

Distribution-Free Predictive Uncertainty Quantification: Strengths and Limits of Conformal Prediction

Aymeric Dieuleveut & Margaux Zaffran

July 15th, 2024

40th Conference on Uncertainty in Artificial Intelligence (UAI)



Inria



Why are we all here today?



Figure 1: us

Why are we all here today?

(Slides available on our webpages)



Why are we all here today?

(Slides available on our webpages)

- Because Conformal Prediction has been a **popular** topic recently.



Vovk et al. (2005) algorithmic learning in a random world cite count.

Why are we all here today?

(Slides available on our webpages)

- Because Conformal Prediction has been a **popular** topic recently.
- Because we believe that conformal methods are **important** tools, whose strengths and limitations are sometimes misunderstood.

Successfully applied to

- Medical applications
- Markets / demand forecasting
- Computer Vision

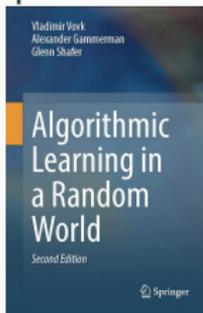


Why are we all here today?

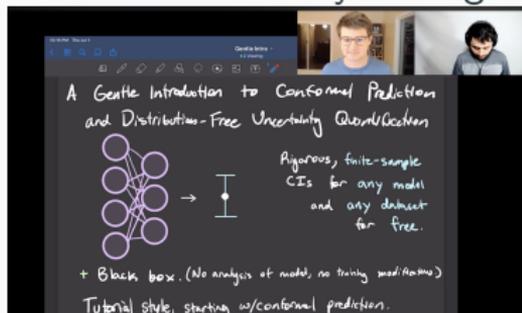


(Slides available on our webpages)

- Because Conformal Prediction has been a **popular** topic recently.
- Because we believe that conformal methods are **important** tools, whose strengths and limitations are sometimes misunderstood.
- To be part of the **diffusion** effort that many colleagues are making.



Book reference: Vovk et al. (2005)
(new edition in 2022)



A gentle tutorial: Angelopoulos and Bates (2023)
+ **Videos playlist**



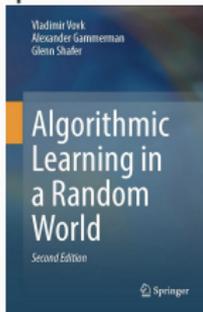
R. J. Tibshirani
introductory lecture's notes

Why are we all here today?

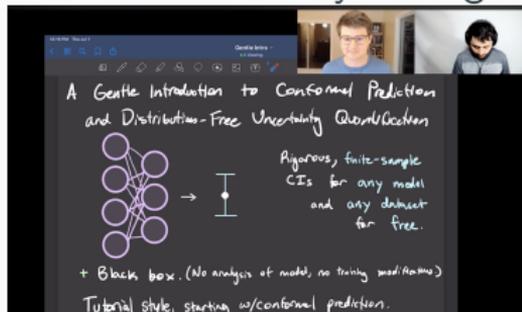


(Slides available on our webpages)

- Because Conformal Prediction has been a **popular** topic recently.
- Because we believe that conformal methods are **important** tools, whose strengths and limitations are sometimes misunderstood.
- To be part of the **diffusion** effort that many colleagues are making.



Book reference: Vovk et al. (2005)
(new edition in 2022)



A gentle tutorial: Angelopoulos and Bates (2023)
+ **Videos playlist**



R. J. Tibshirani

introductory lecture's notes

→ **Based on material freely accessible on this webpage, including sources.**

Feel free to reuse these contents for presentations or teaching!

Goals

- Provide a detailed introduction to the basics
- Demystify the results: fair introduction with limits
- Give you tools to leverage those techniques in your own fields

Disclaimers

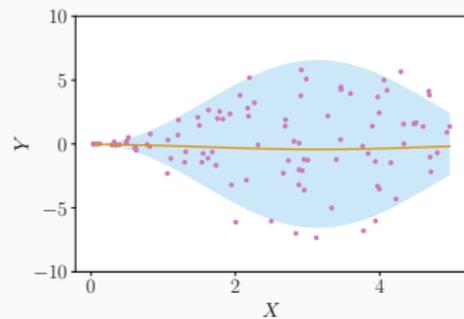
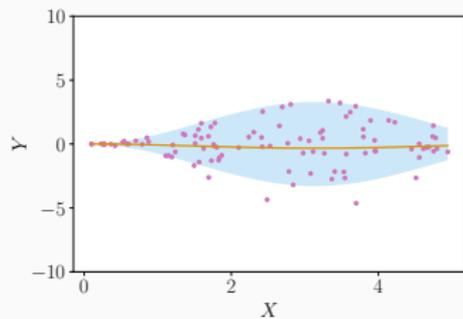
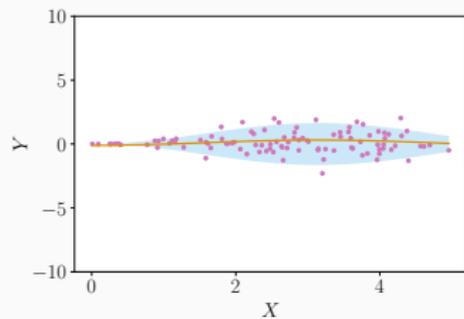
- Many people contributed to the domain - list of references may not be exhaustive
- Multiple other excellent resources

On the importance of quantifying uncertainty

- Obvious in most applications - weather, medical, markets

On the importance of quantifying uncertainty

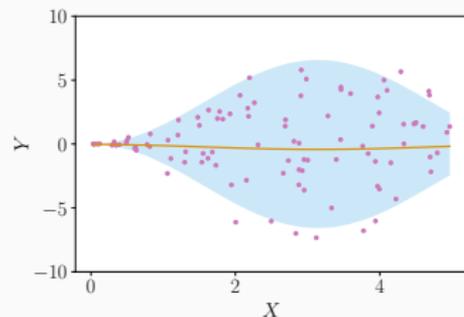
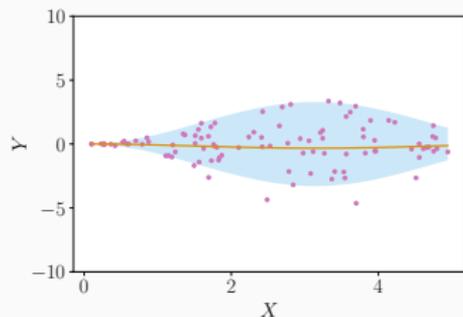
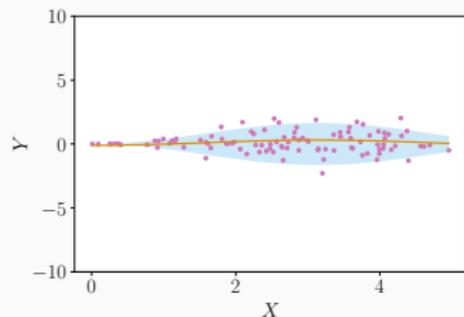
- Obvious in most applications - weather, medical, markets
- Mathematically



↪ Same “best” predictor, yet 3 distinct underlying phenomena!

On the importance of quantifying uncertainty

- Obvious in most applications - weather, medical, markets
- Mathematically



↪ Same “best” predictor, yet 3 distinct underlying phenomena!

⇒ Quantifying uncertainty conveys this information.

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

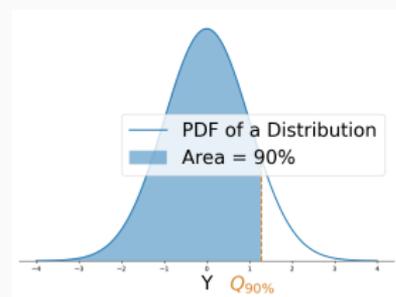
Beyond exchangeability

Some case studies

Concluding remarks

Reminder about quantiles

- Quantile level $\beta \in [0, 1]$
- $Q_Y(\beta) := \inf\{t \in \mathbb{R}, \mathbb{P}(Y \leq t) \geq \beta\}$



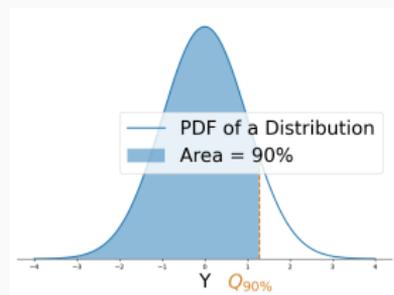
Reminder about quantiles

- Quantile level $\beta \in [0, 1]$
- $Q_Y(\beta) := \inf\{t \in \mathbb{R}, \mathbb{P}(Y \leq t) \geq \beta\}$

- Empirical quantile

$$q_\beta(Y_1, \dots, Y_n)$$

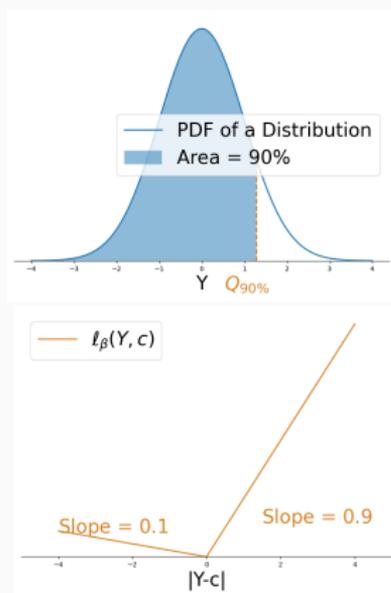
$$:= \lceil \beta \times n \rceil \text{ smallest value of } (Y_1, \dots, Y_n)$$



Reminder about quantiles

- Quantile level $\beta \in [0, 1]$
- $Q_Y(\beta) := \inf\{t \in \mathbb{R}, \mathbb{P}(Y \leq t) \geq \beta\}$
- Empirical quantile
 $q_\beta(Y_1, \dots, Y_n)$
 $:= [\beta \times n]$ smallest value of (Y_1, \dots, Y_n)
- Pinball loss

$$\begin{aligned} \ell_\beta(Y, Y') &= \beta |Y - Y'| \mathbb{1}_{\{Y - Y' \geq 0\}} \\ &\quad + (1 - \beta) |Y - Y'| \mathbb{1}_{\{Y - Y' \leq 0\}} \end{aligned}$$



Reminder about quantiles

- Quantile level $\beta \in [0, 1]$
- $Q_Y(\beta) := \inf\{t \in \mathbb{R}, \mathbb{P}(Y \leq t) \geq \beta\}$

- Empirical quantile

$$q_\beta(Y_1, \dots, Y_n)$$

$:= \lceil \beta \times n \rceil$ smallest value of (Y_1, \dots, Y_n)

- Pinball loss

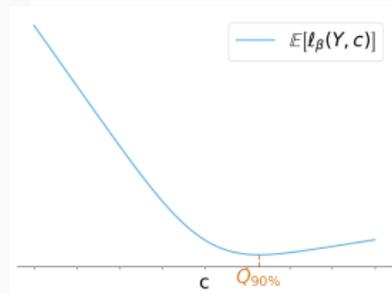
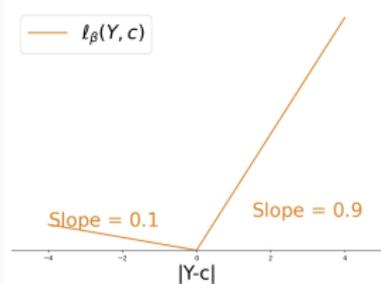
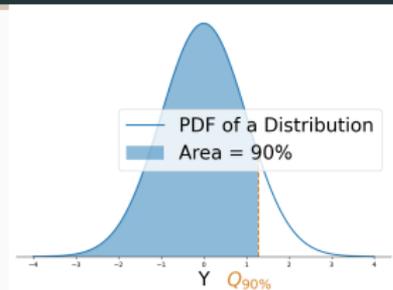
$$\begin{aligned} \ell_\beta(Y, Y') &= \beta |Y - Y'| \mathbb{1}_{\{Y - Y' \geq 0\}} \\ &\quad + (1 - \beta) |Y - Y'| \mathbb{1}_{\{Y - Y' \leq 0\}} \end{aligned}$$

Associated risk:

$$\text{Risk}_{\ell_\beta}(c) = \mathbb{E}[\ell_\beta(Y, c)]$$

Link to quantile:

$$Q_Y(\beta) = \arg \min_{c \in \mathbb{R}} \text{Risk}_{\ell_\beta}(c)$$



Reminder about quantiles

- Quantile level $\beta \in [0, 1]$
- $Q_Y(\beta) := \inf\{t \in \mathbb{R}, \mathbb{P}(Y \leq t) \geq \beta\}$
- Empirical quantile
 $q_\beta(Y_1, \dots, Y_n)$
 $:= \lceil \beta \times n \rceil$ smallest value of (Y_1, \dots, Y_n)
- Pinball loss

$$\begin{aligned} \ell_\beta(Y, Y') &= \beta |Y - Y'| \mathbb{1}_{\{Y - Y' \geq 0\}} \\ &\quad + (1 - \beta) |Y - Y'| \mathbb{1}_{\{Y - Y' \leq 0\}} \end{aligned}$$

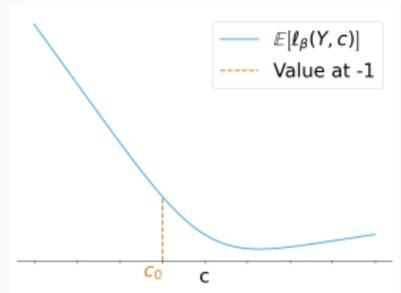
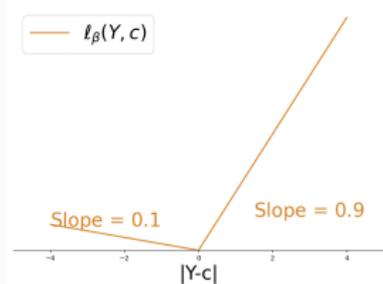
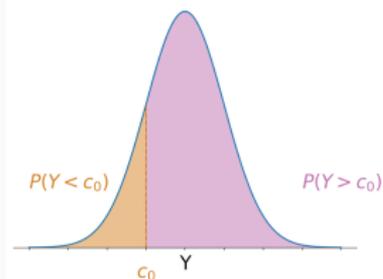
Associated risk:

$$\text{Risk}_{\ell_\beta}(c) = \mathbb{E}[\ell_\beta(Y, c)]$$

Link to quantile:

$$Q_Y(\beta) = \arg \min_{c \in \mathbb{R}} \text{Risk}_{\ell_\beta}(c)$$

Proof: sub-differential



Reminder about quantiles

- Quantile level $\beta \in [0, 1]$
- $Q_Y(\beta) := \inf\{t \in \mathbb{R}, \mathbb{P}(Y \leq t) \geq \beta\}$

- Empirical quantile

$$q_\beta(Y_1, \dots, Y_n)$$

$:= \lceil \beta \times n \rceil$ smallest value of (Y_1, \dots, Y_n)

- Pinball loss

$$\begin{aligned} \ell_\beta(Y, Y') &= \beta |Y - Y'| \mathbb{1}_{\{Y - Y' \geq 0\}} \\ &\quad + (1 - \beta) |Y - Y'| \mathbb{1}_{\{Y - Y' \leq 0\}} \end{aligned}$$

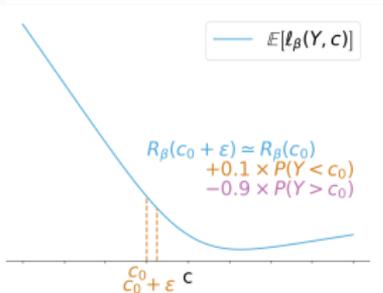
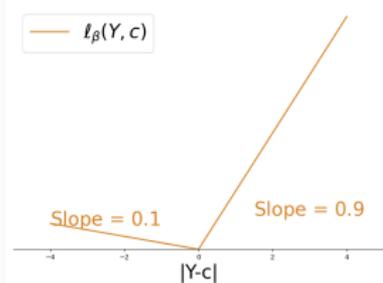
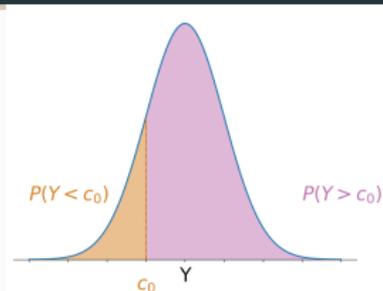
Associated risk:

$$\text{Risk}_{\ell_\beta}(c) = \mathbb{E}[\ell_\beta(Y, c)]$$

Link to quantile:

$$Q_Y(\beta) = \arg \min_{c \in \mathbb{R}} \text{Risk}_{\ell_\beta}(c)$$

Proof: sub-differential



Example (a special quantile: the median).

$$\beta = 0.5$$

$\hookrightarrow Q_Y(0.5)$ represents the median of the distribution of Y .

$$\hookrightarrow Q_Y(0.5) = \operatorname{argmin}_c \mathbb{E}[|Y - c|].$$

Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$

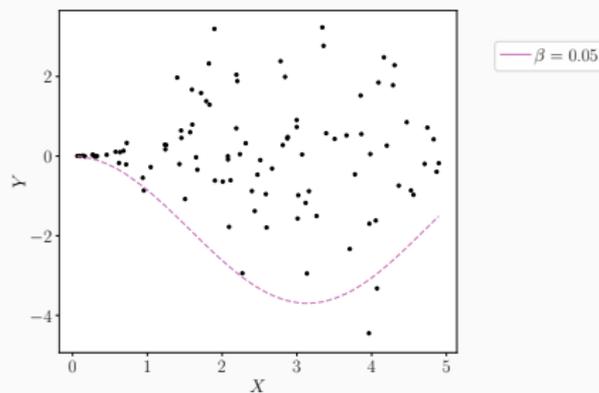
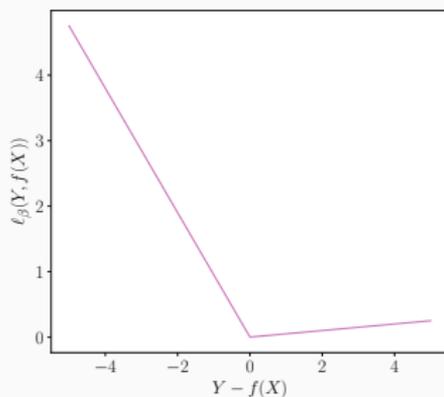
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



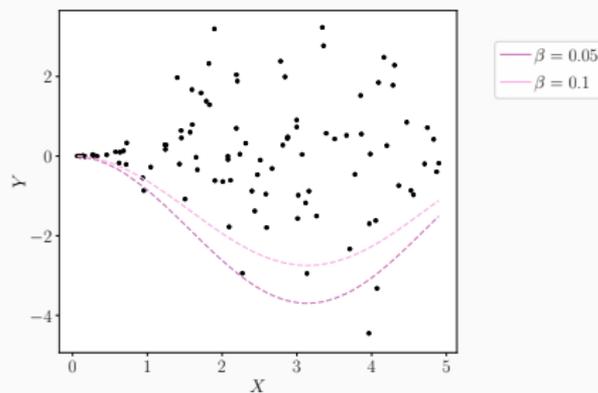
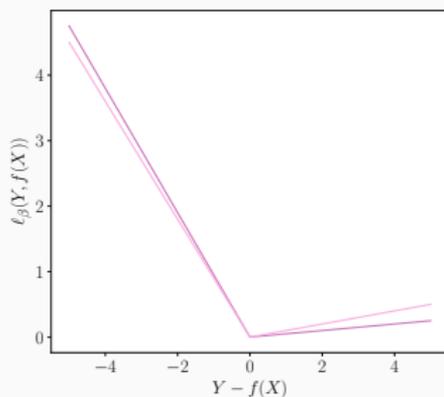
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



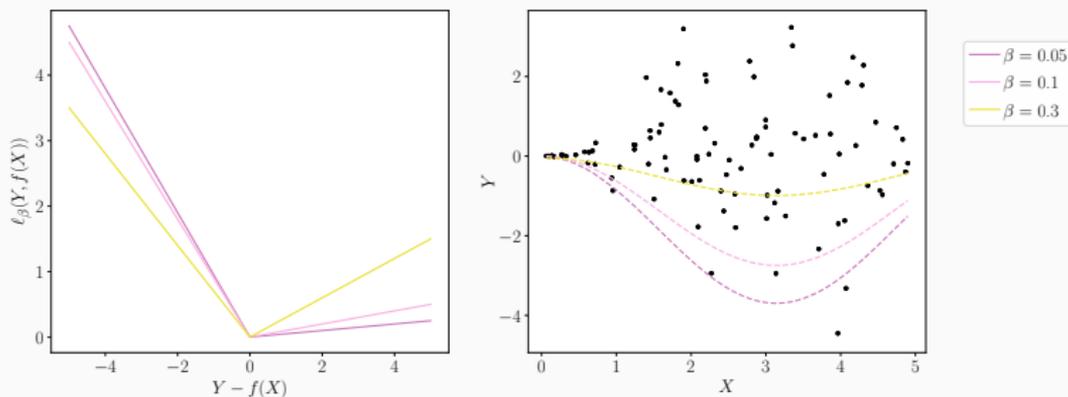
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



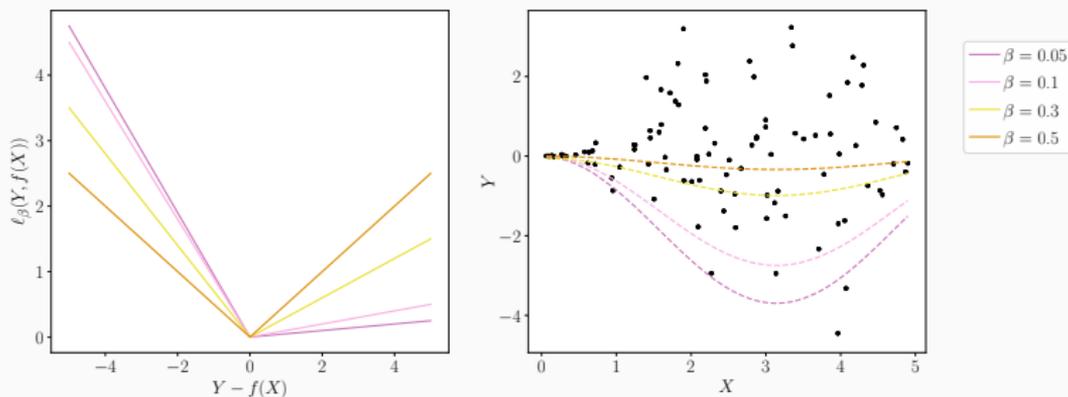
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



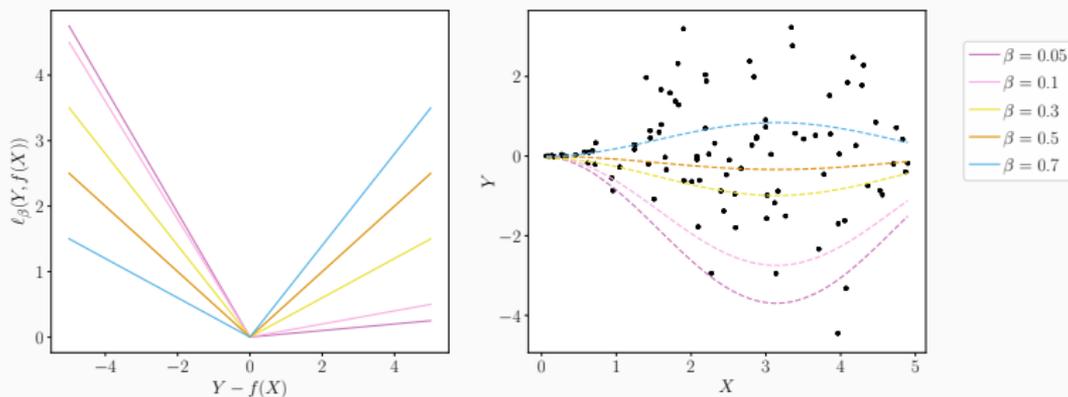
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



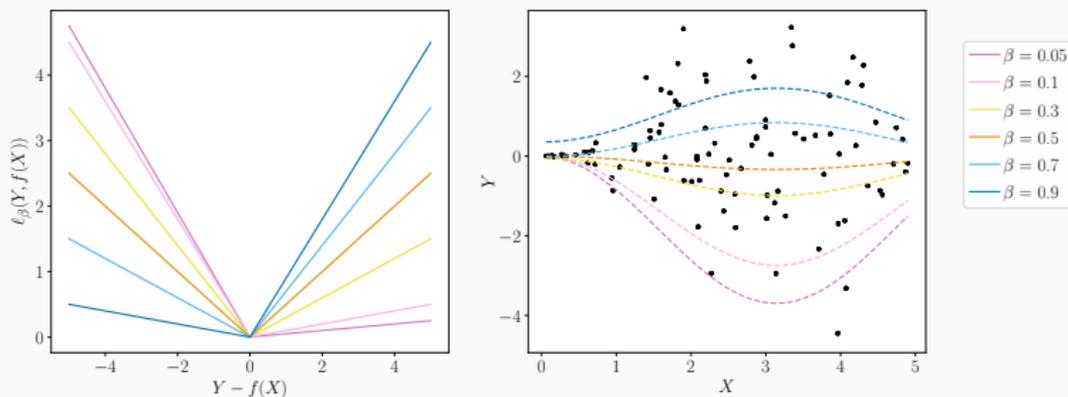
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



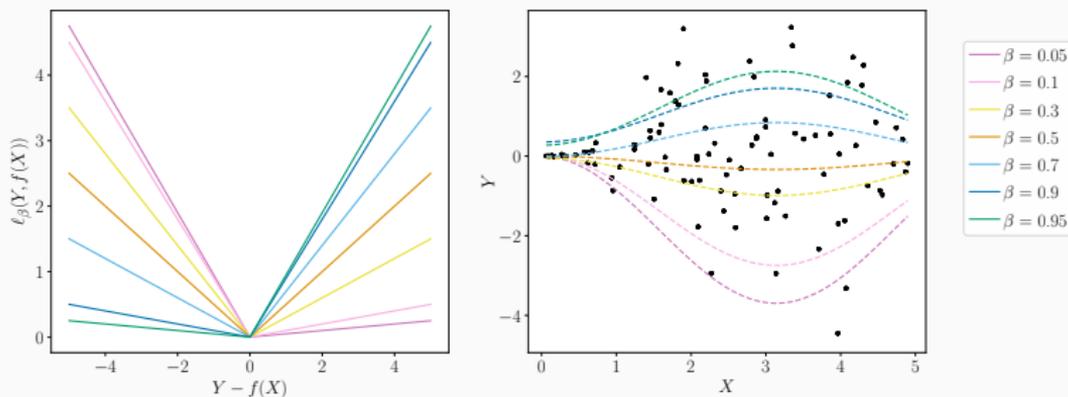
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



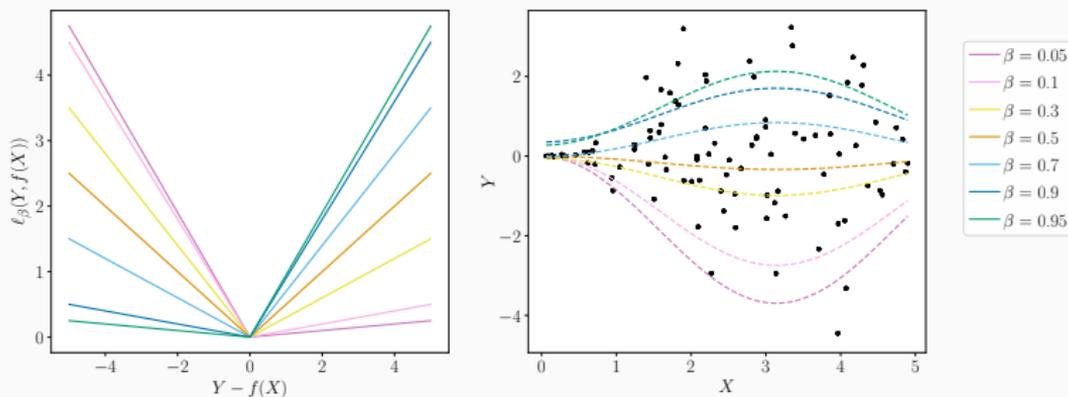
Quantile regression

- Goal : approximate $Q_{Y|X}(\beta)$ – Quantile level β – Pinball loss $\ell_\beta(Y, Y')$.
- Associated risk:

$$\text{Risk}_{\ell_\beta}(f) = \mathbb{E}[\ell_\beta(Y, f(X))]$$

- Bayes predictor:

$$f^* \in \underset{f \in \mathbb{R}^{\mathcal{X}}}{\text{argmin}} \text{Risk}_{\ell_\beta}(f) \quad \Rightarrow \quad f^*(X) = Q_{Y|X}(\beta)$$



