

# Calibration of probabilistic predictions

Johanna Ziegel

joint work with

Sam Allen, Georgios Gavriloopoulos, Alexander Henzi, Gian-Reto Kleger

ETH Zurich

AISTATS Workshop: “Towards Trustworthy Predictions: Theory and Applications of Calibration for Modern AI”

5 May, 2026

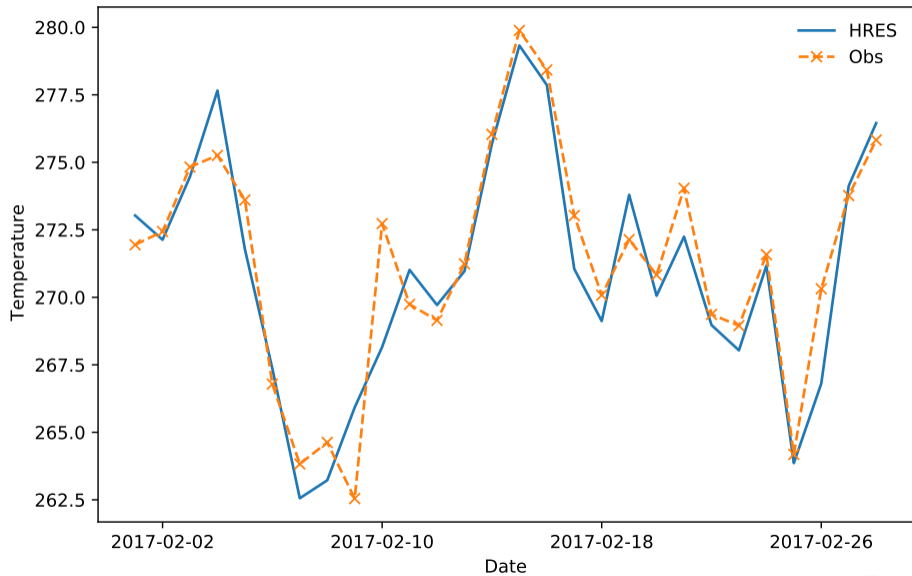
# Introduction: Probabilistic predictions

- ▶ Let  $Y \in \mathbb{R}$  be an unknown future outcome.
  - ▶ Example: Temperature tomorrow at 12:00 in Freiburg.
- ▶ Single valued “best guess”  $z \in \mathbb{R}$  does not quantify uncertainty.
- ▶ Goal: Quantify uncertainty of  $Y$  by a *probabilistic prediction*  $F$ .
  - ▶  $F$  is a distribution on  $\mathbb{R}$  (typically specified as a CDF or density).
- ▶ If  $X$  is the information available for prediction,  $F$  should approximate  $\mathcal{L}(Y | X)$ .

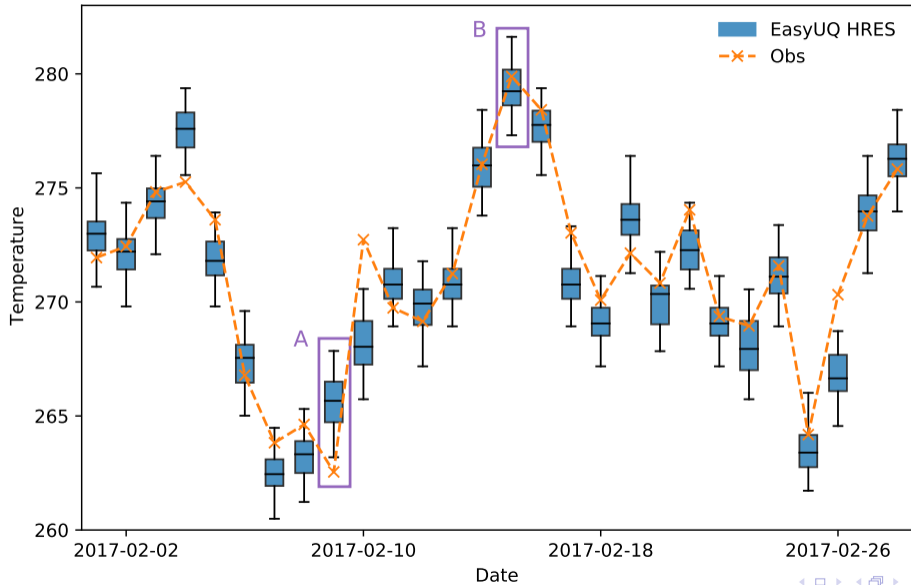
## Introduction: Probabilistic predictions

- ▶ Let  $Y \in \mathbb{R}$  be an unknown future outcome.
  - ▶ Example: Temperature tomorrow at 12:00 in Freiburg.
- ▶ Single valued “best guess”  $z \in \mathbb{R}$  does not quantify uncertainty.
- ▶ Goal: Quantify uncertainty of  $Y$  by a *probabilistic prediction*  $F$ .
  - ▶  $F$  is a distribution on  $\mathbb{R}$  (typically specified as a CDF or density).
- ▶ If  $X$  is the information available for prediction,  $F$  should approximate  $\mathcal{L}(Y | X)$ .

# Illustration



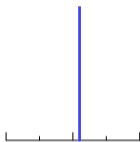
# Illustration



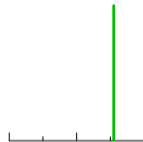
# Probabilistic and point predictions

“Tomorrow at 12:00  
temperature will be 4.5°C.”

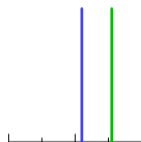
Forecast



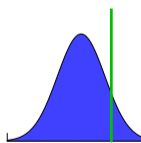
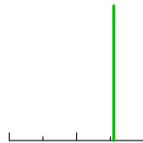
Observation



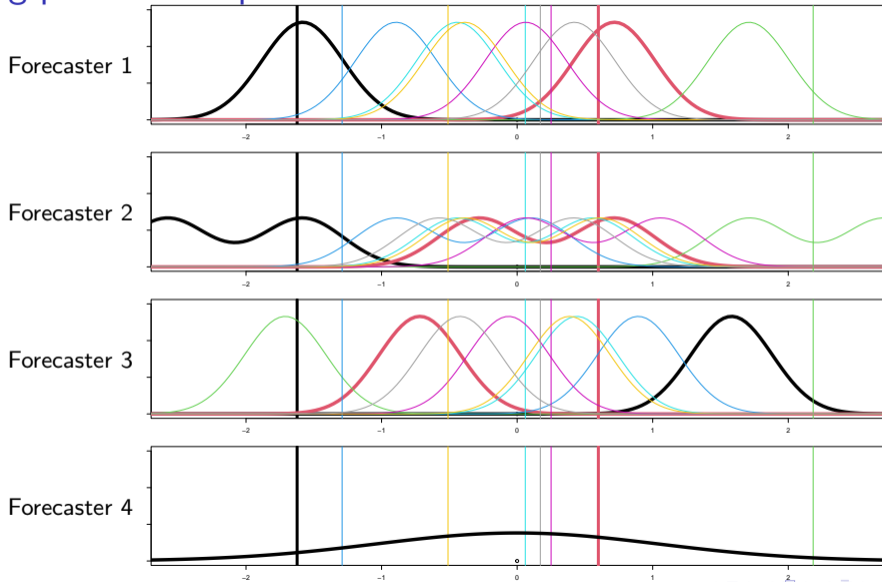
Verification



“Tomorrow at 12:00  
temperature will be  $\mathcal{N}(4.5, \sigma^2)$ .”



