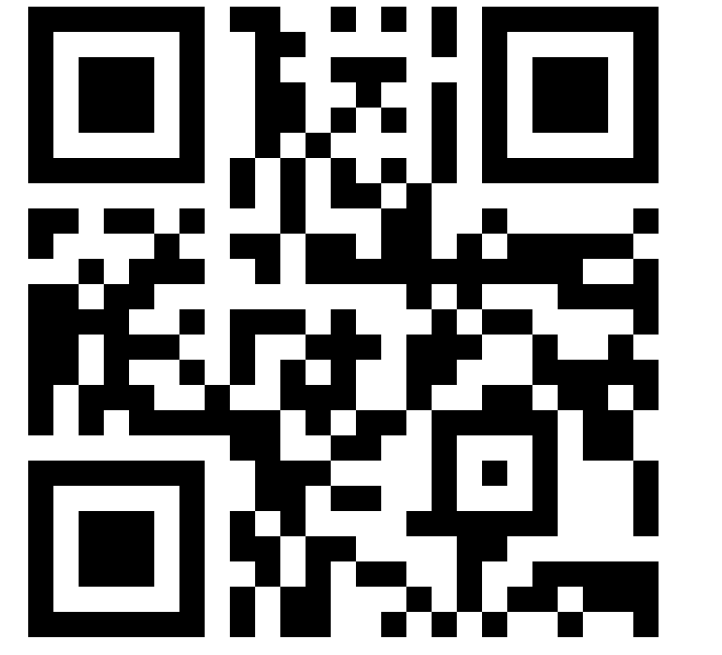


# CONDITIONAL COVERAGE DIAGNOSTICS FOR CONFORMAL PREDICTION

Sacha Braun<sup>1,2</sup>, David Holzmüller<sup>3</sup>, Michael I. Jordan<sup>1,4</sup>, Francis Bach<sup>1,2</sup>

<sup>1</sup>Sierra team, Inria Paris   <sup>2</sup>Ecole Normale Supérieure, PSL   <sup>3</sup>Soda team, Inria Paris-Saclay   <sup>4</sup>UC Berkeley



## 1. The goal: conditional coverage metrics

We want to estimate (among other metrics)

$$\mathbb{E}_X \left[ \left| \mathbb{P}(Y \in C_\alpha(X) \mid X) - (1 - \alpha) \right| \right].$$

## 2. Predicting conditional coverage through probabilistic classification

- Let  $Z = \mathbf{1}\{Y \in C_\alpha(X)\}$ .
- Train a classifier  $h(X)$  to predict  $Z$ .
- Optimal predictor for this problem is  $p(X) = \mathbb{P}(Y \in C_\alpha(X) \mid X)$ .

## 3. Compare validation losses

For a proper loss  $\ell$ , the **Excess Risk of the Target coverage ( $\ell$ -ERT)** is:

$$\ell\text{-ERT}(p) := \mathbb{E}[\ell(1 - \alpha, Z)] - \mathbb{E}[\ell(p(X), Z)].$$

Under conditional coverage, no predictor performs better than the constant  $1 - \alpha$ .

### The lower bound guarantee with cross-validation

For any learned classifier  $h$ , we get a **lower bound** on the true miscoverage:

$$\ell\text{-ERT}(h) \leq \ell\text{-ERT}(p).$$

We estimate this quantity using the empirical risk on held-out data, so the lower bound holds in expectation.