Warning - No theoretical guarantee with a finite sample!

$$\mathbb{P}\left(Y \in \left[\hat{Q}_{Y|X}(\beta/2); \hat{Q}_{Y|X}(1 - \beta/2)\right]\right) \neq 1 - \beta$$

Quantifying predictive uncertainty

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables
- n training samples $(X_i, Y_i)_{i=1}^n$
- **Goal:** predict an unseen point Y_{n+1} at X_{n+1} with **confidence**
- **How?** Given a miscoverage level $\alpha \in [0, 1]$, build a predictive set \mathcal{C}_α such that:

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

and \mathcal{C}_α should be as small as possible, in order to be informative

For example: $\alpha = 0.1$ and obtain a 90% coverage interval

Quantifying predictive uncertainty

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables
- n training samples $(X_i, Y_i)_{i=1}^n$
- **Goal:** predict an unseen point Y_{n+1} at X_{n+1} with **confidence**
- **How?** Given a miscoverage level $\alpha \in [0, 1]$, build a predictive set \mathcal{C}_α such that:

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

and \mathcal{C}_α should be as small as possible, in order to be informative

For example: $\alpha = 0.1$ and obtain a 90% coverage interval

- Construction of the predictive intervals should be
 - **agnostic to the model**
 - **agnostic to the data distribution**
- **Validity** should be ensured
 - in **finite samples**
 - for all **data distribution** and **underlying model**

Quantile Regression

Split Conformal Prediction (SCP)

Standard regression case

Conformalized Quantile Regression (CQR)

Generalization of SCP: going beyond regression

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

Quantile Regression

Split Conformal Prediction (SCP)

Standard regression case

Conformalized Quantile Regression (CQR)

Generalization of SCP: going beyond regression

On the design choices of conformity scores and (empirical) conditional guarantees

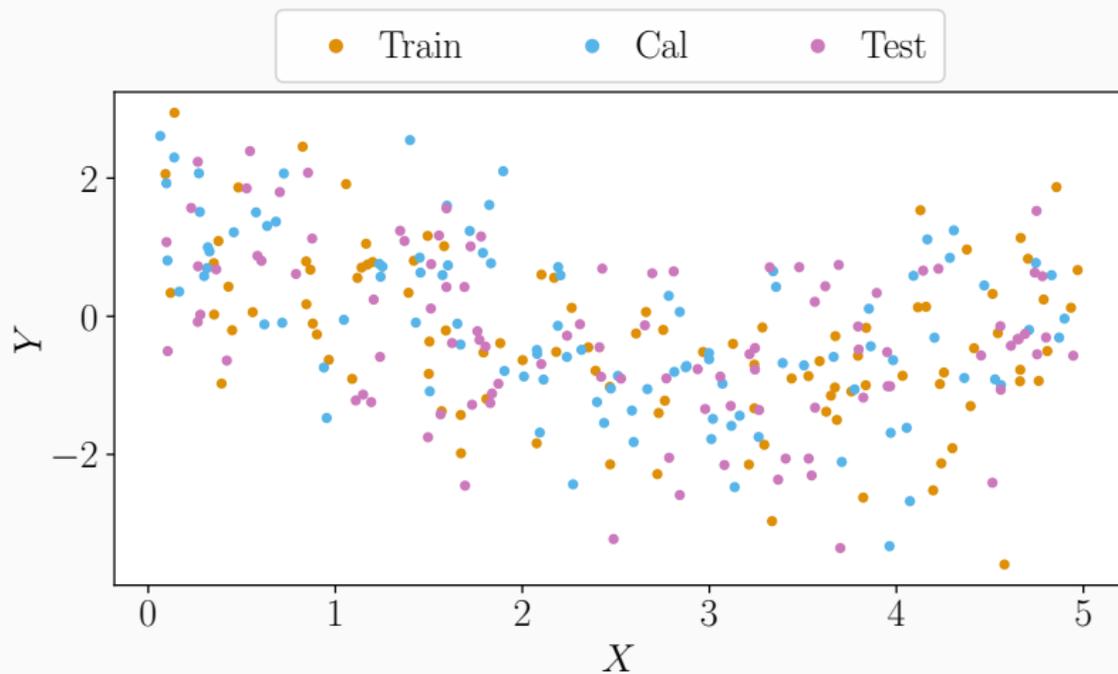
Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

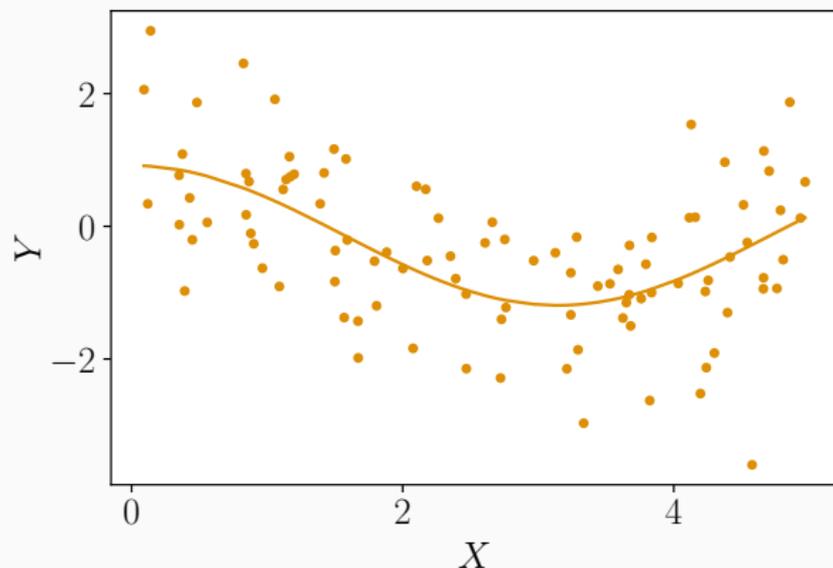
Split Conformal Prediction (SCP)^{1,2,3}: toy example



¹Vovk et al. (2005), *Algorithmic Learning in a Random World*

²Papadopoulos et al. (2002), *Inductive Confidence Machines for Regression*, ECML

³Lei et al. (2018), *Distribution-Free Predictive Inference for Regression*, JRSS B

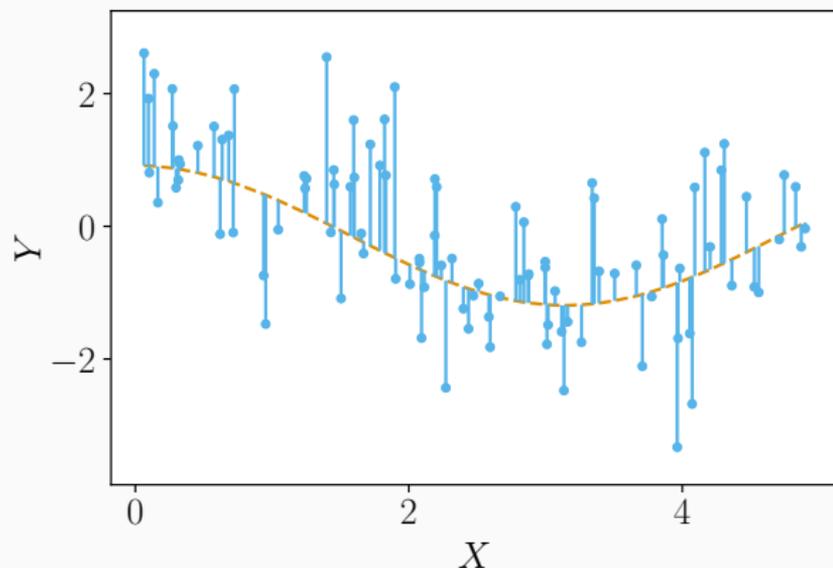


► Learn (or get) $\hat{\mu}$

¹Vovk et al. (2005), *Algorithmic Learning in a Random World*

²Papadopoulos et al. (2002), *Inductive Confidence Machines for Regression*, ECML

³Lei et al. (2018), *Distribution-Free Predictive Inference for Regression*, JRSS B

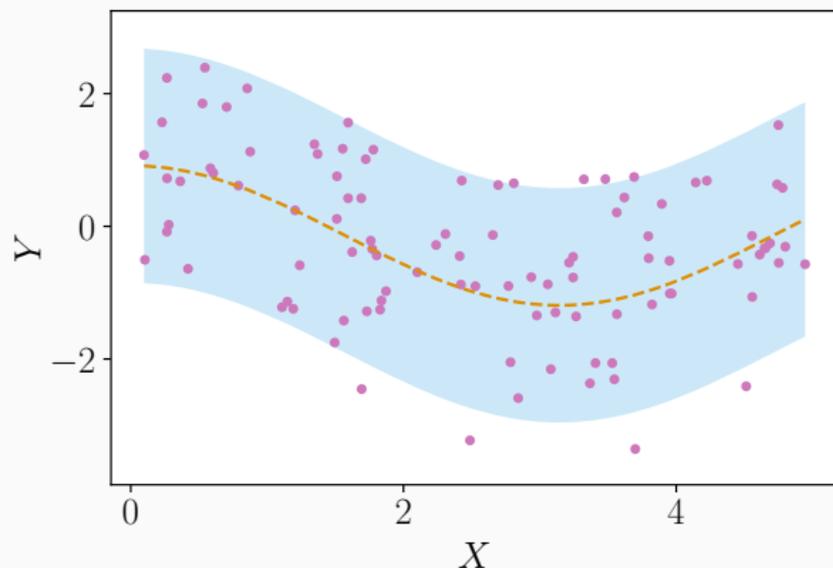


- ▶ Predict with $\hat{\mu}$
- ▶ Get the `|residuals|`, a.k.a. conformity scores
- ▶ Compute the $(1 - \alpha)$ empirical quantile of $\mathcal{S} = \{|\text{residuals}|\}_{\text{Cal}} \cup \{+\infty\}$, noted $q_{1-\alpha}(\mathcal{S})$

¹Vovk et al. (2005), *Algorithmic Learning in a Random World*

²Papadopoulos et al. (2002), *Inductive Confidence Machines for Regression*, ECML

³Lei et al. (2018), *Distribution-Free Predictive Inference for Regression*, JRSS B



- ▶ Predict with $\hat{\mu}$
- ▶ Build $\hat{C}_\alpha(x)$: $[\hat{\mu}(x) \pm q_{1-\alpha}(\mathcal{S})]$

¹Vovk et al. (2005), *Algorithmic Learning in a Random World*

²Papadopoulos et al. (2002), *Inductive Confidence Machines for Regression*, ECML

³Lei et al. (2018), *Distribution-Free Predictive Inference for Regression*, JRSS B



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by *training the algorithm \mathcal{A} on the **proper training set***



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, get prediction values with $\hat{\mu}$



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by training the algorithm \mathcal{A} on the **proper training set**
3. On the **calibration set**, get prediction values with $\hat{\mu}$
4. Obtain a set of $\#Cal + 1$ **conformity scores** :

$$\mathcal{S} = \{S_i = |\hat{\mu}(X_i) - Y_i|, i \in \text{Cal}\} \cup \{+\infty\}$$

(+ worst-case scenario)



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by *training the algorithm \mathcal{A} on the proper training set*
3. On the **calibration set**, get prediction values with $\hat{\mu}$
4. Obtain a set of $\#Cal + 1$ **conformity scores** :

$$\mathcal{S} = \{S_i = |\hat{\mu}(X_i) - Y_i|, i \in \text{Cal}\} \cup \{+\infty\}$$

(+ worst-case scenario)

5. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by training the algorithm \mathcal{A} on the **proper training set**
3. On the **calibration set**, get prediction values with $\hat{\mu}$
4. Obtain a set of $\#Cal + 1$ **conformity scores**:

$$\mathcal{S} = \{S_i = |\hat{\mu}(X_i) - Y_i|, i \in \text{Cal}\} \cup \{+\infty\}$$

(+ worst-case scenario)

5. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$
6. For a new point X_{n+1} , return

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



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by training the algorithm \mathcal{A} on the **proper training set**
3. On the **calibration set**, get prediction values with $\hat{\mu}$
4. Obtain a set of $\#Cal$ **conformity scores**:

$$\mathcal{S} = \{S_i = |\hat{\mu}(X_i) - Y_i|, i \in \text{Cal}\}$$



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by training the algorithm \mathcal{A} on the **proper training set**
3. On the **calibration set**, get prediction values with $\hat{\mu}$
4. Obtain a set of $\#Cal$ **conformity scores**:

$$\mathcal{S} = \{S_i = |\hat{\mu}(X_i) - Y_i|, i \in \text{Cal}\}$$

5. Compute the $(1 - \alpha) \left(\frac{1}{\#Cal} + 1 \right)$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get $\hat{\mu}$ by training the algorithm \mathcal{A} on the **proper training set**
3. On the **calibration set**, get prediction values with $\hat{\mu}$
4. Obtain a set of $\#Cal$ **conformity scores**:

$$\mathcal{S} = \{S_i = |\hat{\mu}(X_i) - Y_i|, i \in \text{Cal}\}$$

5. Compute the $(1 - \alpha) \left(\frac{1}{\#Cal} + 1 \right)$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$
6. For a new point X_{n+1} , return

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

Definition (Exchangeability).

$(X_i, Y_i)_{i=1}^n$ are **exchangeable** if, for any permutation σ of $\llbracket 1, n \rrbracket$:

$$((X_1, Y_1), \dots, (X_n, Y_n)) \stackrel{d}{=} ((X_{\sigma(1)}, Y_{\sigma(1)}), \dots, (X_{\sigma(n)}, Y_{\sigma(n)})).$$

Definition (Exchangeability).

$(X_i, Y_i)_{i=1}^n$ are **exchangeable** if, for any permutation σ of $\llbracket 1, n \rrbracket$:

$$((X_1, Y_1), \dots, (X_n, Y_n)) \stackrel{d}{=} ((X_{\sigma(1)}, Y_{\sigma(1)}), \dots, (X_{\sigma(n)}, Y_{\sigma(n)})).$$

Toy case: Z_1 and Z_2 are exchangeable if $(Z_1, Z_2) \stackrel{d}{=} (Z_2, Z_1)$.

Definition (Exchangeability).

$(X_i, Y_i)_{i=1}^n$ are **exchangeable** if, for any permutation σ of $\llbracket 1, n \rrbracket$:

$$((X_1, Y_1), \dots, (X_n, Y_n)) \stackrel{d}{=} ((X_{\sigma(1)}, Y_{\sigma(1)}), \dots, (X_{\sigma(n)}, Y_{\sigma(n)})) .$$

Example (exchangeable sequences).

- i.i.d. samples

Definition (Exchangeability).

$(X_i, Y_i)_{i=1}^n$ are **exchangeable** if, for any permutation σ of $\llbracket 1, n \rrbracket$:

$$((X_1, Y_1), \dots, (X_n, Y_n)) \stackrel{d}{=} ((X_{\sigma(1)}, Y_{\sigma(1)}), \dots, (X_{\sigma(n)}, Y_{\sigma(n)})).$$

Example (exchangeable sequences).

- i.i.d. samples

- The components of $\mathcal{N} \left(\begin{pmatrix} m \\ \vdots \\ \vdots \\ m \end{pmatrix}, \begin{pmatrix} \sigma^2 & & & \\ & \ddots & \gamma^2 & \\ & \gamma^2 & \ddots & \\ & & & \sigma^2 \end{pmatrix} \right)$

SCP enjoys finite sample guarantees proved in Vovk et al. (2005); Lei et al. (2018).

Theorem (Marginal validity).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^a. SCP applied on $(X_i, Y_i)_{i=1}^n$ outputs $\hat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \hat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

SCP enjoys finite sample guarantees proved in Vovk et al. (2005); Lei et al. (2018).

Theorem (Marginal validity).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^a. SCP applied on $(X_i, Y_i)_{i=1}^n$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

Additionally, if the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are a.s. distinct:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

^aOnly the calibration and test data need to be exchangeable.

Lemma (Quantile lemma).

If $(U_1, \dots, U_n, U_{n+1})$ are **exchangeable**, then for any $\beta \in]0, 1[$:

$$\mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, +\infty)) \geq \beta.$$

Additionally, if U_1, \dots, U_n, U_{n+1} are almost surely distinct, then:

$$\mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, +\infty)) \leq \beta + \frac{1}{n+1}.$$

Lemma (Quantile lemma).

If $(U_1, \dots, U_n, U_{n+1})$ are **exchangeable**, then for any $\beta \in]0, 1[$:

$$\mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, +\infty)) \geq \beta.$$

Additionally, if U_1, \dots, U_n, U_{n+1} are almost surely distinct, then:

$$\mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, +\infty)) \leq \beta + \frac{1}{n+1}.$$

When $(X_i, Y_i)_{i=1}^{n+1}$ are exchangeable, the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are exchangeable.
 \hookrightarrow applying the quantile lemma to the scores concludes the proof.

$$\begin{aligned} \{Y_{n+1} \in \widehat{C}_{n,\alpha}(X_{n+1})\} &= \{\widehat{\mu}(X_{n+1}) - q_{1-\alpha}(S) \leq Y_{n+1} \leq \widehat{\mu}(X_{n+1}) + q_{1-\alpha}(S)\} \\ &= \{|Y_{n+1} - \widehat{\mu}(X_{n+1})| \leq q_{1-\alpha}(S)\} \\ &= \{S_{n+1} \leq q_{1-\alpha}(S)\}. \end{aligned}$$

First note that $U_{n+1} \leq q_\beta(U_1, \dots, U_n, +\infty) \iff U_{n+1} \leq q_\beta(U_1, \dots, U_n, U_{n+1})$.

First note that $U_{n+1} \leq q_\beta(U_1, \dots, U_n, +\infty) \iff U_{n+1} \leq q_\beta(U_1, \dots, U_n, U_{n+1})$.

By exchangeability, for any $i \in \llbracket 1, n+1 \rrbracket$:

$\mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, U_{n+1})) \stackrel{d}{=} \mathbb{P}(U_i \leq q_\beta(U_1, \dots, U_n, U_{n+1}))$. Thus:

$$\begin{aligned} \mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, U_{n+1})) &= \frac{1}{n+1} \sum_{i=1}^{n+1} \mathbb{P}(U_i \leq q_\beta(U_1, \dots, U_n, U_{n+1})) \\ &= \frac{1}{n+1} \mathbb{E} \left[\sum_{i=1}^{n+1} \mathbb{1} \{U_i \leq q_\beta(U_1, \dots, U_n, U_{n+1})\} \right] \\ &\geq \frac{1}{n+1} \mathbb{E} [\lceil \beta(n+1) \rceil] \\ &= \frac{\lceil \beta(n+1) \rceil}{n+1} \\ &\geq \beta, \end{aligned}$$

proving the first statement.

First note that $U_{n+1} \leq q_\beta(U_1, \dots, U_n, +\infty) \iff U_{n+1} \leq q_\beta(U_1, \dots, U_n, U_{n+1})$.

By exchangeability, for any $i \in \llbracket 1, n+1 \rrbracket$:

$\mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, U_{n+1})) \stackrel{d}{=} \mathbb{P}(U_i \leq q_\beta(U_1, \dots, U_n, U_{n+1}))$. Thus:

$$\begin{aligned} \mathbb{P}(U_{n+1} \leq q_\beta(U_1, \dots, U_n, U_{n+1})) &= \frac{1}{n+1} \sum_{i=1}^{n+1} \mathbb{P}(U_i \leq q_\beta(U_1, \dots, U_n, U_{n+1})) \\ &= \frac{1}{n+1} \mathbb{E} \left[\sum_{i=1}^{n+1} \mathbb{1} \{U_i \leq q_\beta(U_1, \dots, U_n, U_{n+1})\} \right] \\ &= \frac{1}{n+1} \mathbb{E} [\lceil \beta(n+1) \rceil] \quad \text{if all } (U_i) \text{ are distinct} \\ &= \frac{\lceil \beta(n+1) \rceil}{n+1} \\ &\leq \beta + \frac{1}{n+1}, \end{aligned}$$

proving the **second** statement.

SCP enjoys finite sample guarantees proved in Vovk et al. (2005); Lei et al. (2018).

Theorem (Marginal validity Vovk et al. (2005)).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^d. SCP applied on $(X_i, Y_i)_{i=1}^n$ outputs $\hat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \hat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

Additionally, if the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are a.s. distinct:

$$\mathbb{P} \left\{ Y_{n+1} \in \hat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

^dOnly the calibration and test data need to be exchangeable.

SCP enjoys finite sample guarantees proved in Vovk et al. (2005); Lei et al. (2018).

Theorem (Marginal validity Vovk et al. (2005)).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^d. SCP applied on $(X_i, Y_i)_{i=1}^n$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

Additionally, if the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are a.s. distinct:

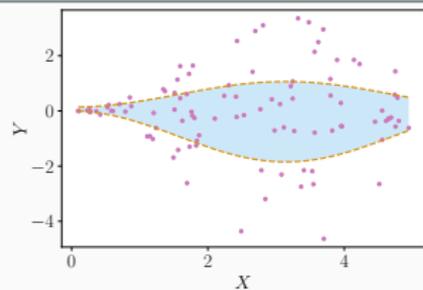
$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

^dOnly the calibration and test data need to be exchangeable.

✗ Marginal coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \mid X_{n+1} = x \right\} \geq 1 - \alpha$

Conditional coverage implies adaptiveness

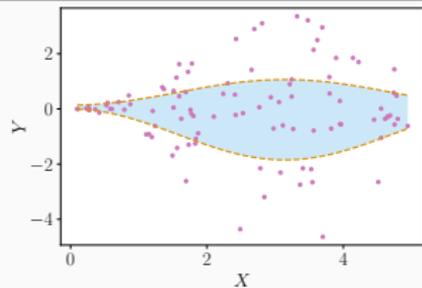
no coverage



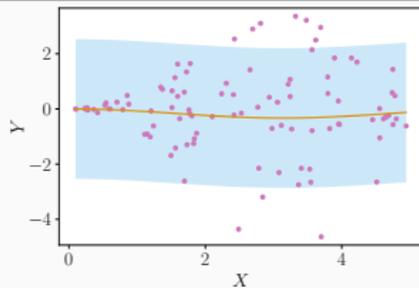
Conditional coverage implies adaptiveness

- **Marginal** coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha (X_{n+1}) \right\}$ the errors may differ across regions of the input space (i.e. non-adaptive)

no coverage

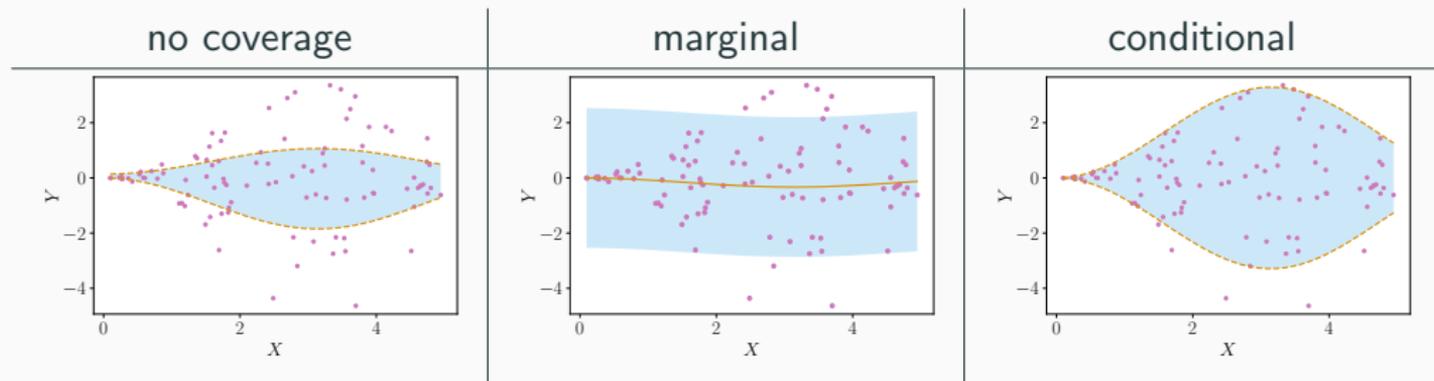


marginal



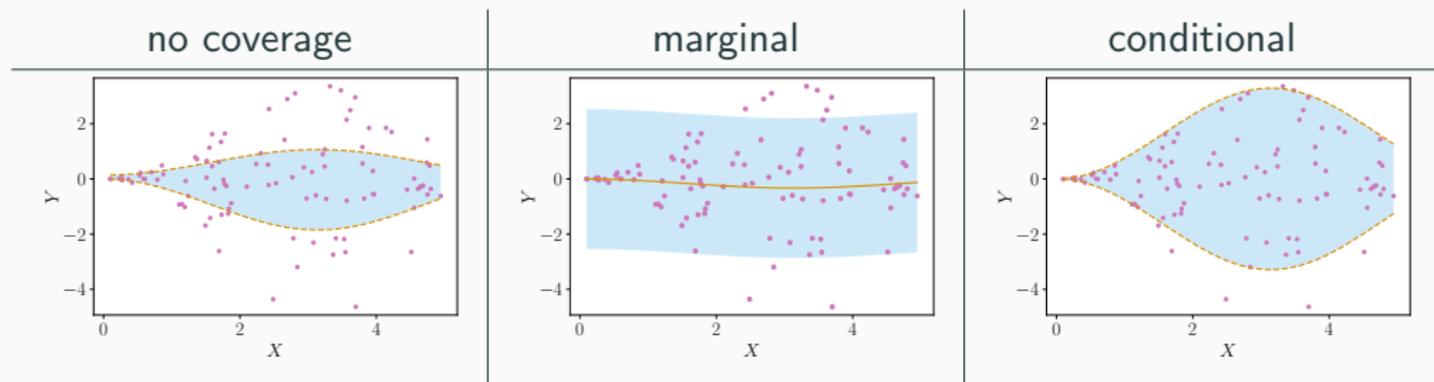
Conditional coverage implies adaptiveness

- **Marginal** coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha (X_{n+1}) \right\}$ the errors may differ across regions of the input space (i.e. non-adaptive)
- **Conditional** coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha (X_{n+1}) \mid X_{n+1} \right\}$ errors are evenly distributed (i.e. fully adaptive)

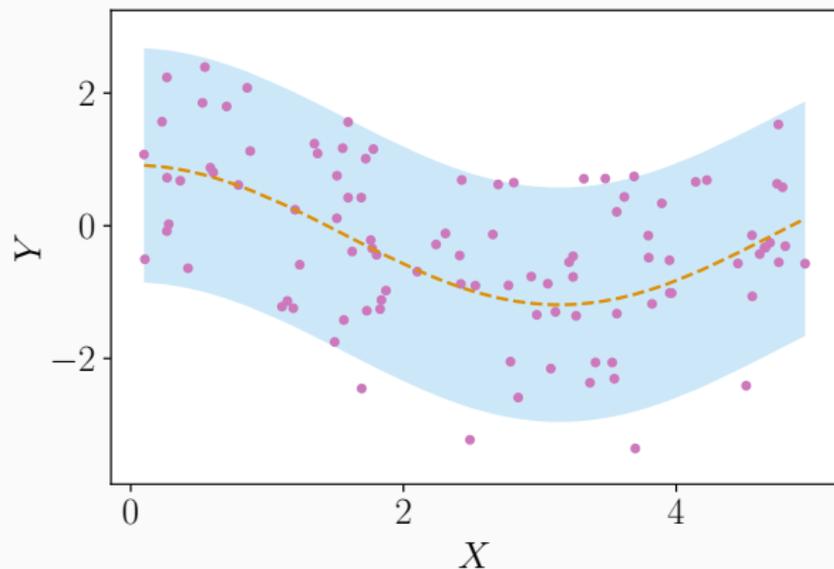


Conditional coverage implies adaptiveness

- **Marginal** coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha (X_{n+1}) \right\}$ the errors may differ across regions of the input space (i.e. non-adaptive)
- **Conditional** coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha (X_{n+1}) \mid X_{n+1} \right\}$ errors are evenly distributed (i.e. fully adaptive)
- Conditional coverage is **stronger** than marginal coverage