# Evaluating probabilistic predictions



## Calibration: Compatibility between forecasts and observations

Probabilities derived from predictive distributions should align with observed frequencies.

Most popular: Probabilistic calibration/“Flat PIT histogram”

$$F_i(Y_i) \sim \text{UNIF}(0, 1) \quad \text{for all } i$$

- ▶  $Y_i \in \mathbb{R}$ ,  $F_i$  predictive CDF for  $Y_i$
- ▶ Suitable randomization if  $F_i$  is not continuous
- ▶ Closely related to validity of conformal predictive systems
- ▶ Ensures correct marginal coverage of prediction intervals
- ▶ **Binary outcomes:**  $Y_i \in \{0, 1\} : \mathbb{P}(Y_i = 1 | p_i) = p_i$
- ▶ *Many* notions of calibration, except for binary outcomes. . .

## Calibration: Compatibility between forecasts and observations

Probabilities derived from predictive distributions should align with observed frequencies.

Most popular: Probabilistic calibration/“Flat PIT histogram”

$$F_i(Y_i) \sim \text{UNIF}(0, 1) \quad \text{for all } i$$

- ▶  $Y_i \in \mathbb{R}$ ,  $F_i$  predictive CDF for  $Y_i$
- ▶ Suitable randomization if  $F_i$  is not continuous
- ▶ Closely related to validity of conformal predictive systems
- ▶ Ensures correct marginal coverage of prediction intervals
- ▶ **Binary outcomes:**  $Y_i \in \{0, 1\} : \mathbb{P}(Y_i = 1 | p_i) = p_i$
- ▶ *Many* notions of calibration, except for binary outcomes...

## Calibration: Compatibility between forecasts and observations

Probabilities derived from predictive distributions should align with observed frequencies.

Most popular: Probabilistic calibration/“Flat PIT histogram”

$$F_i(Y_i) \sim \text{UNIF}(0, 1) \quad \text{for all } i$$

- ▶  $Y_i \in \mathbb{R}$ ,  $F_i$  predictive CDF for  $Y_i$
- ▶ Suitable randomization if  $F_i$  is not continuous
- ▶ Closely related to validity of conformal predictive systems
- ▶ Ensures correct marginal coverage of prediction intervals
- ▶ **Binary outcomes:**  $Y_i \in \{0, 1\} : \mathbb{P}(Y_i = 1 | p_i) = p_i$
- ▶ *Many* notions of calibration, except for binary outcomes...

## Why do we want flat PIT histograms?

Let  $F$  be a (deterministic!) CDF and  $Y$  a random variable,  $V \sim \text{UNIF}[0, 1]$  independent of  $Y$ . Then,

$$\begin{aligned} Y \sim F &\iff F(Y-) + V(F(Y) - F(Y-)) \sim \text{UNIF}[0, 1] \\ &\iff \mathbb{P}(F(Y) \leq \alpha) \leq \alpha \leq \mathbb{P}(F(Y-) < \alpha) \forall \alpha \end{aligned}$$

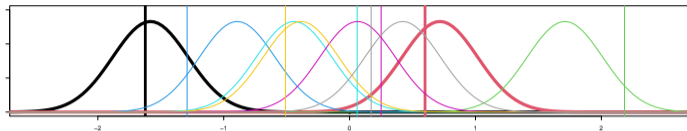
## Why do we want flat PIT histograms?

Let  $F$  be a (deterministic!) CDF and  $Y$  a random variable,  $V \sim \text{UNIF}[0, 1]$  independent of  $Y$ . Then,

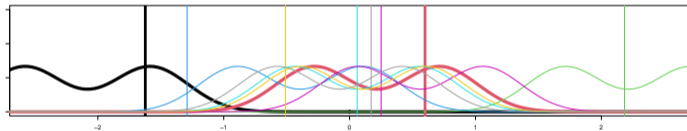
$$\begin{aligned} Y \sim F &\iff F(Y-) + V(F(Y) - F(Y-)) \sim \text{UNIF}[0, 1] \\ &\iff \mathbb{P}(F(Y) \leq \alpha) \leq \alpha \leq \mathbb{P}(F(Y-) < \alpha) \forall \alpha \end{aligned}$$

# Evaluating probabilistic predictions

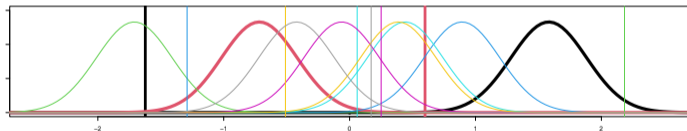
$$\mu \sim \mathcal{N}(0, 1), \quad Y \sim \mathcal{N}(\mu, 0.09)$$



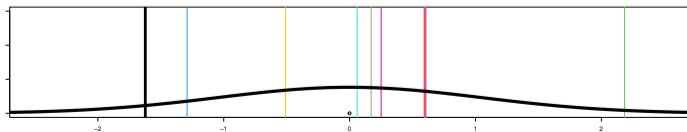
Probabilistic calibration ✓



Probabilistic calibration ✓



Probabilistic calibration ✗



Probabilistic calibration ✓

## Many notions of calibration ...

Auto-calibration:

$$\mathbb{P}(Y_i > y \mid F_i) = 1 - F_i(y) \quad \forall y$$

$$\mathcal{L}(Y_i \mid F_i) = F_i$$



Isotonic calibration:

$$\mathbb{P}(Y_i > y \mid \mathcal{A}(F_i)) = 1 - F_i(y) \quad \forall y$$

$$\mathcal{L}(Y_i \mid \mathcal{A}(F_i)) = F_i$$



Threshold calibration:

$$\mathbb{P}(Y_i > y \mid F_i(y)) = 1 - F_i(y) \quad \forall y$$



Marginal calibration:

$$\mathbb{P}(Y_i > y) = 1 - \mathbb{E}F_i(y) \quad \forall y$$

Quantile calibration:

$$q_\alpha(Y_i \mid F_i^{-1}(\alpha)) = F_i^{-1}(\alpha) \quad \forall \alpha$$



Probabilistic calibration:

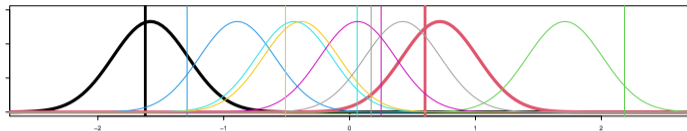
$$F_i(Y_i) \sim \text{UNIF}(0, 1)$$

$$\mathbb{P}(F_i(Y_i) \leq \alpha) \leq \alpha \leq \mathbb{P}(F_i(Y_i) < \alpha) \quad \forall \alpha$$

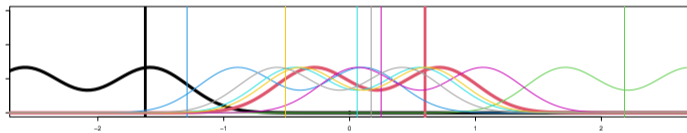
And if we want to focus on tails of  $F_i$  ... (Allen et al., 2025b)

# Evaluating probabilistic predictions

$$\mu \sim \mathcal{N}(0, 1), \quad Y \sim \mathcal{N}(\mu, 0.09)$$



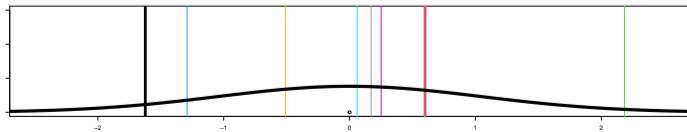
Auto-calibration ✓  
Probabilistic calibration ✓  
Marginal calibration ✓



Auto-calibration ✗  
Probabilistic calibration ✓  
Marginal calibration ✗



Auto-calibration ✗  
Probabilistic calibration ✗  
Marginal calibration ✓



Auto-calibration ✓  
Probabilistic calibration ✓  
Marginal calibration ✓

- ▶ Probabilistic predictions should be calibrated, ideally, *auto-calibrated*.
- ▶ Subject to calibration, they should be *sharp* in order to be informative.
- ▶ Comparison of probabilistic predictions with proper scoring rules:  
Assign a real-valued score assessing calibration and sharpness simultaneously.

**Logarithmic Score (LogS)**  $f$  density of  $F$

$$\text{LogS}(F, y) = -\log f(y)$$

**Continuous Ranked Probability Score (CRPS)**  $F$  CDF, finite mean

$$\text{CRPS}(F, y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{y \leq z\})^2 dz$$

- ▶ Connections of proper scoring rules to calibration through scoring rule decompositions and bias-variance decompositions.

For references see review of Waghamare and Ziegel (2026).

# Conformal prediction

**Goal:** Provide predictions with calibration guarantees out-of-sample.

# What is at the heart of conformal prediction?

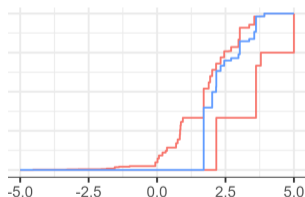
“In-sample calibration yields conformal calibration guarantees.”

Predictive system

A set  $\Pi \subseteq \mathbb{R} \times [0, 1]$  of the form

$$\Pi = \{(y, \tau) \mid \Pi_\ell(y) \leq \tau \leq \Pi_u(y)\}$$

with  $\Pi_\ell \leq \Pi_u$  increasing,  $\lim_{y \rightarrow -\infty} \Pi_\ell(y) = 0$ ,  $\lim_{y \rightarrow \infty} \Pi_u(y) = 1$ .



Conformal calibration guarantee:

We can construct a predictive system that contains a calibrated CDF.

# What is at the heart of conformal prediction?

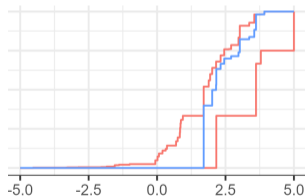
“In-sample calibration yields conformal calibration guarantees.”

## Predictive system

A set  $\Pi \subseteq \mathbb{R} \times [0, 1]$  of the form

$$\Pi = \{(y, \tau) \mid \Pi_\ell(y) \leq \tau \leq \Pi_u(y)\}$$

with  $\Pi_\ell \leq \Pi_u$  increasing,  $\lim_{y \rightarrow -\infty} \Pi_\ell(y) = 0$ ,  $\lim_{y \rightarrow \infty} \Pi_u(y) = 1$ .



## Conformal calibration guarantee:

We can construct a predictive system that contains a calibrated CDF.

# What is at the heart of conformal prediction?

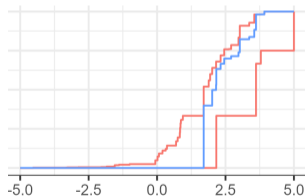
“In-sample calibration yields conformal calibration guarantees.”

## Predictive system

A set  $\Pi \subseteq \mathbb{R} \times [0, 1]$  of the form

$$\Pi = \{(y, \tau) \mid \Pi_\ell(y) \leq \tau \leq \Pi_u(y)\}$$

with  $\Pi_\ell \leq \Pi_u$  increasing,  $\lim_{y \rightarrow -\infty} \Pi_\ell(y) = 0$ ,  $\lim_{y \rightarrow \infty} \Pi_u(y) = 1$ .



## Conformal calibration guarantee:

We can construct a predictive system that contains a calibrated CDF.

## Example of in-sample calibration:

Let  $w_1, \dots, w_m \in \mathbb{R}$ . Define

$$F(y) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{w_i \leq y\}, \quad y \in \mathbb{R}.$$

Draw  $W$  uniformly at random from  $w_1, \dots, w_m$ .

Then  $F$  is *in-sample* probabilistically calibrated, that is,

$$\mathbb{P}(F(W) \leq \alpha) \leq \alpha \leq \mathbb{P}(F(W-) < \alpha), \quad \alpha \in (0, 1).$$

$$F(W) \approx \text{UNIF}(0, 1)$$

Let  $W_1, \dots, W_{n+1} \in \mathbb{R}$  be exchangeable and define for  $w \in \mathbb{R}$

$$F^w(y) = \frac{1}{n+1} \sum_{i=1}^n \mathbb{1}\{W_i \leq y\} + \frac{1}{n+1} \mathbb{1}\{w \leq y\}, \quad y \in \mathbb{R},$$

and

$$\Pi_\ell(y) = \inf\{F^w(y) \mid w \in \mathbb{R}\}, \quad \Pi_u(y) = \sup\{F^w(y) \mid w \in \mathbb{R}\},$$

Then,

$$\Pi_\ell(y) \leq F^{W_{n+1}}(y) \leq \Pi_u(y), \quad \text{and}$$

$$\mathbb{P}(F^{W_{n+1}}(W_{n+1}) \leq \alpha) \leq \alpha \leq \mathbb{P}(F^{W_{n+1}}(W_{n+1}-) < \alpha), \quad \alpha \in (0, 1).$$

