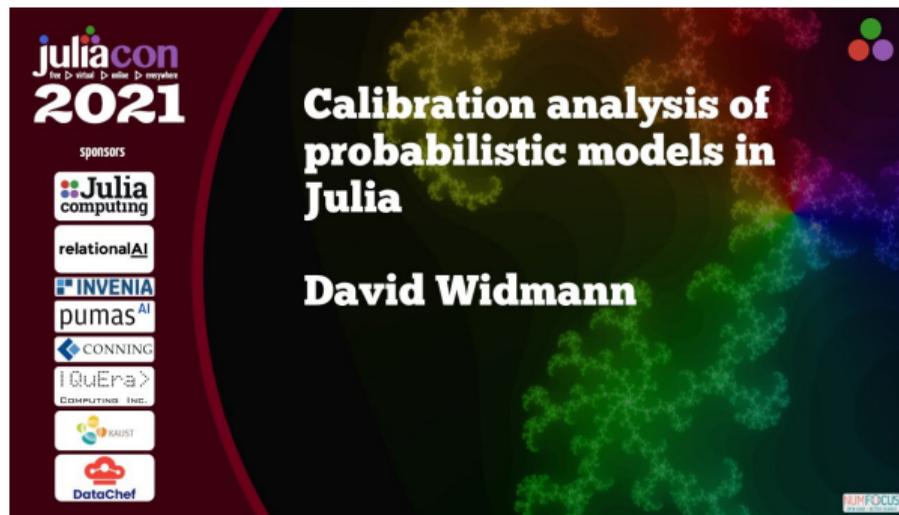
```
julia> test = ConsistencyTest(ece, predictions, observations);
```

```
julia> pvalue(test; bootstrap_iters=10_000)
```

```
0.0188
```

## Additional resources

- ▶ Online documentation: <https://devmotion.github.io/CalibrationErrors.jl/>
- ▶ Talk at JuliaCon 2021: <https://youtu.be/PrLsXFvwzuA>



Slides available at <https://talks.widmann.dev/2021/07/calibration/>

Concluding remarks

## Important takeaways

- ▶ More fine-grained analysis of calibration can be important
- ▶ MMD-like kernel calibration error can be applied to probabilistic models beyond classification
- ▶ Estimators of kernel calibration error have appealing properties
- ▶ Calibration errors and reliability diagrams can be misleading