Standard mean-regression SCP is not adaptive



- ▶ Predict with $\hat{\mu}$
- ▶ Build $\hat{C}_\alpha(x)$: $[\hat{\mu}(x) \pm q_{1-\alpha}(\mathcal{S})]$

Quantile Regression

Split Conformal Prediction (SCP)

Standard regression case

Conformalized Quantile Regression (CQR)

Generalization of SCP: going beyond regression

On the design choices of conformity scores and (empirical) conditional guarantees

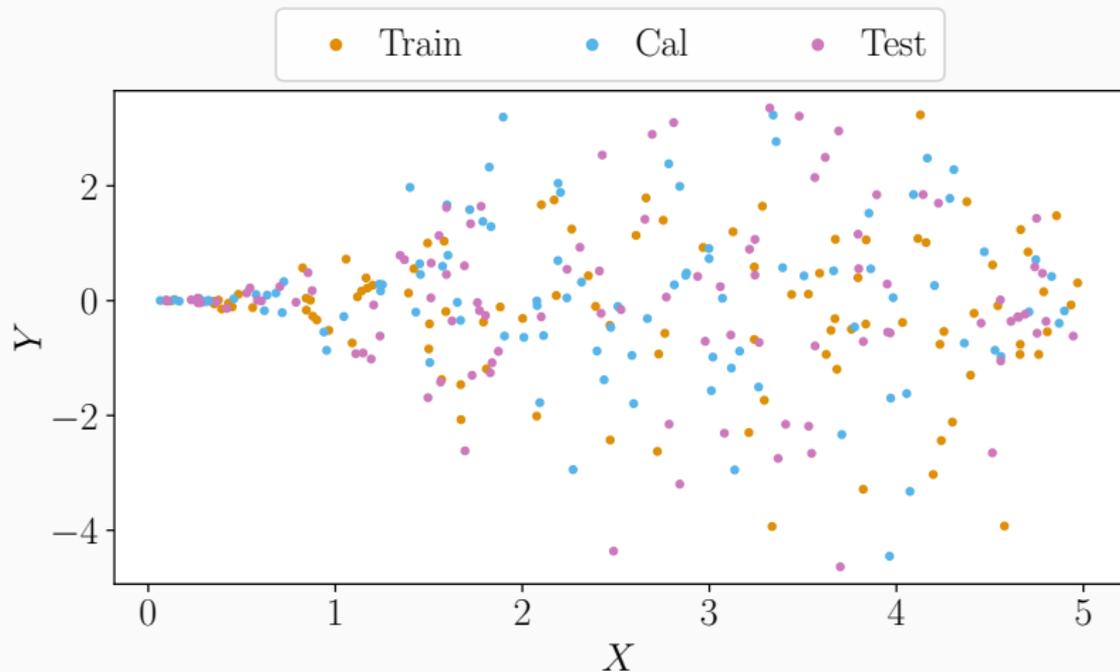
Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

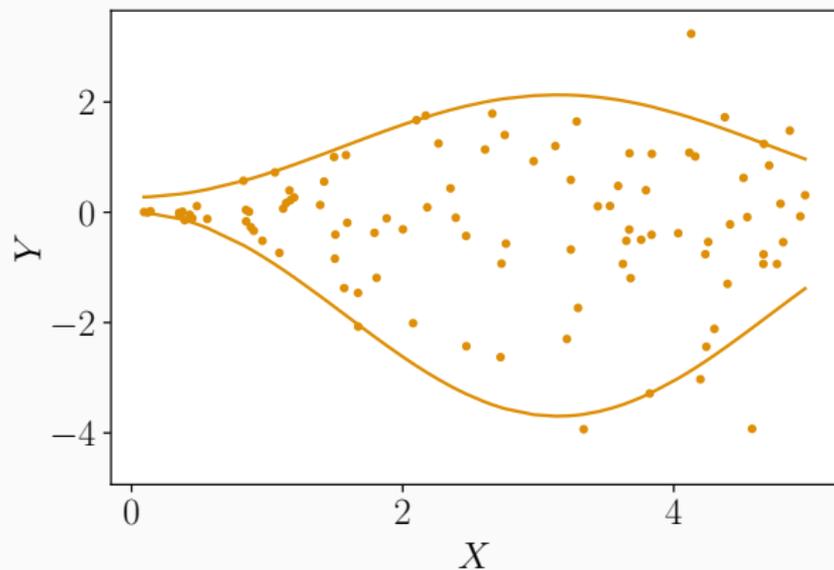
Some case studies

Concluding remarks

Conformalized Quantile Regression (CQR)⁵

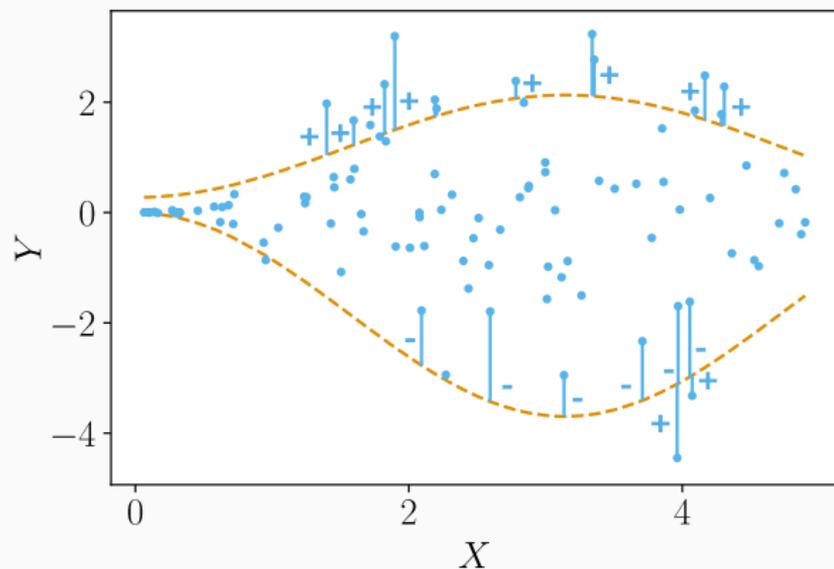


⁵Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS



► Learn (or get) $\widehat{QR}_{\text{lower}}$ and $\widehat{QR}_{\text{upper}}$

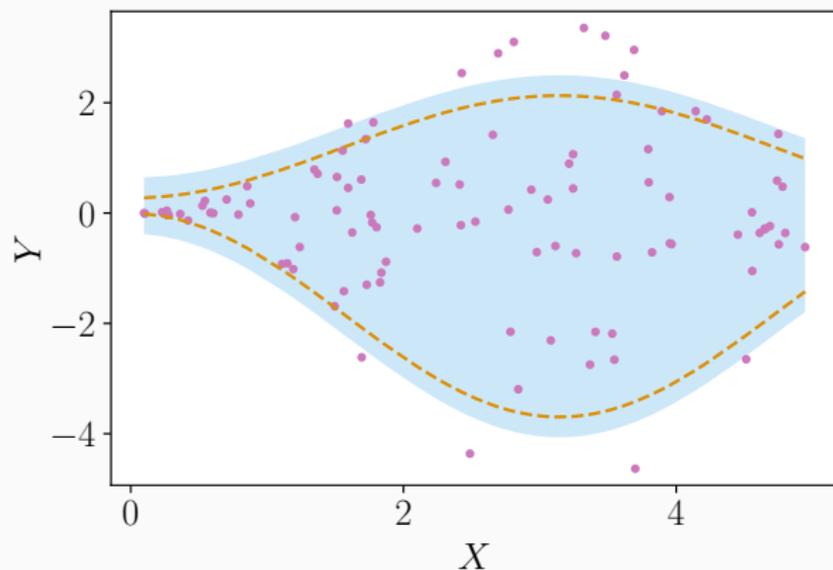
⁵Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS



- ▶ Predict with $\widehat{QR}_{\text{lower}}$ and $\widehat{QR}_{\text{upper}}$
- ▶ Get the scores $\mathcal{S} = \{S_i\}_{\text{Cal}} \cup \{+\infty\}$
- ▶ Compute the $(1 - \alpha)$ empirical quantile of \mathcal{S} , noted $q_{1-\alpha}(\mathcal{S})$

$$\Leftrightarrow S_i := \max \left\{ \widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i) \right\}$$

⁵Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS



► Predict with \widehat{QR}_{lower} and \widehat{QR}_{upper}

► Build

$$\widehat{C}_\alpha(x) = [\widehat{QR}_{lower}(x) - q_{1-\alpha}(\mathcal{S}); \widehat{QR}_{upper}(x) + q_{1-\alpha}(\mathcal{S})]$$

⁵Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \widehat{QR}_{lower} and \widehat{QR}_{upper} by training the algorithm \mathcal{A} on the **proper training set**



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \widehat{QR}_{lower} and \widehat{QR}_{upper} by training the algorithm \mathcal{A} on the **proper training set**
3. Obtain a set of $\#Cal + 1$ **conformity scores** \mathcal{S} :

$$\mathcal{S} = \{S_i = \max(\widehat{QR}_{lower}(X_i) - Y_i, Y_i - \widehat{QR}_{upper}(X_i)), i \in \text{Cal}\} \cup \{+\infty\}$$



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \widehat{QR}_{lower} and \widehat{QR}_{upper} by training the algorithm \mathcal{A} on the **proper training set**
3. Obtain a set of $\#Cal + 1$ **conformity scores** \mathcal{S} :

$$\mathcal{S} = \{S_i = \max(\widehat{QR}_{lower}(X_i) - Y_i, Y_i - \widehat{QR}_{upper}(X_i)), i \in \text{Cal}\} \cup \{+\infty\}$$
4. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \widehat{QR}_{lower} and \widehat{QR}_{upper} by training the algorithm \mathcal{A} on the **proper training set**
3. Obtain a set of $\#Cal + 1$ **conformity scores** \mathcal{S} :

$$\mathcal{S} = \{S_i = \max\left(\widehat{QR}_{lower}(X_i) - Y_i, Y_i - \widehat{QR}_{upper}(X_i)\right), i \in \text{Cal}\} \cup \{+\infty\}$$

4. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$

5. For a new point X_{n+1} , return

$$\widehat{C}_\alpha(X_{n+1}) = [\widehat{QR}_{lower}(X_{n+1}) - q_{1-\alpha}(\mathcal{S}); \widehat{QR}_{upper}(X_{n+1}) + q_{1-\alpha}(\mathcal{S})]$$

1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \widehat{QR}_{lower} and \widehat{QR}_{upper} by training the algorithm \mathcal{A} on the **proper training set**

3. Obtain a set of $\#Cal$ **conformity scores** \mathcal{S} :

$$\mathcal{S} = \{S_i = \max(\widehat{QR}_{lower}(X_i) - Y_i, Y_i - \widehat{QR}_{upper}(X_i)), i \in \text{Cal}\}$$

4. Compute the $(1 - \alpha) \left(\frac{1}{\#Cal} + 1 \right)$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$
5. For a new point X_{n+1} , return

$$\widehat{C}_\alpha(X_{n+1}) = [\widehat{QR}_{lower}(X_{n+1}) - q_{1-\alpha}(\mathcal{S}); \widehat{QR}_{upper}(X_{n+1}) + q_{1-\alpha}(\mathcal{S})]$$

This procedure enjoys the finite sample guarantee proposed and proved in Romano et al. (2019).

Theorem (Marginal validity of CQR Romano et al. (2019)).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^a. CQR on $(X_i, Y_i)_{i=1}^n$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

If, in addition, the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are almost surely distinct, then

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

^aOnly the calibration and test data need to be exchangeable.

Proof: quantile lemma again $Y_{n+1} \in \widehat{C}_{n,\alpha}(X_{n+1}) \Leftrightarrow S_{n+1} \leq q_{1-\alpha}(S)$.

CQR: theoretical guarantees

This procedure enjoys the finite sample guarantee proposed and proved in Romano et al. (2019).

Theorem (Marginal validity of CQR Romano et al. (2019)).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^a. CQR on $(X_i, Y_i)_{i=1}^n$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

If, in addition, the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are almost surely distinct, then

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

^aOnly the calibration and test data need to be exchangeable.

Proof: quantile lemma again $Y_{n+1} \in \widehat{C}_{n,\alpha}(X_{n+1}) \Leftrightarrow S_{n+1} \leq q_{1-\alpha}(S)$.

X Marginal coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \mid X_{n+1} = x \right\} \geq 1 - \alpha$

Quantile Regression

Split Conformal Prediction (SCP)

Standard regression case

Conformalized Quantile Regression (CQR)

Generalization of SCP: going beyond regression

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = s(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = s(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

Ex 1: $s(\hat{A}(X_i), Y_i) := |\hat{\mu}(X_i) - Y_i|$ in regression with standard scores

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = s(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

Ex 1: $s(\hat{A}(X_i), Y_i) := |\hat{\mu}(X_i) - Y_i|$ in regression with standard scores

Ex 2: $s(\hat{A}(X_i), Y_i) := \max\left(\widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i)\right)$ in CQR

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = \mathbf{s}(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

Ex 1: $\mathbf{s}(\hat{A}(X_i), Y_i) := |\hat{\mu}(X_i) - Y_i|$ in regression with standard scores

Ex 2: $\mathbf{s}(\hat{A}(X_i), Y_i) := \max\left(\widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i)\right)$ in CQR

4. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = \mathbf{s}(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

Ex 1: $\mathbf{s}(\hat{A}(X_i), Y_i) := |\hat{\mu}(X_i) - Y_i|$ in regression with standard scores

Ex 2: $\mathbf{s}(\hat{A}(X_i), Y_i) := \max\left(\widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i)\right)$ in CQR

4. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$
5. For a new point X_{n+1} , return

$$\hat{C}_\alpha(X_{n+1}) = \{y \text{ such that } \mathbf{s}(\hat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the proper training set*
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = \mathbf{s}(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

Ex 1: $\mathbf{s}(\hat{A}(X_i), Y_i) := |\hat{\mu}(X_i) - Y_i|$ in regression with standard scores

Ex 2: $\mathbf{s}(\hat{A}(X_i), Y_i) := \max(\widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i))$ in CQR

4. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$
5. For a new point X_{n+1} , return

$$\hat{C}_\alpha(X_{n+1}) = \{y \text{ such that } \mathbf{s}(\hat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

Ex 1: $\hat{C}_\alpha(X_{n+1}) = [\hat{\mu}(X_{n+1}) \pm q_{1-\alpha}(\mathcal{S})]$

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the proper training set*
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = \mathbf{s}(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

Ex 1: $\mathbf{s}(\hat{A}(X_i), Y_i) := |\hat{\mu}(X_i) - Y_i|$ in regression with standard scores

Ex 2: $\mathbf{s}(\hat{A}(X_i), Y_i) := \max(\widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i))$ in CQR

4. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$
5. For a new point X_{n+1} , return

$$\hat{C}_\alpha(X_{n+1}) = \{y \text{ such that } \mathbf{s}(\hat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

Ex 2: $\hat{C}_\alpha(X_{n+1}) = [\widehat{QR}_{\text{lower}}(X_{n+1}) - q_{1-\alpha}(\mathcal{S});$

$$\widehat{QR}_{\text{upper}}(X_{n+1}) + q_{1-\alpha}(\mathcal{S})]$$

SCP is defined by the conformity score function



1. Randomly split the training data into a **proper training set** (size $\#Tr$) and a **calibration set** (size $\#Cal$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#Cal + 1$ **conformity scores**

$$\mathcal{S} = \{S_i = \mathbf{s}(\hat{A}(X_i), Y_i), i \in \text{Cal}\} \cup \{+\infty\}$$

Ex 1: $\mathbf{s}(\hat{A}(X_i), Y_i) := |\hat{\mu}(X_i) - Y_i|$ in regression with standard scores

Ex 2: $\mathbf{s}(\hat{A}(X_i), Y_i) := \max\left(\widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i)\right)$ in CQR

4. Compute the $1 - \alpha$ quantile of these scores, noted $q_{1-\alpha}(\mathcal{S})$
5. For a new point X_{n+1} , return

$$\hat{C}_\alpha(X_{n+1}) = \{y \text{ such that } \mathbf{s}(\hat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

\hookrightarrow The definition of the **conformity scores** is crucial, as they incorporate almost all the information: data + underlying model

This procedure enjoys the finite sample guarantee proposed and proved in Vovk et al. (2005).

Theorem (Marginal validity of SCP Vovk et al. (2005)).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^a. SCP on $(X_i, Y_i)_{i=1}^n$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

If, in addition, the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are almost surely distinct, then

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

^aOnly the calibration and test data need to be exchangeable.

Proof: application of the quantile lemma.

SCP: theoretical guarantees

This procedure enjoys the finite sample guarantee proposed and proved in Vovk et al. (2005).

Theorem (Marginal validity of SCP Vovk et al. (2005)).

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**^a. SCP on $(X_i, Y_i)_{i=1}^n$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

If, in addition, the scores $\{S_i\}_{i \in \text{Cal}} \cup \{S_{n+1}\}$ are almost surely distinct, then

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

^aOnly the calibration and test data need to be exchangeable.

Proof: application of the quantile lemma.

x Marginal coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \mid X_{n+1} = x \right\} \geq 1 - \alpha$

- $Y \in \{1, \dots, C\}$ (C classes)
- $\hat{A}(X) = (\hat{p}_1(X), \dots, \hat{p}_C(X))$ (estimated probabilities)

- $Y \in \{1, \dots, C\}$ (C classes)
- $\hat{A}(X) = (\hat{p}_1(X), \dots, \hat{p}_C(X))$ (estimated probabilities)
- $s(\hat{A}(X), Y) := 1 - (\hat{A}(X))_Y$

- $Y \in \{1, \dots, C\}$ (C classes)
- $\hat{A}(X) = (\hat{p}_1(X), \dots, \hat{p}_C(X))$ (estimated probabilities)
- $s(\hat{A}(X), Y) := 1 - (\hat{A}(X))_Y$
- For a new point X_{n+1} , return
$$\hat{C}_\alpha(X_{n+1}) = \{y \text{ such that } s(\hat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

Ex: $Y_i \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal_i	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.15	0.15	0.20	0.15	0.15	0.25	0.20
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.60	0.55	0.50	0.45	0.40	0.35	0.45
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.25	0.30	0.30	0.40	0.45	0.40	0.35
S_j	0.05	0.1	0.15	0.40	0.45	0.50	0.55	0.55	0.6	0.65

Ex: $Y_i \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal_i	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.15	0.15	0.20	0.15	0.15	0.25	0.20
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.60	0.55	0.50	0.45	0.40	0.35	0.45
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.25	0.30	0.30	0.40	0.45	0.40	0.35
S_j	0.05	0.1	0.15	0.40	0.45	0.50	0.55	0.55	0.6	0.65

- $q_{1-\alpha}(S) = 0.65$

Ex: $Y_i \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal_i	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.15	0.15	0.20	0.15	0.15	0.25	0.20
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.60	0.55	0.50	0.45	0.40	0.35	0.45
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.25	0.30	0.30	0.40	0.45	0.40	0.35
S_i	0.05	0.1	0.15	0.40	0.45	0.50	0.55	0.55	0.6	0.65

- $q_{1-\alpha}(S) = 0.65$
- $\hat{A}(X_{n+1}) = (0.05, 0.60, 0.35)$

SCP: standard classification in practice

Ex: $Y_i \in \{ \text{“dog”}, \text{“tiger”}, \text{“cat”} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal_i	“dog”	“dog”	“dog”	“tiger”	“tiger”	“tiger”	“tiger”	“cat”	“cat”	“cat”
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.15	0.15	0.20	0.15	0.15	0.25	0.20
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.60	0.55	0.50	0.45	0.40	0.35	0.45
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.25	0.30	0.30	0.40	0.45	0.40	0.35
S_j	0.05	0.1	0.15	0.40	0.45	0.50	0.55	0.55	0.6	0.65

- $q_{1-\alpha}(S) = 0.65$
- $\hat{A}(X_{n+1}) = (0.05, 0.60, 0.35)$
 $\hookrightarrow s(\hat{A}(X_{n+1}), \text{“dog”}) = 0.95$

“dog” $\notin \hat{C}_\alpha(X_{n+1})$

SCP: standard classification in practice

Ex: $Y_i \in \{ \text{“dog”}, \text{“tiger”}, \text{“cat”} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal_i	“dog”	“dog”	“dog”	“tiger”	“tiger”	“tiger”	“tiger”	“cat”	“cat”	“cat”
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.15	0.15	0.20	0.15	0.15	0.25	0.20
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.60	0.55	0.50	0.45	0.40	0.35	0.45
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.25	0.30	0.30	0.40	0.45	0.40	0.35
S_j	0.05	0.1	0.15	0.40	0.45	0.50	0.55	0.55	0.6	0.65

- $q_{1-\alpha}(\mathcal{S}) = 0.65$
- $\hat{A}(X_{n+1}) = (0.05, \mathbf{0.60}, 0.35)$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{“dog”}) = 0.95$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \mathbf{\text{“tiger”}}) = 0.40 \leq q_{1-\alpha}(\mathcal{S})$

“dog” $\notin \hat{C}_\alpha(X_{n+1})$
“tiger” $\in \hat{C}_\alpha(X_{n+1})$

SCP: standard classification in practice

Ex: $Y_i \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal_i	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.15	0.15	0.20	0.15	0.15	0.25	0.20
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.60	0.55	0.50	0.45	0.40	0.35	0.45
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.25	0.30	0.30	0.40	0.45	0.40	0.35
S_j	0.05	0.1	0.15	0.40	0.45	0.50	0.55	0.55	0.6	0.65

- $q_{1-\alpha}(\mathcal{S}) = 0.65$
- $\hat{A}(X_{n+1}) = (0.05, 0.60, \mathbf{0.35})$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"dog"}) = 0.95$ "dog" $\notin \hat{C}_\alpha(X_{n+1})$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"tiger"}) = 0.40 \leq q_{1-\alpha}(\mathcal{S})$ "tiger" $\in \hat{C}_\alpha(X_{n+1})$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"cat"}) = 0.65 \leq q_{1-\alpha}(\mathcal{S})$ "cat" $\in \hat{C}_\alpha(X_{n+1})$

SCP: standard classification in practice

Ex: $Y_i \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal_i	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.15	0.15	0.20	0.15	0.15	0.25	0.20
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.60	0.55	0.50	0.45	0.40	0.35	0.45
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.25	0.30	0.30	0.40	0.45	0.40	0.35
S_j	0.05	0.1	0.15	0.40	0.45	0.50	0.55	0.55	0.6	0.65

- $q_{1-\alpha}(\mathcal{S}) = 0.65$
- $\hat{A}(X_{n+1}) = (0.05, 0.60, 0.35)$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"dog"}) = 0.95$ "dog" $\notin \hat{C}_\alpha(X_{n+1})$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"tiger"}) = 0.40 \leq q_{1-\alpha}(\mathcal{S})$ "tiger" $\in \hat{C}_\alpha(X_{n+1})$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"cat"}) = 0.65 \leq q_{1-\alpha}(\mathcal{S})$ "cat" $\in \hat{C}_\alpha(X_{n+1})$
- $\hat{C}_\alpha(X_{n+1}) = \{ \text{"tiger"}, \text{"cat"} \}$

Ex: $Y \in \{\text{"dog"}, \text{"tiger"}, \text{"cat"}\}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.05	0.10	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.70	0.25	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.25	0.65	0.60	0.55
S_i	0.05	0.1	0.15	0.15	0.20	0.25	0.30	0.35	0.40	0.45

- $q_{1-\alpha}(\mathcal{S}) = 0.45$

Ex: $Y \in \{\text{"dog"}, \text{"tiger"}, \text{"cat"}\}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _i	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.05	0.10	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.70	0.25	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.25	0.65	0.60	0.55
S_i	0.05	0.1	0.15	0.15	0.20	0.25	0.30	0.35	0.40	0.45

- $q_{1-\alpha}(\mathcal{S}) = 0.45$
- $\hat{A}(X_{n+1}) = (0.05, 0.60, 0.35)$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"dog"}) = 0.95$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"tiger"}) = 0.40$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"cat"}) = 0.65$

Ex: $Y \in \{\text{"dog"}, \text{"tiger"}, \text{"cat"}\}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _i	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.05	0.10	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.70	0.25	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.25	0.65	0.60	0.55
S_i	0.05	0.1	0.15	0.15	0.20	0.25	0.30	0.35	0.40	0.45

- $q_{1-\alpha}(\mathcal{S}) = 0.45$
- $\hat{A}(X_{n+1}) = (0.05, 0.60, 0.35)$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"dog"}) = 0.95$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"tiger"}) = 0.40 \leq q_{1-\alpha}(\mathcal{S})$
 - $\hookrightarrow s(\hat{A}(X_{n+1}), \text{"cat"}) = 0.65$

Ex: $Y \in \{\text{"dog"}, \text{"tiger"}, \text{"cat"}\}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.05	0.10	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.70	0.25	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.25	0.65	0.60	0.55
S_i	0.05	0.1	0.15	0.15	0.20	0.25	0.30	0.35	0.40	0.45

- $q_{1-\alpha}(\mathcal{S}) = 0.45$

- $\hat{A}(X_{n+1}) = (0.05, 0.60, 0.35)$

$$\hookrightarrow s(\hat{A}(X_{n+1}), \text{"dog"}) = 0.95$$

$$\hookrightarrow s(\hat{A}(X_{n+1}), \text{"tiger"}) = 0.40 \leq q_{1-\alpha}(\mathcal{S})$$

$$\hookrightarrow s(\hat{A}(X_{n+1}), \text{"cat"}) = 0.65$$

$$\text{"dog"} \notin \hat{C}_\alpha(X_{n+1})$$

$$\text{"tiger"} \in \hat{C}_\alpha(X_{n+1})$$

$$\text{"cat"} \notin \hat{C}_\alpha(X_{n+1})$$

Ex: $Y \in \{\text{"dog"}, \text{"tiger"}, \text{"cat"}\}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.05	0.10	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.70	0.25	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.25	0.65	0.60	0.55
S_i	0.05	0.1	0.15	0.15	0.20	0.25	0.30	0.35	0.40	0.45

- $q_{1-\alpha}(\mathcal{S}) = 0.45$

- $\hat{A}(X_{n+1}) = (0.05, 0.60, 0.35)$

$$\hookrightarrow s(\hat{A}(X_{n+1}), \text{"dog"}) = 0.95$$

$$\hookrightarrow s(\hat{A}(X_{n+1}), \text{"tiger"}) = 0.40 \leq q_{1-\alpha}(\mathcal{S})$$

$$\hookrightarrow s(\hat{A}(X_{n+1}), \text{"cat"}) = 0.65$$

$$\text{"dog"} \notin \hat{C}_\alpha(X_{n+1})$$

$$\text{"tiger"} \in \hat{C}_\alpha(X_{n+1})$$

$$\text{"cat"} \notin \hat{C}_\alpha(X_{n+1})$$

- $\hat{C}_\alpha(X_{n+1}) = \{\text{"tiger"}\}$

efficiency yet non-adaptivity of the simplest classification scores

- ✓ Outputs the most efficient set possible (i.e. achieving the smallest average set size, Sadinle et al., 2018),
- ✗ Does not allow to discriminate between “easy” and “hard” test point. In practice, it leads to predictive sets that under-cover (resp. over-cover) on “hard” (resp. “easy”) subgroups. This is due to the fact that the same threshold $q_{1-\alpha}(\mathcal{S})$ is applied to any test point.