Proof: Conditional on empirical distribution  $\hat{\mathbb{P}}_{n+1}$  of  $(W_i)_{i=1}^{n+1}$ ,  $W_{n+1}$  is a random draw from  $W_1, \dots, W_{n+1}$ . By in-sample probabilistic calibration:

$$\mathbb{P}(F^{W_{n+1}}(W_{n+1}) \leq \alpha \mid \hat{\mathbb{P}}_{n+1}) \leq \alpha \leq \mathbb{P}(F^{W_{n+1}}(W_{n+1}-) < \alpha \mid \hat{\mathbb{P}}_{n+1}) \dots$$

Let  $W_1, \dots, W_{n+1} \in \mathbb{R}$  be exchangeable and define for  $w \in \mathbb{R}$

$$F^w(y) = \frac{1}{n+1} \sum_{i=1}^n \mathbb{1}\{W_i \leq y\} + \frac{1}{n+1} \mathbb{1}\{w \leq y\}, \quad y \in \mathbb{R},$$

and

$$\Pi_\ell(y) = \inf\{F^w(y) \mid w \in \mathbb{R}\}, \quad \Pi_u(y) = \sup\{F^w(y) \mid w \in \mathbb{R}\},$$

Then,

$$\Pi_\ell(y) \leq F^{W_{n+1}}(y) \leq \Pi_u(y), \quad \text{and}$$

$$\mathbb{P}(F^{W_{n+1}}(W_{n+1}) \leq \alpha) \leq \alpha \leq \mathbb{P}(F^{W_{n+1}}(W_{n+1}-) < \alpha), \quad \alpha \in (0, 1).$$

Proof: Conditional on empirical distribution  $\hat{\mathbb{P}}_{n+1}$  of  $(W_i)_{i=1}^{n+1}$ ,  $W_{n+1}$  is a random draw from  $W_1, \dots, W_{n+1}$ . By in-sample probabilistic calibration:

$$\mathbb{P}(F^{W_{n+1}}(W_{n+1}) \leq \alpha \mid \hat{\mathbb{P}}_{n+1}) \leq \alpha \leq \mathbb{P}(F^{W_{n+1}}(W_{n+1}-) < \alpha \mid \hat{\mathbb{P}}_{n+1}) \dots$$

## (Classical) conformal prediction trick

Use conformity measure  $A(\hat{\mathbb{P}}, (x, y))$  to lift the one-dimensional result to general spaces  $\mathcal{X} \times \mathcal{Y}$ .

- ▶ Let  $(X_1, Y_1), \dots, (X_{n+1}, Y_{n+1}) \in \mathcal{X} \times \mathbb{R}$  be exchangeable.
- ▶ Apply one-dimensional result to

$$W_1 = A(\hat{\mathbb{P}}^y, (X_1, Y_1)), \dots, W_n = A(\hat{\mathbb{P}}^y, (X_n, Y_n)), w(y) = A(\hat{\mathbb{P}}^y, (X_{n+1}, y))$$

to obtain valid prediction sets.

- ▶ Predictive CDFs only available for specific conformity measures  $A$ .  
(Classical) conformal predictive systems

## (Classical) conformal prediction trick

Use conformity measure  $A(\hat{\mathbb{P}}, (x, y))$  to lift the one-dimensional result to general spaces  $\mathcal{X} \times \mathcal{Y}$ .

- ▶ Let  $(X_1, Y_1), \dots, (X_{n+1}, Y_{n+1}) \in \mathcal{X} \times \mathbb{R}$  be exchangeable.
- ▶ Apply one-dimensional result to

$$W_1 = A(\hat{\mathbb{P}}^y, (X_1, Y_1)), \dots, W_n = A(\hat{\mathbb{P}}^y, (X_n, Y_n)), w(y) = A(\hat{\mathbb{P}}^y, (X_{n+1}, y))$$

to obtain valid prediction sets.

- ▶ Predictive CDFs only available for specific conformity measures  $A$ .  
(Classical) conformal predictive systems

## Alternative

Use other **in-sample** calibrated procedures.

# Isotonic calibration

- ▶ Middle ground between probabilistic and auto-calibration
- ▶ Based on Isotonic Distributional Regression (IDR) (Henzi, Ziegel, and Gneiting, 2021)

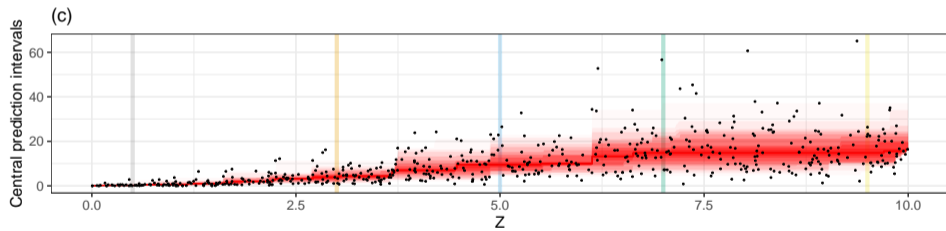
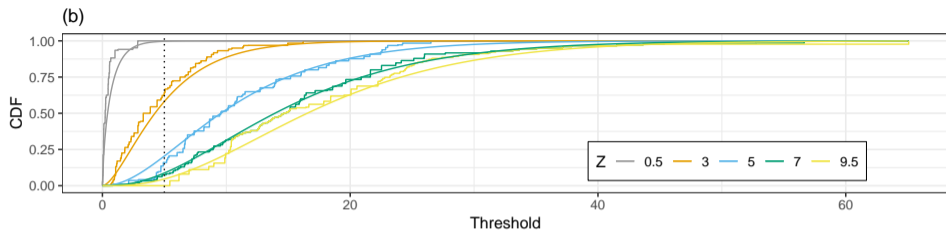
**IDR estimator** Let  $\leq$  be a partial order on  $\mathcal{X}$ .

Define  $\hat{\mathbf{F}} = (F_{x_k})_{k=1}^m$  as

$$\hat{\mathbf{F}} = \underset{F_i \preceq_{\text{st}} F_j \text{ if } x_i \leq x_j}{\operatorname{argmin}} \sum_{\ell=1}^m \operatorname{CRPS}(F_{\ell}, y_{\ell}).$$

Continuous ranked probability score (CRPS)

$$\operatorname{CRPS}(F, y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{y \leq z\})^2 dz$$



**Illustration of IDR:**  $n = 600$  draws of  $Z \sim \text{Unif}(0, 10)$  and  $Y \sim \text{Gamma}(\text{sh} = \sqrt{Z}, \text{sc} = \min(\max(Z, 1), 6))$ .

