

# Conformal Prediction with Missing Values

---

Margaux Zaffran

54èmes Journées de Statistiques, Bruxelles, 2023

Session MALIA





**Aymeric Dieuleveut**

Ecole

Polytechnique

*Paris - France*



**Julie Josse**

INRIA

IDESP

*Montpellier - France*



**Yaniv Romano**

Technion - Israel Institute  
of Technology

*Haifa - Israel*

What about splitting the data?

Standard Split Conformal Prediction for Mean-Regression

Conformalized Quantile Regression

Predictive uncertainty quantification with missing values

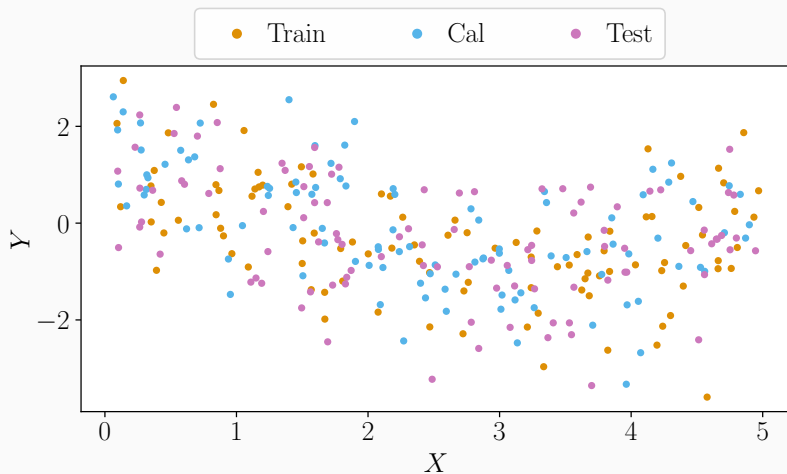
What about splitting the data?

Standard Split Conformal Prediction for Mean-Regression

Conformalized Quantile Regression

Predictive uncertainty quantification with missing values

# Split Conformal Prediction (SCP)<sup>1,2,3</sup>: toy example

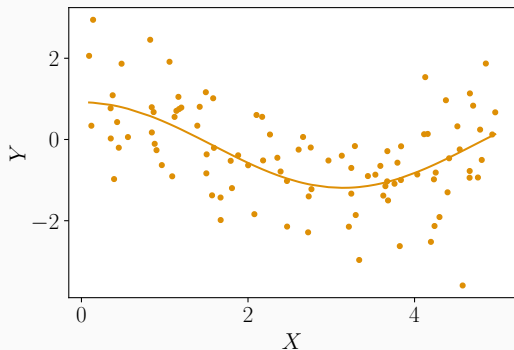


<sup>1</sup>Vovk et al. (2005), *Algorithmic Learning in a Random World*

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

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

# Split Conformal Prediction (SCP)<sup>1,2,3</sup>: training step



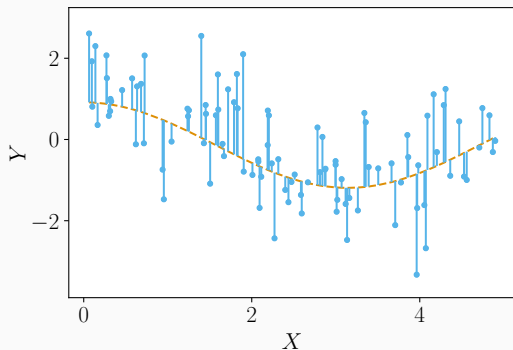
► Learn (or get)  $\hat{\mu}$

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

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

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

# Split Conformal Prediction (SCP)<sup>1,2,3</sup>: calibration step



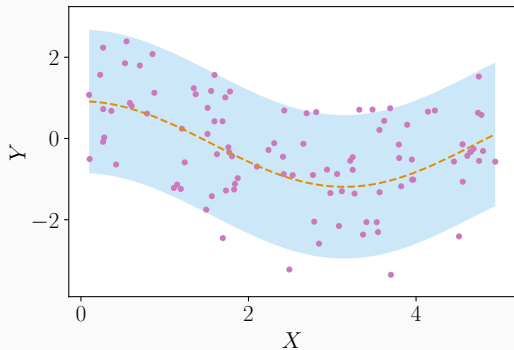
- ▶ Predict with  $\hat{\mu}$
- ▶ Get the **|residuals|**, a.k.a. scores  $\{\mathcal{S}^{(k)}\}_{k \in \text{Cal}}$
- ▶ Compute the  $(1 - \alpha)$  empirical quantile of  $\mathcal{S} = \{|\text{residuals}|\}_{\text{Cal}} \cup \{+\infty\}$ , noted  $q_{1-\alpha}(\mathcal{S})$

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

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

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

# Split Conformal Prediction (SCP)<sup>1,2,3</sup>: prediction step



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

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

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

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





## Definition (Exchangeability)

$(X^{(k)}, Y^{(k)})_{k=1}^n$  are **exchangeable** if for any permutation  $\sigma$  of  $\llbracket 1, n \rrbracket$  we have:

$$\begin{aligned} & \mathcal{L}((X^{(1)}, Y^{(1)}), \dots, (X^{(n)}, Y^{(n)})) \\ = & \mathcal{L}((X^{(\sigma(1))}, Y^{(\sigma(1))}), \dots, (X^{(\sigma(n))}, Y^{(\sigma(n))})), \end{aligned}$$

where  $\mathcal{L}$  designates the joint distribution.

### Definition (Exchangeability)

$(X^{(k)}, Y^{(k)})_{k=1}^n$  are **exchangeable** if for any permutation  $\sigma$  of  $\llbracket 1, n \rrbracket$  we have:

$$\begin{aligned} & \mathcal{L}((X^{(1)}, Y^{(1)}), \dots, (X^{(n)}, Y^{(n)})) \\ = & \mathcal{L}((X^{(\sigma(1))}, Y^{(\sigma(1))}), \dots, (X^{(\sigma(n))}, Y^{(\sigma(n))})), \end{aligned}$$

where  $\mathcal{L}$  designates the joint distribution.

**Toy case:**  $Z^{(1)}$  and  $Z^{(2)}$  are exchangeable if  $(Z^{(1)}, Z^{(2)}) \stackrel{\mathcal{L}}{=} (Z^{(2)}, Z^{(1)})$ .

## Definition (Exchangeability)

$(X^{(k)}, Y^{(k)})_{k=1}^n$  are **exchangeable** if for any permutation  $\sigma$  of  $\llbracket 1, n \rrbracket$  we have:

$$\begin{aligned} & \mathcal{L}((X^{(1)}, Y^{(1)}), \dots, (X^{(n)}, Y^{(n)})) \\ &= \mathcal{L}((X^{(\sigma(1))}, Y^{(\sigma(1))}), \dots, (X^{(\sigma(n))}, Y^{(\sigma(n))})), \end{aligned}$$

where  $\mathcal{L}$  designates the joint distribution.

