

# Conformal prediction and beyond

## Uncertainty quantification for regression & time-series problems

Gerard Castro

Universitat de Barcelona

July 11, 2024



UNIVERSITAT DE  
BARCELONA

Facultat de Matemàtiques  
i Informàtica

1 Introduction

2 Conformal prediction

3 Beyond exchangeability

4 Results

5 Conclusions

# Motivation

- We extract  $n$  samples from  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$  random variables with unknown marginal & joint distributions.
- Given a new sample  $X_{n+1}$  & miscoverage level  $\alpha \in [0, 1]$ :
  - We want to **estimate** a predictive **interval**  $\mathcal{C}_\alpha$  such that the probability of  $Y_{n+1}$  falling into  $\mathcal{C}_\alpha$  is at least  $1 - \alpha$ , *i.e.*

$$\mathbb{P}\{Y_{n+1} \in \mathcal{C}_\alpha(X_{n+1})\} \geq 1 - \alpha$$

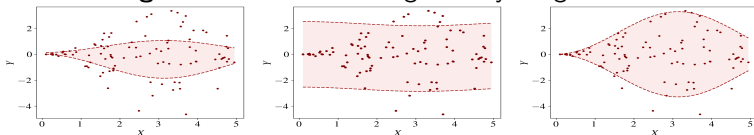
- The interval should be the **smallest** possible while **keeping coverage**. **Conditional** coverage ideally sought.

# Motivation

- We extract  $n$  samples from  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$  random variables with unknown marginal & joint distributions.
- Given a new sample  $X_{n+1}$  & miscoverage level  $\alpha \in [0, 1]$ :
  - We want to **estimate** a predictive **interval**  $\mathcal{C}_\alpha$  such that the probability of  $Y_{n+1}$  falling into  $\mathcal{C}_\alpha$  is at least  $1 - \alpha$ , *i.e.*

$$\mathbb{P}\{Y_{n+1} \in \mathcal{C}_\alpha(X_{n+1})\} \geq 1 - \alpha$$

- The interval should be the **smallest** possible while **keeping coverage**. **Conditional** coverage ideally sought.



(a) No coverage

(b) Marginal coverage

(c) Conditional cov.

**Figure 1:** Different types of coverage.

- ① Introduction
- ② Conformal prediction**
- ③ Beyond exchangeability
- ④ Results
- ⑤ Conclusions

# Split Conformal Prediction (SCP)

We need to use out-of-training data to understand how errors distribute: we need to "*conformalize*" the predictions to the data using a "*conformity score*". SCP proposes:

- 1 Split data into **training**  $\text{Tr}$  & **calibration**  $\text{Cal}$ .
- 2 Obtain  $\hat{\mu}$  by training it in  $\text{Tr}$ .
- 3 Obtain a set  $\mathcal{S}$  of conformity scores by using the  $\text{Cal}$  set:  
$$\mathcal{S}_{\text{Cal}} := \{|Y_i - \hat{\mu}(X_i)|, i \in \text{Cal}\}.$$
- 4 Compute the  $1 - \alpha$  "*empirical quantile*" of  $\mathcal{S}_{\text{Cal}}$ :  $q_{1-\alpha}(\mathcal{S}_{\text{Cal}})$ .

# Split Conformal Prediction (SCP)

We need to use out-of-training data to understand how errors distribute: we need to "*conformalize*" the predictions to the data using a "*conformity score*". SCP proposes:

- 1 Split data into **training**  $\text{Tr}$  & **calibration**  $\text{Cal}$ .
- 2 Obtain  $\hat{\mu}$  by training it in  $\text{Tr}$ .
- 3 Obtain a set  $\mathcal{S}$  of conformity scores by using the  $\text{Cal}$  set:  
$$\mathcal{S}_{\text{Cal}} := \{|Y_i - \hat{\mu}(X_i)|, i \in \text{Cal}\}.$$
- 4 Compute  $(1 - \alpha) \left( \frac{1}{\#\text{Cal}} + 1 \right)$  quantile of  $\mathcal{S}_{\text{Cal}}$ :  $q_{1-\alpha}(\mathcal{S}_{\text{Cal}})$ .