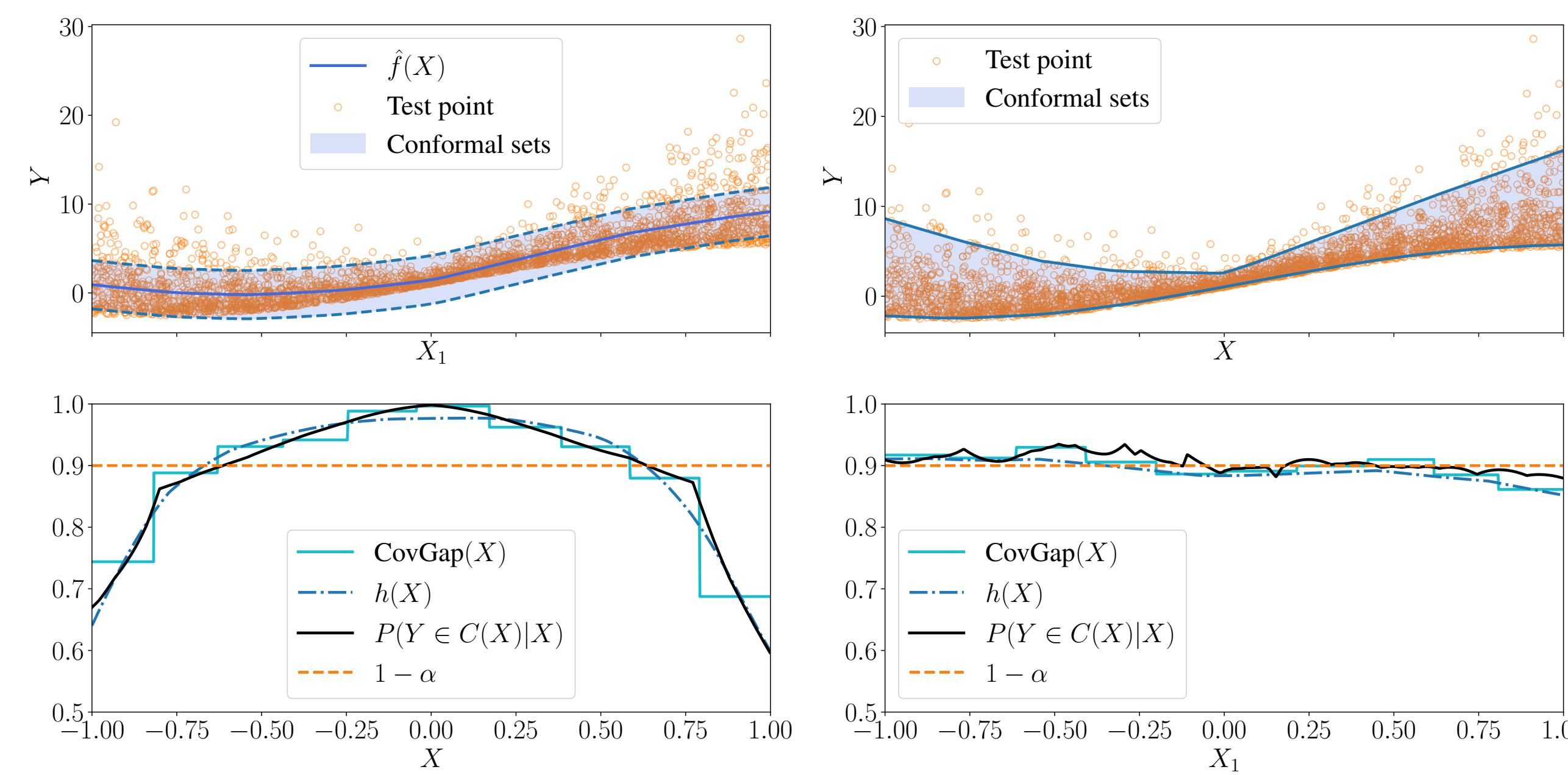
## 4. One loss for every distance function

Different proper scoring rules  $\ell$  allow us to estimate specific, interpretable distances to perfect coverage:

Metric	Proper score $\ell(p, z)$	Estimated target quantity ( $\ell$ -ERT)
$L_1$ -ERT	$\text{sgn}(p - (1 - \alpha))(1 - \alpha - z)$	$\mathbb{E}_X[ 1 - \alpha - p(X) ]$
$L_2$ -ERT	$(z - p)^2$	$\mathbb{E}_X[(1 - \alpha - p(X))^2]$
KL-ERT	$-z \log p - (1 - z) \log(1 - p)$	$\mathbb{E}_X[D_{\text{KL}}(p(X) \parallel 1 - \alpha)]$
$d$ -ERT	$-f_\alpha(p) - (y - p)f'_\alpha(p)$	$\mathbb{E}_X[d(1 - \alpha, p(X))]$

where  $f_\alpha(q) = d(1 - \alpha, q)$  is convex with  $f_\alpha(1 - \alpha) = 0$  and  $f'_\alpha$  is a subderivative of  $f_\alpha$  satisfying  $f'_\alpha(1 - \alpha) = 0$ .