## Examples of exchangeable sequences

- i.i.d. samples

## Definition (Exchangeability)

$(X^{(k)}, Y^{(k)})_{k=1}^n$  are **exchangeable** if for any permutation  $\sigma$  of  $\llbracket 1, n \rrbracket$  we have:

$$\begin{aligned} & \mathcal{L}((X^{(1)}, Y^{(1)}), \dots, (X^{(n)}, Y^{(n)})) \\ &= \mathcal{L}((X^{(\sigma(1))}, Y^{(\sigma(1))}), \dots, (X^{(\sigma(n))}, Y^{(\sigma(n))})), \end{aligned}$$

where  $\mathcal{L}$  designates the joint distribution.

## Examples of exchangeable sequences

- i.i.d. samples

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

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

### Theorem

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

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

Additionally, if the scores  $\{S^{(k)}\}_{k \in \text{Cal}}$  are a.s. distinct:

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

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

## Theorem

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

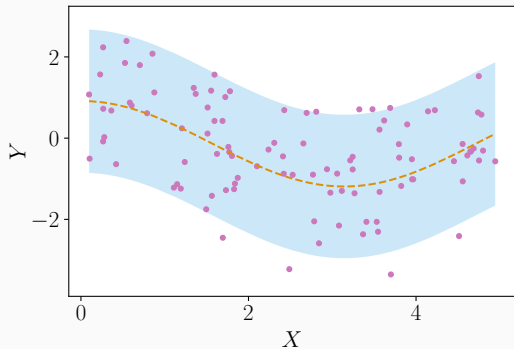
$$\mathbb{P} \left\{ Y^{(n+1)} \in \widehat{C}_\alpha \left( X^{(n+1)} \right) \right\} \geq 1 - \alpha.$$

Additionally, if the scores  $\{S^{(k)}\}_{k \in \text{Cal}}$  are a.s. distinct:

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

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

## Standard mean-regression SCP is not adaptive



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



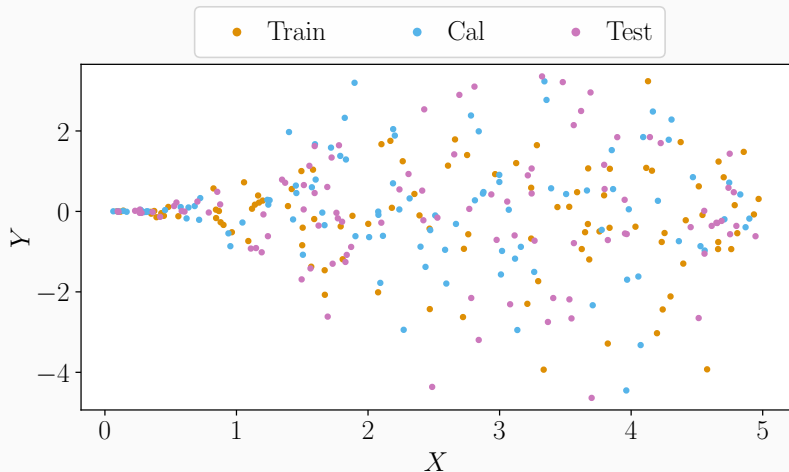
What about splitting the data?

Standard Split Conformal Prediction for Mean-Regression

Conformalized Quantile Regression

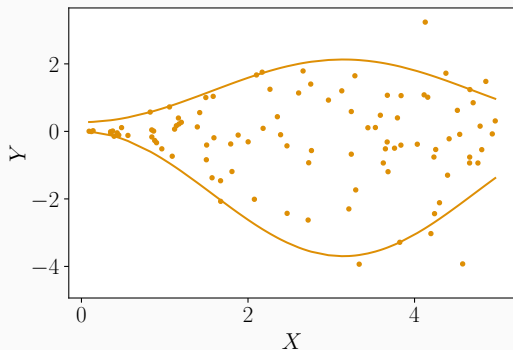
Predictive uncertainty quantification with missing values

# Conformalized Quantile Regression (CQR)<sup>4</sup>



<sup>4</sup>Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

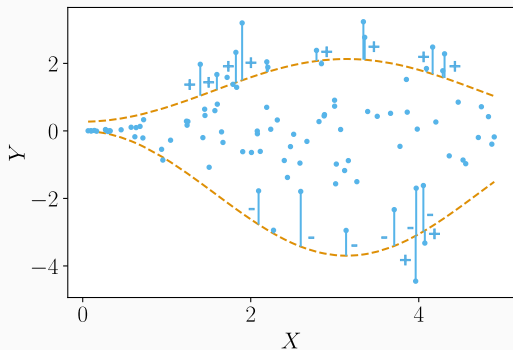
# Conformalized Quantile Regression (CQR)<sup>4</sup>: training step



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

<sup>4</sup>Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

# Conformalized Quantile Regression (CQR)<sup>4</sup>: calibration step

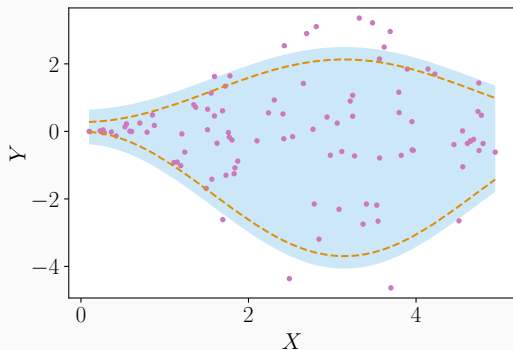


- ▶ Predict with  $\widehat{QR}_{\text{lower}}$  and  $\widehat{QR}_{\text{upper}}$
- ▶ Get the scores  $\mathcal{S} = \{S^{(k)}\}_{\text{Cal}} \cup \{+\infty\}$
- ▶ Compute the  $(1 - \alpha)$  empirical quantile of  $\mathcal{S}$ , noted  $q_{1-\alpha}(\mathcal{S})$

$$\hookrightarrow S^{(k)} := \max \left\{ \widehat{QR}_{\text{lower}}(X^{(k)}) - Y^{(k)}, Y^{(k)} - \widehat{QR}_{\text{upper}}(X^{(k)}) \right\}$$

<sup>4</sup>Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

# Conformalized Quantile Regression (CQR)<sup>4</sup>: prediction step



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

► Build

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

<sup>4</sup>Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

## CQR: theoretical guarantees

CQR is a **particular case of SCP**.

CQR is a **particular case of SCP**.