# Split Conformal Prediction (SCP)

We need to use out-of-training data to understand how errors distribute: we need to "*conformalize*" the predictions to the data using a "*conformity score*". SCP proposes:

- 1 Split data into **training** Tr & **calibration** Cal.
- 2 Obtain  $\hat{\mu}$  by training it in Tr.
- 3 Obtain a set  $\mathcal{S}$  of conformity scores by using the Cal set:  
$$\mathcal{S}_{\text{Cal}} := \{|Y_i - \hat{\mu}(X_i)|, i \in \text{Cal}\}.$$
- 4 Compute  $(1 - \alpha) \left( \frac{1}{\#\text{Cal}} + 1 \right)$  quantile of  $\mathcal{S}_{\text{Cal}}$ :  $q_{1-\alpha}(\mathcal{S}_{\text{Cal}})$ .
- 5 For a new sample  $X_{n+1}$ , return

$$\hat{\mathcal{C}}_\alpha = [\hat{\mu}(X_{n+1}) - q_{1-\alpha}(\mathcal{S}_{\text{Cal}}), \hat{\mu}(X_{n+1}) + q_{1-\alpha}(\mathcal{S}_{\text{Cal}})]$$

## Note

The only hypothesis required is **data exchangeability**.



# Conformalized Quantile Regression (CQR)

We need to use out-of-training data to understand how errors distribute: we need to "conformalize" the predictions to the data using a "conformity score".  $\mathcal{SCP}$  CQR proposes:

- 1 Split data into **training**  $\text{Tr}$  & **calibration**  $\text{Cal}$ .
- 2 Obtain  $\hat{\mu}_{\text{down}}$  &  $\hat{\mu}_{\text{up}}$  trained in  $\text{Tr}$ .
- 3 Obtain a set  $\mathcal{S}$  of conformity scores by using the  $\text{Cal}$  set:  
$$\mathcal{S}_{\text{Cal}} := \{\max(\hat{\mu}_{\text{down}}(X_i) - Y_i, Y_i - \hat{\mu}_{\text{up}}(X_i)), i \in \text{Cal}\}.$$
- 4 Compute  $(1 - \alpha) \left( \frac{1}{\#\text{Cal}} + 1 \right)$  quantile of  $\mathcal{S}_{\text{Cal}}$ :  $q_{1-\alpha}(\mathcal{S}_{\text{Cal}})$ .
- 5 For a new sample  $X_{n+1}$ , return

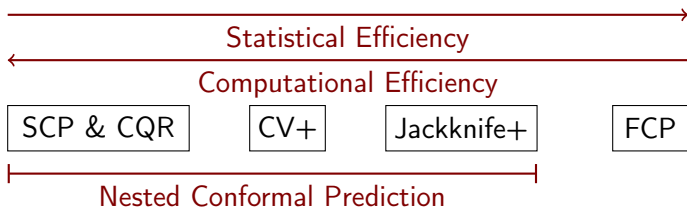
$$\hat{C}_\alpha(X_{n+1}) = [\hat{\mu}_{\text{down}}(X_{n+1}) - q_{1-\alpha}(\mathcal{S}_{\text{Cal}}), \hat{\mu}_{\text{up}}(X_{n+1}) + q_{1-\alpha}(\mathcal{S}_{\text{Cal}})]$$

## Note

The only hypothesis required is **data exchangeability**.

## Other flavours

Then, other flavours are proposed to reconcile this trade-off between statistical and computational efficiency, for instance CV+ & J+aB:



**Figure 2:** Trade-off between statistical & computational efficiency.

- Both CV+ & J+aB are based on defining multiple folds to apply a similar methodology as SCP: cross-validation & leave-one-out (LOO) folds, respectively.