SCP: classification with Adaptive Prediction Sets⁸

1. Sort in decreasing order $\hat{p}_{\sigma_x(1)}(x) \geq \dots \geq \hat{p}_{\sigma_x(C)}(x)$

⁸Romano et al. (2020b), *Classification with Valid and Adaptive Coverage*, NeurIPS

SCP: classification with Adaptive Prediction Sets⁸

1. Sort in decreasing order $\hat{p}_{\sigma_x(1)}(x) \geq \dots \geq \hat{p}_{\sigma_x(C)}(x)$

2. $s(x, y; \hat{p}) := \sum_{k=1}^{\sigma_x^{-1}(y)} \hat{p}_{\sigma_x(k)}(x)$ (sum of the estimated probabilities associated to classes at least as large as that of the true class Y)

⁸Romano et al. (2020b), *Classification with Valid and Adaptive Coverage*, NeurIPS

SCP: classification with Adaptive Prediction Sets⁸

1. Sort in decreasing order $\hat{p}_{\sigma_x(1)}(x) \geq \dots \geq \hat{p}_{\sigma_x(C)}(x)$

2. $s(x, y; \hat{p}) := \sum_{k=1}^{\sigma_x^{-1}(y)} \hat{p}_{\sigma_x(k)}(x)$ (sum of the estimated probabilities associated to classes at least as large as that of the true class Y)

3. Return the set of classes $\{\sigma_{X_{n+1}}(1), \dots, \sigma_{X_{n+1}}(r^*)\}$, where

$$r^* = \arg \max_{1 \leq r \leq C} \left\{ \sum_{k=1}^r \hat{p}_{\sigma_{X_{n+1}}(k)}(X_{n+1}) < q_{1-\alpha}(\mathcal{S}) \right\} + 1$$

⁸Romano et al. (2020b), *Classification with Valid and Adaptive Coverage*, NeurIPS

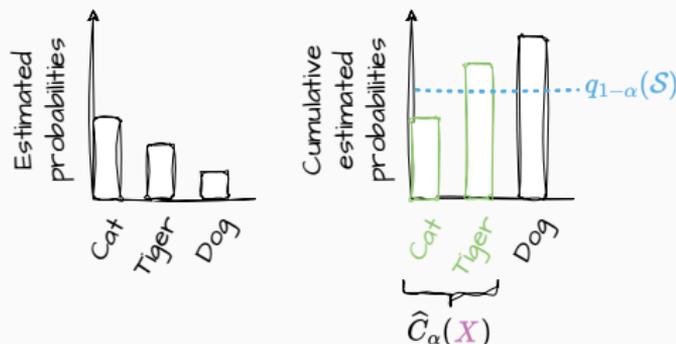
SCP: classification with Adaptive Prediction Sets⁸

1. Sort in decreasing order $\hat{p}_{\sigma_x(1)}(x) \geq \dots \geq \hat{p}_{\sigma_x(C)}(x)$

2. $s(x, y; \hat{p}) := \sum_{k=1}^{\sigma_x^{-1}(y)} \hat{p}_{\sigma_x(k)}(x)$ (sum of the estimated probabilities associated to classes at least as large as that of the true class Y)

3. Return the set of classes $\{\sigma_{X_{n+1}}(1), \dots, \sigma_{X_{n+1}}(r^*)\}$, where

$$r^* = \arg \max_{1 \leq r \leq C} \left\{ \sum_{k=1}^r \hat{p}_{\sigma_{X_{n+1}}(k)}(X_{n+1}) < q_{1-\alpha}(\mathcal{S}) \right\} + 1$$



⁸Romano et al. (2020b), *Classification with Valid and Adaptive Coverage*, NeurIPS
Figure highly inspired by Angelopoulos and Bates (2023).

Ex: $Y \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.10	0.25	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.75	0.40	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.15	0.35	0.60	0.55
S_i	0.95	0.90	0.85	0.85	0.80	0.75	0.75	0.75	0.60	0.55

Ex: $Y \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.10	0.25	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.75	0.40	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.15	0.35	0.60	0.55
S_i	0.95	0.90	0.85	0.85	0.80	0.75	0.75	0.75	0.60	0.55

- $q_{1-\alpha}(\mathcal{S}) = 0.95$

Ex: $Y \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.10	0.25	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.75	0.40	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.15	0.35	0.60	0.55
S_i	0.95	0.90	0.85	0.85	0.80	0.75	0.75	0.75	0.60	0.55

- $q_{1-\alpha}(\mathcal{S}) = 0.95$
 \hookrightarrow Ex 1: $\hat{A}(X_{n+1}) = (0.05, 0.45, 0.5)$

Ex: $Y \in \{\text{"dog"}, \text{"tiger"}, \text{"cat"}\}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.10	0.25	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.75	0.40	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.15	0.35	0.60	0.55
S_i	0.95	0.90	0.85	0.85	0.80	0.75	0.75	0.75	0.60	0.55

- $q_{1-\alpha}(\mathcal{S}) = 0.95$

\hookrightarrow Ex 1: $\hat{A}(X_{n+1}) = (0.05, 0.45, 0.5), r^* = 2$

$$\hat{C}_\alpha(X_{n+1}) = \{\text{"tiger"}, \text{"cat"}\}$$

Ex: $Y \in \{ \text{"dog"}, \text{"tiger"}, \text{"cat"} \}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.10	0.25	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.75	0.40	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.15	0.35	0.60	0.55
S_i	0.95	0.90	0.85	0.85	0.80	0.75	0.75	0.75	0.60	0.55

- $q_{1-\alpha}(\mathcal{S}) = 0.95$

\hookrightarrow Ex 1: $\hat{A}(X_{n+1}) = (0.05, 0.45, 0.5), r^* = 2$

$$\hat{C}_\alpha(X_{n+1}) = \{ \text{"tiger"}, \text{"cat"} \}$$

\hookrightarrow Ex 2: $\hat{A}(X_{n+1}) = (0.03, 0.95, 0.02)$

Ex: $Y \in \{\text{"dog"}, \text{"tiger"}, \text{"cat"}\}$, with $\alpha = 0.1$

- Scores on the calibration set

Cal _{<i>i</i>}	"dog"	"dog"	"dog"	"tiger"	"tiger"	"tiger"	"tiger"	"cat"	"cat"	"cat"
$\hat{p}_{\text{dog}}(X_i)$	0.95	0.90	0.85	0.05	0.05	0.05	0.10	0.25	0.10	0.15
$\hat{p}_{\text{tiger}}(X_i)$	0.02	0.05	0.10	0.85	0.80	0.75	0.75	0.40	0.30	0.30
$\hat{p}_{\text{cat}}(X_i)$	0.03	0.05	0.05	0.10	0.15	0.20	0.15	0.35	0.60	0.55
S_i	0.95	0.90	0.85	0.85	0.80	0.75	0.75	0.75	0.60	0.55

- $q_{1-\alpha}(\mathcal{S}) = 0.95$

\hookrightarrow Ex 1: $\hat{A}(X_{n+1}) = (0.05, 0.45, 0.5), r^* = 2$

$$\hat{C}_\alpha(X_{n+1}) = \{\text{"tiger"}, \text{"cat"}\}$$

\hookrightarrow Ex 2: $\hat{A}(X_{n+1}) = (0.03, 0.95, 0.02), r^* = 1$

$$\hat{C}_\alpha(X_{n+1}) = \{\text{"tiger"}\}$$

- **Simple** procedure which quantifies the uncertainty of **any** predictive model \hat{A} by returning predictive regions
- **Finite-sample** guarantees
- **Distribution-free** as long as the data are **exchangeable** (and so are the scores)

- **Simple** procedure which quantifies the uncertainty of **any** predictive model \hat{A} by returning predictive regions
- **Finite-sample** guarantees
- **Distribution-free** as long as the data are **exchangeable** (and so are the scores)
- **Marginal** theoretical guarantee over the joint (X, Y) distribution, and **not conditional**, i.e., no guarantee that for any x :

$$\mathbb{P} \left\{ Y_{n+1} \in \hat{C}_\alpha(X_{n+1}) \mid X_{n+1} = x \right\} \geq 1 - \alpha.$$

- **Simple** procedure which quantifies the uncertainty of **any** predictive model \hat{A} by returning predictive regions
- **Finite-sample** guarantees
- **Distribution-free** as long as the data are **exchangeable** (and so are the scores)
- **Marginal** theoretical guarantee over the joint (X, Y) distribution, and **not conditional**, i.e., no guarantee that for any x :

$$\mathbb{P} \left\{ Y_{n+1} \in \hat{C}_\alpha(X_{n+1}) \mid X_{n+1} = x \right\} \geq 1 - \alpha.$$

↪ marginal also over the whole calibration set and the test point!

Challenges: open questions (non exhaustive!)

- Conditional coverage
- Computational cost vs statistical power
- Exchangeability

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

On distribution-free X -conditional validity

Y -conditional validity

Impact of the calibration set on the coverage

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

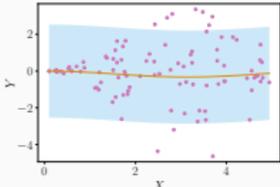
Concluding remarks

SCP: what choices for the regression scores?

$$\hat{C}_\alpha(X_{n+1}) = \{y \text{ such that } s(\hat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

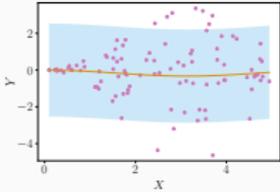
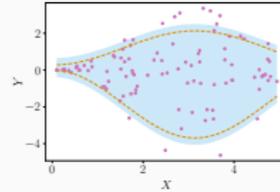
SCP: what choices for the regression scores?

$$\hat{C}_\alpha(X_{n+1}) = \{y \text{ such that } s(\hat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

	Standard SCP Vovk et al. (2005)		
$s(\hat{A}(X), Y)$	$ \hat{\mu}(X) - Y $		
$\hat{C}_\alpha(x)$	$[\hat{\mu}(x) \pm q_{1-\alpha}(\mathcal{S})]$		
Visu.			
✓	black-box around a “usable” prediction		
✗	not adaptive		

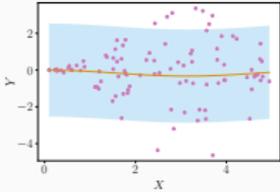
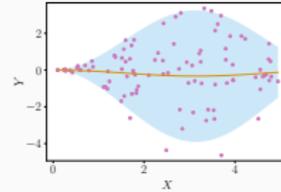
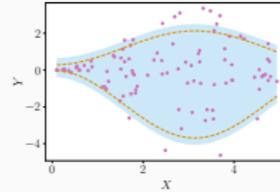
SCP: what choices for the regression scores?

$$\widehat{C}_\alpha(X_{n+1}) = \{y \text{ such that } s(\widehat{A}(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S})\}$$

	Standard SCP Vovk et al. (2005)	CQR Romano et al. (2019)
$s(\widehat{A}(X), Y)$	$ \widehat{\mu}(X) - Y $	$\max(\widehat{Q}R_{\text{lower}}(X) - Y, Y - \widehat{Q}R_{\text{upper}}(X))$
$\widehat{C}_\alpha(x)$	$[\widehat{\mu}(x) \pm q_{1-\alpha}(\mathcal{S})]$	$[\widehat{Q}R_{\text{lower}}(x) - q_{1-\alpha}(\mathcal{S}); \widehat{Q}R_{\text{upper}}(x) + q_{1-\alpha}(\mathcal{S})]$
Visu.		
✓	black-box around a “usable” prediction	adaptive
✗	not adaptive	no black-box around a “usable” prediction

SCP: what choices for the regression scores?

$$\widehat{C}_\alpha(X_{n+1}) = \{y \text{ such that } s(\widehat{A}(X_{n+1}), y) \leq q_{1-\alpha}(S)\}$$

	Standard SCP Vovk et al. (2005)	Locally weighted SCP Lei et al. (2018)	CQR Romano et al. (2019)
$s(\widehat{A}(X), Y)$	$ \widehat{\mu}(X) - Y $	$\frac{ \widehat{\mu}(X) - Y }{\widehat{\rho}(X)}$	$\max(\widehat{Q}R_{\text{lower}}(X) - Y, Y - \widehat{Q}R_{\text{upper}}(X))$
$\widehat{C}_\alpha(x)$	$[\widehat{\mu}(x) \pm q_{1-\alpha}(S)]$	$[\widehat{\mu}(x) \pm q_{1-\alpha}(S)\widehat{\rho}(x)]$	$[\widehat{Q}R_{\text{lower}}(x) - q_{1-\alpha}(S); \widehat{Q}R_{\text{upper}}(x) + q_{1-\alpha}(S)]$
Visu.			
✓	black-box around a “usable” prediction	black-box around a “usable” prediction	adaptive
✗	not adaptive	limited adaptiveness	no black-box around a “usable” prediction

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

On distribution-free X -conditional validity

Y -conditional validity

Impact of the calibration set on the coverage

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

Definition of distribution-free features conditional validity

\hat{C}_α = **estimated** predictive set based on n data points.

Definition (Distribution-free X -conditional validity).

\hat{C}_α achieves **distribution-free X -conditional validity** if:

- for any distribution \mathcal{D} ,
- for any associated exchangeable joint distribution $\mathcal{D}^{\text{exch}(n+1)}$,

we have that:

$$\mathbb{P}_{\mathcal{D}^{\text{exch}(n+1)}} \left(Y^{(n+1)} \in \hat{C}_\alpha \left(X^{(n+1)} \right) \mid X^{(n+1)} \right) \stackrel{\text{a.s.}}{\geq} 1 - \alpha.$$

Theorem (Impossibility results Vovk (2012); Lei and Wasserman (2014)).

If \widehat{C}_α is distribution-free X -conditionally valid, then, for any \mathcal{D} , for \mathcal{D}_X -almost all \mathcal{D}_X -non-atoms $x \in \mathcal{X}$, it holds:

- ▶ *Regression*: $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(\text{mes} \left(\widehat{C}_\alpha(x) \right) = \infty \right) \geq 1 - \alpha,$
- ▶ *Classification*: for any $y \in \mathcal{Y}$, $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(y \in \widehat{C}_\alpha(x) \right) \geq 1 - \alpha.$

Theorem (Impossibility results Vovk (2012); Lei and Wasserman (2014)).

If \hat{C}_α is **distribution-free X -conditionally valid**, then, for any \mathcal{D} , for \mathcal{D}_X -almost all \mathcal{D}_X -non-atoms $x \in \mathcal{X}$, it holds:

- ▶ *Regression*: $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(\text{mes} \left(\hat{C}_\alpha(x) \right) = \infty \right) \geq 1 - \alpha,$
- ▶ *Classification*: for any $y \in \mathcal{Y}$, $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(y \in \hat{C}_\alpha(x) \right) \geq 1 - \alpha.$

↔ distribution-free X -conditional hardness result apply beyond CP

Theorem (Impossibility results Vovk (2012); Lei and Wasserman (2014)).

If \widehat{C}_α is distribution-free X -conditionally valid, then, **for any** \mathcal{D} , for \mathcal{D}_X -almost all \mathcal{D}_X -non-atoms $x \in \mathcal{X}$, it holds:

- ▶ *Regression*: $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(\text{mes} \left(\widehat{C}_\alpha(x) \right) = \infty \right) \geq 1 - \alpha,$
- ▶ *Classification*: for any $y \in \mathcal{Y}$, $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(y \in \widehat{C}_\alpha(x) \right) \geq 1 - \alpha.$

↔ distribution-free X -conditional hardness result apply beyond CP

↔ X -conditional estimators are overly large even on easy cases

Theorem (Impossibility results Vovk (2012); Lei and Wasserman (2014)).

If \widehat{C}_α is distribution-free X -conditionally valid, then, for any \mathcal{D} , for \mathcal{D}_X -almost all \mathcal{D}_X -non-atoms $x \in \mathcal{X}$, it holds:

- ▶ *Regression*: $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(\text{mes} \left(\widehat{C}_\alpha(x) \right) = \infty \right) \geq 1 - \alpha,$
- ▶ *Classification*: for any $y \in \mathcal{Y}$, $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(y \in \widehat{C}_\alpha(x) \right) \geq 1 - \alpha.$

↪ distribution-free X -conditional hardness result apply beyond CP

↪ X -conditional estimators are overly large even on easy cases

↪ the lower bounds are tight

Example (Naive estimator).

$C_\alpha(\cdot; \xi) \equiv \mathcal{Y} \mathbb{1} \{ \xi \leq 1 - \alpha \} + \emptyset \mathbb{1} \{ \xi > \alpha \},$ where $\xi \sim \mathcal{U}([0, 1]).$

Theorem (Impossibility results Vovk (2012); Lei and Wasserman (2014)).

If \widehat{C}_α is distribution-free X -conditionally valid, then, for any \mathcal{D} , for \mathcal{D}_X -almost all \mathcal{D}_X -non-atoms $x \in \mathcal{X}$, it holds:

- ▶ **Regression:** $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(\text{mes} \left(\widehat{C}_\alpha(x) \right) = \infty \right) \geq 1 - \alpha,$
- ▶ **Classification:** for any $y \in \mathcal{Y}$, $\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left(y \in \widehat{C}_\alpha(x) \right) \geq 1 - \alpha.$

↔ distribution-free X -conditional hardness result apply beyond CP

↔ X -conditional estimators are overly large even on easy cases

↔ the lower bounds are tight

↔ Classification: every label is likely to be included in \widehat{C}_α .

\widehat{C}_α is likely to be large: for any \mathcal{D} , for \mathcal{D}_X -almost all \mathcal{D}_X -non-atoms $x \in \mathcal{X}$,
 $\mathbb{E}_{\mathcal{D}^{\otimes(n)}} \left[\#\widehat{C}_\alpha(x) \right] \geq (1 - \alpha)\#\mathcal{Y}.$

Definition (distribution-free $(1 - \alpha, \delta)$ - \mathcal{X} -conditional validity).

Let $\delta > 0$ be a tolerance level.

An estimator \widehat{C}_α achieves distribution-free $(1 - \alpha, \delta)$ - \mathcal{X} -conditional validity if for any distribution \mathcal{D} , for any $\mathcal{X} \subseteq \mathcal{X}$ such that $\mathbb{P}_{\mathcal{D}_X}(X \in \mathcal{X}) \geq \delta$, and for any associated exchangeable joint distribution $\mathcal{D}^{\text{exch}(n+1)}$, we have:

$$\mathbb{P}_{\mathcal{D}^{\text{exch}(n+1)}} \left(Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \mid X_{n+1} \in \mathcal{X} \right) \geq 1 - \alpha.$$

Weaker notion of \mathcal{X} -conditional validity (Barber et al., 2021a)

Definition (distribution-free $(1 - \alpha, \delta)$ - \mathcal{X} -conditional validity).

Let $\delta > 0$ be a tolerance level.

An estimator \hat{C}_α achieves distribution-free $(1 - \alpha, \delta)$ - \mathcal{X} -conditional validity if for any distribution \mathcal{D} , for any $\mathcal{X} \subseteq \mathcal{X}$ such that $\mathbb{P}_{\mathcal{D}_X}(X \in \mathcal{X}) \geq \delta$, and for any associated exchangeable joint distribution $\mathcal{D}^{\text{exch}(n+1)}$, we have:

$$\mathbb{P}_{\mathcal{D}^{\text{exch}(n+1)}} \left(Y_{n+1} \in \hat{C}_\alpha(X_{n+1}) \mid X_{n+1} \in \mathcal{X} \right) \geq 1 - \alpha.$$

Informal theorem (lower bound on $(1 - \alpha, \delta)$ - \mathcal{X} -cond. valid efficiency)

An estimator achieving $(1 - \alpha, \delta)$ - \mathcal{X} -conditional validity can not be more efficient than an estimator achieving **distribution-free marginal validity at the level $1 - \alpha\delta$** .

\Leftrightarrow In practice, consider small $\delta \rightarrow$ unefficient predictive sets.

Definition (distribution-free group-features-conditional validity).

Let $G := \left(G^{(k)}\right)_{k=1}^K$ represents groups on the features space (possibly overlapping).

An estimator \hat{C}_α achieves distribution-free G -conditional validity if for any distribution \mathcal{D} , and for any associated exchangeable joint distribution $\mathcal{D}^{\text{exch}(n+1)}$, we have:

$$\mathbb{P}_{\mathcal{D}^{\text{exch}(n+1)}} \left(Y_{n+1} \in \hat{C}_\alpha(X_{n+1}, G_{n+1}) \mid G_{n+1} \right) \stackrel{\text{a.s.}}{\geq} 1 - \alpha.$$

Theorem (General MCV hardness result).

If \widehat{C}_α is distribution-free group-features-conditionally valid then for any distribution \mathcal{D} , for any group g such that $\mathcal{D}_G(g) := \mathbb{P}_{\mathcal{D}}(G = g) > 0$, it holds:

► *Regression*

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha(X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1},$$

► *Classification*

$$\text{for any } y \in \mathcal{Y}, \mathbb{P}_{\mathcal{P}^{\otimes(n+1)}} \left(\text{alert} < 2 | \text{handout} : 0 > y \in \widehat{C}_\alpha(X_{n+1}, g) \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1}.$$

Theorem (General MCV hardness result).

If \widehat{C}_α is distribution-free group-features-conditionally valid then for any distribution \mathcal{D} , for any group g such that $\mathcal{D}_G(g) := \mathbb{P}_{\mathcal{D}}(G = g) > 0$, it holds:

► *Regression*

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha(X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1},$$

► *Classification*

$$\text{for any } y \in \mathcal{Y}, \mathbb{P}_{\mathcal{P}^{\otimes(n+1)}} \left(\text{alert} < 2 | \text{handout} : 0 > y \in \widehat{C}_\alpha(X_{n+1}, g) \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1}.$$

Theorem (General MCV hardness result).

If \widehat{C}_α is distribution-free group-features-conditionally valid then for any distribution \mathcal{D} , for any group g such that $\mathcal{D}_G(g) := \mathbb{P}_{\mathcal{D}}(G = g) > 0$, it holds:

► *Regression*

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha(X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1},$$

► *Classification*

$$\text{for any } y \in \mathcal{Y}, \mathbb{P}_{\mathcal{P}^{\otimes(n+1)}} \left(\text{alert} < 2 | \text{handout} : 0 > y \in \widehat{C}_\alpha(X_{n+1}, g) \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1}.$$

Theorem (General MCV hardness result).

If \widehat{C}_α is distribution-free group-features-conditionally valid then for any distribution \mathcal{D} , for any group g such that $\mathcal{D}_G(g) := \mathbb{P}_{\mathcal{D}}(G = g) > 0$, it holds:

► *Regression*

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha(X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1},$$

► *Classification*

$$\text{for any } y \in \mathcal{Y}, \mathbb{P}_{\mathcal{P}^{\otimes(n+1)}} \left(\text{alert} < 2 | \text{handout} : 0 > y \in \widehat{C}_\alpha(X_{n+1}, g) \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1}.$$

Irreducible term: consider \widehat{C}_α outputting \mathcal{Y} with probability $1 - \alpha$ and \emptyset otherwise.

Theorem (General MCV hardness result).

If \widehat{C}_α is distribution-free group-features-conditionally valid then for any distribution \mathcal{D} , for any group g such that $\mathcal{D}_G(g) := \mathbb{P}_{\mathcal{D}}(G = g) > 0$, it holds:

► *Regression*

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha(X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g)\sqrt{n+1},$$

► *Classification*

$$\text{for any } y \in \mathcal{Y}, \mathbb{P}_{\mathcal{P}^{\otimes(n+1)}} \left(\text{alert} < 2 | \text{handout} : 0 > y \in \widehat{C}_\alpha(X_{n+1}, g) \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g)\sqrt{n+1}.$$

Irreducible term: consider \widehat{C}_α outputting \mathcal{Y} with probability $1 - \alpha$ and \emptyset otherwise.

$\Delta_{g,n}$ term: smaller than $\mathcal{D}_G(g)\sqrt{n+1}$

Theorem (General MCV hardness result).

If \widehat{C}_α is distribution-free group-features-conditionally valid then for any distribution \mathcal{D} , for any group g such that $\mathcal{D}_G(g) := \mathbb{P}_{\mathcal{D}}(G = g) > 0$, it holds:

► *Regression*

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha(X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g)\sqrt{n+1},$$

► *Classification*

$$\text{for any } y \in \mathcal{Y}, \mathbb{P}_{\mathcal{P}^{\otimes(n+1)}} \left(\text{alert} < 2 | \text{handout} : 0 > y \in \widehat{C}_\alpha(X_{n+1}, g) \right) \geq 1 - \alpha - \Delta_{g,n} \geq 1 - \alpha - \mathcal{D}_G(g)\sqrt{n+1}.$$

Irreducible term: consider \widehat{C}_α outputting \mathcal{Y} with probability $1 - \alpha$ and \emptyset otherwise.

$\Delta_{g,n}$ term: smaller than $\mathcal{D}_G(g)\sqrt{n+1}$

↪ gets negligible (making the lower bound nearly $1 - \alpha$) **only** for low probability groups compared to n .

Restricting the link between G and (X or Y) does not allow informative G -conditional-coverage (Zaffran et al., 2024)

Analogous statements are also available for the classification framework.

Restricting the link between G and $(X$ or $Y)$ does not allow informative G -conditional-coverage (Zaffran et al., 2024)

Theorem ($G \perp\!\!\!\perp X$ hardness result).

If any \hat{C}_α is G -conditionally-valid under $G \perp\!\!\!\perp X$, then for any distribution \mathcal{D} such that $G \perp\!\!\!\perp X$, for any group g such that $\mathcal{D}_G(g) > 0$, it holds:

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\hat{C}_\alpha (X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1}.$$

Restricting the link between G and $(X$ or $Y)$ does not allow informative G -conditional-coverage (Zaffran et al., 2024)

Theorem ($G \perp\!\!\!\perp X$ hardness result).

If any \widehat{C}_α is G -conditionally-valid under $G \perp\!\!\!\perp X$, then for any distribution \mathcal{D} such that $G \perp\!\!\!\perp X$, for any group g such that $\mathcal{D}_G(g) > 0$, it holds:

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha (X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1}.$$

Theorem ($Y \perp\!\!\!\perp G \mid X$ hardness result).

If any \widehat{C}_α is G -conditionally-valid under $Y \perp\!\!\!\perp G \mid X$, then for any distribution \mathcal{D} such that $Y \perp\!\!\!\perp G \mid X$, for any mask g such that $\frac{1}{\sqrt{2}} \geq \mathcal{D}_G(g) > 0$, it holds:

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\widehat{C}_\alpha (X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - 2\mathcal{D}_G(g) \sqrt{n+1}.$$

Analogous statements are also available for the classification framework.

Restricting the link between G and $(X$ or $Y)$ does not allow informative G -conditional-coverage (Zaffran et al., 2024)

Theorem ($G \perp\!\!\!\perp X$ hardness result).

If any \hat{C}_α is G -conditionally-valid under $G \perp\!\!\!\perp X$, then for any distribution \mathcal{D} such that $G \perp\!\!\!\perp X$, for any group g such that $\mathcal{D}_G(g) > 0$, it holds:

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\hat{C}_\alpha (X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - \mathcal{D}_G(g) \sqrt{n+1}.$$

Theorem ($Y \perp\!\!\!\perp G \mid X$ hardness result).

If any \hat{C}_α is G -conditionally-valid under $Y \perp\!\!\!\perp G \mid X$, then for any distribution \mathcal{D} such that $Y \perp\!\!\!\perp G \mid X$, for any mask g such that $\frac{1}{\sqrt{2}} \geq \mathcal{D}_G(g) > 0$, it holds:

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\text{mes} \left(\hat{C}_\alpha (X_{n+1}, g) \right) = \infty \right) \geq 1 - \alpha - 2\mathcal{D}_G(g) \sqrt{n+1}.$$

\Rightarrow Need to restrict **both** the link between G and X , as well as between G and Y .

Analogous statements are also available for the classification framework.

- Approximate conditional coverage

↪ Romano et al. (2020a); Guan (2022); Jung et al. (2023); Gibbs et al. (2023)

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

- Approximate conditional coverage
↪ Romano et al. (2020a); Guan (2022); Jung et al. (2023); Gibbs et al. (2023)
Target $\mathbb{P}(Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) | X_{n+1} \in \mathcal{R}(x)) \geq 1 - \alpha$
- Asymptotic (with the sample size) conditional coverage
↪ Romano et al. (2019); Kivaranovic et al. (2020); Chernozhukov et al. (2021); Sesia and Romano (2021); Izbicki et al. (2022)

- Approximate conditional coverage
↪ Romano et al. (2020a); Guan (2022); Jung et al. (2023); Gibbs et al. (2023)
Target $\mathbb{P}(Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) | X_{n+1} \in \mathcal{R}(x)) \geq 1 - \alpha$
- Asymptotic (with the sample size) conditional coverage
↪ **Romano et al. (2019)**; Kivaranovic et al. (2020); Chernozhukov et al. (2021); Sesia and Romano (2021); Izbicki et al. (2022)

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

On distribution-free X -conditional validity

Y -conditional validity

Impact of the calibration set on the coverage

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

Achieving Y -conditional validity in classification

1. Randomly split the training data into a **proper training set** (size $\#\text{Tr}$) and a **calibration set** (size $\#\text{Cal}$)
2. Get \hat{A} (by training \mathcal{A} on the **proper training set** $(X_i, Y_i)_{i \in \text{Tr}}$)

3. **For** any candidate $y \in \mathcal{Y}$:

On the **calibration set**, obtain a set of $\#\text{Cal}_y + 1$ **conformity scores**:

$$\mathcal{S}_y = \{S_i = \mathbf{s}(X_i, y; \hat{A}), i \in \text{Cal such that } Y_i = y\} \cup \{+\infty\}$$

4. For a new point X_{n+1} , return $\hat{C}_{n,\alpha}(X_{n+1}) \{y \text{ such that } \mathbf{s}(X_{n+1}, y; \hat{A}) \leq q_{1-\alpha}(\mathcal{S}_y)\}$

Achieving Y -conditional validity in classification

1. Randomly split the training data into a **proper training set** (size $\#\text{Tr}$) and a **calibration set** (size $\#\text{Cal}$)
2. Get \hat{A} (by training \mathcal{A} on the **proper training set** $(X_i, Y_i)_{i \in \text{Tr}}$)

3. **For** any candidate $y \in \mathcal{Y}$:

On the **calibration set**, obtain a set of $\#\text{Cal}_y + 1$ **conformity scores** :

$$\mathcal{S}_y = \{S_i = \mathbf{s}(X_i, y; \hat{A}), i \in \text{Cal such that } Y_i = y\} \cup \{+\infty\}$$

4. For a new point X_{n+1} , return $\hat{C}_{n,\alpha}(X_{n+1}) \{y \text{ such that } \mathbf{s}(X_{n+1}, y; \hat{A}) \leq q_{1-\alpha}(\mathcal{S}_y)\}$

\hookrightarrow What if there is a high class imbalance?