Therefore, it enjoys finite sample guarantees proved in Romano et al. (2019).

### Theorem

Suppose  $(X^{(k)}, Y^{(k)})_{k=1}^{n+1}$  are *exchangeable (or i.i.d.)*. CQR applied on  $(X^{(k)}, Y^{(k)})_{k=1}^n$  outputs  $\widehat{C}_\alpha(\cdot)$  such that:

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

Additionally, if the scores  $\{S^{(k)}\}_{k \in \text{Cal}}$  are *a.s. distinct*:

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

CQR is a **particular case of SCP**.

Therefore, it enjoys finite sample guarantees proved in Romano et al. (2019).

### Theorem

Suppose  $(X^{(k)}, Y^{(k)})_{k=1}^{n+1}$  are *exchangeable (or i.i.d.)*. CQR applied on  $(X^{(k)}, Y^{(k)})_{k=1}^n$  outputs  $\widehat{C}_\alpha(\cdot)$  such that:

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

Additionally, if the scores  $\{S^{(k)}\}_{k \in \text{Cal}}$  are a.s. distinct:

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

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



What about splitting the data?

## Predictive uncertainty quantification with missing values

Learning with Missing Data

Conformal Prediction with Missing Values

Missing Data Augmentation

Experimental Results

Conclusions

What about splitting the data?

Predictive uncertainty quantification with missing values

Learning with Missing Data

Conformal Prediction with Missing Values

Missing Data Augmentation

Experimental Results

Conclusions

## Missing values are ubiquitous and challenging

Data:  $(X^{(k)}, Y^{(k)})_{k=1}^n$

Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>
22.42	0.55	0.67	0.03
8.26	0.72	0.18	0.55
19.41	0.60	0.58	NA
19.75	0.54	0.43	0.96
7.32	NA	0.19	NA
13.55	0.65	0.69	0.50
20.75	NA	NA	0.61
9.26	0.89	NA	0.84
9.68	0.963	0.45	0.65

## Missing values are ubiquitous and challenging

Data:  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^n$

Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	Mask M =		
				(M <sub>1</sub>	M <sub>2</sub>	M <sub>3</sub> )
22.42	0.55	0.67	0.03	0	0	0
8.26	0.72	0.18	0.55	0	0	0
19.41	0.60	0.58	NA	0	0	1
19.75	0.54	0.43	0.96	0	0	0
7.32	NA	0.19	NA	1	0	1
13.55	0.65	0.69	0.50	0	0	0
20.75	NA	NA	0.61	1	1	0
9.26	0.89	NA	0.84	0	1	0
9.68	0.963	0.45	0.65	0	0	0

## Missing values are ubiquitous and challenging

Data:  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^n$

Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	Mask M =		
				(M <sub>1</sub>	M <sub>2</sub>	M <sub>3</sub> )
22.42	0.55	0.67	0.03	0	0	0
8.26	0.72	0.18	0.55	0	0	0
19.41	0.60	0.58	NA	0	0	1
19.75	0.54	0.43	0.96	0	0	0
7.32	NA	0.19	NA	1	0	1
13.55	0.65	0.69	0.50	0	0	0
20.75	NA	NA	0.61	1	1	0
9.26	0.89	NA	0.84	0	1	0
9.68	0.963	0.45	0.65	0	0	0

↔  $2^d$  potential masks.

## Missing values are ubiquitous and challenging

Data:  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^n$

Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	Mask M =		
				(M <sub>1</sub>	M <sub>2</sub>	M <sub>3</sub> )
22.42	0.55	0.67	0.03	0	0	0
8.26	0.72	0.18	0.55	0	0	0
19.41	0.60	0.58	NA	0	0	1
19.75	0.54	0.43	0.96	0	0	0
7.32	NA	0.19	NA	1	0	1
13.55	0.65	0.69	0.50	0	0	0
20.75	NA	NA	0.61	1	1	0
9.26	0.89	NA	0.84	0	1	0
9.68	0.963	0.45	0.65	0	0	0

↔  $2^d$  potential masks.

↔  $M$  can depend on  $X$  or  $Y$ .

## Missing values are ubiquitous and challenging

Data:  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^n$

Y	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	Mask M =		
				(M <sub>1</sub>	M <sub>2</sub>	M <sub>3</sub> )
22.42	0.55	0.67	0.03	0	0	0
8.26	0.72	0.18	0.55	0	0	0
19.41	0.60	0.58	NA	0	0	1
19.75	0.54	0.43	0.96	0	0	0
7.32	NA	0.19	NA	1	0	1
13.55	0.65	0.69	0.50	0	0	0
20.75	NA	NA	0.61	1	1	0
9.26	0.89	NA	0.84	0	1	0
9.68	0.963	0.45	0.65	0	0	0

↔  $2^d$  potential masks.

↔  $M$  can depend on  $X$  or  $Y$ .

⇒ Statistical and computational challenges.

## Supervised learning with missing values: impute-then-regress

Impute-then-regress procedures are widely used.



# Supervised learning with missing values: impute-then-regress

Impute-then-regress procedures are widely used.

1. Replace NA using an **imputation function** (e.g. the mean), noted  $\phi$ .

$x^{(1)}$	-1	-10	6	0
$x^{(2)}$	4	NA	-2	2
$x^{(3)}$	5	1	2	NA
$x^{(4)}$	0	NA	NA	1

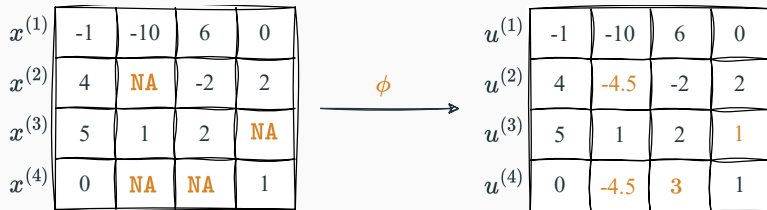
$\xrightarrow{\phi}$

$u^{(1)}$	-1	-10	6	0
$u^{(2)}$	4	-4.5	-2	2
$u^{(3)}$	5	1	2	1
$u^{(4)}$	0	-4.5	3	1

# Supervised learning with missing values: impute-then-regress

Impute-then-regress procedures are widely used.

1. Replace NA using an **imputation function** (e.g. the mean), noted  $\phi$ .



2. Train your algorithm (Random Forest, Neural Nets, etc.) on the imputed

$$\text{data: } \left\{ \underbrace{\phi \left( X_{\text{obs}(M^{(k)})}^{(k)}, M^{(k)} \right)}_{U^{(k)} = \text{imputed } X^{(k)}}, Y^{(k)} \right\}_{k=1}^n .$$

# Supervised learning with missing values: impute-then-regress

Impute-then-regress procedures are widely used.