- ① Introduction
- ② Conformal prediction
- ③ Beyond exchangeability**
- ④ Results
- ⑤ Conclusions

# Covariate shift: changes in features' distribution

- $\{(X_i, Y_i)\}_{i=1}^n \stackrel{\text{exch.}}{\sim} P_X \times P_{Y|X}$
- $(X_{n+1}, Y_{n+1}) \sim \tilde{P}_X \times P_{Y|X}$
- *Tibshirani et al. (2019)* heuristic idea:
  - 1 Estimate how "close" a sample  $X_i (\sim P_X)$  is *w.r.t.* to the test point  $(\sim \tilde{P}_X)$  using the likelihood ratio:  $w(X_i) := \frac{d\tilde{P}_X(X_i)}{dP_X(X_i)}$ .
  - 2 Normalize the weights:  $\omega_i := \frac{w(X_i)}{\sum_{j=1}^{n+1} w(X_j)}$ .
  - 3 Build the predictive interval  $\mathcal{C}_\alpha$  using the weighted calibration samples:

$$\hat{\mathcal{C}}_\alpha(X_{n+1}) = \{Y : s_{\hat{\mu}}(X_{n+1}, Y) \leq q_{1-\alpha}(\{\omega_i S_i\}_{i \in \text{Cal}})\}$$

## Label shift: changes in target's distribution

- $\{(X_i, Y_i)\}_{i=1}^n \stackrel{\text{exch.}}{\sim} P_{X|Y} \times P_Y$
- $(X_{n+1}, Y_{n+1}) \sim P_{X|Y} \times \tilde{P}_Y$
- *A. Podkopaev & A. Ramdas (2021)* adapts former idea letting weights as function of  $Y$ ,  $\omega_i^Y$ :
  - 1 Estimate how "close" a label  $Y_i$  ( $\sim P_Y$ ) is *w.r.t.* to the hypothetical point ( $\sim \tilde{P}_Y$ ) using the likelihood ratio:  

$$w(Y_i) := \frac{d\tilde{P}_Y(Y_i)}{dP_Y(Y_i)}.$$
  - 2 Normalize the weights:  $\omega_i^Y := \frac{w(Y_i)}{\sum_{j=1}^n w(Y_j) + w(Y)}$ .
  - 3 Build the predictive interval  $\mathcal{C}_\alpha$  traversing all the variable output's space and using the weighted calibration samples:

$$\hat{\mathcal{C}}_\alpha(X_{n+1}) = \{Y : s_{\hat{\mu}}(X_{n+1}, Y) \leq q_{1-\alpha}(\{\omega_i^Y S_i\}_{i \in \text{Cal}})\}$$

# Time series data: samples (temporal) auto-correlation

- Assume a setup like  $Y_t = \mu(X_t) + \epsilon_t$ , where  $\epsilon_t$  are *i.i.d.* according to a cumulative distribution function  $F$ .
- Let the first  $T$  sample points  $\mathcal{D} := \{(X_t, Y_t)_{t=1}^T\}$  be training data: we **want** a **sequence** of  $s \geq 1$  **intervals** of  $\alpha$  miscoverage level,  $\{C_{T, T+i}^\alpha\}_{i=1}^s$  (for the unknown labels  $\{Y_{T+i}^\alpha\}_{i=1}^s$ ).
  - $s$  is the batch size ( $n^{\text{Q}}$  steps to look ahead)
- Also, once **new samples**  $\{(X_{T+i}, Y_{T+i})\}_{i=1}^s$  become **available**, we would like to also **leverage them**.
  - We want to use the most recent  $T + s$  points for the  $\{C_{T+s, j}^\alpha\}_{j=T+s+1}^{T+2s}$  intervals.

# Time series data: samples (temporal) auto-correlation

- Assume a setup like  $Y_t = \mu(X_t) + \epsilon_t$ , where  $\epsilon_t$  are *i.i.d.* according to a cumulative distribution function  $F$ .
- Let the first  $T$  sample points  $\mathcal{D} := \{(X_t, Y_t)_{t=1}^T\}$  be training data: we **want** a **sequence** of  $s \geq 1$  **intervals** of  $\alpha$  miscoverage level,  $\{C_{T, T+i}^\alpha\}_{i=1}^s$  (for the unknown labels  $\{Y_{T+i}^\alpha\}_{i=1}^s$ ).
  - $s$  is the batch size ( $n^{\text{Q}}$  steps to look ahead)
- Also, once **new samples**  $\{(X_{T+i}, Y_{T+i})\}_{i=1}^s$  become **available**, we would like to also **leverage them**.
  - We want to use the most recent  $T + s$  points for the  $\{C_{T+s, j}^\alpha\}_{j=T+s+1}^{T+2s}$  intervals.