Ding et al. (2023) proposed to instead obtain **cluster**-conditional coverage.

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

On distribution-free X -conditional validity

Y -conditional validity

Impact of the calibration set on the coverage

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

Probably Approximately Correct bounds on calibration-conditional coverage (Vovk, 2012; Bian and Barber, 2023)

Theorem (calibration conditional validity of SCP).

SCP outputs \hat{C}_α such that for any distribution \mathcal{D} and any $0 < \delta \leq 0.5$:

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left(\mathbb{P}_{\mathcal{D}} \left(Y_{n+1} \notin \hat{C}_{n,\alpha}(X_{n+1}) \mid (X_i, Y_i)_{i=1}^n \right) \leq \alpha + \sqrt{\frac{\log(1/\delta)}{2\#\text{Cal}}} \right) \geq 1 - \delta.$$

\leftrightarrow controls the deviation of miscoverage with respect to the nominal level of a predictive set built on a given calibration set.

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Full Conformal Prediction

Jackknife+

Beyond exchangeability

Some case studies

Concluding remarks

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Full Conformal Prediction

Jackknife+

Beyond exchangeability

Some case studies

Concluding remarks

SCP suffers from data splitting:

- lower statistical efficiency (lower model accuracy and higher predictive set size)
- higher statistical variability

SCP suffers from data splitting:

- lower statistical efficiency (lower model accuracy and higher predictive set size)
- higher statistical variability

Can we avoid splitting the data set?

The naive idea does not enjoy valid coverage (even empirically)

- A naive idea:
 - Get \hat{A} by training the algorithm \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$.

The naive idea does not enjoy valid coverage (even empirically)

- A naive idea:
 - Get \hat{A} by training the algorithm \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$.
 - compute the empirical quantile $q_{1-\alpha}(\mathcal{S})$ of the set of scores

$$\mathcal{S} = \left\{ \mathbf{s} \left(\hat{A}(X_i), Y_i \right) \right\}_{i=1}^n \cup \{\infty\}.$$

The naive idea does not enjoy valid coverage (even empirically)

- A naive idea:
 - Get \hat{A} by training the algorithm \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$.
 - compute the empirical quantile $q_{1-\alpha}(\mathcal{S})$ of the set of scores

$$\mathcal{S} = \left\{ \mathbf{s} \left(\hat{A}(X_i), Y_i \right) \right\}_{i=1}^n \cup \{\infty\}.$$

- output the set $\left\{ y \text{ such that } \mathbf{s} \left(\hat{A}(X_{n+1}), y \right) \leq q_{1-\alpha}(\mathcal{S}) \right\}$.

The naive idea does not enjoy valid coverage (even empirically)

- A naive idea:
 - Get \hat{A} by training the algorithm \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$.
 - compute the empirical quantile $q_{1-\alpha}(\mathcal{S})$ of the set of scores

$$\mathcal{S} = \left\{ \mathbf{s} \left(\hat{A}(X_i), Y_i \right) \right\}_{i=1}^n \cup \{\infty\}.$$

- output the set $\{y \text{ such that } \mathbf{s} \left(\hat{A}(X_{n+1}), y \right) \leq q_{1-\alpha}(\mathcal{S})\}$.

✗ \hat{A} obtained w. the training set $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ but not X_{n+1} .

Example (“Naive Idea” sets with an interpolating algorithm).

Assume \mathcal{A} interpolates:

- $\hat{A} = \mathcal{A}((x_1, y_1), \dots, (x_n, y_n))$
- $\hat{A}(x_k) - y_k = 0$ for any $k \in \llbracket 1, n \rrbracket$

\Rightarrow Naive method above (with MAE score functions) outputs $\{\hat{A}(X_{n+1})\}$
(a single point) for any new test point!

Full Conformal Prediction⁹ does not discard training points!

- Full (or transductive) Conformal Prediction
 - avoids data splitting

⁹Vovk et al. (2005), *Algorithmic Learning in a Random World*

Full Conformal Prediction⁹ does not discard training points!

- Full (or transductive) Conformal Prediction
 - avoids data splitting
 - at the cost of many more model fits

⁹Vovk et al. (2005), *Algorithmic Learning in a Random World*

Full Conformal Prediction⁹ does not discard training points!

- Full (or transductive) Conformal Prediction
 - avoids data splitting
 - at the cost of many more model fits
- **Idea:** the most probable labels Y_{n+1} live in \mathcal{Y} , and have a low enough conformity score. By looping over all possible $y \in \mathcal{Y}$, the ones leading to the smallest conformity scores will be found.

⁹Vovk et al. (2005), *Algorithmic Learning in a Random World*

For any candidate (X_{n+1}, y) :

1. Get \hat{A}_y by training \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \cup \{(X_{n+1}, y)\}$

Full Conformal Prediction (CP): recovering exchangeability

For any candidate (X_{n+1}, y) :

1. Get \hat{A}_y by training \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \cup \{(X_{n+1}, y)\}$
2. Obtain a set of training scores

$$\mathcal{S}_y^{(\text{train})} = \left\{ \mathbf{s}(\hat{A}_y(X_i), Y_i) \right\}_{i=1}^n \cup \left\{ \mathbf{s}(\hat{A}_y(X_{n+1}), y) \right\}$$

and compute their $1 - \alpha$ empirical quantile $q_{1-\alpha}(\mathcal{S}_y^{(\text{train})})$

Full Conformal Prediction (CP): recovering exchangeability

For any candidate (X_{n+1}, y) :

1. Get \hat{A}_y by training \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \cup \{(X_{n+1}, y)\}$
2. Obtain a set of training scores

$$\mathcal{S}_y^{(\text{train})} = \left\{ \mathbf{s}(\hat{A}_y(X_i), Y_i) \right\}_{i=1}^n \cup \left\{ \mathbf{s}(\hat{A}_y(X_{n+1}), y) \right\}$$

and compute their $1 - \alpha$ empirical quantile $q_{1-\alpha}(\mathcal{S}_y^{(\text{train})})$

Output the set $\left\{ y \text{ such that } \mathbf{s}(\hat{A}_y(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S}_y^{(\text{train})}) \right\}$.

Full Conformal Prediction (CP): recovering exchangeability

For any candidate (X_{n+1}, y) :

1. Get \hat{A}_y by training \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \cup \{(X_{n+1}, y)\}$
2. Obtain a set of training scores

$$\mathcal{S}_y^{(\text{train})} = \left\{ \mathbf{s}(\hat{A}_y(X_i), Y_i) \right\}_{i=1}^n \cup \left\{ \mathbf{s}(\hat{A}_y(X_{n+1}), y) \right\}$$

and compute their $1 - \alpha$ empirical quantile $q_{1-\alpha}(\mathcal{S}_y^{(\text{train})})$

Output the set $\left\{ y \text{ such that } \mathbf{s}(\hat{A}_y(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S}_y^{(\text{train})}) \right\}$.

- ✓ Test point treated in the same way than train points

Full Conformal Prediction (CP): recovering exchangeability

For any candidate (X_{n+1}, y) :

1. Get \hat{A}_y by training \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \cup \{(X_{n+1}, y)\}$
2. Obtain a set of training scores

$$\mathcal{S}_y^{(\text{train})} = \left\{ \mathbf{s}(\hat{A}_y(X_i), Y_i) \right\}_{i=1}^n \cup \left\{ \mathbf{s}(\hat{A}_y(X_{n+1}), y) \right\}$$

and compute their $1 - \alpha$ empirical quantile $q_{1-\alpha}(\mathcal{S}_y^{(\text{train})})$

Output the set $\left\{ y \text{ such that } \mathbf{s}(\hat{A}_y(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S}_y^{(\text{train})}) \right\}$.

- ✓ Test point treated in the same way than train points
- ✓ Any score works

Full Conformal Prediction (CP): recovering exchangeability

For any candidate (X_{n+1}, y) :

1. Get \hat{A}_y by training \mathcal{A} on $\{(X_1, Y_1), \dots, (X_n, Y_n)\} \cup \{(X_{n+1}, y)\}$
2. Obtain a set of training scores

$$\mathcal{S}_y^{(\text{train})} = \left\{ \mathbf{s}(\hat{A}_y(X_i), Y_i) \right\}_{i=1}^n \cup \left\{ \mathbf{s}(\hat{A}_y(X_{n+1}), y) \right\}$$

and compute their $1 - \alpha$ empirical quantile $q_{1-\alpha}(\mathcal{S}_y^{(\text{train})})$

Output the set $\left\{ y \text{ such that } \mathbf{s}(\hat{A}_y(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S}_y^{(\text{train})}) \right\}$.

- ✓ Test point treated in the same way than train points
- ✓ Any score works
- ✗ Computationally costly

Definition (Symmetrical algorithm).

A deterministic algorithm $\mathcal{A} : (U_1, \dots, U_n) \mapsto \hat{A}$ is **symmetric** if for any permutation σ of $\llbracket 1, n \rrbracket$: $\mathcal{A}(U_1, \dots, U_n) \stackrel{\text{a.s.}}{=} \mathcal{A}(U_{\sigma(1)}, \dots, U_{\sigma(n)})$.

Definition (Symmetrical algorithm).

A deterministic algorithm $\mathcal{A} : (U_1, \dots, U_n) \mapsto \hat{A}$ is **symmetric** if for any permutation σ of $\llbracket 1, n \rrbracket$: $\mathcal{A}(U_1, \dots, U_n) \stackrel{\text{a.s.}}{=} \mathcal{A}(U_{\sigma(1)}, \dots, U_{\sigma(n)})$.

Lemma (Exchangeable scores).

If the algorithm $\mathcal{A} : (U_1, \dots, U_n) \mapsto \hat{A}$ is **symmetric**, and $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**, then S_1, \dots, S_{n+1} are exchangeable, with

$$S_i := \mathbf{s}(\hat{A}_{Y_{n+1}}(X_i), Y_i).$$

Definition (Symmetrical algorithm).

A deterministic algorithm $\mathcal{A} : (U_1, \dots, U_n) \mapsto \hat{A}$ is **symmetric** if for any permutation σ of $\llbracket 1, n \rrbracket$: $\mathcal{A}(U_1, \dots, U_n) \stackrel{\text{a.s.}}{=} \mathcal{A}(U_{\sigma(1)}, \dots, U_{\sigma(n)})$.

Lemma (Exchangeable scores).

If the algorithm $\mathcal{A} : (U_1, \dots, U_n) \mapsto \hat{A}$ is **symmetric**, and $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**, then S_1, \dots, S_{n+1} are exchangeable, with

$$S_i := \mathbf{s}(\hat{A}_{Y_{n+1}}(X_i), Y_i).$$

Moreover

$$Y_{n+1} \in \widehat{C}_\alpha^{\text{Full}}(X_{n+1}) := \left\{ y \text{ such that } \mathbf{s}(\hat{A}_y(X_{n+1}), y) \leq q_{1-\alpha}(\mathcal{S}_y^{(\text{train})}) \right\}$$

$$\Leftrightarrow \mathbf{s}(\hat{A}_{Y_{n+1}}(X_{n+1}), Y_{n+1}) \leq q_{1-\alpha}(\mathcal{S}_{Y_{n+1}}^{(\text{train})})$$

$$\Leftrightarrow S_{n+1} \leq q_{1-\alpha}(S_1, \dots, S_n, S_{n+1}) !$$

Full CP enjoys finite sample guarantees proved in Vovk et al. (2005).

Theorem (Marginal validity of Full CP Vovk et al. (2005)).

Suppose that

- (i) $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**,
- (ii) the algorithm \mathcal{A} is **symmetric**.

Full CP applied on $(X_i, Y_i)_{i=1}^n \cup \{X_{n+1}\}$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

Additionally, if the scores are a.s. distinct:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{n+1}.$$

Full CP enjoys finite sample guarantees proved in Vovk et al. (2005).

Theorem (Marginal validity of Full CP Vovk et al. (2005)).

Suppose that

- (i) $(X_i, Y_i)_{i=1}^{n+1}$ are **exchangeable**,
- (ii) the algorithm \mathcal{A} is **symmetric**.

Full CP applied on $(X_i, Y_i)_{i=1}^n \cup \{X_{n+1}\}$ outputs $\widehat{C}_\alpha(\cdot)$ such that:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \geq 1 - \alpha.$$

Additionally, if the scores are a.s. distinct:

$$\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \right\} \leq 1 - \alpha + \frac{1}{n+1}.$$

x Marginal coverage: $\mathbb{P} \left\{ Y_{n+1} \in \widehat{C}_\alpha(X_{n+1}) \mid X_{n+1} = x \right\} \geq 1 - \alpha$

Example (FCP sets with an interpolating algorithm).

Assume \mathcal{A} interpolates:

- $\hat{A} = \mathcal{A}((x_1, y_1), \dots, (x_{n+1}, y_{n+1}))$
- $\hat{A}(x_k) - y_k = 0$ for any $k \in \llbracket 1, n+1 \rrbracket$

Example (FCP sets with an interpolating algorithm).

Assume \mathcal{A} interpolates:

- $\hat{A} = \mathcal{A}((x_1, y_1), \dots, (x_{n+1}, y_{n+1}))$
- $\hat{A}(x_k) - y_k = 0$ for any $k \in \llbracket 1, n+1 \rrbracket$

\Rightarrow Full Conformal Prediction (*with standard score functions*) outputs \mathcal{Y} (the whole label space) for any new test point!

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Full Conformal Prediction

Jackknife+

Beyond exchangeability

Some case studies

Concluding remarks

Jackknife: the naive idea does not enjoy valid coverage

- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$



Jackknife: the naive idea does not enjoy valid coverage

- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$
- **LOO scores** $\mathcal{S} = \left\{ |\hat{A}_{-i}(X_i) - Y_i| \right\}_i \cup \{+\infty\}$ (in standard mean regression)



Jackknife: the naive idea does not enjoy valid coverage

- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$
- **LOO scores** $\mathcal{S} = \left\{ |\hat{A}_{-i}(X_i) - Y_i| \right\}_i \cup \{+\infty\}$ (in standard mean regression)
- Get \hat{A} by training \mathcal{A} on \mathcal{D}_n



Jackknife: the naive idea does not enjoy valid coverage

- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$
- **LOO scores** $\mathcal{S} = \left\{ |\hat{A}_{-i}(X_i) - Y_i| \right\}_i \cup \{+\infty\}$ (in standard mean regression)
- Get \hat{A} by training \mathcal{A} on \mathcal{D}_n
- Build the predictive interval: $\left[\hat{A}(X_{n+1}) \pm q_{1-\alpha}(\mathcal{S}) \right]$



Jackknife: the naive idea does not enjoy valid coverage



- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$
- **LOO scores** $\mathcal{S} = \left\{ |\hat{A}_{-i}(X_i) - Y_i| \right\}_i \cup \{+\infty\}$ (in standard mean regression)
- Get \hat{A} by training \mathcal{A} on \mathcal{D}_n
- Build the predictive interval: $\left[\hat{A}(X_{n+1}) \pm q_{1-\alpha}(\mathcal{S}) \right]$

Warning

No guarantee on the prediction of \hat{A} with scores based on $(\hat{A}_{-i})_i$, without assuming a form of **stability** on \mathcal{A} .

- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$





- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$
- **LOO predictions / predictive intervals**

$$\mathcal{S}_{\text{up/down}} = \left\{ \hat{A}_{-i}(X_{n+1}) \pm |\hat{A}_{-i}(X_i) - Y_i| \right\}_i \cup \{\pm\infty\}$$

(in standard mean regression)



- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$

- **LOO predictions / predictive intervals**

$$\mathcal{S}_{\text{up/down}} = \left\{ \hat{A}_{-i}(X_{n+1}) \pm |\hat{A}_{-i}(X_i) - Y_i| \right\}_i \cup \{\pm\infty\}$$

(in standard mean regression)

- Build the predictive interval: $[q_{\alpha, \text{inf}}(\mathcal{S}_{\text{down}}); q_{1-\alpha}(\mathcal{S}_{\text{up}})]$



- Based on **leave-one-out (LOO) residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Get \hat{A}_{-i} by training \mathcal{A} on $\mathcal{D}_n \setminus (X_i, Y_i)$

- **LOO predictions / predictive intervals**

$$\mathcal{S}_{\text{up/down}} = \left\{ \hat{A}_{-i}(X_{n+1}) \pm |\hat{A}_{-i}(X_i) - Y_i| \right\}_i \cup \{\pm\infty\}$$

(in standard mean regression)

- Build the predictive interval: $[q_{\alpha, \text{inf}}(\mathcal{S}_{\text{down}}); q_{1-\alpha}(\mathcal{S}_{\text{up}})]$

Theorem (Marginal validity of Jackknife+ Barber et al. (2021b)).

If $\mathcal{D}_n \cup (X_{n+1}, Y_{n+1})$ are exchangeable and \mathcal{A} is symmetric:
 $\mathbb{P}(Y_{n+1} \in \hat{C}_\alpha(X_{n+1})) \geq 1 - 2\alpha.$

Recall $q_{\beta, \text{inf}}(X_1, \dots, X_n) := \lfloor \beta \times n \rfloor$ smallest value of (X_1, \dots, X_n)

- Based on cross-validation residuals
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Split \mathcal{D}_n into K folds F_1, \dots, F_K
- Get \hat{A}_{-F_k} by training \mathcal{A} on $\mathcal{D}_n \setminus F_k$





- Based on **cross-validation residuals**
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Split \mathcal{D}_n into K folds F_1, \dots, F_K
- Get \hat{A}_{-F_k} by training \mathcal{A} on $\mathcal{D}_n \setminus F_k$
- **Cross-val predictions / predictive intervals**

$$\mathcal{S}_{\text{up/down}} = \left\{ \left\{ \hat{A}_{-F_k}(X_{n+1}) \pm |\hat{A}_{-F_k}(X_i) - Y_i| \right\}_{i \in F_k} \right\}_k \cup \{\pm\infty\}$$

(in standard mean regression)



- Based on cross-validation residuals
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Split \mathcal{D}_n into K folds F_1, \dots, F_K
- Get \hat{A}_{-F_k} by training \mathcal{A} on $\mathcal{D}_n \setminus F_k$
- Cross-val predictions / predictive intervals

$$\mathcal{S}_{\text{up/down}} = \left\{ \left\{ \hat{A}_{-F_k}(X_{n+1}) \pm |\hat{A}_{-F_k}(X_i) - Y_i| \right\}_{i \in F_k} \right\}_k \cup \{\pm\infty\}$$

(in standard mean regression)

- Build the predictive interval: $[q_{\alpha, \text{inf}}(\mathcal{S}_{\text{down}}); q_{1-\alpha}(\mathcal{S}_{\text{up}})]$



- Based on cross-validation residuals
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ training data
- Split \mathcal{D}_n into K folds F_1, \dots, F_K
- Get \hat{A}_{-F_k} by training \mathcal{A} on $\mathcal{D}_n \setminus F_k$
- Cross-val predictions / predictive intervals

$$\mathcal{S}_{\text{up/down}} = \left\{ \left\{ \hat{A}_{-F_k}(X_{n+1}) \pm |\hat{A}_{-F_k}(X_i) - Y_i| \right\}_{i \in F_k} \right\}_k \cup \{\pm\infty\}$$

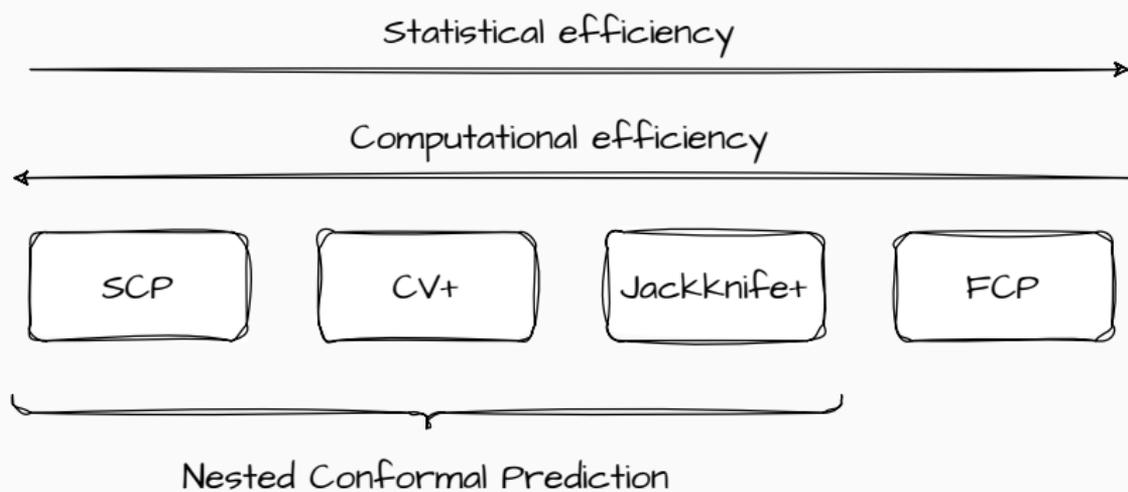
(in standard mean regression)

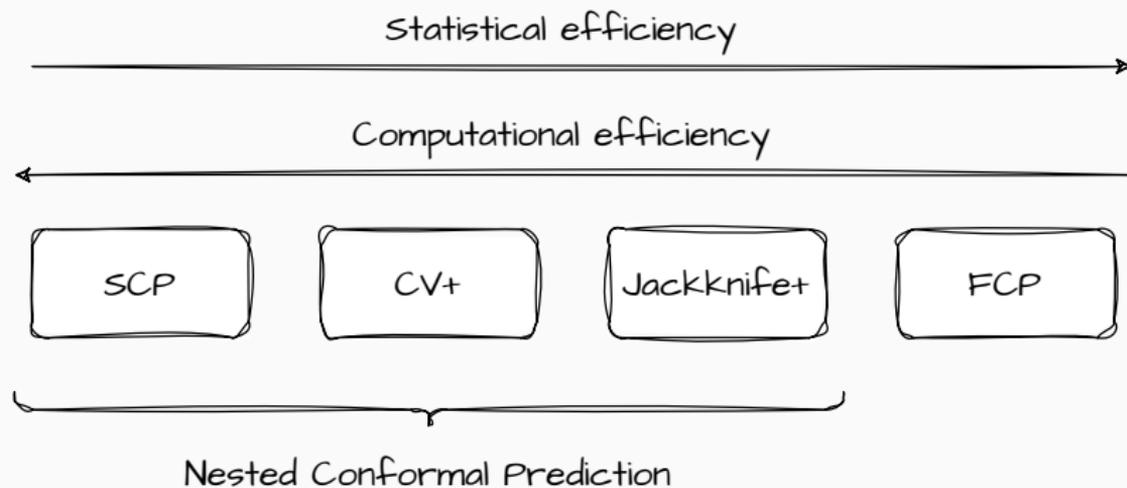
- Build the predictive interval: $[q_{\alpha, \text{inf}}(\mathcal{S}_{\text{down}}); q_{1-\alpha}(\mathcal{S}_{\text{up}})]$

Theorem (Marginal validity of CV+ Barber et al. (2021b)).

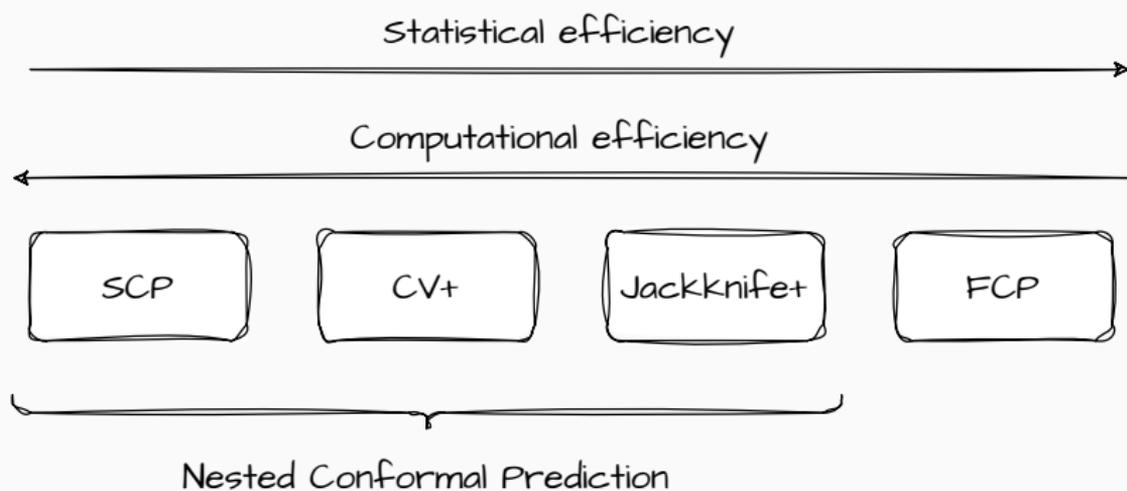
If $\mathcal{D}_n \cup (X_{n+1}, Y_{n+1})$ are exchangeable and \mathcal{A} is symmetric: $\mathbb{P}(Y_{n+1} \in \hat{C}_\alpha(X_{n+1})) \geq 1 - 2\alpha - \min\left(\frac{2(1-1/K)}{n/K+1}, \frac{1-K/n}{K+1}\right) \geq 1 - 2\alpha - \sqrt{2/n}$.

General overview





- Generalized framework encapsulating out-of-sample methods: Nested CP (Gupta et al., 2022) → extends $JK+/CV+$ for any score.