## Why IDR?

- ▶ Non-parametric distributional regression procedure under order constraints
- ▶ Explicit expression for estimator available
- ▶ Implementations available (R and Python)
- ▶ Consistency results available (under regularity conditions)

### Theorem (In-sample isotonic calibration of IDR)

*IDR is in-sample isotonically calibrated, more specifically,*

$$\hat{\mathbb{P}}_m(Y > y \mid \mathcal{A}(X)) = 1 - F_X^Y(y), \quad y \in \mathbb{R},$$

*and hence, in particular, threshold calibrated, quantile calibrated, and probabilistically calibrated. Here,  $(X, Y) \sim \hat{\mathbb{P}}_m$ , and  $\hat{\mathbb{P}}_m$  is the empirical distribution of  $(x_j, y_j)_{j=1}^m$ .*

Henzi, Ziegel, and Gneiting (2021); Arnold and Ziegel (2025)

- ▶ Choice: How is the partial order on  $\mathcal{X}$  constructed?

## Why IDR?

- ▶ Non-parametric distributional regression procedure under order constraints
- ▶ Explicit expression for estimator available
- ▶ Implementations available (R and Python)
- ▶ Consistency results available (under regularity conditions)

### Theorem (In-sample isotonic calibration of IDR)

*IDR is in-sample isotonically calibrated, more specifically,*

$$\hat{\mathbb{P}}_m(Y > y \mid \mathcal{A}(X)) = 1 - F_X^Y(y), \quad y \in \mathbb{R},$$

*and hence, in particular, threshold calibrated, quantile calibrated, and probabilistically calibrated. Here,  $(X, Y) \sim \hat{\mathbb{P}}_m$ , and  $\hat{\mathbb{P}}_m$  is the empirical distribution of  $(x_j, y_j)_{j=1}^m$ .*

Henzi, Ziegel, and Gneiting (2021); Arnold and Ziegel (2025)

- ▶ Choice: How is the partial order on  $\mathcal{X}$  constructed?

## Why IDR?

- ▶ Non-parametric distributional regression procedure under order constraints
- ▶ Explicit expression for estimator available
- ▶ Implementations available (R and Python)
- ▶ Consistency results available (under regularity conditions)

### Theorem (In-sample isotonic calibration of IDR)

*IDR is in-sample isotonically calibrated, more specifically,*

$$\hat{\mathbb{P}}_m(Y > y \mid \mathcal{A}(X)) = 1 - F_X^Y(y), \quad y \in \mathbb{R},$$

*and hence, in particular, threshold calibrated, quantile calibrated, and probabilistically calibrated. Here,  $(X, Y) \sim \hat{\mathbb{P}}_m$ , and  $\hat{\mathbb{P}}_m$  is the empirical distribution of  $(x_j, y_j)_{j=1}^m$ .*

Henzi, Ziegel, and Gneiting (2021); Arnold and Ziegel (2025)

- ▶ Choice: How is the partial order on  $\mathcal{X}$  constructed?

Let  $(X_1, Y_1), \dots, (X_{n+1}, Y_{n+1}) \in \mathcal{X} \times \mathbb{R}$  be exchangeable.

Let  $\Pi$  be constructed with IDR (*conformal IDR*):

- ▶ Let  $F_{X_k}^z$  be the IDR CDF computed from  $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, z)$ .
- ▶ Define

$$\Pi_{\ell, X_{n+1}}(y) = \inf\{F_{X_{n+1}}^z(y) \mid z \in \mathbb{R}\}, \quad \Pi_{u, X_{n+1}}(z) = \sup\{F_{X_{n+1}}^z(y) \mid z \in \mathbb{R}\},$$

### Theorem (Conformal calibration guarantee)

*Predictive system contains an isotonically calibrated CDF:*

$$F_{X_{n+1}}^{Y_{n+1}}(y) = 1 - \mathbb{P}(Y_{n+1} > y \mid \mathcal{A}(F_{X_{n+1}}^{Y_{n+1}})), \quad y \in \mathbb{R},$$

and

$$\Pi_{\ell, X_{n+1}}(y) \leq F_{X_{n+1}}^{Y_{n+1}}(y) \leq \Pi_{u, X_{n+1}}(y), \quad y \in \mathbb{R}$$

## Comments

- ▶ Conformal guarantee does not depend of any isotonicity assumption.
- ▶ The partial order on  $\mathcal{X}$  can be estimated on the same sample (computational challenge! “full conformal”) or on an independent sample (“split conformal”).
- ▶ Asymptotically,

$$F_{X_{n+1}}^{Y_{n+1}}(y) \rightarrow 1 - \mathbb{P}(Y_{n+1} > y \mid \mathcal{A}(X_{n+1})), \quad y \in \mathbb{R},$$

only depends on  $\mathcal{L}(Y \mid X)$  and not on  $\mathcal{L}(X)$ .  
Not true in general for conformal prediction.

## Comments

- ▶ Conformal guarantee does not depend of any isotonicity assumption.
- ▶ The partial order on  $\mathcal{X}$  can be estimated on the same sample (computational challenge! “full conformal”) or on an independent sample (“split conformal”).
- ▶ Asymptotically,

$$F_{X_{n+1}}^{Y_{n+1}}(y) \rightarrow 1 - \mathbb{P}(Y_{n+1} > y \mid \mathcal{A}(X_{n+1})), \quad y \in \mathbb{R},$$

only depends on  $\mathcal{L}(Y \mid X)$  and not on  $\mathcal{L}(X)$ .

Not true in general for conformal prediction.

## Auto-calibration

Let  $(x_1, y_1), \dots, (x_m, y_m) \in \mathcal{X} \times \mathbb{R}$ .

▶ Let  $B_1, \dots, B_{m'}$  be a partition of  $\{1, \dots, m\}$ .

▶

$$F_{x_k}(y) = \frac{1}{|B_i|} \sum_{j \in B_i} \mathbb{1}\{y_j \leq y\}, \quad k \in B_i, y \in \mathbb{R}$$

is in-sample auto-calibrated, that is,

$$\hat{\mathbb{P}}_m(Y \leq y | F_X) = F_X(y), \quad y \in \mathbb{R},$$

hence, in particular, isotonically calibrated, threshold calibrated, quantile calibrated, and probabilistically calibrated.

Here,  $(X, Y) \sim \hat{\mathbb{P}}_m$ , and  $\hat{\mathbb{P}}_m$  is the empirical distribution of  $(x_j, y_j)_{j=1}^m$ .

▶ We call this *conformal binning*. Closely related to Venn prediction.

▶ All in-sample auto-calibrated procedures are of this form.

▶ Choice: How is the partition constructed?

## Auto-calibration

Let  $(x_1, y_1), \dots, (x_m, y_m) \in \mathcal{X} \times \mathbb{R}$ .

▶ Let  $B_1, \dots, B_{m'}$  be a partition of  $\{1, \dots, m\}$ .



$$F_{x_k}(y) = \frac{1}{|B_i|} \sum_{j \in B_i} \mathbb{1}\{y_j \leq y\}, \quad k \in B_i, y \in \mathbb{R}$$

is **in-sample auto-calibrated**, that is,

$$\hat{\mathbb{P}}_m(Y \leq y | F_X) = F_X(y), \quad y \in \mathbb{R},$$

hence, in particular, isotonically calibrated, threshold calibrated, quantile calibrated, and probabilistically calibrated.