1. Replace NA using an imputation function (e.g. the mean), noted  $\phi$ .

$x^{(1)}$	-1	-10	6	0
$x^{(2)}$	4	NA	-2	2
$x^{(3)}$	5	1	2	NA
$x^{(4)}$	0	NA	NA	1

$\xrightarrow{\phi}$

$u^{(1)}$	-1	-10	6	0
$u^{(2)}$	4	-4.5	-2	2
$u^{(3)}$	5	1	2	1
$u^{(4)}$	0	-4.5	3	1

2. Train your algorithm (Random Forest, Neural Nets, etc.) on the imputed

$$\text{data: } \left\{ \underbrace{\phi \left( X_{\text{obs}(M^{(k)})}^{(k)}, M^{(k)} \right)}_{U^{(k)} = \text{imputed } X^{(k)}}, Y^{(k)} \right\}_{k=1}^n .$$

↪ we consider an **impute-then-regress** pipeline in this work.

## Predictive uncertainty quantification with missing values

**Goal:** predict  $Y^{(n+1)}$  with **confidence**  $1 - \alpha$ , i.e. build the smallest  $\mathcal{C}_\alpha$  such that:

## Predictive uncertainty quantification with missing values

**Goal:** predict  $Y^{(n+1)}$  with **confidence**  $1 - \alpha$ , i.e. build the smallest  $\mathcal{C}_\alpha$  such that:

### 1. Marginal Validity (MV)

$$\mathbb{P} \left\{ Y^{(n+1)} \in \mathcal{C}_\alpha \left( X^{(n+1)}, M^{(n+1)} \right) \right\} \geq 1 - \alpha. \quad (\text{MV})$$

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

# Predictive uncertainty quantification with missing values

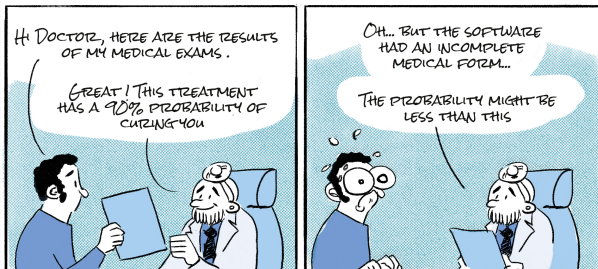
**Goal:** predict  $Y^{(n+1)}$  with **confidence**  $1 - \alpha$ , i.e. build the smallest  $C_\alpha$  such that:

## 1. Marginal Validity (MV)

$$\mathbb{P} \left\{ Y^{(n+1)} \in C_\alpha \left( X^{(n+1)}, M^{(n+1)} \right) \right\} \geq 1 - \alpha. \quad (\text{MV})$$

## 2. Mask-Conditional-Validity (MCV)

$$\forall m \in \{0, 1\}^d : \mathbb{P} \left\{ Y^{(n+1)} \in C_\alpha \left( X^{(n+1)}, m \right) \mid M^{(n+1)} = m \right\} \geq 1 - \alpha. \quad (\text{MCV})$$



Illustrations @theoreminger

# Predictive uncertainty quantification with missing values

**Goal:** predict  $Y^{(n+1)}$  with **confidence**  $1 - \alpha$ , i.e. build the smallest  $\mathcal{C}_\alpha$  such that:

## 1. Marginal Validity (MV)

$$\mathbb{P} \left\{ Y^{(n+1)} \in \mathcal{C}_\alpha \left( X^{(n+1)}, M^{(n+1)} \right) \right\} \geq 1 - \alpha. \quad (\text{MV})$$

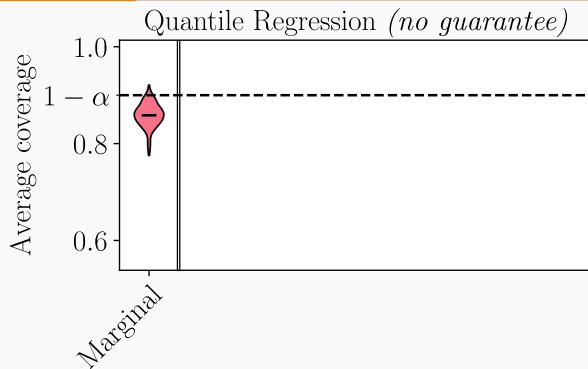
## 2. Mask-Conditional-Validity (MCV)

$$\forall m \in \{0, 1\}^d : \mathbb{P} \left\{ Y^{(n+1)} \in \mathcal{C}_\alpha \left( X^{(n+1)}, m \right) \mid M^{(n+1)} = m \right\} \geq 1 - \alpha. \quad (\text{MCV})$$

3 considered approaches to reach these goals.

	Quantile Regression (QR)		
(MV)	?		
(MCV)	?		

## Quantile Regression (QR) intervals

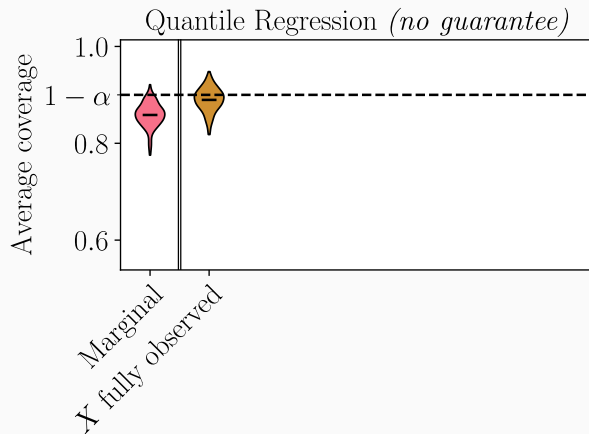


- Marginal validity (eq. (MV), i.e. on average) is not reached!