- Generalized framework encapsulating out-of-sample methods: Nested CP (Gupta et al., 2022) \rightarrow extends $JK+/CV+$ for any score.
- Accelerating FCP: Nourtdinov et al. (2001); Lei (2019); Ndiaye and Takeuchi (2019); Cherubin et al. (2021); Ndiaye and Takeuchi (2022); Ndiaye (2022)

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

Exchangeability does not hold in many practical applications

- CP requires **exchangeable** data points to ensure validity

Exchangeability does not hold in many practical applications

- CP requires **exchangeable** data points to ensure validity
- ✗ Covariate shift, i.e. \mathcal{L}_X changes but $\mathcal{L}_{Y|X}$ stays constant

Exchangeability does not hold in many practical applications

- CP requires **exchangeable** data points to ensure validity
- ✗ Covariate shift, i.e. \mathcal{L}_X changes but $\mathcal{L}_{Y|X}$ stays constant
- ✗ Label shift, i.e. \mathcal{L}_Y changes but $\mathcal{L}_{X|Y}$ stays constant

Exchangeability does not hold in many practical applications

- CP requires **exchangeable** data points to ensure validity
- ✗ Covariate shift, i.e. \mathcal{L}_X changes but $\mathcal{L}_{Y|X}$ stays constant
- ✗ Label shift, i.e. \mathcal{L}_Y changes but $\mathcal{L}_{X|Y}$ stays constant
- ✗ Arbitrary distribution shift

Exchangeability does not hold in many practical applications

- CP requires **exchangeable** data points to ensure validity
- ✗ Covariate shift, i.e. \mathcal{L}_X changes but $\mathcal{L}_{Y|X}$ stays constant
- ✗ Label shift, i.e. \mathcal{L}_Y changes but $\mathcal{L}_{X|Y}$ stays constant
- ✗ Arbitrary distribution shift
- ✗ Possibly many shifts, not only one

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_X \times P_{Y|X}$
 - $(X_{n+1}, Y_{n+1}) \sim \tilde{P}_X \times P_{Y|X}$

¹²Tibshirani et al. (2019), *Conformal Prediction Under Covariate Shift*, NeurIPS

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_X \times P_{Y|X}$
 - $(X_{n+1}, Y_{n+1}) \sim \tilde{P}_X \times P_{Y|X}$
- **Idea:** give more importance to calibration points that are closer in distribution to the test point

¹²Tibshirani et al. (2019), *Conformal Prediction Under Covariate Shift*, NeurIPS

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_X \times P_{Y|X}$
 - $(X_{n+1}, Y_{n+1}) \sim \tilde{P}_X \times P_{Y|X}$
- **Idea:** give more importance to calibration points that are closer in distribution to the test point
- **In practice:**
 1. estimate the **likelihood ratio** $w(X_i) = \frac{d\tilde{P}_X(X_i)}{dP_X(X_i)}$

¹²Tibshirani et al. (2019), *Conformal Prediction Under Covariate Shift*, NeurIPS

- **Setting:**

- $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_X \times P_{Y|X}$
- $(X_{n+1}, Y_{n+1}) \sim \tilde{P}_X \times P_{Y|X}$

- **Idea:** give more importance to calibration points that are closer in distribution to the test point

- **In practice:**

1. estimate the **likelihood ratio** $w(X_i) = \frac{d\tilde{P}_X(X_i)}{dP_X(X_i)}$
2. normalize the weights, i.e. $\omega_i = \omega(X_i) = \frac{w(X_i)}{\sum_{j=1}^{n+1} w(X_j)}$

¹²Tibshirani et al. (2019), *Conformal Prediction Under Covariate Shift*, NeurIPS

- **Setting:**

- $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_X \times P_{Y|X}$
- $(X_{n+1}, Y_{n+1}) \sim \tilde{P}_X \times P_{Y|X}$

- **Idea:** give more importance to calibration points that are closer in distribution to the test point

- **In practice:**

1. estimate the **likelihood ratio** $w(X_i) = \frac{d\tilde{P}_X(X_i)}{dP_X(X_i)}$
2. normalize the weights, i.e. $\omega_i = \omega(X_i) = \frac{w(X_i)}{\sum_{j=1}^{n+1} w(X_j)}$

3. outputs $\hat{C}_\alpha(X_{n+1}) =$

$$\left\{ y : \mathbf{s}(\hat{A}(X_{n+1}), y) \leq Q_{1-\alpha} \left(\sum_{i \in \text{Cal}} \omega_i \delta_{S_i} + \omega_{n+1} \delta_\infty \right) \right\}$$

¹²Tibshirani et al. (2019), *Conformal Prediction Under Covariate Shift*, NeurIPS

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_{X|Y} \times P_Y$
 - $(X_{n+1}, Y_{n+1}) \sim P_{X|Y} \times \tilde{P}_Y$
 - **Classification**

¹³Podkopaev and Ramdas (2021), *Distribution-free uncertainty quantification for classification under label shift*, 60 / 78

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_{X|Y} \times P_Y$
 - $(X_{n+1}, Y_{n+1}) \sim P_{X|Y} \times \tilde{P}_Y$
 - **Classification**
- **Idea:** give more importance to calibration points that are closer in distribution to the test point

¹³Podkopaev and Ramdas (2021), *Distribution-free uncertainty quantification for classification under label shift* / 78

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_{X|Y} \times P_Y$
 - $(X_{n+1}, Y_{n+1}) \sim P_{X|Y} \times \tilde{P}_Y$
 - **Classification**
- **Idea:** give more importance to calibration points that are closer in distribution to the test point
- **Trouble:** the actual test labels are **unknown**

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_{X|Y} \times P_Y$
 - $(X_{n+1}, Y_{n+1}) \sim P_{X|Y} \times \tilde{P}_Y$
 - **Classification**
- **Idea:** give more importance to calibration points that are closer in distribution to the test point
- **Trouble:** the actual test labels are **unknown**
- **In practice:**
 1. estimate the **likelihood ratio** $w(Y_i) = \frac{d\tilde{P}_Y(Y_i)}{dP_Y(Y_i)}$ using algorithms from the existing label shift literature

¹³Podkopaev and Ramdas (2021), *Distribution-free uncertainty quantification for classification under label shift* 60 / 78

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_{X|Y} \times P_Y$
 - $(X_{n+1}, Y_{n+1}) \sim P_{X|Y} \times \tilde{P}_Y$
 - **Classification**
- **Idea:** give more importance to calibration points that are closer in distribution to the test point
- **Trouble:** the actual test labels are **unknown**
- **In practice:**
 1. estimate the **likelihood ratio** $w(Y_i) = \frac{d\tilde{P}_Y(Y_i)}{dP_Y(Y_i)}$ using algorithms from the existing label shift literature
 2. normalize the weights, i.e. $\omega_i^y = \omega^y(X_i) = \frac{w(Y_i)}{\sum_{j=1}^n w(Y_j) + w(y)}$

- **Setting:**
 - $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} P_{X|Y} \times P_Y$
 - $(X_{n+1}, Y_{n+1}) \sim P_{X|Y} \times \tilde{P}_Y$
 - **Classification**
- **Idea:** give more importance to calibration points that are closer in distribution to the test point
- **Trouble:** the actual test labels are **unknown**
- **In practice:**

1. estimate the **likelihood ratio** $w(Y_i) = \frac{d\tilde{P}_Y(Y_i)}{dP_Y(Y_i)}$ using algorithms from the existing label shift literature

2. normalize the weights, i.e. $\omega_i^y = \omega^y(X_i) = \frac{w(Y_i)}{\sum_{j=1}^n w(Y_j) + w(y)}$

3. outputs $\hat{C}_\alpha(X_{n+1}) =$

$$\left\{ y : \mathfrak{s}(\hat{A}(X_{n+1}), y) \leq Q_{1-\alpha} \left(\sum_{i \in \text{Cal}} \omega_i^y \delta_{S_i} + \omega_{n+1}^y \delta_\infty \right) \right\}$$

- Arbitrary distribution shift: Cauchois et al. (2020) leverages ideas from the distributionally robust optimization literature
- Two major **general theoretical results** beyond exchangeability:

- Arbitrary distribution shift: Cauchois et al. (2020) leverages ideas from the distributionally robust optimization literature
- Two major **general theoretical results** beyond exchangeability:
 - Chernozhukov et al. (2018)
 - ↔ If the learnt model is accurate and the data noise is strongly mixing, then CP is valid asymptotically ✓

- Arbitrary distribution shift: Cauchois et al. (2020) leverages ideas from the distributionally robust optimization literature
- Two major **general theoretical results** beyond exchangeability:
 - Chernozhukov et al. (2018)
 - ↪ If the learnt model is accurate and the data noise is strongly mixing, then CP is valid asymptotically ✓
 - Barber et al. (2022)
 - ↪ Quantifies the coverage loss depending on the strength of exchangeability violation
 - $$\mathbb{P}(Y_{n+1} \in \hat{C}_\alpha(X_{n+1})) \geq 1 - \alpha - \frac{\text{average violation of exchangeability}}{\text{by each calibration point}}$$
 - ↪ proposed algorithm: **reweighting** again!
 - e.g., in a temporal setting, give higher weights to more recent points.

- **Data:** T_0 random variables $(X_1, Y_1), \dots, (X_{T_0}, Y_{T_0})$ in $\mathbb{R}^d \times \mathbb{R}$
- **Aim:** predict the response values as well as predictive intervals for T_1 subsequent observations $X_{T_0+1}, \dots, X_{T_0+T_1}$ sequentially: at any prediction step $t \in \llbracket T_0 + 1, T_0 + T_1 \rrbracket$, $Y_{t-T_0}, \dots, Y_{t-1}$ have been revealed
- Build the smallest interval \widehat{C}_α^t such that:

$$\mathbb{P} \left\{ Y_t \in \widehat{C}_\alpha^t(X_t) \right\} \geq 1 - \alpha, \text{ for } t \in \llbracket T_0 + 1, T_0 + T_1 \rrbracket,$$

often relaxed in:

$$\frac{1}{T_1} \sum_{t=T_0+1}^{T_0+T_1} \mathbb{1} \left\{ Y_t \in \widehat{C}_\alpha^t(X_t) \right\} \approx 1 - \alpha.$$

- Consider splitting strategies that respect the temporal structure

Recent developments

- Consider splitting strategies that respect the temporal structure
- Gibbs and Candès (2021) propose a method which reacts faster to temporal evolution
 - **Idea**: track the previous coverages of the predictive intervals ($\mathbb{1}\{Y_t \in \hat{C}_\alpha(X_t)\}$)
 - **Tool**: update the empirical quantile level with a learning rate γ
 - Asymptotic guarantee (on average) for **any distribution** (even adversarial)

Recent developments

- Consider splitting strategies that respect the temporal structure
- Gibbs and Candès (2021) propose a method which reacts faster to temporal evolution
 - **Idea**: track the previous coverages of the predictive intervals ($\mathbb{1}\{Y_t \in \widehat{C}_\alpha(X_t)\}$)
 - **Tool**: update the empirical quantile level with a learning rate γ
 - Asymptotic guarantee (on average) for **any distribution** (even adversarial)
- Zaffran et al. (2022) studies the influence of this learning rate γ and proposes, along with Gibbs and Candès (2022), a method not requiring to choose γ

Recent developments

- Consider splitting strategies that respect the temporal structure
- Gibbs and Candès (2021) propose a method which reacts faster to temporal evolution
 - **Idea**: track the previous coverages of the predictive intervals ($\mathbb{1}\{Y_t \in \widehat{C}_\alpha(X_t)\}$)
 - **Tool**: update the empirical quantile level with a learning rate γ
 - Asymptotic guarantee (on average) for **any distribution** (even adversarial)
- Zaffran et al. (2022) studies the influence of this learning rate γ and proposes, along with Gibbs and Candès (2022), a method not requiring to choose γ
- Bhatnagar et al. (2023) enjoys **anytime** regret bound, by leveraging tools from the strongly adaptive regret minimization literature

Recent developments

- Consider splitting strategies that respect the temporal structure
- Gibbs and Candès (2021) propose a method which reacts faster to temporal evolution
 - **Idea**: track the previous coverages of the predictive intervals ($\mathbb{1}\{Y_t \in \widehat{C}_\alpha(X_t)\}$)
 - **Tool**: update the empirical quantile level with a learning rate γ
 - Asymptotic guarantee (on average) for **any distribution** (even adversarial)
- Zaffran et al. (2022) studies the influence of this learning rate γ and proposes, along with Gibbs and Candès (2022), a method not requiring to choose γ
- Bhatnagar et al. (2023) enjoys anytime regret bound, by leveraging tools from the strongly adaptive regret minimization literature
- Bastani et al. (2022) proposes an algorithm achieving stronger coverage guarantees (conditional on specified overlapping subsets, and threshold calibrated) without hold-out set

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Healthcare

Electricity

Concluding remarks

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Healthcare

Electricity

Concluding remarks

- Medical application
- Image based task
- Pixel by pixel analysis \rightsquigarrow
applications to segmentation
for self-driving cars

Image-to-Image Regression with Distribution-Free Uncertainty Quantification and Applications in Imaging

Anastasios N. Angelopoulos^{*1} Amit Kohli^{*1} Stephen Bates¹ Michael I. Jordan¹ Jitendra Malik¹
Thayer Alshaabi² Srigokul Upadhyayula^{2,3} Yaniv Romano⁴

- Medical application
- Image based task
- Pixel by pixel analysis \rightsquigarrow
applications to segmentation
for self-driving cars

1. **Task:** *Image to Image regression* - for each pixel of an image, predict a real valued output from the entire image.
2. **UQ Goal:** provide a predictive interval for each pixel, such that the output is in the interval at least 90% of the time.

Image-to-Image Regression with Distribution-Free Uncertainty Quantification and Applications in Imaging

Anastasios N. Angelopoulos^{*1} Amit Kohli^{*1} Stephen Bates¹ Michael I. Jordan¹ Jitendra Malik¹
Thayer Alshaabi² Srigokul Upadhyayula^{2,3} Yaniv Romano⁴

Image to Image regression with DF-UQ - Angelopoulos et al. (2022b)

- Medical application
- Image based task
- Pixel by pixel analysis \rightsquigarrow applications to segmentation for self-driving cars

1. **Task:** *Image to Image regression* - for each pixel of an image, predict a real valued output from the entire image.
2. **UQ Goal:** provide a predictive interval for each pixel, such that the output is in the interval at least 90% of the time.

Image-to-Image Regression with Distribution-Free Uncertainty Quantification and Applications in Imaging

Anastasios N. Angelopoulos^{*1} Amit Kohli^{*1} Stephen Bates¹ Michael I. Jordan¹ Jitendra Malik¹
Thayer Alshaabi² Srigokul Upadhyayula^{2,3} Yaniv Romano⁴

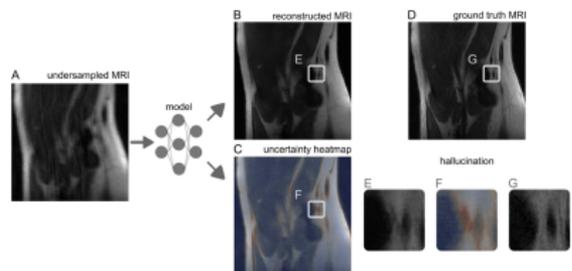


Figure 1. An algorithmic MRI reconstruction with uncertainty. A rapidly acquired but undersampled MR image of a knee (A) is fed into a model that predicts a sharp reconstruction (B) with calibrated uncertainty (C). In (C), red means high uncertainty and blue means low uncertainty. Wherever the reconstruction contains hallucinations, the uncertainty is high; see the hallucination in the image patch (E), which has high uncertainty in (F), and does not exist in the ground truth (G). For experimental details, see Section 3.4.

Figure 2: Image from Angelopoulos et al. (2022b)

Method:

1. Split conformal prediction method - isolate calibration set
2. On the **proper training set**, learn:
 - Mean regressor - $\hat{\mu} : \mathbb{R}^{NM} \rightarrow [0; 1]$

Method:

1. Split conformal prediction method - isolate calibration set
2. On the **proper training set**, learn:
 - Mean regressor - $\hat{\mu} : \mathbb{R}^{NM} \rightarrow [0; 1]$
 - Heuristic notion of uncertainty: $\tilde{u}, \tilde{\ell} : \mathbb{R}^{NM} \rightarrow [0; 1]$, such that

$$[\hat{\mu}(X) - \tilde{\ell}(X); \hat{\mu}(X) + \tilde{u}(X)]$$

→ 3 regressors are used

4 techniques are experimented for these regressors, including QR.

Method:

1. Split conformal prediction method - isolate calibration set
2. On the **proper training set**, learn:
 - Mean regressor - $\hat{\mu} : \mathbb{R}^{NM} \rightarrow [0; 1]$
 - Heuristic notion of uncertainty: $\tilde{u}, \tilde{\ell} : \mathbb{R}^{NM} \rightarrow [0; 1]$, such that

$$[\hat{\mu}(X) - \tilde{\ell}(X); \hat{\mu}(X) + \tilde{u}(X)]$$

→ 3 regressors are used

4 techniques are experimented for these regressors, including QR.

3. Calibration step: leverage the **calibration set**.
 - In spirit, almost equivalent to CQR but with a multiplicative form.
 - Precisely, relies on RCPS (Bates et al., 2021a)

Method:

1. Split conformal prediction method - isolate calibration set
2. On the **proper training set**, learn:
 - Mean regressor - $\hat{\mu} : \mathbb{R}^{NM} \rightarrow [0; 1]$
 - Heuristic notion of uncertainty: $\tilde{u}, \tilde{\ell} : \mathbb{R}^{NM} \rightarrow [0; 1]$, such that

$$[\hat{\mu}(X) - \tilde{\ell}(X); \hat{\mu}(X) + \tilde{u}(X)]$$

→ 3 regressors are used

4 techniques are experimented for these regressors, including QR.

3. Calibration step: leverage the **calibration set**.
 - In spirit, almost equivalent to CQR but with a multiplicative form.
 - Precisely, relies on RCPS (Bates et al., 2021a)

Guarantee:

$$\mathbb{P} [\mathbb{E} [\text{Average miscoverage on all pixels of a test image} \geq \alpha | \text{Cal}]] \leq \delta$$

→ Marginal validity on the **test**, with high probability w.r.t. the **calibration set**.

Abstract

Image-to-image regression is an important learning task, used frequently in biological imaging. Current algorithms, however, do not generally offer statistical guarantees that protect against a model's mistakes and hallucinations. To address this, we develop uncertainty quantification techniques with rigorous statistical guarantees for image-to-image regression problems. In particular, we show how to derive uncertainty intervals around each pixel that are guaranteed to contain the true value with a user-specified confidence probability. Our methods work in conjunction

2. Methods

We now formally describe the method for constructing uncertainty intervals. Each pixel in the image will get its own uncertainty interval, as in (1), that is statistically guaranteed to contain the true value with high probability.

Abstract

Image-to-image regression is an important learning task, used frequently in biological imaging. Current algorithms, however, do not generally offer statistical guarantees that protect against a model's mistakes and hallucinations. To address this, we develop uncertainty quantification techniques with rigorous statistical guarantees for image-to-image regression problems. In particular, we show how to derive uncertainty intervals around each pixel that are guaranteed to contain the true value with a user-specified confidence probability. Our methods work in conjunction

2. Methods

We now formally describe the method for constructing uncertainty intervals. Each pixel in the image will get its own uncertainty interval, as in (1), that is statistically guaranteed to contain the true value with high probability.

- Not a conditional coverage claim!
- The statement is on-average on the test point - easy or hard.

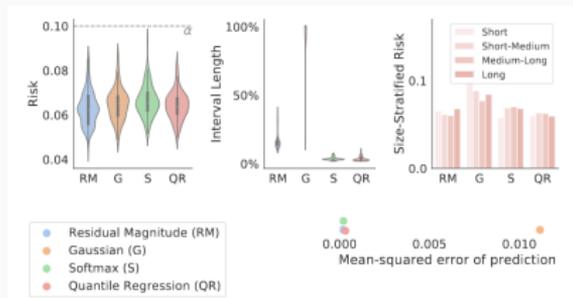
Size-stratified risk. Next, we seek prediction sets that do not systematically make mistakes in difficult parts of the image. Our risk control requirement in Definition 2.1 may be satisfied even if the prediction sets systematically fail to contain the most difficult pixels. For example, if $\alpha = 0.1$ and 90% of pixels are covered by fixed-width intervals of size 0.01, then the requirement is satisfied—however, the sets no longer serve as useful notions of uncertainty. To

- Hard problem (impossibility results!)
- Introduce metrics to see *if* and *on which underlying regressors* such problem happens.

Image to Image regression with DF-UQ - Angelopoulos et al. (2022b)

Example of such metrics (see also Feldman et al., 2021) :

- Link between the size of the PI and the coverage level \rightarrow



Example of such metrics (see also

Feldman et al., 2021) :

- Link between the size of the PI and the coverage level \rightarrow
- Localization of the errors \downarrow



Figure 3. Examples of quantitative phase reconstructions of leukocytes with uncertainty shown in the following order: input (we only show one of the two illuminations), prediction, uncertainty visualization (produced with quantile regression), absolute difference between prediction and ground truth (renormalized for visualization), ground truth.

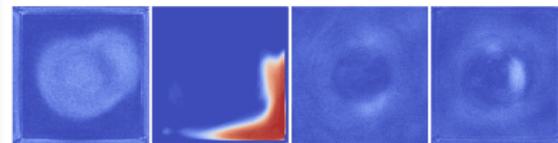
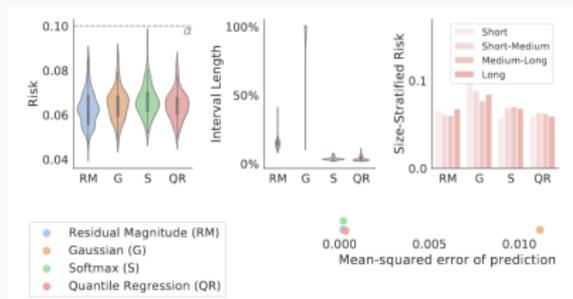


Figure 8. Spatial variations in microcoverage in the BSCCM dataset are shown for each of the four methods as a heatmap. Blue represents 0% microcoverage and red represents 100%. The methods are, in order, residual magnitude, gaussian, softmax, and quantile regression.

Figure 3: All images from Angelopoulos et al. (2022b)

Image to Image regression with DF-UQ - Angelopoulos et al. (2022b)

Example of such metrics (see also

Feldman et al., 2021) :

- Link between the size of the PI and the coverage level \rightarrow
- Localization of the errors \downarrow



Figure 3. Examples of quantitative phase reconstructions of leukocytes with uncertainty shown in the following order: input (we only show one of the two illuminations), prediction, uncertainty visualization (produced with quantile regression), absolute difference between prediction and ground truth (renormalized for visualization), ground truth.

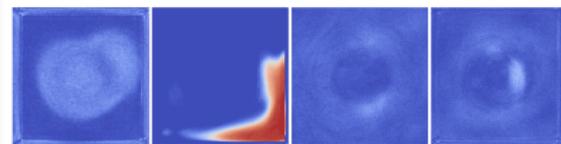
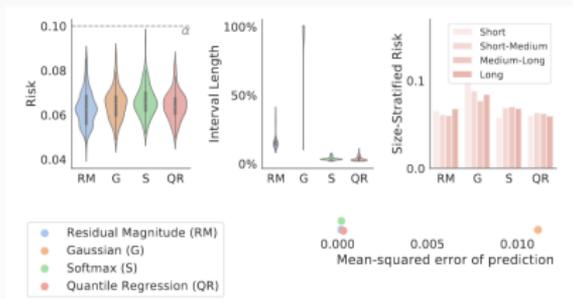


Figure 8. Spatial variations in miscoverage in the BSCCM dataset are shown for each of the four methods as a heatmap. Blue represents 0% miscoverage and red represents 100%. The methods are, in order, residual magnitude, gaussian, softmax, and quantile regression.

Figure 3: All images from Angelopoulos et al. (2022b)

Take aways:

- Elegant application of SCP with CQR type score
- **Test marginal** and **calibration** + **train** conditional validity guarantees with HP
- Main problem is Test conditionality \rightarrow look at metrics to evaluate which methods performs best!

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Healthcare

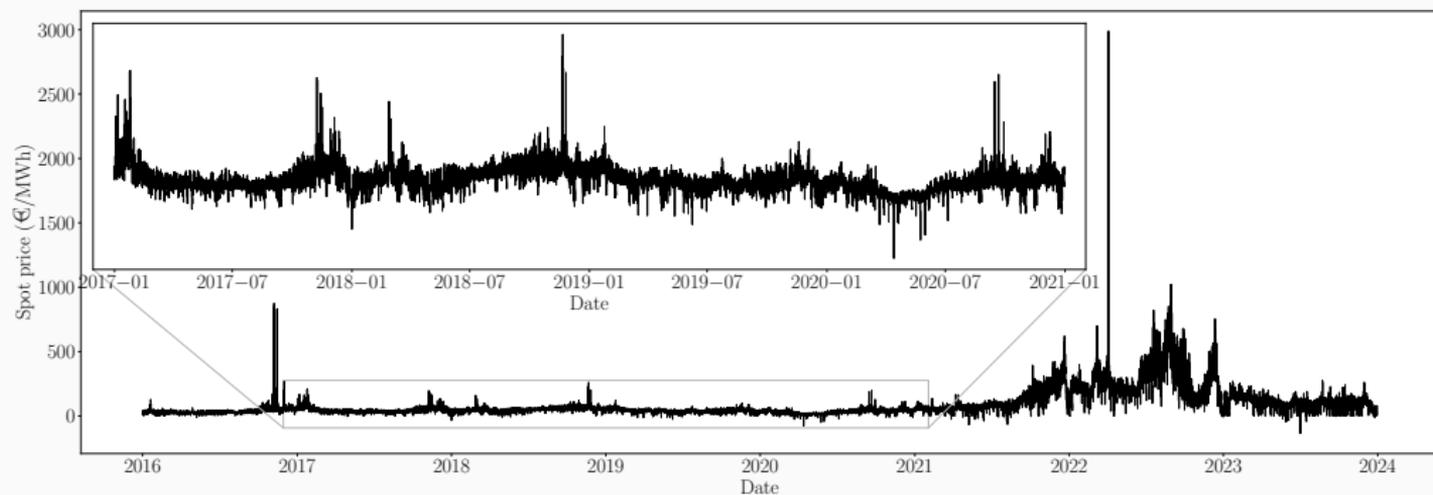
Electricity

Concluding remarks

Hourly day-ahead market prices (between producers and suppliers)

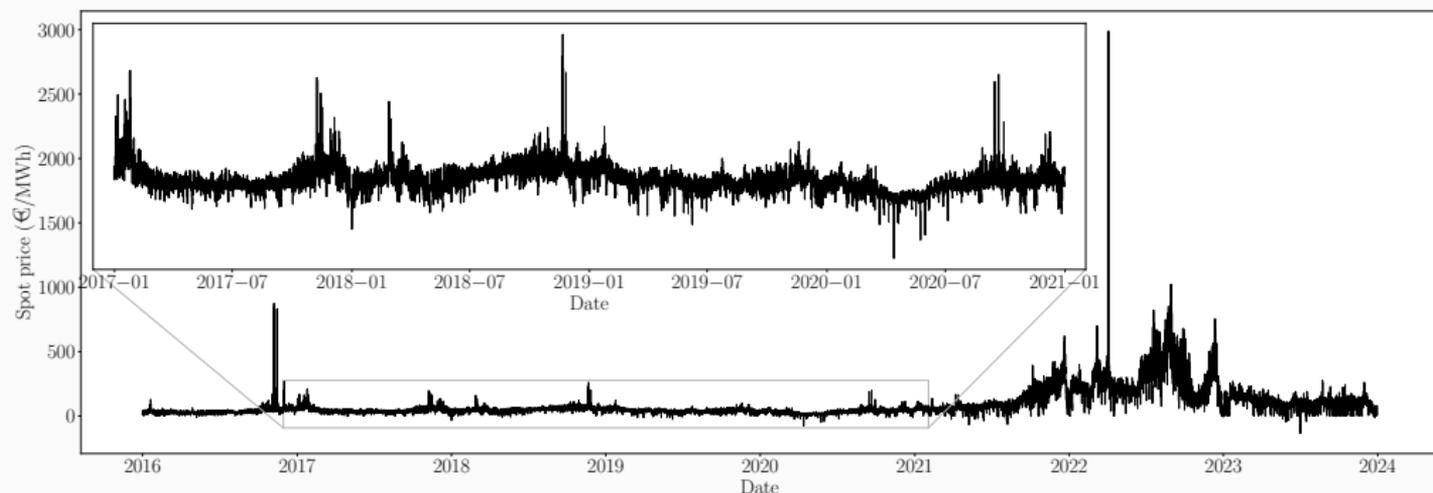
Forecasting French spot electricity prices

Hourly day-ahead market prices (between producers and suppliers)



Forecasting French spot electricity prices

Hourly day-ahead market prices (between producers and suppliers)

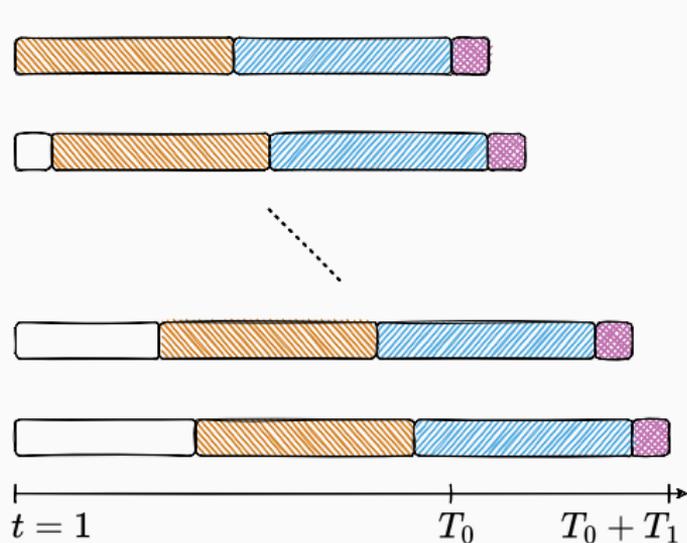


To which extent are they forecastable?