⇒ C. Xu & Y. Xie (2021) proposes the "EnbPI" methodology:

- It uses no data-splitting but LOO estimators ( $\hat{\mu}_{-i}$  model trained with  $\mathcal{D} \setminus \{(X_i, Y_i)\}$ ).
- Models not refitted during test time, but newest samples' residuals used to further *conformalize* predictions.

## EnbPI idea

There are  $T$  training samples and we build  $T_1$  intervals (indices  $T + 1, \dots, T + T_1$ ):

- Obtain  $B$  bootstrapped models  $\mu^b$  by:
  - Sampling, with replacement, an index set  $S_b := (i_1, \dots, i_T)$
  - Fitting the bootstrapped model with  $S_b$
- For  $i = 1, \dots, T$ :
  - Aggregate  $\mu^b$  with any function  $\phi$ : obtaining  $\hat{\mu}_{-i}^\phi$ .
  - Compute conformity scores:  $\epsilon_i^\phi := |Y_i - \hat{\mu}_{-i}^\phi(X_i)|$ .
- For each  $t = T + 1, \dots, T + T_1$  timestamps, return in batches of  $s$  size:

$$\hat{C}_{T,t}^\alpha(X_t) = \left[ \hat{\mu}_{-t}^\phi(X_t) \pm w_t^\phi \right], \text{ where } \begin{cases} \hat{\mu}_{-t}^\phi(X_t) : 1 - \alpha \text{ quant. } \{ \hat{\mu}_{-i}^\phi(X_t) \}_{i=1}^T \\ w_t^\phi : 1 - \alpha \text{ quantile of } \{ \epsilon_i^\phi \}_{i=1}^T \end{cases}$$

- "*Partial fit*" step: for each  $s$  returned intervals, conformity score  $w_t^\phi$  is re-computed with the most recent observations.



- ① Introduction
- ② Conformal prediction
- ③ Beyond exchangeability
- ④ Results**
  - Assessment
  - Regression problem
  - Time series problem
- ⑤ Conclusions

- ① Introduction
- ② Conformal prediction
- ③ Beyond exchangeability
- ④ **Results**
  - Assessment
  - Regression problem
  - Time series problem
- ⑤ Conclusions

## Metrics definition

The following metrics will be used:

- **Coverage level:** *i.e.* fraction of true labels lying within the prediction intervals (the closer to  $1 - \alpha$ , the better)
- **Interval width:** intervals' mean width (the smaller, the better)
- **"Informativeness":** best width-coverage ratio, assessed through CWC score (the higher, the better):

$w$  mean width

$$\text{CWC} = (1 - w) * \exp(-\eta(c - (1 - \alpha))^2), \text{ with } \begin{cases} c \text{ attained coverage} \\ \eta \text{ balancing term} \end{cases}$$

- **Adaptability:** ability of achieving conditional coverage, assessed through SSC score (the closer to  $1 - \alpha$ , the better).
  - Maximum coverage violation along all width groups.
  - Only usable for non-constant width intervals.
- **Computational efficiency:** measured by CPU time.

1 Introduction

2 Conformal prediction

3 Beyond exchangeability

**4 Results**

Assessment

Regression problem

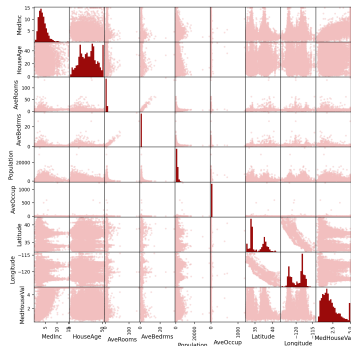
Time series problem

5 Conclusions

# Data & modeling

A tabular regression problem is considered with:

- The sklearn built-in California Housing dataset (20,640 samples, 8 features).
- A (light) gradient boosting regressor, LGBM, automatically fine-tuned through grid-search.
- A 5-fold cross-validation assessment for  $\alpha = 0.20$  miscoverage level.