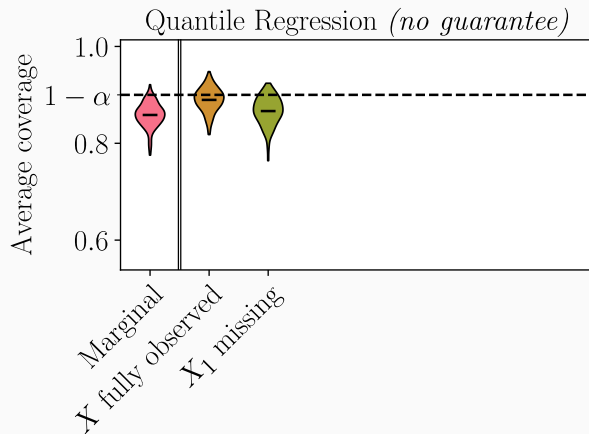
	QR		
(MV)	$\times$		
(MCV)			



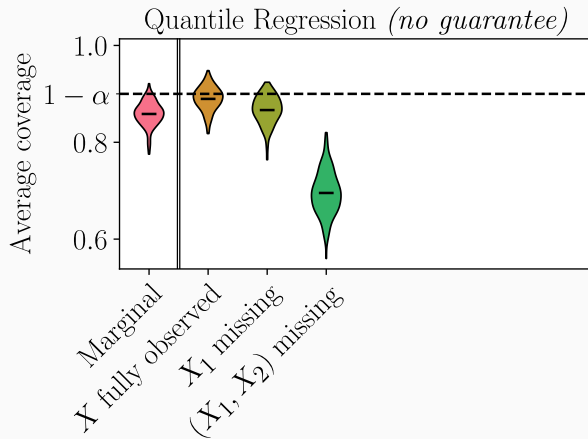
## Quantile Regression (QR) intervals



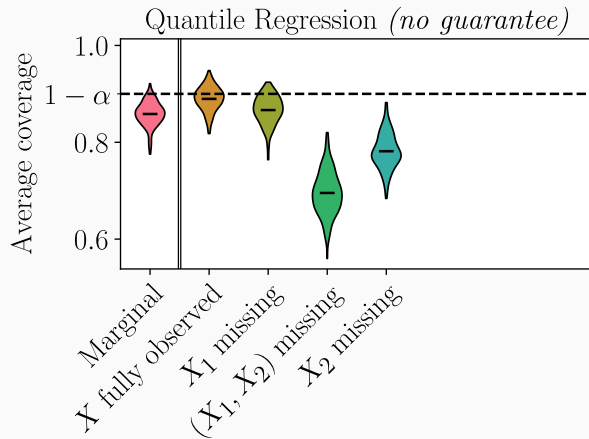
# Quantile Regression (QR) intervals



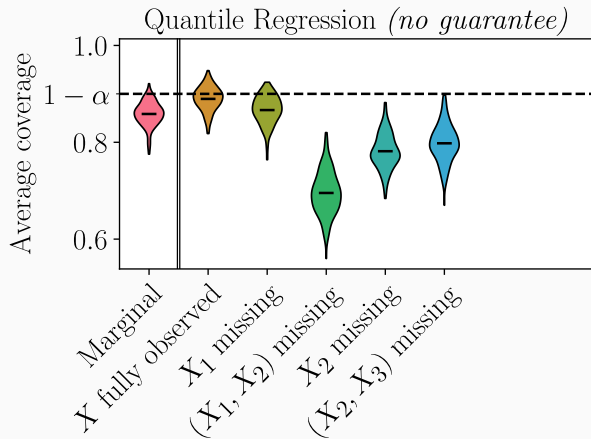
# Quantile Regression (QR) intervals



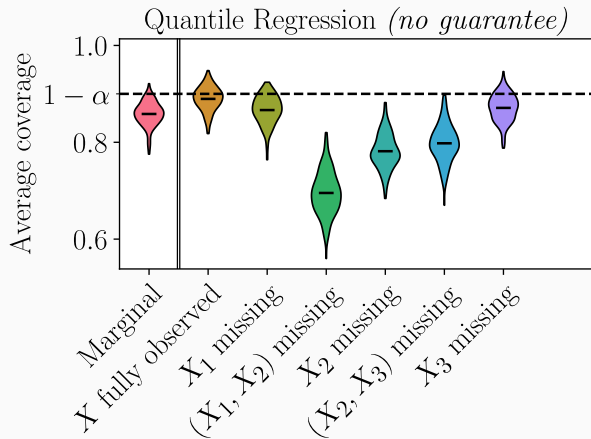
# Quantile Regression (QR) intervals



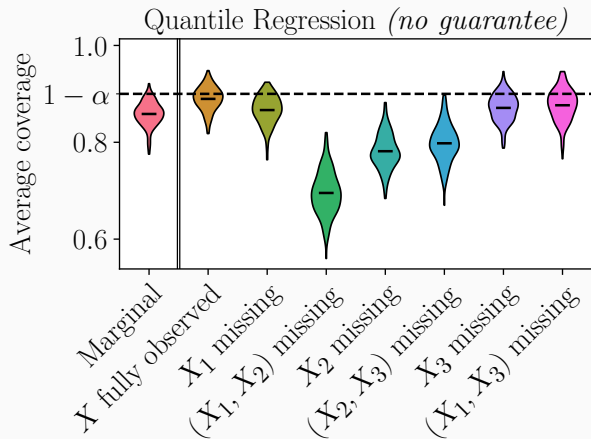
# Quantile Regression (QR) intervals



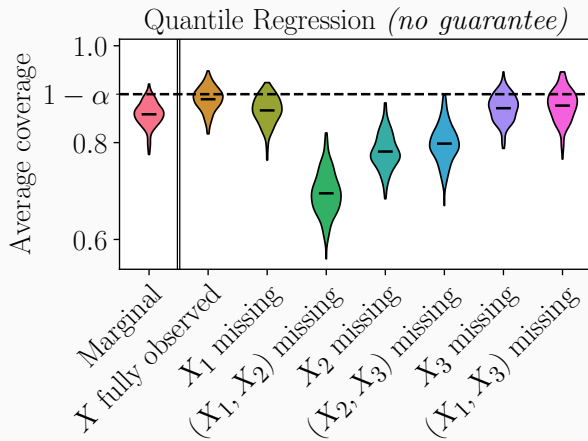
# Quantile Regression (QR) intervals



# Quantile Regression (QR) intervals



# Quantile Regression (QR) intervals

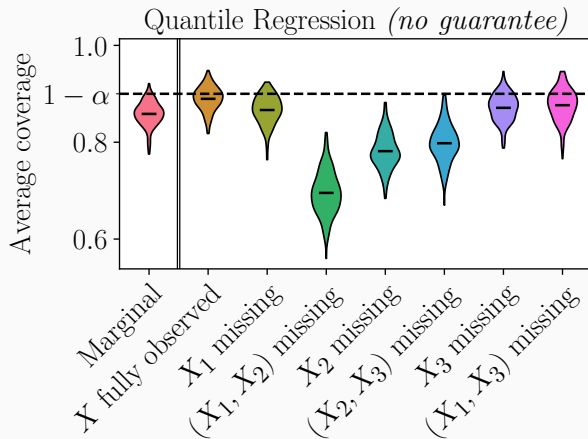


- The predictive uncertainty strongly depends on the mask

	QR		
(MV)	✗		
(MCV)	✗		



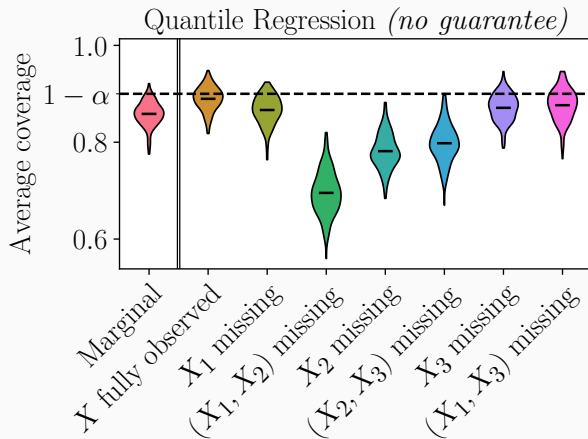
# Quantile Regression (QR) intervals



- The predictive uncertainty strongly depends on the mask

↔ missing values induce heteroskedasticity

## Quantile Regression (QR) intervals



- The predictive uncertainty strongly depends on the mask

↔ missing values induce heteroskedasticity

↔ supported by theory on the Gaussian Linear Model

## Missing values induce heteroskedasticity

Theoretical study of the Gaussian linear model ( $Y = \beta^T X + \varepsilon$ ) generalizes

↔ **oracle** intervals: smallest predictive interval when the distribution of  $Y|(X, M)$  is known

## Missing values induce heteroskedasticity

Theoretical study of the Gaussian linear model ( $Y = \beta^T X + \varepsilon$ ) generalizes  
 $\hookrightarrow$  **oracle** intervals: smallest predictive interval when the distribution of  $Y|(X, M)$   
is known

### Proposition (Oracle intervals under the Gaussian lin. mod.)

$$\mathcal{L}_\alpha^*(m) = 2 \times q_{1-\alpha/2}^{\mathcal{N}(0,1)} \times \sqrt{\beta_{\text{mis}(m)}^T \Sigma_{\text{mis}|\text{obs}}^m \beta_{\text{mis}(m)} + \sigma_\varepsilon^2}.$$

## Missing values induce heteroskedasticity

Theoretical study of the Gaussian linear model ( $Y = \beta^T X + \varepsilon$ ) generalizes  
 $\hookrightarrow$  **oracle** intervals: smallest predictive interval when the distribution of  $Y|(X, M)$   
is known

### Proposition (Oracle intervals under the Gaussian lin. mod.)

$$\mathcal{L}_\alpha^*(m) = 2 \times q_{1-\alpha/2}^{\mathcal{N}(0,1)} \times \sqrt{\beta_{\text{mis}(m)}^T \Sigma_{\text{mis}|\text{obs}}^m \beta_{\text{mis}(m)} + \sigma_\varepsilon^2}.$$