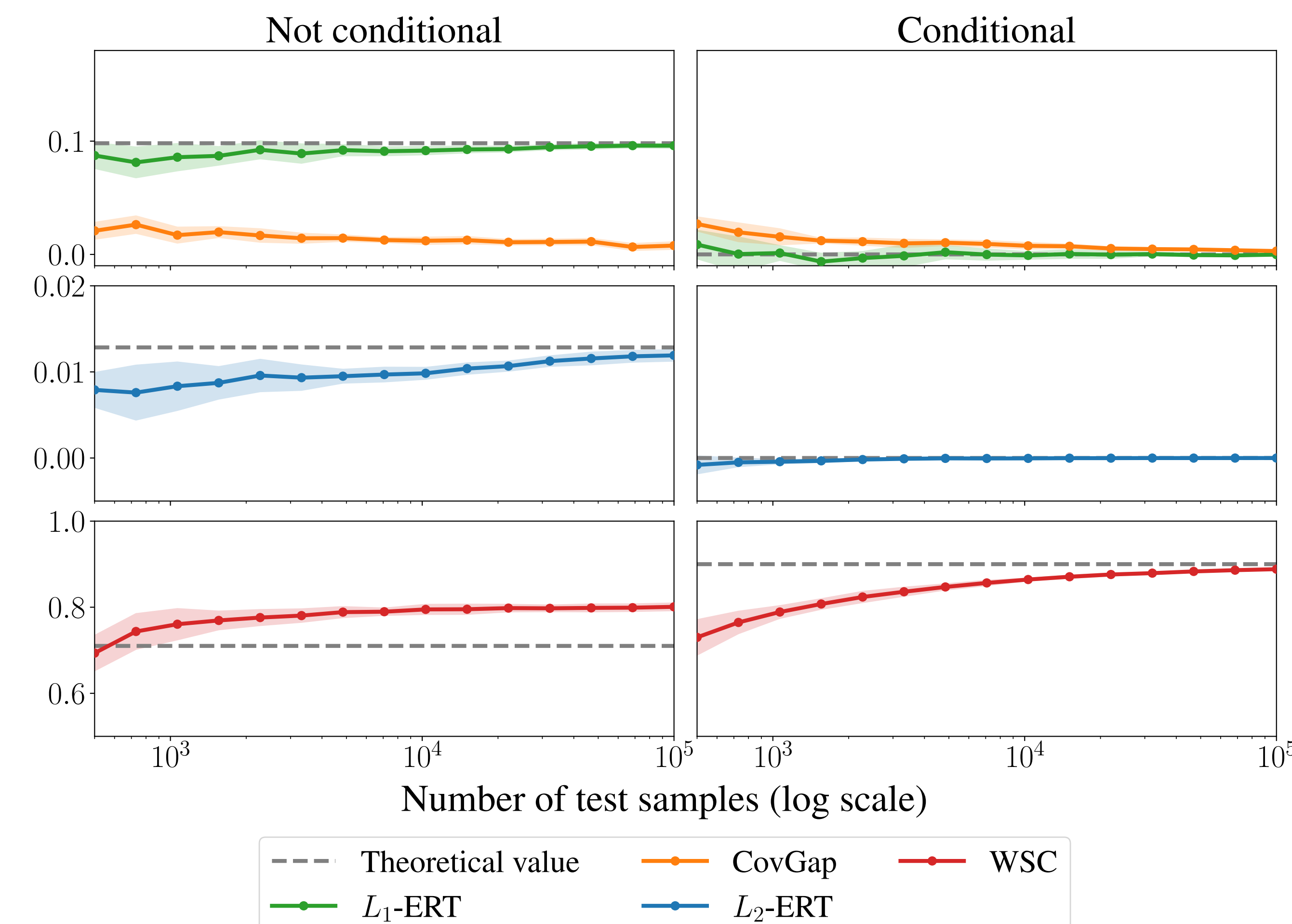
## 5. Visualizing ERT in 1D



$h(\cdot)$  is a neural network. Left: Non-conditional CP fails locally;  $h(X)$  captures the variance ( $L_1$ -ERT  $\approx 0.075 > 0$ ). Right: More conditional sets;  $h(X)$  stays flat at  $1 - \alpha$  ( $L_1$ -ERT  $\approx 0$ ).

## 6. Sample efficiency vs. baselines

ERT stabilizes with far fewer samples than WSC or CovGap, providing a reliable diagnostic rapidly without false positives.



$X \sim \mathcal{U}([-1, 1]^8)$  and  $Y \sim \mathcal{N}(0, \sigma(X_1))$ .

Left: Marginal but not conditional setting. ERT detects failure instantly. Right: Perfectly conditional setting (oracle). WSC hallucinated failures; ERT correctly stays near zero.

## 7. Estimating under- and over-coverage

By strategically clipping the predicted probability  $p$  within the proper loss function, we can decompose the miscoverage to isolate specific failures:

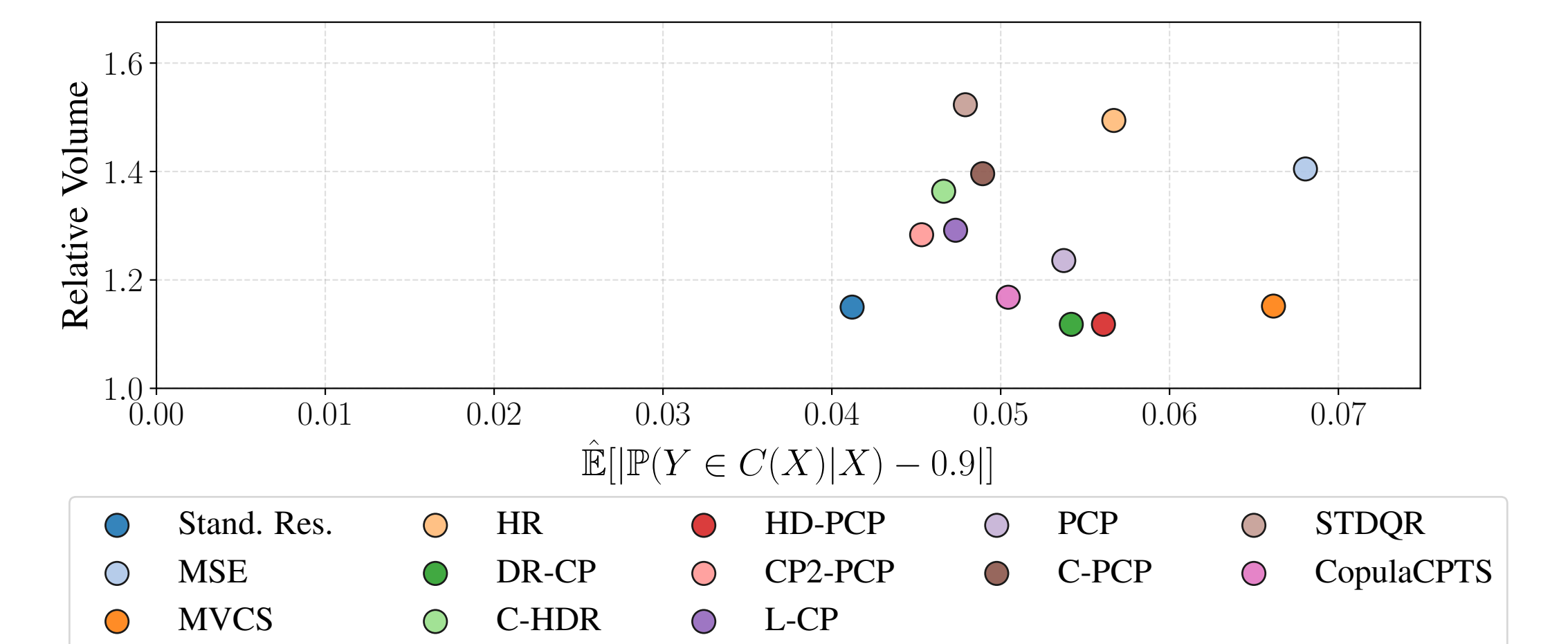
- $\ell_+$ -ERT (over-coverage): Computed using  $\ell_+(p, y) = \ell(\max\{p, 1 - \alpha\}, y)$ .
- $\ell_-$ -ERT (under-coverage): Computed using  $\ell_-(p, y) = \ell(\min\{p, 1 - \alpha\}, y)$ .

## 8. Modern classifiers beat simple ones

Classifier	$L_1$ -ERT	$L_2$ -ERT	Time / 1K (s)
TabICLv1.1 (GPU)	<b>72.7</b> <sub>1.6</sub>	<b>54.0</b> <sub>1.4</sub>	16.0 <sub>0.0</sub>
RealTabPFN-2.5 (GPU)	72.1 <sub>1.4</sub>	49.7 <sub>1.1</sub>	9.2 <sub>0.0</sub>
LightGBM (CPU)	68.9 <sub>1.8</sub>	46.4 <sub>1.0</sub>	2.4 <sub>0.0</sub>
RandomForest (CPU)	67.3 <sub>2.0</sub>	36.9 <sub>2.4</sub>	4.4 <sub>0.0</sub>
PartitionWise (CPU)	33.7 <sub>1.9</sub>	10.5 <sub>0.7</sub>	<b>0.2</b> <sub>0.0</sub>

Average % of the maximum ERT recovered across datasets, for different number of test samples. Traditional partition-wise estimators (like CovGap) capture less than 40% of the actual miscoverage compared to gradient-boosted trees and foundation models.

## 9. Conformal prediction scores benchmark



Comparison of conditional coverage deviation and normalized prediction set volumes across various conformal prediction strategies ( $\alpha = 0.1$ ). Volumes are scaled by the power  $1/\text{dim}_{\text{output}}$  and normalized by the minimum observed volume. The x-axis displays the  $L_1$ -ERT metric, which estimates the conditional coverage deviation  $\mathbb{E}[|\mathbb{P}(Y \in C(X) \mid X) - (1 - \alpha)|]$ .

## 10. Code & package :

```
pip install covmetrics
from covmetrics import ERT
ERT_value = ERT().evaluate(x, z, alpha)
```