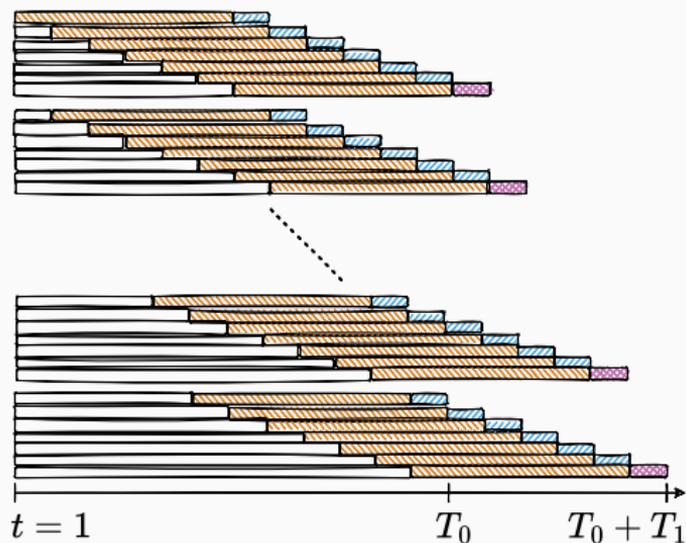
↪ forecasts errors **no lower than 10%** of the realized price!

Temporal splitting strategies: Online Sequential Split Conformal Prediction (OSSCP, Zaffran et al., 2022; Dutot et al., 2024)

□ Unused data ▨ Proper training set Tr_t ▨ Calibration set Cal_t ▨ Test point



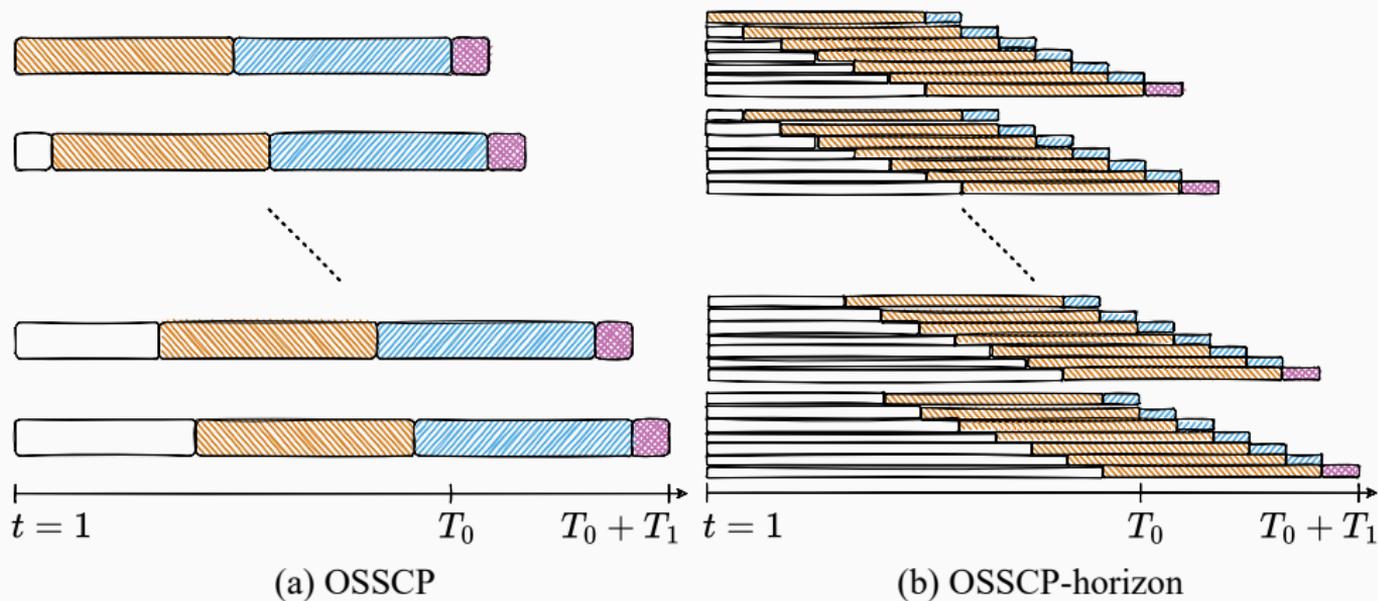
(a) OSSCP



(b) OSSCP-horizon

Temporal splitting strategies: Online Sequential Split Conformal Prediction (OSSCP, Zaffran et al., 2022; Dutot et al., 2024)

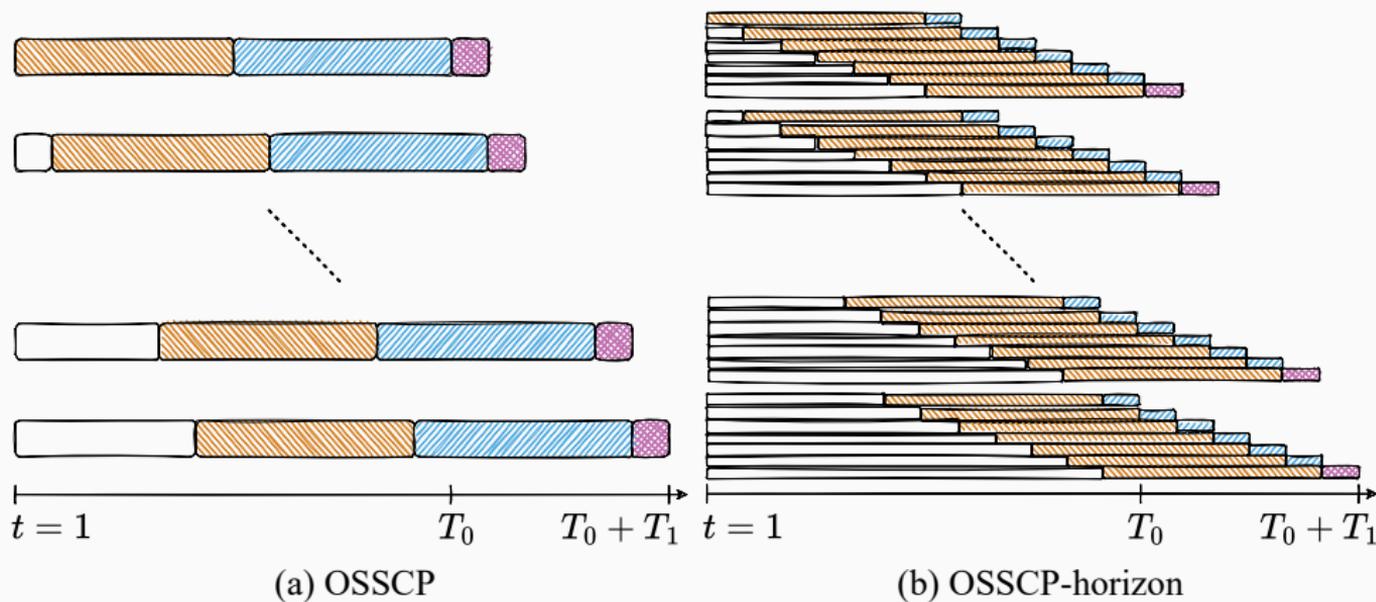
Unused data
 Proper training set Tr_t
 Calibration set Cal_t
 Test point



↪ OSSCP improves robustness in temporal settings;

Temporal splitting strategies: Online Sequential Split Conformal Prediction (OSSCP, Zaffran et al., 2022; Dutot et al., 2024)

□ Unused data ▨ Proper training set Tr_t ▨ Calibration set Cal_t ▨ Test point



↪ OSSCP improves robustness in temporal settings;

↪ OSSCP-horizon drastically improves robustness in non-stationary temporal settings.

Adaptive Conformal Inference (ACI) was initially proposed to handle distribution shift.

Adaptive Conformal Inference (ACI) was initially proposed to handle distribution shift.

It relies on updating online an *effective miscoverage rate* α_t , with the scheme

$$\alpha_{t+1} := \alpha_t + \gamma \left(\alpha - \mathbb{1} \left\{ Y^{(t)} \notin \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} \right),$$

and $\alpha_1 = \alpha$, $\gamma \geq 0$.

Adaptive Conformal Inference (ACI) was initially proposed to handle distribution shift.

It relies on updating online an *effective miscoverage rate* α_t , with the scheme

$$\alpha_{t+1} := \alpha_t + \gamma \left(\alpha - \mathbb{1} \left\{ Y^{(t)} \notin \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} \right),$$

and $\alpha_1 = \alpha$, $\gamma \geq 0$.

Intuition: if we did make an error, the interval was too small so we want to increase its length by taking a higher quantile (a smaller α_t). Reversely if we included the point.

Adaptive Conformal Inference (ACI) was initially proposed to handle distribution shift.

It relies on updating online an *effective miscoverage rate* α_t , with the scheme

$$\alpha_{t+1} := \alpha_t + \gamma \left(\alpha - \mathbb{1} \left\{ Y^{(t)} \notin \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} \right),$$

and $\alpha_1 = \alpha$, $\gamma \geq 0$.

Intuition: if we did make an **error**, the interval was **too small** so we want to **increase its length** by taking a **higher quantile** (a **smaller** α_t). Reversely if we included the point.

Adaptive Conformal Inference (ACI) was initially proposed to handle distribution shift.

It relies on updating online an *effective miscoverage rate* α_t , with the scheme

$$\alpha_{t+1} := \alpha_t + \gamma \left(\alpha - \mathbb{1} \left\{ Y^{(t)} \notin \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} \right),$$

and $\alpha_1 = \alpha$, $\gamma \geq 0$.

Intuition: if we did make an error, the interval was too small so we want to increase its length by taking a higher quantile (a smaller α_t). Reversely if we included the point.

Guarantee: *Asymptotic validity* result for *any sequence of observations*.

$$\frac{1}{T_1} \sum_{t=T_0+1}^{T_0+T_1} \mathbb{1} \left\{ Y^{(t)} \in \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} \xrightarrow{T_1 \rightarrow +\infty} 1 - \alpha$$

Adaptive Conformal Inference (ACI) was initially proposed to handle distribution shift.

It relies on updating online an *effective miscoverage rate* α_t , with the scheme

$$\alpha_{t+1} := \alpha_t + \gamma \left(\alpha - \mathbb{1} \left\{ Y^{(t)} \notin \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} \right),$$

and $\alpha_1 = \alpha$, $\gamma \geq 0$.

Intuition: if we did make an error, the interval was too small so we want to increase its length by taking a higher quantile (a smaller α_t). Reversely if we included the point.

Guarantee: *Asymptotic validity* result for *any sequence of observations*.

$$\left| \frac{1}{T_1} \sum_{t=T_0+1}^{T_0+T_1} \mathbb{1} \left\{ Y^{(t)} \in \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} - (1 - \alpha) \right| \leq \frac{2}{\gamma T_1}$$

Adaptive Conformal Inference (ACI) was initially proposed to handle distribution shift.

It relies on updating online an *effective miscoverage rate* α_t , with the scheme

$$\alpha_{t+1} := \alpha_t + \gamma \left(\alpha - \mathbb{1} \left\{ Y^{(t)} \notin \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} \right),$$

and $\alpha_1 = \alpha$, $\gamma \geq 0$.

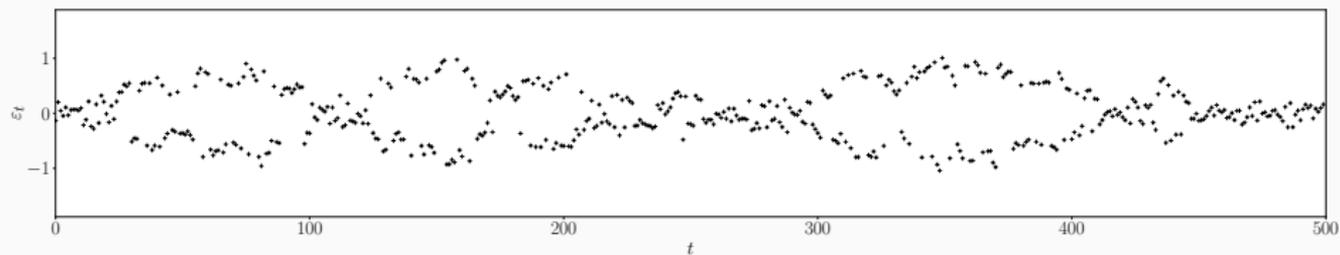
Intuition: if we did make an error, the interval was too small so we want to increase its length by taking a higher quantile (a smaller α_t). Reversely if we included the point.

Guarantee: *Asymptotic validity* result for *any sequence of observations*.

$$\left| \frac{1}{T_1} \sum_{t=T_0+1}^{T_0+T_1} \mathbb{1} \left\{ Y^{(t)} \in \widehat{C}_{\alpha_t} \left(X^{(t)} \right) \right\} - (1 - \alpha) \right| \leq \frac{2}{\gamma T_1}$$

\Rightarrow favors large γ .

Visualisation of ACI procedure



Visualisation of ACI procedure

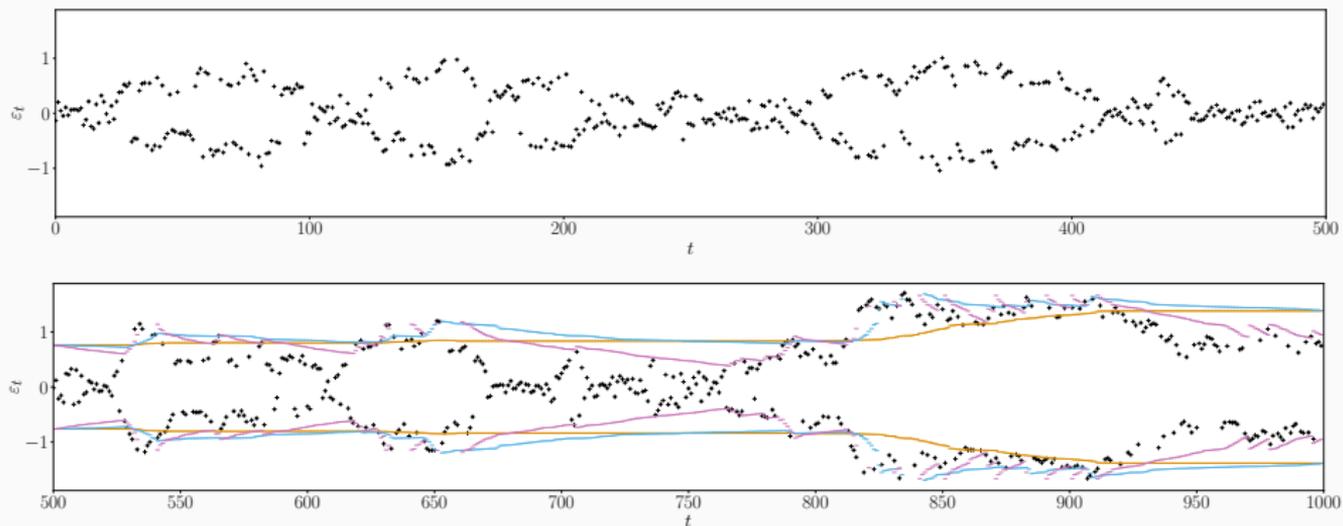
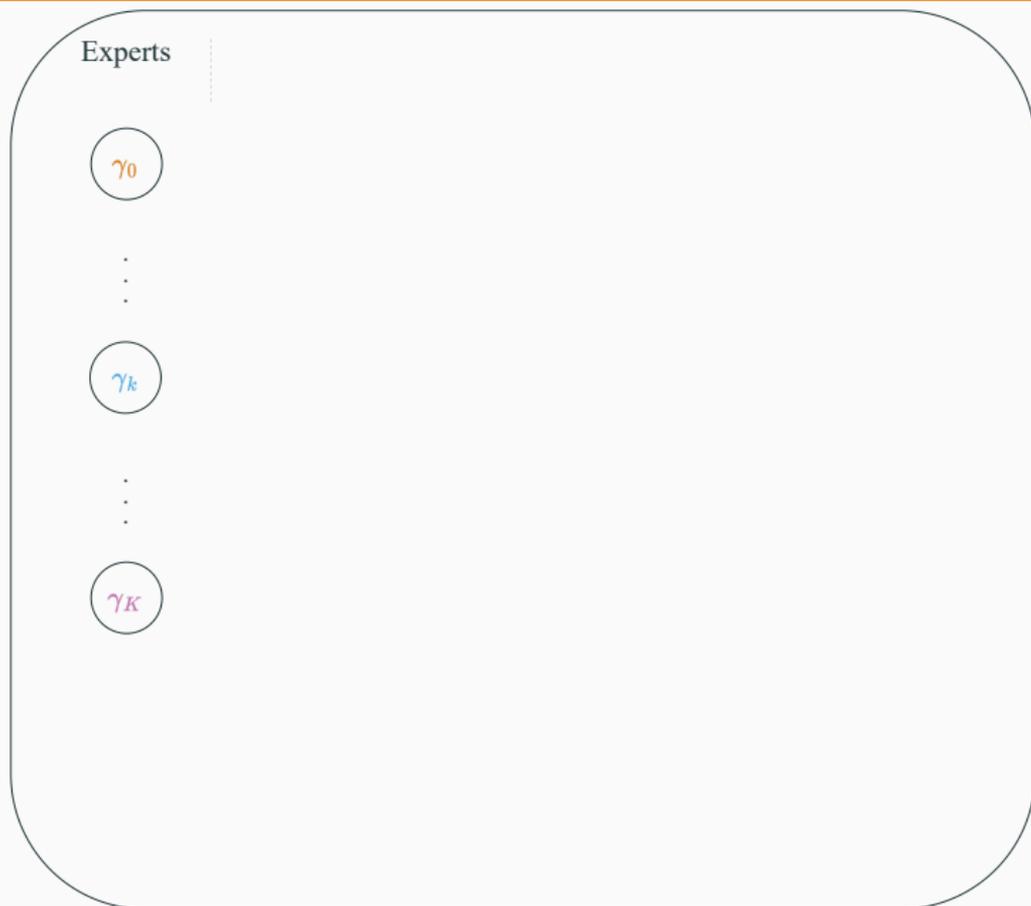
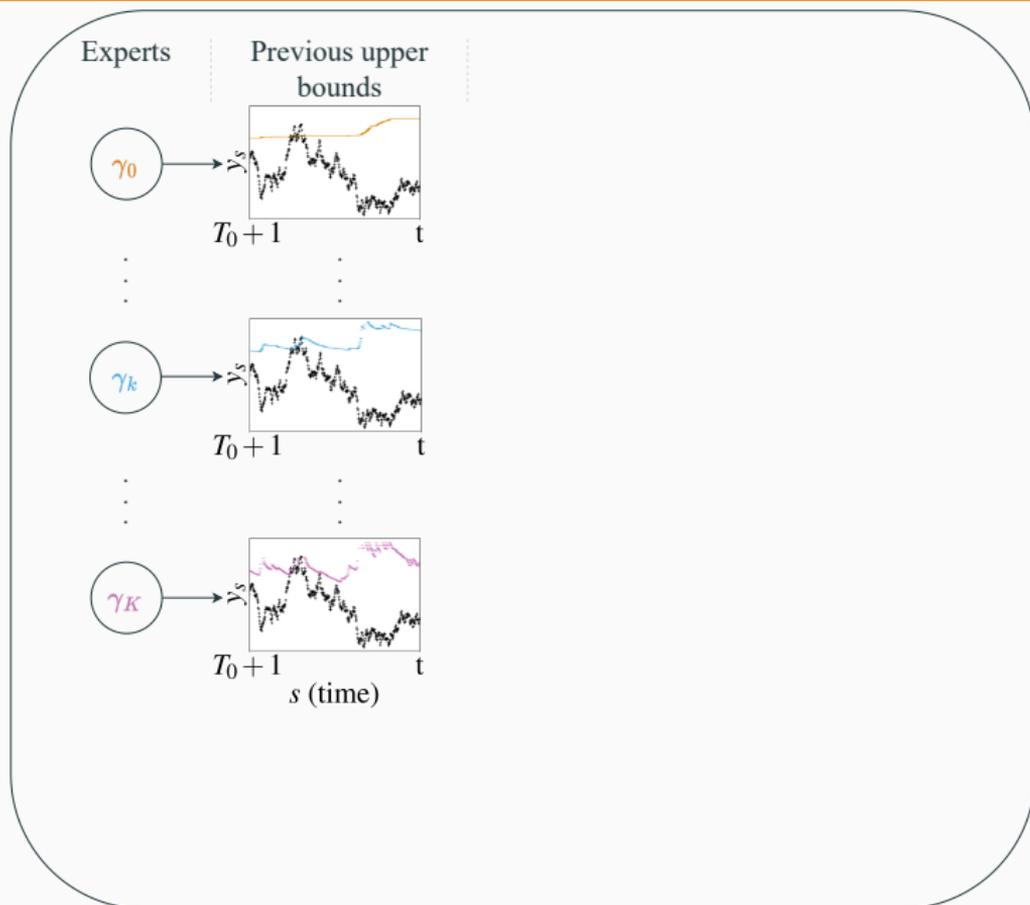


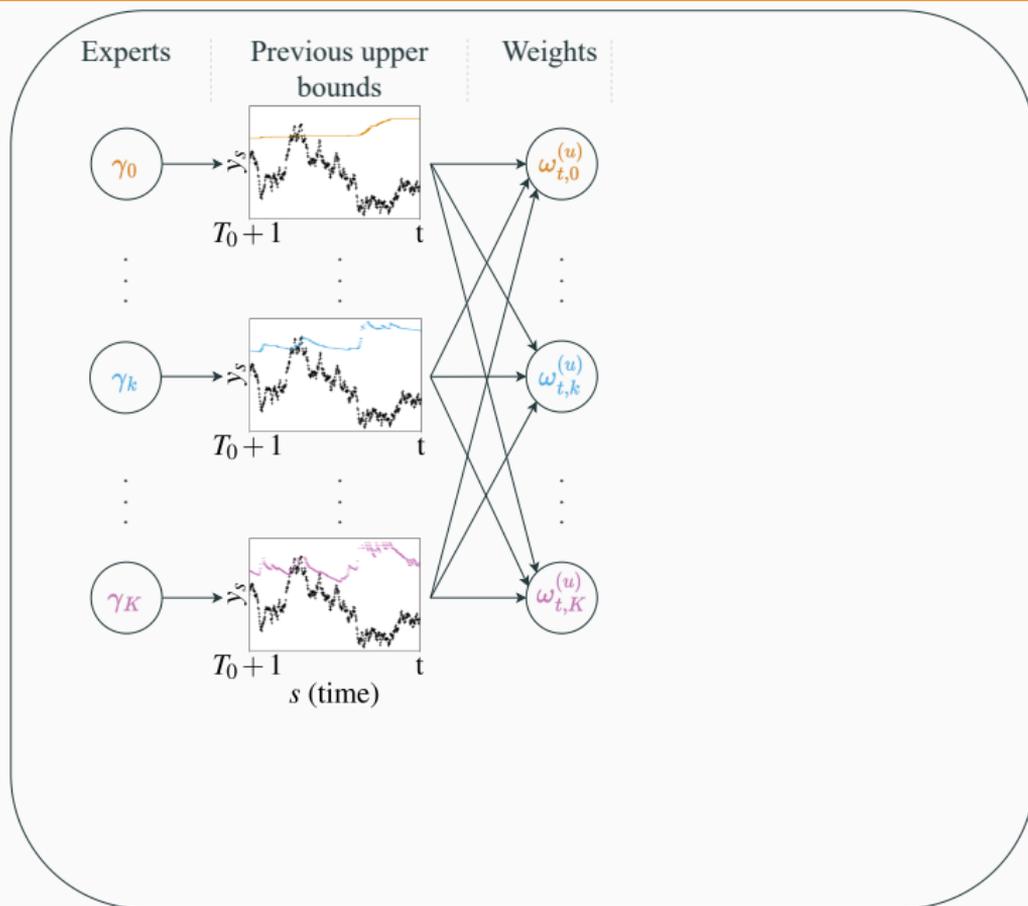
Figure 4: Visualisation of ACI with different values of γ ($\gamma = 0$, $\gamma = 0.01$, $\gamma = 0.05$)



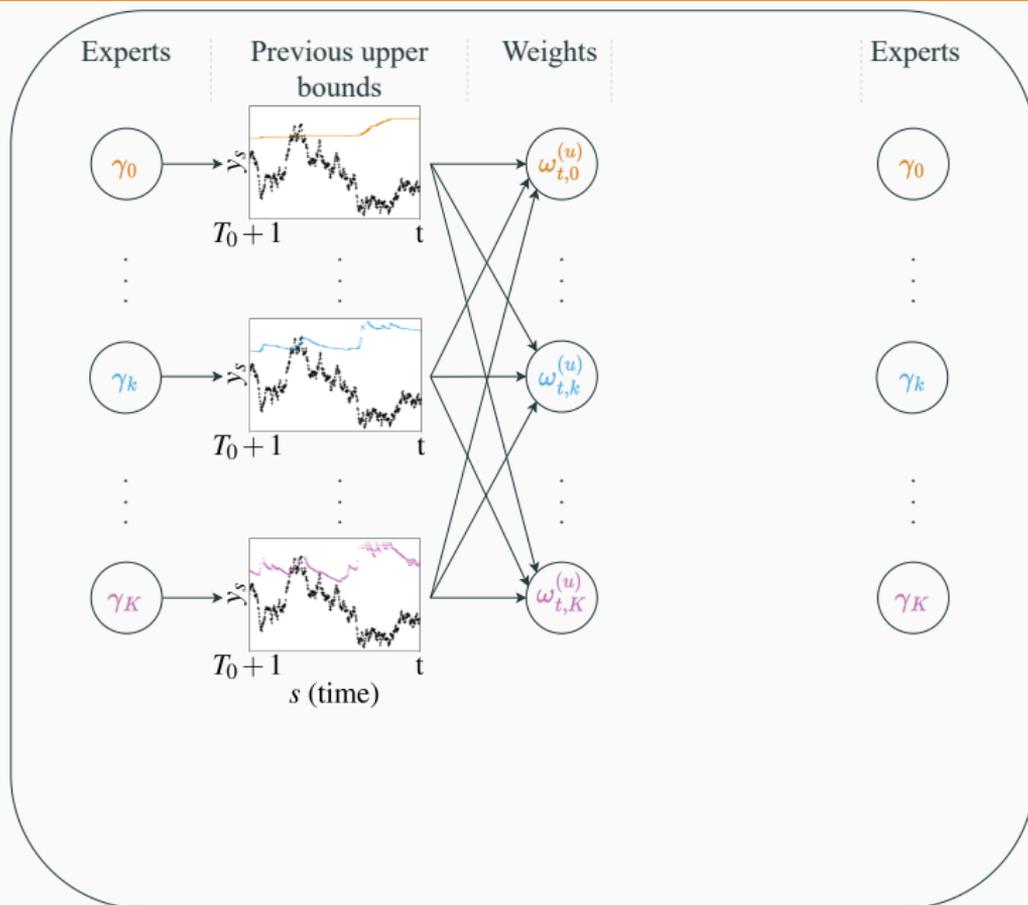
AgACI: adaptive wrapper around ACI, upper bound (Zaffran et al., 2022)