- Even with an homoskedastic noise, missingness generates **heteroskedasticity**

## Missing values induce heteroskedasticity

Theoretical study of the Gaussian linear model ( $Y = \beta^T X + \varepsilon$ ) generalizes  
 $\hookrightarrow$  **oracle** intervals: smallest predictive interval when the distribution of  $Y|(X, M)$   
is known

### Proposition (Oracle intervals under the Gaussian lin. mod.)

$$\mathcal{L}_\alpha^*(m) = 2 \times q_{1-\alpha/2}^{\mathcal{N}(0,1)} \times \sqrt{\beta_{\text{mis}(m)}^T \Sigma_{\text{mis}|\text{obs}}^m \beta_{\text{mis}(m)} + \sigma_\varepsilon^2}.$$

- Even with an homoskedastic noise, missingness generates heteroskedasticity
- **The uncertainty increases when missing values are associated with larger regression coefficients (i.e. the most predictive variables)**

What about splitting the data?

**Predictive uncertainty quantification with missing values**

Learning with Missing Data

**Conformal Prediction with Missing Values**

Missing Data Augmentation

Experimental Results

Conclusions

### Lemma

Assume  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^n$  are i.i.d. (or exchangeable).

Then, for any missing mechanism, for almost all imputation function<sup>1</sup>  $\phi$ :

$(\phi(X_{obs(M^{(k)})}^{(k)}, M^{(k)}), Y^{(k)})_{k=1}^n$  are **exchangeable**.

---

<sup>1</sup>Even if the imputation is not accurate, the guarantee will hold.



### Lemma

Assume  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^n$  are i.i.d. (or exchangeable).

Then, for any missing mechanism, for almost all imputation function<sup>1</sup>  $\phi$ :

$\left(\phi\left(X_{\text{obs}(M^{(k)})}^{(k)}, M^{(k)}\right), Y^{(k)}\right)_{k=1}^n$  are **exchangeable**.

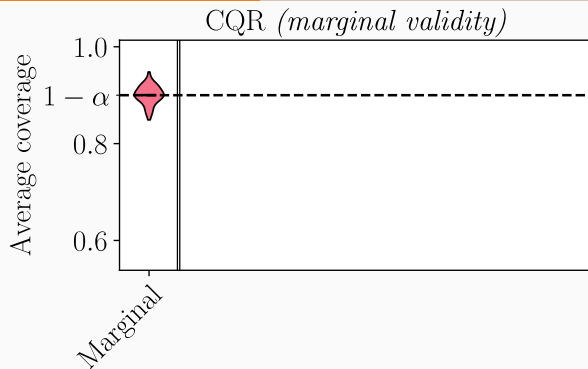
$\Rightarrow$  CQR, and Conformal Prediction, applied on an imputed data set still enjoys marginal guarantees:

$$\mathbb{P}\left\{Y^{(n+1)} \in \widehat{C}_\alpha\left(X^{(n+1)}, M^{(n+1)}\right)\right\} \geq 1 - \alpha.$$

---

<sup>1</sup>Even if the imputation is not accurate, the guarantee will hold.

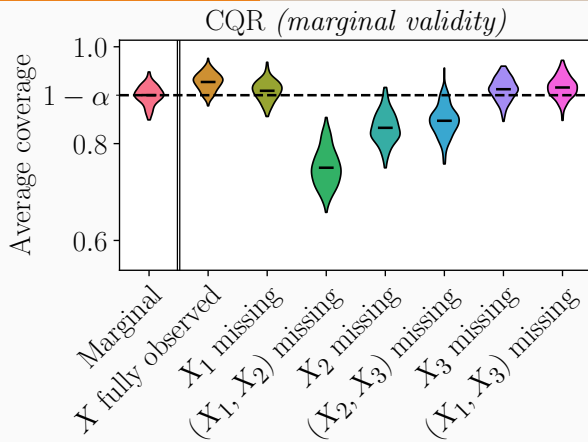
## CQR is marginally valid on imputed data sets



- Marginal (i.e. average) coverage is indeed recovered!