Here,  $(X, Y) \sim \hat{\mathbb{P}}_m$ , and  $\hat{\mathbb{P}}_m$  is the empirical distribution of  $(x_j, y_j)_{j=1}^m$ .

▶ We call this *conformal binning*. Closely related to Venn prediction.

▶ All in-sample auto-calibrated procedures are of this form.

▶ Choice: How is the partition constructed?

## Auto-calibration

Let  $(x_1, y_1), \dots, (x_m, y_m) \in \mathcal{X} \times \mathbb{R}$ .

▶ Let  $B_1, \dots, B_{m'}$  be a partition of  $\{1, \dots, m\}$ .

▶

$$F_{x_k}(y) = \frac{1}{|B_i|} \sum_{j \in B_i} \mathbb{1}\{y_j \leq y\}, \quad k \in B_i, y \in \mathbb{R}$$

is **in-sample auto-calibrated**, that is,

$$\hat{\mathbb{P}}_m(Y \leq y | F_X) = F_X(y), \quad y \in \mathbb{R},$$

hence, in particular, isotonically calibrated, threshold calibrated, quantile calibrated, and probabilistically calibrated.

Here,  $(X, Y) \sim \hat{\mathbb{P}}_m$ , and  $\hat{\mathbb{P}}_m$  is the empirical distribution of  $(x_j, y_j)_{j=1}^m$ .

▶ We call this *conformal binning*. Closely related to Venn prediction.

▶ All in-sample auto-calibrated procedures are of this form.

▶ Choice: How is the partition constructed?

## Auto-calibration

Let  $(x_1, y_1), \dots, (x_m, y_m) \in \mathcal{X} \times \mathbb{R}$ .

▶ Let  $B_1, \dots, B_{m'}$  be a partition of  $\{1, \dots, m\}$ .



$$F_{x_k}(y) = \frac{1}{|B_i|} \sum_{j \in B_i} \mathbb{1}\{y_j \leq y\}, \quad k \in B_i, y \in \mathbb{R}$$

is **in-sample auto-calibrated**, that is,

$$\hat{\mathbb{P}}_m(Y \leq y | F_X) = F_X(y), \quad y \in \mathbb{R},$$

hence, in particular, isotonically calibrated, threshold calibrated, quantile calibrated, and probabilistically calibrated.

Here,  $(X, Y) \sim \hat{\mathbb{P}}_m$ , and  $\hat{\mathbb{P}}_m$  is the empirical distribution of  $(x_j, y_j)_{j=1}^m$ .

- ▶ We call this *conformal binning*. Closely related to Venn prediction.
- ▶ All in-sample auto-calibrated procedures are of this form.
- ▶ Choice: How is the partition constructed?

# Thickness of predictive systems

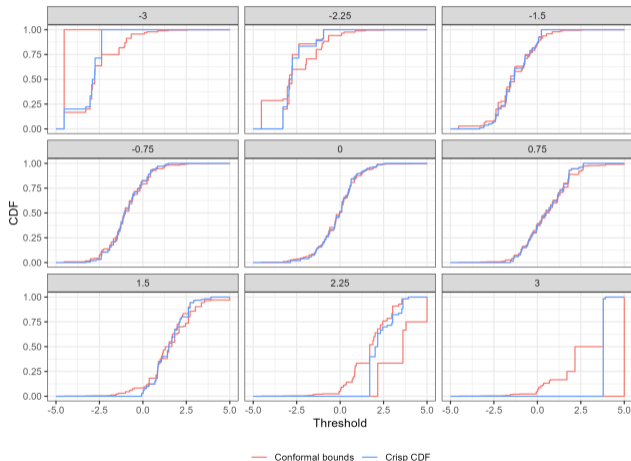
- ▶ Predictive systems are only useful if they are thin.
- ▶ Classical conformal predictive systems
  - ▶ Thickness is  $1/(n + 1)$ .
- ▶ Auto-calibration

If partition is determined based on  $X_1, \dots, X_{n+1}$  (example:  $k$ -means clustering),

  - ▶ Thickness is  $1/(\text{size of partition element containing } n + 1)$ .
- ▶ Isotonic calibration with IDR
  - ▶ Expected thickness is less or equal to  $14n^{-1/6}$ .

# Tiny simulation example for conformal IDR

$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(X, 1), n = 512.$$



- ▶ Principled approach to choose a crisp conformal IDR.
- ▶ Expected thickness goes to zero asymptotically.
- ▶ Thickness of conformal IDR informs about epistemic uncertainty.

# Aleatoric and epistemic uncertainty

## Aleatoric uncertainty

Aleatoric uncertainty of future outcome  $Y$  is fully described by

$$\mathcal{L}(Y | X).$$

Uncertainty remains even with infinite amounts of data  $(X_i, Y_i)$ .

## Epistemic uncertainty (second order probabilities, ambiguity, ...)

Uncertainty due to our approximation of  $\mathcal{L}(Y | X)$  based on limited data, limited knowledge of data generating process, parameter estimation, ...

Uncertainty goes away if we have infinite amounts of data.

- ▶ With IDR we recover  $\mathcal{L}(Y | \mathcal{A}(X))$ .

# Aleatoric and epistemic uncertainty

## Aleatoric uncertainty

Aleatoric uncertainty of future outcome  $Y$  is fully described by

$$\mathcal{L}(Y | X).$$

Uncertainty remains even with infinite amounts of data  $(X_i, Y_i)$ .

## Epistemic uncertainty (second order probabilities, ambiguity, ...)

Uncertainty due to our approximation of  $\mathcal{L}(Y | X)$  based on limited data, limited knowledge of data generating process, parameter estimation, ...

Uncertainty goes away if we have infinite amounts of data.

- ▶ With IDR we recover  $\mathcal{L}(Y | \mathcal{A}(X))$ .

# Aleatoric and epistemic uncertainty

## Aleatoric uncertainty

Aleatoric uncertainty of future outcome  $Y$  is fully described by

$$\mathcal{L}(Y | X).$$

Uncertainty remains even with infinite amounts of data  $(X_i, Y_i)$ .

## Epistemic uncertainty (second order probabilities, ambiguity, ...)

Uncertainty due to our approximation of  $\mathcal{L}(Y | X)$  based on limited data, limited knowledge of data generating process, parameter estimation, ...

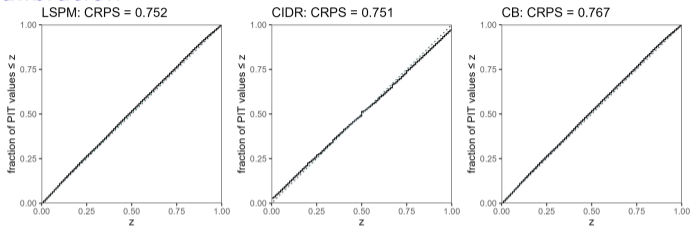
Uncertainty goes away if we have infinite amounts of data.