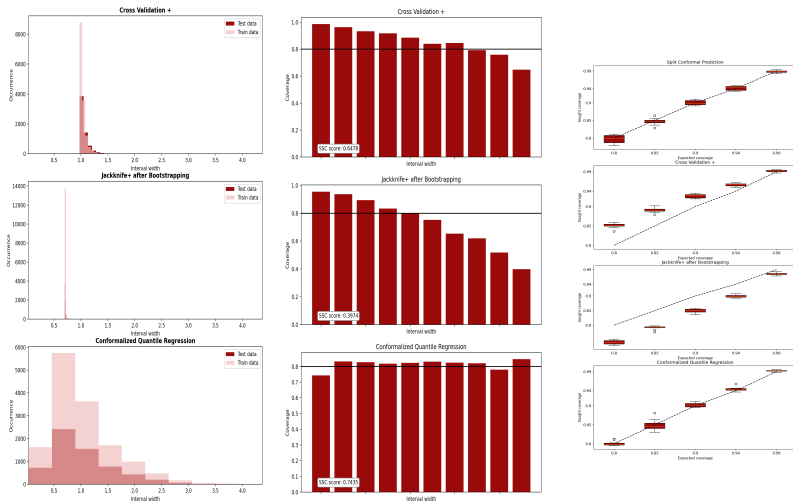
**Figure 3:** Marginal distributions.

## Metrics table

Strategy	Coverage	RMSE	Train. time	Inf. time
SCP	$0.806 \pm 0.008$	$0.472 \pm 0.007$	$1.6 \pm 0.2$	$0.07 \pm 0.05$
CV+	$0.853 \pm 0.004$	$0.467 \pm 0.009$	$9 \pm 3$	$8.0 \pm 0.3$
J+aB	$0.734 \pm 0.007$	$0.467 \pm 0.009$	$51 \pm 5$	$9.7 \pm 0.4$
CQR	$0.805 \pm 0.010$	$0.494 \pm 0.013$	$2.6 \pm 0.1$	$0.10 \pm 0.04$

Strategy	Coverage	Width	CWC	SSC
SCP	$0.806 \pm 0.008$	$0.971 \pm 0.015$	$0.798 \pm 0.004$	—
CV+	$0.853 \pm 0.004$	$1.042 \pm 0.005$	$0.784 \pm 0.002$	$0.65 \pm 0.01$
J+aB	$0.734 \pm 0.007$	$0.710 \pm 0.003$	$0.853 \pm 0.001$	—
CQR	$0.805 \pm 0.010$	$1.013 \pm 0.013$	$0.790 \pm 0.004$	$0.75 \pm 0.04$

## Other visualizations



**Figure 4:** Width histograms & coverage vs. width &  $\alpha$ .

1 Introduction

2 Conformal prediction

3 Beyond exchangeability

**4 Results**

Assessment

Regression problem

**Time series problem**

Original dataset

Change point in test

5 Conclusions



# Dataset

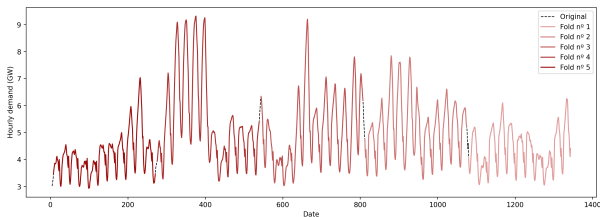
A time series forecasting problem is considered with:

- Victoria electricity demand dataset (1340 samples, features: time, demand lagged up to 7 days & temperature).
- A `sklearn` random forest regressor automatically fine-tuned through grid-search.
- A 5-fold cross-validation for  $\alpha = 0.20$ :

# Dataset

A time series forecasting problem is considered with:

- Victoria electricity demand dataset (1340 samples, features: time, demand lagged up to 7 days & temperature).
- A `sklearn` random forest regressor automatically fine-tuned through grid-search.
- A 5-fold cross-validation for  $\alpha = 0.20$ :



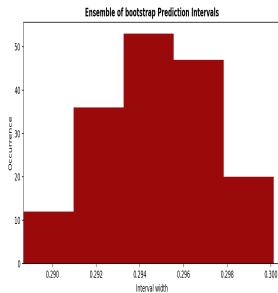
**Figure 5:** 5-fold CV splits.