	QR	CQR	
(MV)	✗	✓	
(MCV)	✗		

## CQR is marginally valid on imputed data sets

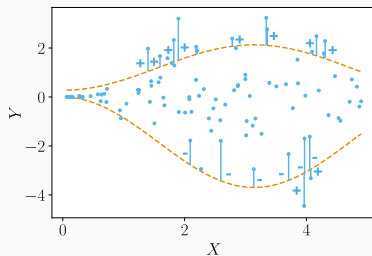


- Disparities between masks is not corrected by the conformalization step.

	QR	CQR
(MV)	✗	✓
(MCV)	✗	✗

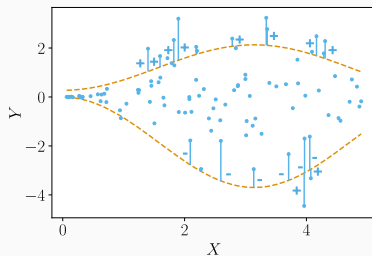
# Conformalization step is independent of the important variable: the mask!

**Observation:** the  $\alpha$ -correction term is computed among all the data points, regardless of their mask!



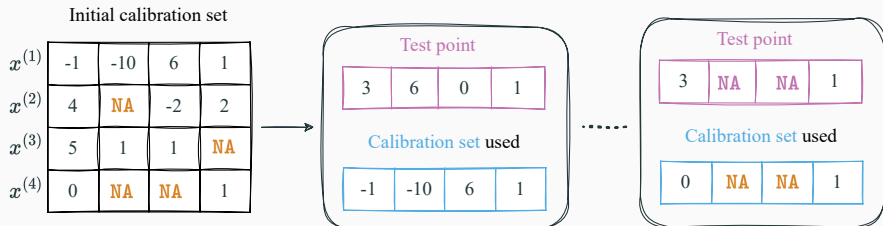
# Conformalization step is independent of the important variable: the mask!

**Observation:** the  $\alpha$ -correction term is computed among all the data points, regardless of their mask!



**Warning:**  $2^d$  possible masks

⇒ Splitting the calibration set by mask is infeasible (lack of data)!



What about splitting the data?

**Predictive uncertainty quantification with missing values**

Learning with Missing Data

Conformal Prediction with Missing Values

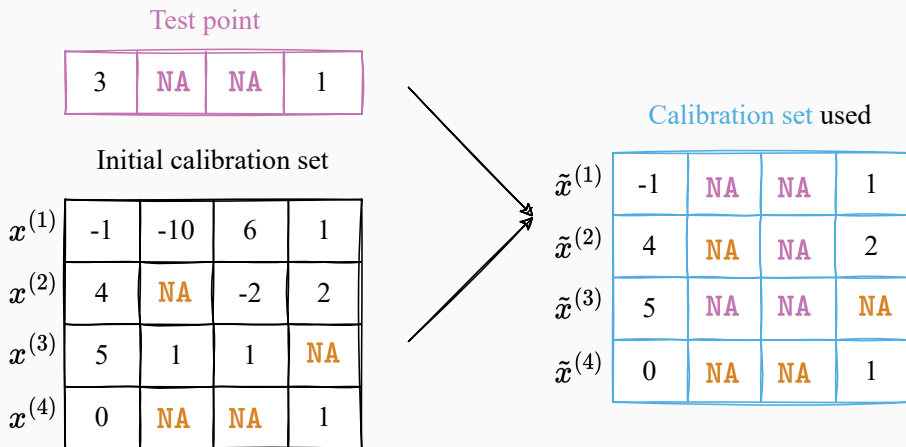
**Missing Data Augmentation**

Experimental Results

Conclusions

# Missing Data Augmentation (MDA)

**Idea:** for each **test point**, modify the **calibration points** to mimic the **test mask**



**Algorithms:** MDA with **Exact** masking or with **Nested** masking.

Test point

3	NA	NA	1
---	----	----	---

Initial calibration set

$x^{(1)}$	-1	-10	6	1
$x^{(2)}$	4	NA	-2	2
$x^{(3)}$	5	1	1	NA
$x^{(4)}$	0	NA	NA	1

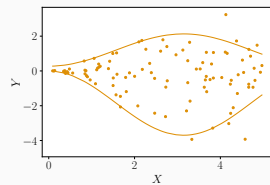
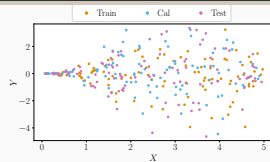
Calibration set used

$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$				
$\tilde{x}^{(4)}$	0	NA	NA	1



# CQR-MDA with exact masking in words

1. Split the training set into a **proper training set** and **calibration set**
2. Train the imputation function on the **proper training set**
3. Impute the **proper training set**
4. Train the quantile regressors on the imputed **proper training set**



## CQR-MDA with exact masking in words

1. Split the training set into a **proper training set** and **calibration set**
2. Train the imputation function on the proper training set
3. Impute the proper training set
4. Train the **quantile regressors** on the imputed proper training set
5. For a test point  $(X^{(n+1)}, M^{(n+1)})$ :

3	NA	NA	1
---	----	----	---











### Theorem (CP-MDA-Exact achieves MCV)

If the data is exchangeable and  $M \perp\!\!\!\perp (X, Y)$ , then for almost all imputation function CP-MDA-Exact is such that for any  $m \in \{0, 1\}^d$ :

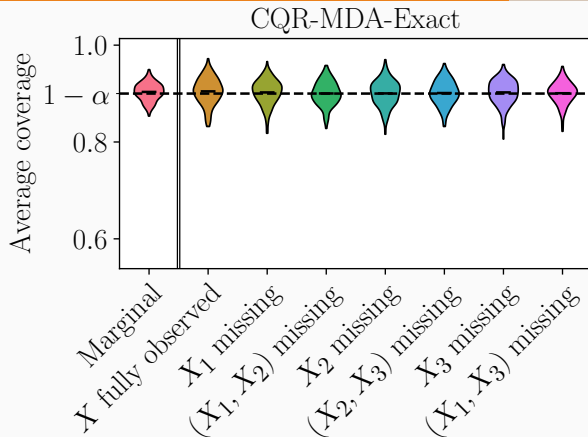
$$\mathbb{P}\left(Y \in \widehat{C}_\alpha(X, m) \mid M = m\right) \geq 1 - \alpha,$$

and if additionally the scores are almost surely distinct:

$$\mathbb{P}\left(Y \in \widehat{C}_\alpha(X, m) \mid M = m\right) \leq 1 - \alpha + \frac{1}{1 + \#\text{Cal}^m}.$$