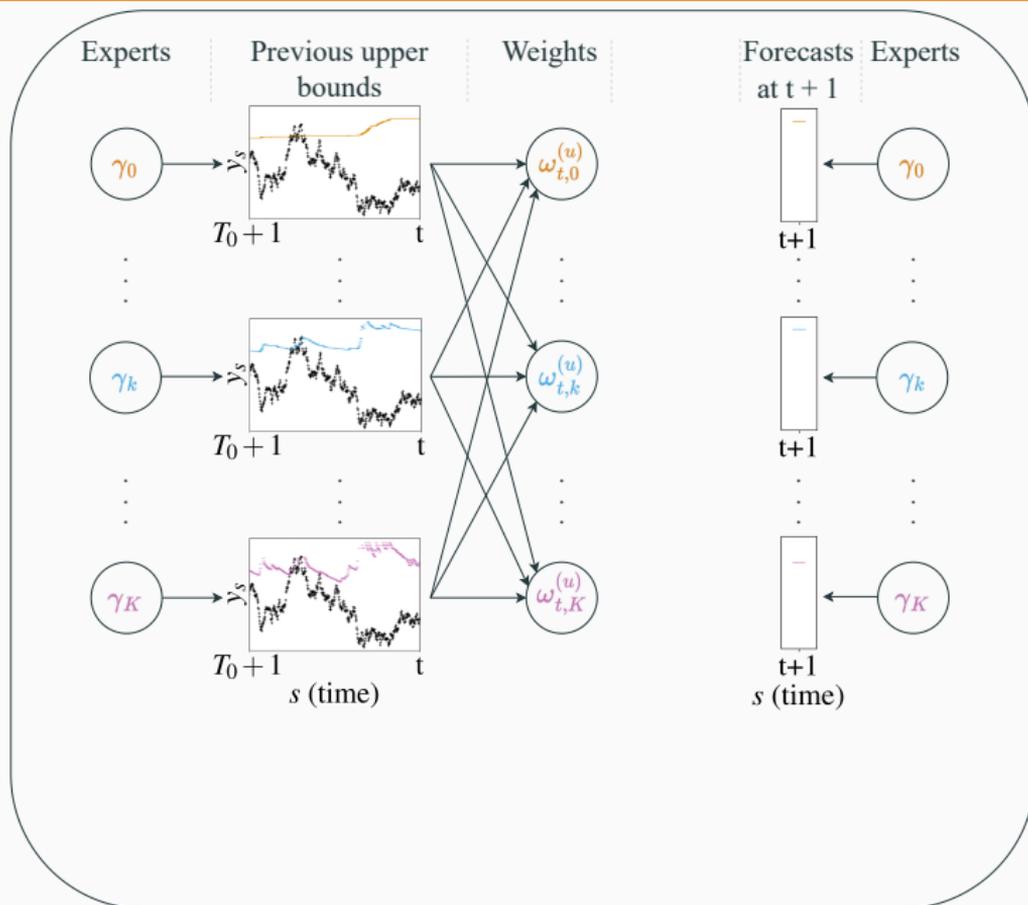
AgACI: adaptive wrapper around ACI, upper bound (Zaffran et al., 2022)



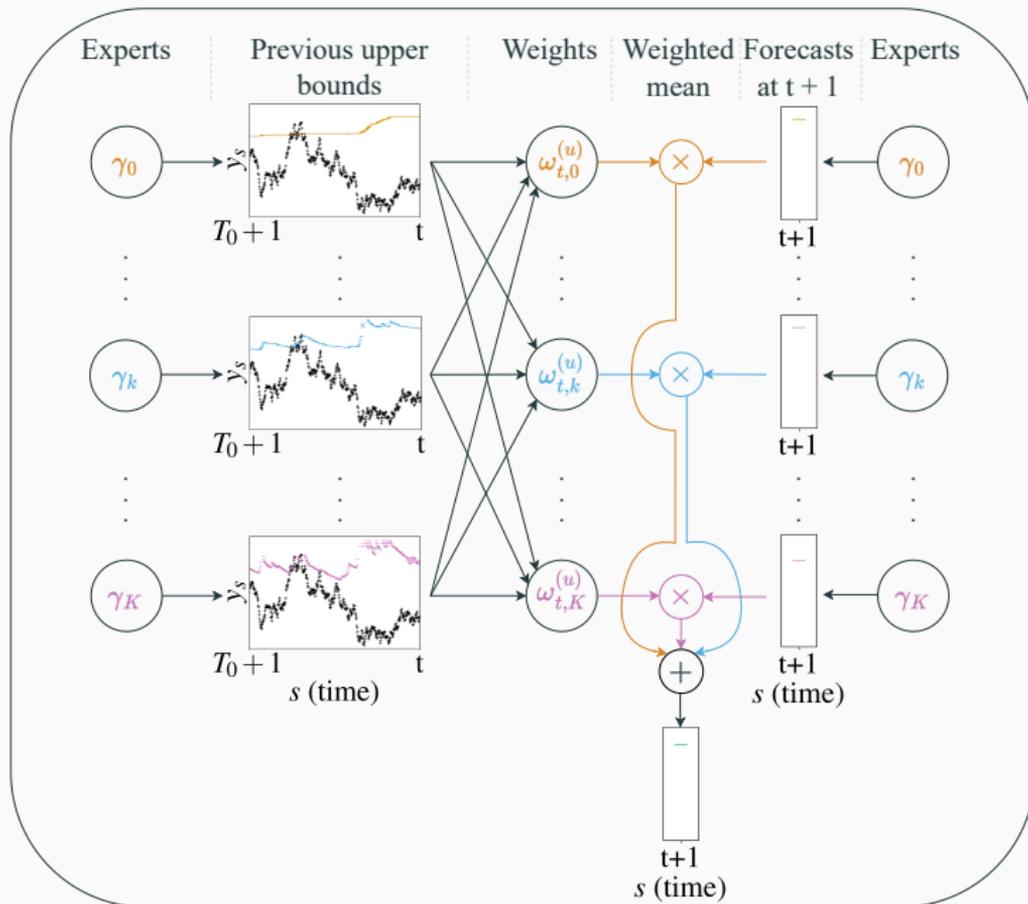
AgACI: adaptive wrapper around ACI, upper bound (Zaffran et al., 2022)



AgACI: adaptive wrapper around ACI, upper bound (Zaffran et al., 2022)



AgACI: adaptive wrapper around ACI, upper bound (Zaffran et al., 2022)



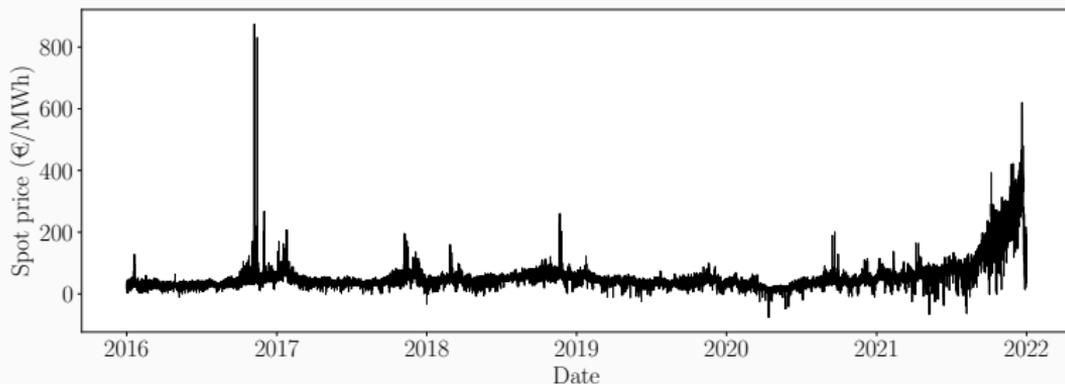
- Synthetic data with ARMA noise

- **Synthetic data with ARMA noise**
 - Benchmarks are not robust to the increase in the temporal dependence;

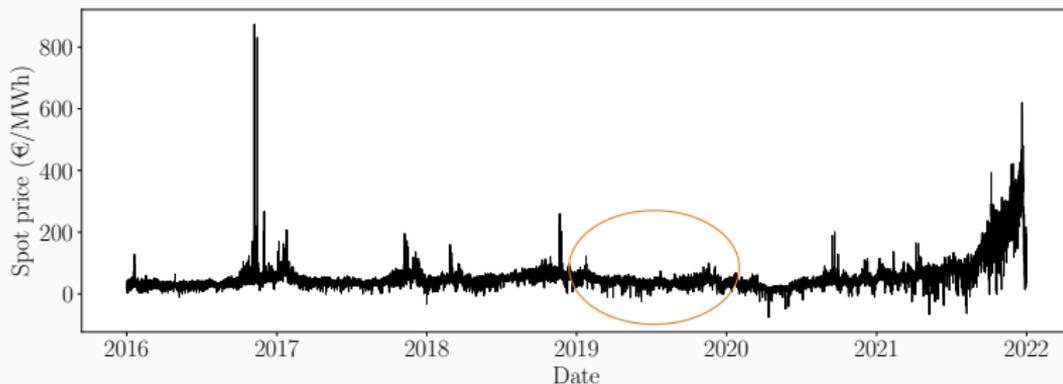
- **Synthetic data with ARMA noise**
 - Benchmarks are not robust to the increase in the temporal dependence;
 - ACI is robust, maintaining validity, with an appropriate γ ;

- **Synthetic data with ARMA noise**
 - Benchmarks are not robust to the increase in the temporal dependence;
 - ACI is robust, maintaining validity, with an appropriate γ ;
 - AgACI is robust, maintaining validity, not the smallest.

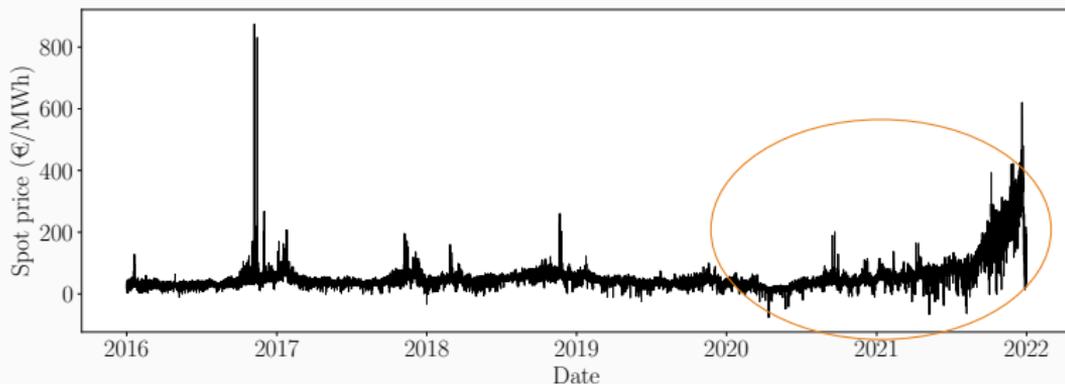
- **Synthetic data with ARMA noise**
 - Benchmarks are not robust to the increase in the temporal dependence;
 - ACI is robust, maintaining validity, with an appropriate γ ;
 - AgACI is robust, maintaining validity, not the smallest.
- **French electricity spot prices**



- **Synthetic data with ARMA noise**
 - Benchmarks are not robust to the increase in the temporal dependence;
 - ACI is robust, maintaining validity, with an appropriate γ ;
 - AgACI is robust, maintaining validity, not the smallest.
- **French electricity spot prices**
 - 2019: AgACI provides validity with a reasonable efficiency;



- **Synthetic data with ARMA noise**
 - Benchmarks are not robust to the increase in the temporal dependence;
 - ACI is robust, maintaining validity, with an appropriate γ ;
 - AgACI is robust, maintaining validity, not the smallest.
- **French electricity spot prices**
 - 2019: AgACI provides validity with a reasonable efficiency;
 - 2020 and 2021: AgACI fails to ensure validity, and the various forecasting models considered behave differently.



Improving adaptiveness for high non-stationarity (Dutot et al., 2024)

Online aggregation of various AgACI, each of them being trained with different underlying forecasting models, for each bound independently.

Improving adaptiveness for high non-stationarity (Dutot et al., 2024)

Online aggregation of various AgACI, each of them being trained with different underlying forecasting models, for each bound independently.

- ✓ Retrieves validity even in the most hazardous period of 2020 and 2021.

Improving adaptiveness for high non-stationarity (Dutot et al., 2024)

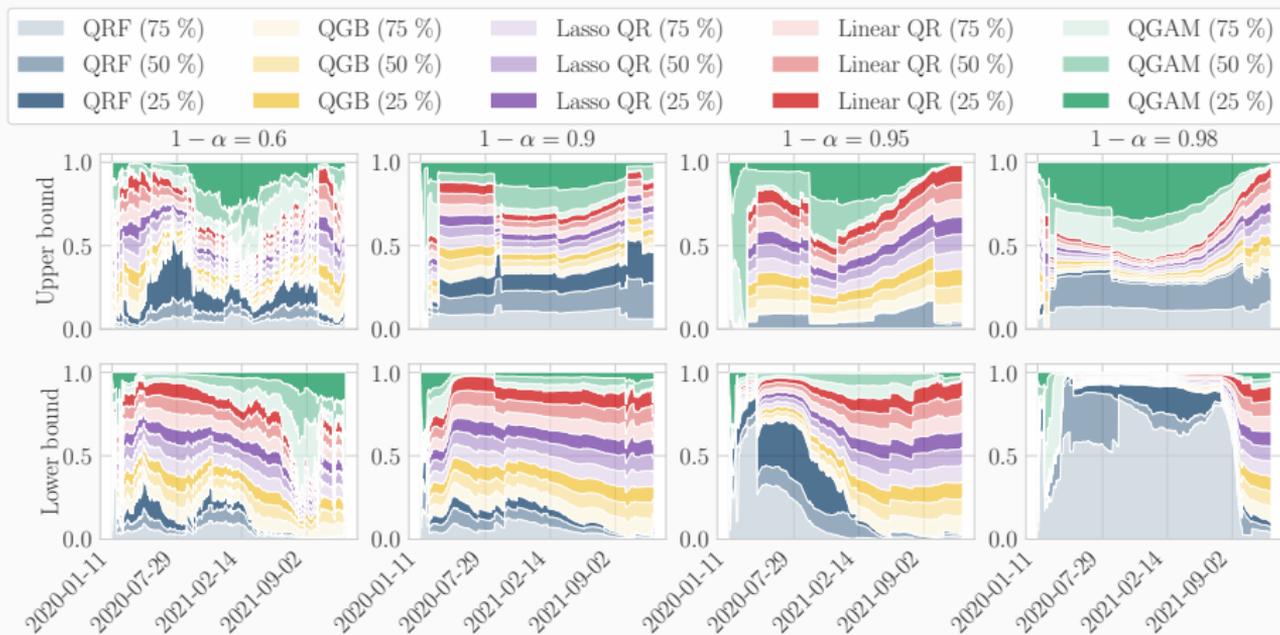
Online aggregation of various AgACI, each of them being trained with different underlying forecasting models, for each bound independently.

- ✓ Retrieves validity even in the most hazardous period of 2020 and 2021.
- ✓ Analyzing its weights provides interpretability.

Improving adaptiveness for high non-stationarity (Dutot et al., 2024)

Online aggregation of various AgACI, each of them being trained with different underlying forecasting models, for each bound independently.

- ✓ Retrieves validity even in the most hazardous period of 2020 and 2021.
- ✓ Analyzing its weights provides interpretability.



Aggregating the two bounds independently (as in AgACI and beyond):

Aggregating the two bounds independently (as in AgACI and beyond):

- ✓ Allows more flexible and adaptive behavior in practice, catching the varying nature of the predictive distribution tails

Aggregating the two bounds independently (as in AgACI and beyond):

- ✓ Allows more flexible and adaptive behavior in practice, catching the varying nature of the predictive distribution tails
- ✗ Prevents from obtaining theoretical guarantees (by opposition to Gibbs and Candès, 2022)

Aggregating the two bounds independently (as in AgACI and beyond):

- ✓ Allows more flexible and adaptive behavior in practice, catching the varying nature of the predictive distribution tails
- ✗ Prevents from obtaining theoretical guarantees (by opposition to Gibbs and Candès, 2022)

↔ Weaken the objective and consider a more practical theoretical aim?

Quantile Regression

Split Conformal Prediction (SCP)

On the design choices of conformity scores and (empirical) conditional guarantees

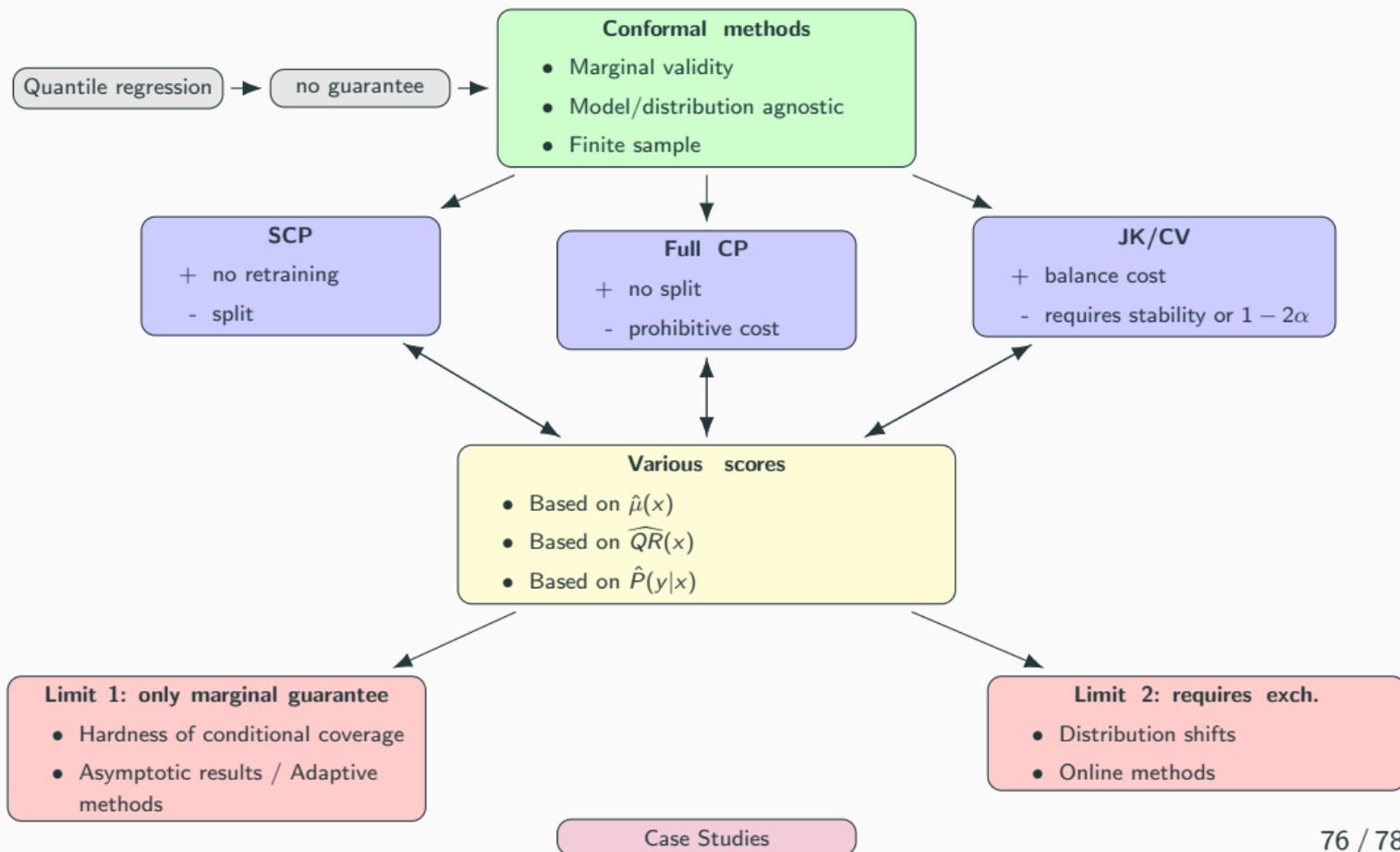
Avoiding data splitting: full conformal and out-of-bags approaches

Beyond exchangeability

Some case studies

Concluding remarks

Summary: Uncertainty quantification through conformal methods



Some (other, non-exhaustives) current open directions

- Outlier detection (Vovk et al., 2003; Bates et al., 2023)
- Selective inference, false discovery rate guarantees (Marandon et al., 2024; Gazin et al., 2024)
- Beyond the indicator loss (Angelopoulos et al., 2022a; Bates et al., 2021b; Angelopoulos et al., 2023; Lekeufack et al., 2024)
- Aggregating predictive sets (Gasparin and Ramdas, 2024b,a; Gasparin et al., 2024)

For discussion and feedback, thanks to

- Julie Josse
- Claire Boyer
- Étienne Roquain

Questions?

- Angelopoulos, A. N. and Bates, S. (2023). Conformal prediction: A gentle introduction. *Foundations and Trends® in Machine Learning*, 16(4).
- Angelopoulos, A. N., Bates, S., Candès, E. J., Jordan, M. I., and Lei, L. (2022a). Learn then test: Calibrating predictive algorithms to achieve risk control.
- Angelopoulos, A. N., Bates, S., Fisch, A., Lei, L., and Schuster, T. (2023). Conformal risk control.
- Angelopoulos, A. N., Kohli, A. P., Bates, S., Jordan, M., Malik, J., Alshaabi, T., Upadhyayula, S., and Romano, Y. (2022b). Image-to-image regression with distribution-free uncertainty quantification and applications in imaging. In *International Conference on Machine Learning*, pages 717–730. PMLR.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. (2021a). The limits of distribution-free conditional predictive inference. *Information and Inference: A Journal of the IMA*, 10(2).

- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. (2021b). Predictive inference with the jackknife+. *The Annals of Statistics*, 49(1).
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. (2022). Conformal prediction beyond exchangeability. To appear in *Annals of Statistics (2023)*.
- Bastani, O., Gupta, V., Jung, C., Noarov, G., Ramalingam, R., and Roth, A. (2022). Practical adversarial multivalid conformal prediction. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Bates, S., Angelopoulos, A., Lei, L., Malik, J., and Jordan, M. (2021a). Distribution-free, risk-controlling prediction sets. *Journal of the ACM (JACM)*, 68(6):1–34.
- Bates, S., Angelopoulos, A., Lei, L., Malik, J., and Jordan, M. (2021b). Distribution-free, risk-controlling prediction sets. *J. ACM*, 68(6).

- Bates, S., Candès, E., Lei, L., Romano, Y., and Sesia, M. (2023). Testing for outliers with conformal p-values. *The Annals of Statistics*, 51(1):149 – 178.
- Bhatnagar, A., Wang, H., Xiong, C., and Bai, Y. (2023). Improved online conformal prediction via strongly adaptive online learning. In *Proceedings of the 40th International Conference on Machine Learning*. PMLR.
- Bian, M. and Barber, R. F. (2023). Training-conditional coverage for distribution-free predictive inference. *Electronic Journal of Statistics*, 17(2):2044 – 2066.
- Cauchois, M., Gupta, S., Ali, A., and Duchi, J. C. (2020). Robust Validation: Confident Predictions Even When Distributions Shift. arXiv: 2008.04267.
- Chernozhukov, V., Wüthrich, K., and Yinchu, Z. (2018). Exact and Robust Conformal Inference Methods for Predictive Machine Learning with Dependent Data. In *Conference On Learning Theory*. PMLR.

- Chernozhukov, V., Wüthrich, K., and Zhu, Y. (2021). Distributional conformal prediction. *Proceedings of the National Academy of Sciences*, 118(48).
- Cherubin, G., Chatzikokolakis, K., and Jaggi, M. (2021). Exact optimization of conformal predictors via incremental and decremental learning. In *Proceedings of the 38th International Conference on Machine Learning*. PMLR.
- Ding, T., Angelopoulos, A., Bates, S., Jordan, M., and Tibshirani, R. J. (2023). Class-conditional conformal prediction with many classes. In Oh, A., Neumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S., editors, *Advances in Neural Information Processing Systems*, volume 36, pages 64555–64576. Curran Associates, Inc.
- Dutot, G., Zaffran, M., Féron, O., and Goude, Y. (2024). Adaptive probabilistic forecasting of french electricity spot prices.

- Feldman, S., Bates, S., and Romano, Y. (2021). Improving Conditional Coverage via Orthogonal Quantile Regression. *arXiv:2106.00394 [cs]*. arXiv: 2106.00394.
- Gasparin, M. and Ramdas, A. (2024a). Conformal online model aggregation.
- Gasparin, M. and Ramdas, A. (2024b). Merging uncertainty sets via majority vote.
- Gasparin, M., Wang, R., and Ramdas, A. (2024). Combining exchangeable p-values.
- Gazin, U., Blanchard, G., and Roquain, E. (2024). Transductive conformal inference with adaptive scores.
- Gibbs, I. and Candès, E. (2021). Adaptive conformal inference under distribution shift. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Gibbs, I. and Candès, E. (2022). Conformal inference for online prediction with arbitrary distribution shifts. arXiv: 2208.08401.

- Gibbs, I., Cherian, J. J., and Candès, E. J. (2023). Conformal prediction with conditional guarantees. arXiv: 2305.12616.
- Guan, L. (2022). Localized conformal prediction: a generalized inference framework for conformal prediction. *Biometrika*, 110(1).
- Gupta, C., Kuchibhotla, A. K., and Ramdas, A. (2022). Nested conformal prediction and quantile out-of-bag ensemble methods. *Pattern Recognition*, 127.
- Izbicki, R., Shimizu, G., and Stern, R. B. (2022). CD-split and HPD-split: Efficient conformal regions in high dimensions. *Journal of Machine Learning Research*, 23(87).
- Jung, C., Noarov, G., Ramalingam, R., and Roth, A. (2023). Batch multivald conformal prediction. In *International Conference on Learning Representations*.

- Kivaranovic, D., Johnson, K. D., and Leeb, H. (2020). Adaptive, Distribution-Free Prediction Intervals for Deep Networks. In *International Conference on Artificial Intelligence and Statistics*. PMLR.
- Lei, J. (2019). Fast exact conformalization of the lasso using piecewise linear homotopy. *Biometrika*, 106(4).
- Lei, J., G'Sell, M., Rinaldo, A., Tibshirani, R. J., and Wasserman, L. (2018). Distribution-Free Predictive Inference for Regression. *Journal of the American Statistical Association*.
- Lei, J. and Wasserman, L. (2014). Distribution-free prediction bands for non-parametric regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1).

- Lekeufack, J., Angelopoulos, A. N., Bajcsy, A., Jordan, M. I., and Malik, J. (2024). Conformal decision theory: Safe autonomous decisions from imperfect predictions.
- Marandon, A., Lei, L., Mary, D., and Roquain, E. (2024). Adaptive novelty detection with false discovery rate guarantee. *The Annals of Statistics*, 52(1):157 – 183.
- Ndiaye, E. (2022). Stable conformal prediction sets. In *Proceedings of the 39th International Conference on Machine Learning*. PMLR.
- Ndiaye, E. and Takeuchi, I. (2019). Computing full conformal prediction set with approximate homotopy. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Ndiaye, E. and Takeuchi, I. (2022). Root-finding approaches for computing conformal prediction set. *Machine Learning*, 112(1).

- Nouretdinov, I., Melluish, T., and Vovk, V. (2001). Ridge regression confidence machine. In *Proceedings of the 18th International Conference on Machine Learning*.
- Papadopoulos, H., Proedrou, K., Vovk, V., and Gammerman, A. (2002). Inductive Confidence Machines for Regression. In *Machine Learning: ECML*. Springer.
- Podkopaev, A. and Ramdas, A. (2021). Distribution-free uncertainty quantification for classification under label shift. In *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*. PMLR.
- Romano, Y., Barber, R. F., Sabatti, C., and Candès, E. (2020a). With Malice Toward None: Assessing Uncertainty via Equalized Coverage. *Harvard Data Science Review*, 2(2).

- Romano, Y., Patterson, E., and Candès, E. (2019). Conformalized Quantile Regression. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Romano, Y., Sesia, M., and Candès, E. (2020b). Classification with valid and adaptive coverage. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Sadinle, M., Lei, J., and Wasserman, L. (2018). Least ambiguous set-valued classifiers with bounded error levels. *Journal of the American Statistical Association*, 114(525):223–234.
- Sesia, M. and Romano, Y. (2021). Conformal prediction using conditional histograms. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.

- Tibshirani, R. J., Barber, R. F., Candes, E., and Ramdas, A. (2019). Conformal Prediction Under Covariate Shift. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- Vovk, V. (2012). Conditional Validity of Inductive Conformal Predictors. In *Asian Conference on Machine Learning*. PMLR.
- Vovk, V. (2015). Cross-conformal predictors. *Annals of Mathematics and Artificial Intelligence*, 74(1-2).
- Vovk, V., Gammerman, A., and Shafer, G. (2005). *Algorithmic Learning in a Random World*. Springer US.
- Vovk, V., Nouretdinov, I., and Gammerman, A. (2003). Testing exchangeability on-line. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning, ICML'03*, page 768–775. AAAI Press.

- Zaffran, M., Dieuleveut, A., Josse, J., and Romano, Y. (2024). Predictive uncertainty quantification with missing values. Preprint submitted to *Journal of Machine Learning Research*, arXiv arXiv:2405.15641.
- Zaffran, M., Féron, O., Goude, Y., Josse, J., and Dieuleveut, A. (2022). Adaptive conformal predictions for time series. In *Proceedings of the 39th International Conference on Machine Learning*. PMLR.