## Metrics table

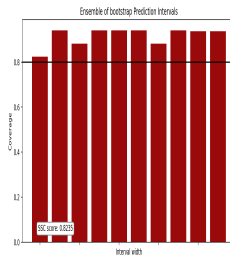
Strategy	Coverage	RMSE	Total time
EnbPI_nP	$0.780 \pm 0.069$	$0.165 \pm 0.067$	$6.2 \pm 0.3$
EnbPI	$0.789 \pm 0.058$	$0.165 \pm 0.067$	$528.3 \pm 0.4$

Strategy	Coverage	Width	CWC	SSC
EnbPI_nP	$0.780 \pm 0.069$	$0.293 \pm 0.013$	$0.935 \pm 0.018$	—
EnbPI	$0.789 \pm 0.058$	$0.300 \pm 0.007$	$0.93 \pm 0.02$	$0.5 \pm 0.2$

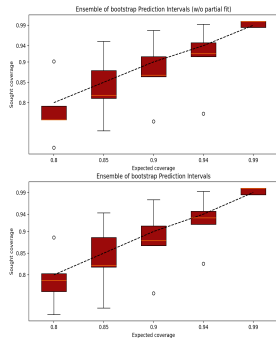
## Other visualizations



(a) EnbPI intervals' width histograms



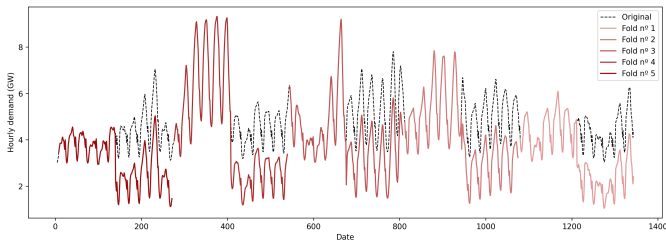
(b) EnbPI coverage in function of intervals' width

(c) EnbPI\_nP & EnbPI coverage in function of  $\alpha$ Figure 6: Width histograms & coverage vs. width &  $\alpha$ .

# Idea

The consistency of former time series may outshine the benefits of the "*partial fit*" EnbPI feature. Thus:

- A change point is added in test to mock off a distribution shift.
- The same random forest regressor will be applied to a 5-fold cross-validation, now for  $\alpha = 0.05$ :

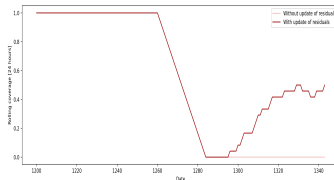
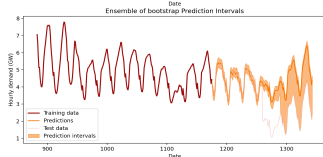
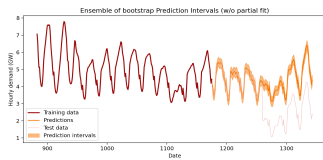


**Figure 7:** 5-fold CV splits with change points in each test's.

## Metrics &amp; plot

Strategy	Coverage	RMSE	Total time
EnbPI_nP	$0.439 \pm 0.075$	$1.431 \pm 0.024$	$6.0 \pm 0.3$
EnbPI	$0.696 \pm 0.042$	$1.431 \pm 0.024$	$530 \pm 1$

Strategy	Coverage	Width	SSC
EnbPI_nP	$0.439 \pm 0.075$	$0.569 \pm 0.043$	—
EnbPI	$0.696 \pm 0.042$	$1.300 \pm 0.034$	$0.07 \pm 0.12$



- ① Introduction
- ② Conformal prediction
- ③ Beyond exchangeability
- ④ Results
- ⑤ Conclusions**

# Regression & time series

The best strategies for exchangeable data are, **decreasingly ordered** by:

- *Statistical efficiency*: CQR, SCP, CV+, J+aB.
  - This is fulfilled independently of  $\alpha$ .
- *Computational efficiency*: SCP, CQR, CV+, J+aB.
- *Predictive power* are: CV+ & J+aB, SCP, CQR.
- "*Informativeness*": J+aB, SCP, CQR, CV+.
- *Adaptability*: CQR, CV+, J+aB (slight to none). Contrarily, SCP intervals are not adaptive at all.

Regarding the time series case, **EnbPI** is a **suitable option** to provide valid intervals.

- EnbPI's adjustment using test residuals is necessary.
- This option also allows all the issued **intervals** to be **adaptive**.



# Thank you for your attention!

Questions?