- ▶ With IDR we recover  $\mathcal{L}(Y | \mathcal{A}(X))$ .

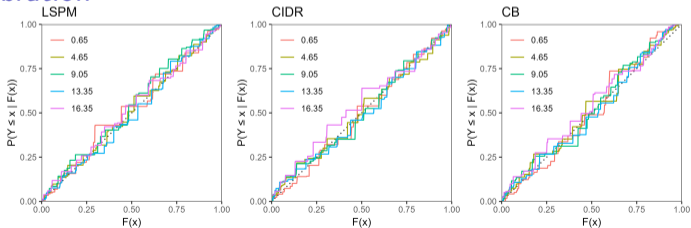
## Case study 1: Temperature forecasts

- ▶ Forecasts of daily mean temperature
- ▶ Post-processing of predictions issued by numerical weather prediction models
- ▶ EUPPBench dataset: Benchmark data set for comparing statistical post-processing methods.
- ▶ 20 years of reforecasts used for training (calibration data)
- ▶ 2017-2018 daily forecasts and observations as test data
- ▶ Forecasts are ensemble mean forecasts of the ECMWF; additional covariates station and season

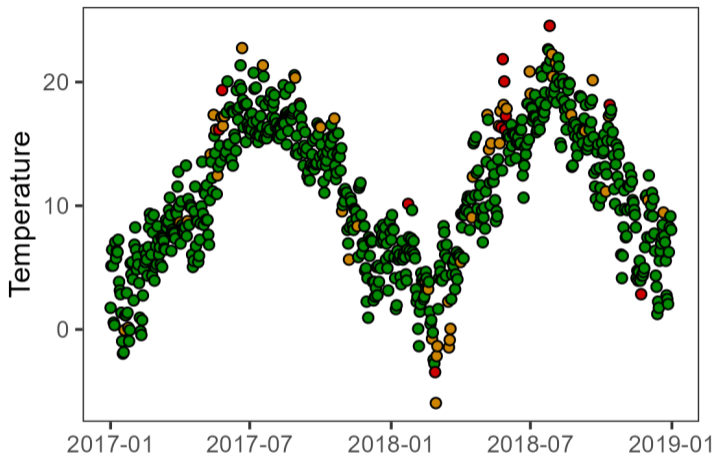
## Probabilistic calibration



## Threshold calibration

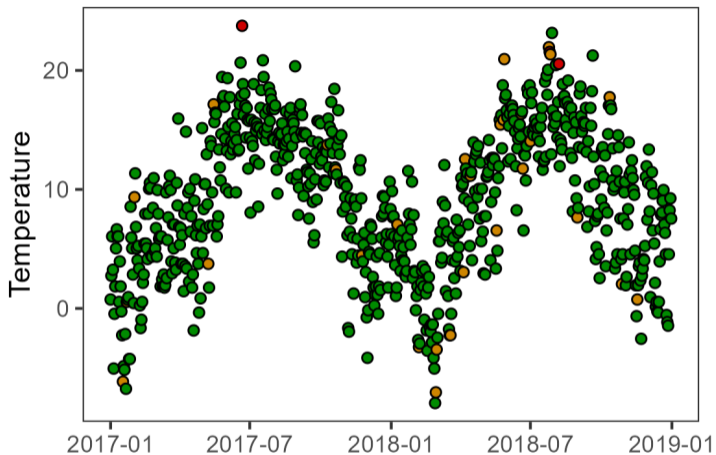


## Epsitemic uncertainty (CIDR) – First attempt



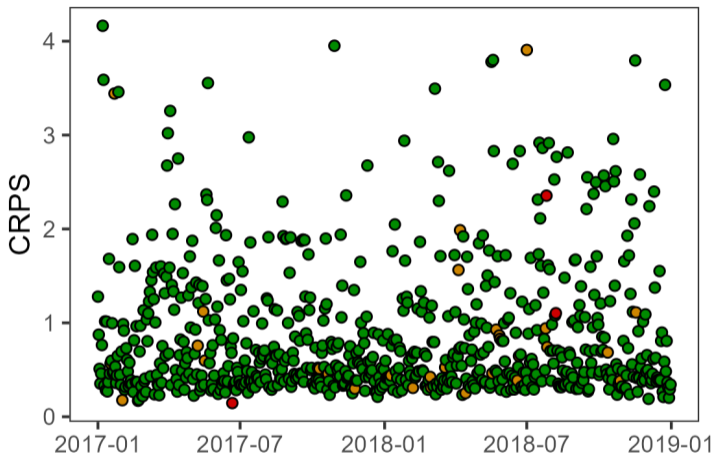
Thickness ● Low ● Medium ● High

## Epsitemic uncertainty (CIDR)



Thickness ● Low ● Medium ● High

## Epsitemic uncertainty (CIDR)

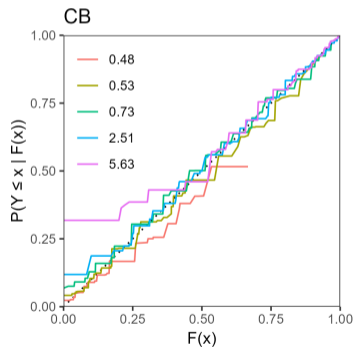
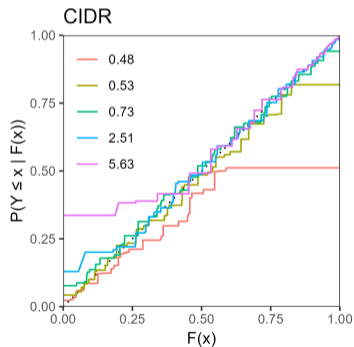
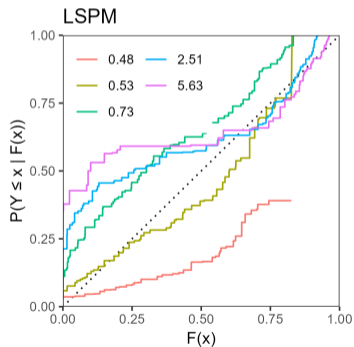


Thickness ● Low ● Medium ● High

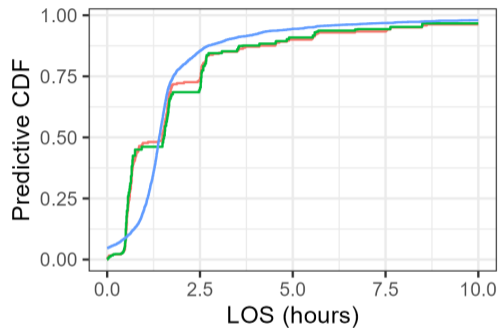
## Case study 2: Length of stay in intensive care units

- Predictions for individual patients' length of stay in ICU's in Switzerland 24h after admission

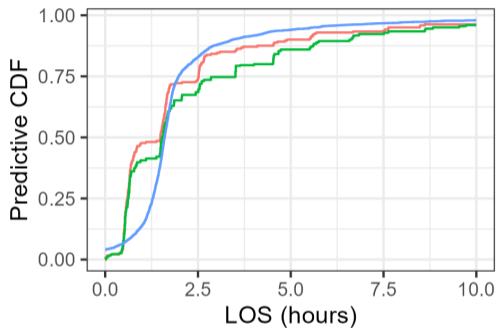
### Threshold calibration



## Examples of predictive cdfs



— CB — CIDR — LSPM



— CB — CIDR — LSPM

# Summary

- ▶ In-sample calibration yields conformal calibration guarantees.
- ▶ Strong out-of-sample calibration guarantees are possible.
- ▶ Arguments can be extended to distribution shifts (work in progress).
- ▶ Conformal binning is simple but works well.  
We used  $k$ -means clustering.
- ▶ Conformal IDR allows to quantify epistemic uncertainty, since IDR converges to a well-understood limiting object.

