

What is conformal prediction?

- Conformal prediction is a generic method for finite-sample valid distribution-free prediction and can be applied with most machine learning algorithms to yield valid prediction regions.
- Given i.i.d. pairs $(X_i, Y_i) \sim P, i = 1, \dots, N$, for a distribution P on $\mathcal{X} \times \mathbb{R}$ (e.g., $\mathcal{X} = \mathbb{R}^d$)

Goal. Build a prediction set $\hat{C}_N : \mathcal{X} \rightarrow \mathcal{P}(\mathbb{R})$, such that for new i.i.d. pair (X_{N+1}, Y_{N+1}) :

$$\mathbb{P}\left(Y_{N+1} \in \hat{C}_N(X_{N+1})\right) \geq 1 - \alpha,$$

(where the probability is over all $N + 1$ pairs).

Split Conformal Prediction

- Split the sample $(X_i, Y_i), 1 \leq i \leq N$ into two parts each with $n = N/2$ observations.
- Compute the estimator $\hat{\mu}_m(\cdot)$ based on the **first** split.
- Let $\tilde{q}_{n,\alpha}$ denote the $(1 - \alpha)(1 + 1/n)$ -th quantile of the residuals $|Y_i - \hat{\mu}_m(X_i)|$, on the **second** split.
- $\tilde{C}_N := \{(x, y) : y \in [\hat{\mu}_m(x) - \tilde{q}_{n,\alpha}, \hat{\mu}_m(x) + \tilde{q}_{n,\alpha}]\}$

By exchangeability, we have finite-sample coverage:

$$\mathbb{P}\left((X_{N+1}, Y_{N+1}) \in \tilde{C}_N\right) \geq 1 - \alpha.$$

This is valid regardless of whether $\hat{\mu}_m$ is a consistent estimator. In practice, the coverage is close to $1 - \alpha$.

Problem Formulation

- Suppose $\hat{\mu}_m^1(\cdot), \dots, \hat{\mu}_m^K(\cdot)$ are K different estimators with different tuning parameters. They can be from different estimation methods such as LASSO, random forest.
- ★ The question we discuss is how to select $\hat{k} \in \{1, 2, \dots, K\}$ and construct a **valid** prediction region with **width close to the smallest**.
- We will call coverage guarantee as **validity** and smallest width property as **efficiency**.

Efficient First Conformal Prediction (EFCP)

Input: Data $(X_i, Y_i), 1 \leq i \leq N$, coverage probability $1 - \alpha$, and K estimation methods for $\mu(\cdot)$.

- Randomly split the data into three parts each with $n := N/3$ observations.
- Fit the estimators $\hat{\mu}_1(\cdot), \dots, \hat{\mu}_K(\cdot)$ on the **first** split of the data.
- $T_{\alpha,k} := (1 - \alpha)(1 + 1/n)$ -th quantile of $|Y_i - \hat{\mu}_k(X_i)|$, residuals in the **second** split of the data. The corresponding conformal prediction region is

$$\hat{C}_k := \{(x, y) : y \in [\hat{\mu}_k(x) \pm T_{\alpha,k}]\}.$$
- Set \hat{k} as the minimizer of $T_{\alpha,k}$ over $1 \leq k \leq K$.
- $T_{\alpha,\hat{k}}^* := (1 - \alpha)(1 + 1/n)$ -th quantile of $|Y_i - \hat{\mu}_{\hat{k}}(X_i)|$, residuals in the **third** split of the data.

6. **Output:**

$$\hat{C}_\alpha^{\text{VFPCP}} := \{(x, y) : y \in [\hat{\mu}_{\hat{k}}(x) \pm T_{\alpha,\hat{k}}^*]\}.$$

Validity First Conformal Prediction (VFPCP)

Input: Data $(X_i, Y_i), 1 \leq i \leq N$, coverage probability $1 - \alpha$, and K estimation methods for $\mu(\cdot)$.

- Randomly split the data into two parts each with $n := N/2$ observations.
- Fit the estimators $\hat{\mu}_1(\cdot), \dots, \hat{\mu}_K(\cdot)$ on the **first** split.
- $T_{\alpha,k} := (1 - \alpha)(1 + 1/n)$ -th quantile of $|Y_i - \hat{\mu}_k(X_i)|$, residuals in the **second** split of the data. The corresponding conformal prediction region is

$$\hat{C}_k := \{(x, y) : y \in [\hat{\mu}_k(x) \pm T_{\alpha,k}]\}.$$
- Set \hat{k} as the minimizer of $T_{\alpha,k}$ over $1 \leq k \leq K$.

5. **Output:**

$$\hat{C}_\alpha^{\text{EFCP}} := \hat{C}_{\hat{k}} = \{(x, y) : y \in [\hat{\mu}_{\hat{k}}(x) \pm T_{\alpha,\hat{k}}]\}.$$

Results and comparison of VFPCP and EFCP

	Coverage	Width
VFPCP	$1 - \alpha - 0$	min-width + $\mathfrak{c} \sqrt{\frac{\log(K/\delta)}{N}}$
EFCP	$1 - \alpha - \mathfrak{c} \sqrt{\frac{\log(K)}{N}}$	min-width + 0

- VFPCP requires **three** splits of the data while EFCP only requires **two** splits of the data.
- Validity and efficiency of VFPCP requires **some continuity assumptions** on the distributions, while EFCP only requires **i.i.d. data**.

Application: tuning-free ridge regression

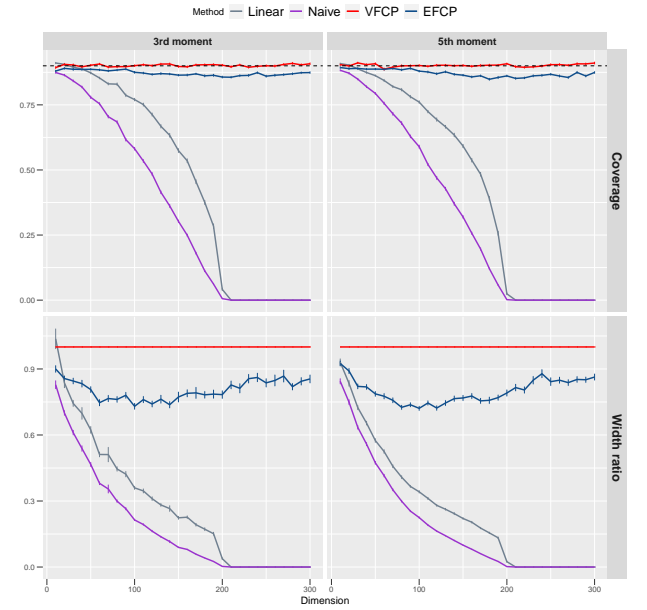


Figure 1: Ridge Regression with a Linear Model

Full paper and supplement (ArXiv):