## References

- S. Allen, G. Gavrilopoulos, A. Henzi, G.-R. Kleger, and J. Ziegel. In-sample calibration yields conformal calibration guarantees. *Preprint, arXiv: 2503. 03841*, 2025a.
- S. Allen, J. Koh, J. Segers, and J. Ziegel. Tail calibration of probabilistic forecasts. *Journal of the American Statistical Association*, 120:2796–2808, 2025b.
- S. Arnold and J. Ziegel. Isotonic conditional laws. *Bernoulli*, 31:1140–1159, 2025.
- A. Henzi, J. F. Ziegel, and T. Gneiting. Isotonic distributional regression. *Journal of the Royal Statistical Society: Series B*, 85:963–993, 2021.
- K. Waghamare and J. Ziegel. Proper scoring rules for estimation and forecast evaluation. *Annual Review of Statistics and Its Application*, 13:271–296, 2026.

Thank you!

# Why the CRPS?

It is a strictly proper scoring rule.

If  $Y \sim F$  and  $G$  is any other CDF, then  $S(F, y)$  is *strictly proper* if

$$\mathbb{E}_F S(F, Y) \leq \mathbb{E}_F S(G, Y)$$

with equality if and only if  $F = G$ .

## Example 1

If  $F, G$  have finite mean, then the CRPS

$$\text{CRPS}(F, Y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{Y \leq z\})^2 dz$$

is strictly proper.

## Example 2

If  $F, G$  have densities  $f, g$ , then the logarithmic score

$$S_{\log}(F, y) = -\log f(y)$$

is strictly proper.

## Why the CRPS?

It is a strictly proper scoring rule.

If  $Y \sim F$  and  $G$  is any other CDF, then  $S(F, y)$  is *strictly proper* if

$$\mathbb{E}_F S(F, Y) \leq \mathbb{E}_F S(G, Y)$$

with equality if and only if  $F = G$ .

### Example 1

If  $F, G$  have finite mean, then the CRPS

$$\text{CRPS}(F, Y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{Y \leq z\})^2 dz$$

is strictly proper.

### Example 2

If  $F, G$  have densities  $f, g$ , then the logarithmic score

$$S_{\log}(F, y) = -\log f(y)$$

is strictly proper.

## Why the CRPS?

It is a strictly proper scoring rule.

If  $Y \sim F$  and  $G$  is any other CDF, then  $S(F, y)$  is *strictly proper* if

$$\mathbb{E}_F S(F, Y) \leq \mathbb{E}_F S(G, Y)$$

with equality if and only if  $F = G$ .

### Example 1

If  $F, G$  have finite mean, then the CRPS

$$\text{CRPS}(F, Y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{Y \leq z\})^2 dz$$

is strictly proper.

### Example 2

If  $F, G$  have densities  $f, g$ , then the logarithmic score

$$S_{\log}(F, y) = -\log f(y)$$

is strictly proper.

## Mathematical setup

“If the **covariate increases** we expect an **increase of the outcome**.”

$$x \leq x' \implies \mathcal{L}(Y | X = x) \preceq_{\text{st}} \mathcal{L}(Y | X = x')$$

$$\iff F_{Y|X=x}(y) \geq F_{Y|X=x'}(y), \quad y \in \mathbb{R}$$

$$\iff q_{\alpha}(Y|X = x) \leq q_{\alpha}(Y|X = x'), \quad \alpha \in (0, 1)$$

**IDR estimator** (for  $x \in \mathbb{R}$ ): Data  $(x_i, y_i)_{i=1}^n$ ,  $x_1 < \dots < x_n$

Define  $\hat{\mathbf{F}} = (\hat{F}_i)_{i=1}^n = (\hat{F}_{Y|X=x_i})_{i=1}^n$  as

$$\hat{\mathbf{F}} = \underset{F_1 \preceq_{\text{st}} \dots \preceq_{\text{st}} F_n}{\operatorname{argmin}} \sum_{\ell=1}^n \text{CRPS}(F_{\ell}, y_{\ell}).$$

Continuous ranked probability score (CRPS)

$$\text{CRPS}(F, Y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{Y \leq z\})^2 dz$$

## Mathematical setup

“If the **covariate increases** we expect an **increase of the outcome**.”

$$x \leq x' \implies \mathcal{L}(Y | X = x) \preceq_{\text{st}} \mathcal{L}(Y | X = x')$$

$$\iff F_{Y|X=x}(y) \geq F_{Y|X=x'}(y), \quad y \in \mathbb{R}$$

$$\iff q_{\alpha}(Y|X = x) \leq q_{\alpha}(Y|X = x'), \quad \alpha \in (0, 1)$$

**IDR estimator** (for  $x \in \mathbb{R}$ ): Data  $(x_i, y_i)_{i=1}^n$ ,  $x_1 < \dots < x_n$

Define  $\hat{\mathbf{F}} = (\hat{F}_i)_{i=1}^n = (\hat{F}_{Y|X=x_i})_{i=1}^n$  as

$$\hat{\mathbf{F}} = \underset{F_1 \preceq_{\text{st}} \dots \preceq_{\text{st}} F_n}{\operatorname{argmin}} \sum_{\ell=1}^n \operatorname{CRPS}(F_{\ell}, y_{\ell}).$$

Continuous ranked probability score (CRPS)

$$\operatorname{CRPS}(F, Y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{Y \leq z\})^2 dz$$

Then

$$\hat{F}_{Y|X=x_i} = \hat{F}_i(y) = \max_{j=i, \dots, n} \min_{k=1, \dots, j} \frac{1}{j - k + 1} \sum_{\ell=k}^j \mathbb{1}\{y_\ell \leq y\}.$$

- ▶  $\hat{F}_1(y), \dots, \hat{F}_n(y)$  is the antitonic regression of the binary outcomes  $\mathbb{1}\{y_1 \leq y\}, \dots, \mathbb{1}\{y_n \leq y\}$ .