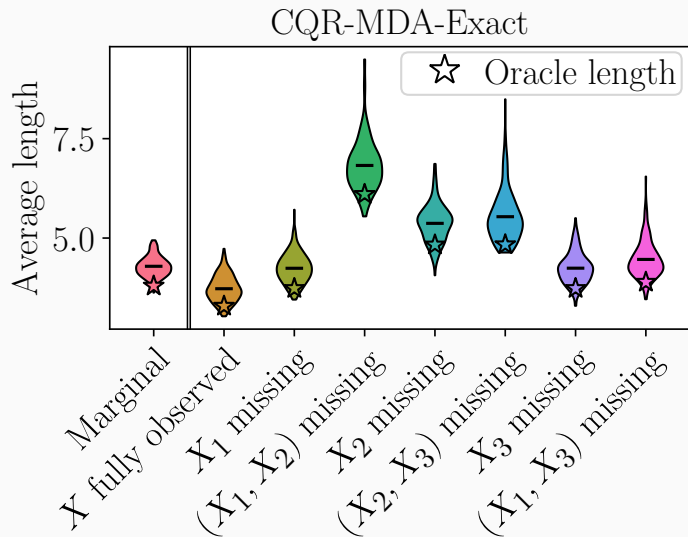


# MDA achieves Mask-Conditional-Validity (MCV), cont'd



	QR	CQR	CQR-MDA
(MV)	✗	✓	✓
(MCV)	✗	✗	✓

# MDA achieves Mask-Conditional-Validity in an informative way



What about splitting the data?

**Predictive uncertainty quantification with missing values**

Learning with Missing Data

Conformal Prediction with Missing Values

Missing Data Augmentation

**Experimental Results**

Conclusions

- Imputation by iterative ridge ( $\sim$  conditional expectation)

- Imputation by iterative ridge ( $\sim$  conditional expectation)
- **Concatenate the mask in the features**

- Imputation by iterative ridge ( $\sim$  conditional expectation)
- **Concatenate the mask in the features**
- Neural network, fitted to minimize the pinball loss

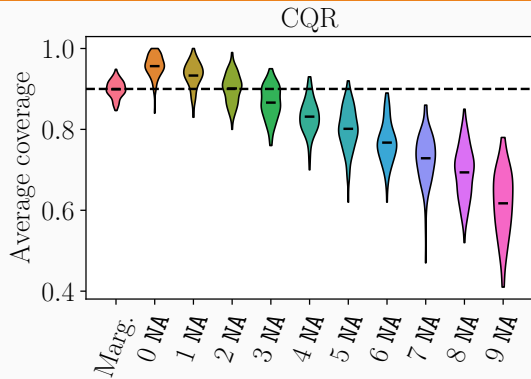
- Imputation by iterative ridge ( $\sim$  conditional expectation)
- **Concatenate the mask in the features**
- Neural network, fitted to minimize the pinball loss
- (Semi)-synthetic experiments:

- Imputation by iterative ridge ( $\sim$  conditional expectation)
- **Concatenate the mask in the features**
- Neural network, fitted to minimize the pinball loss
- (Semi)-synthetic experiments:
  - MCAR missing values, with probability 20%

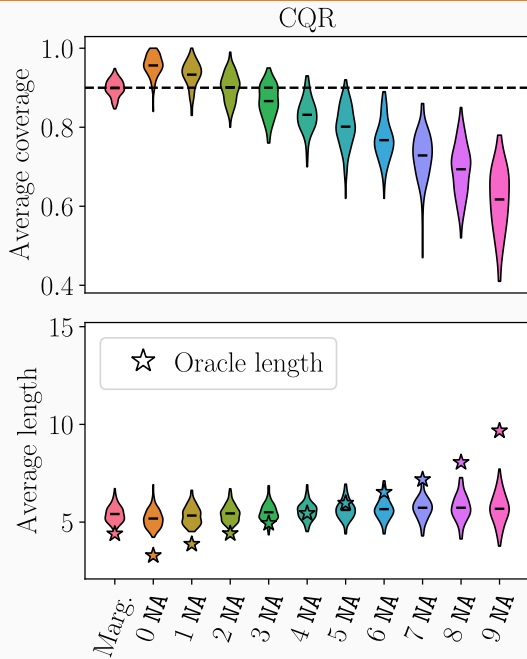


- Imputation by iterative ridge ( $\sim$  conditional expectation)
- **Concatenate the mask in the features**
- Neural network, fitted to minimize the pinball loss
- (Semi)-synthetic experiments:
  - MCAR missing values, with probability 20%
  - 100 repetitions

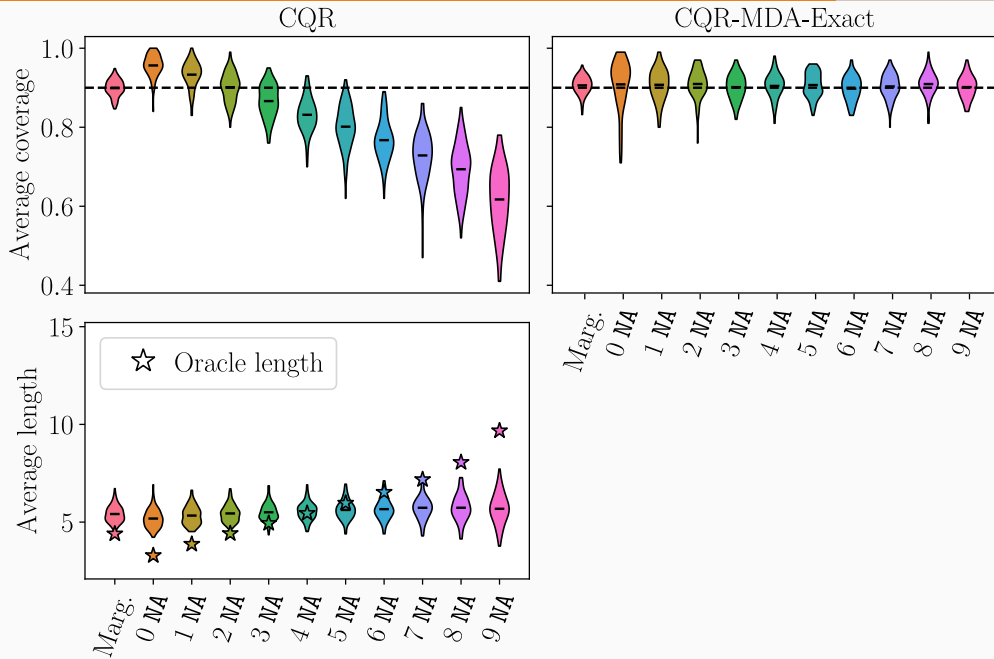
# Synthetic experiments (Gaussian linear model, $d = 10$ )



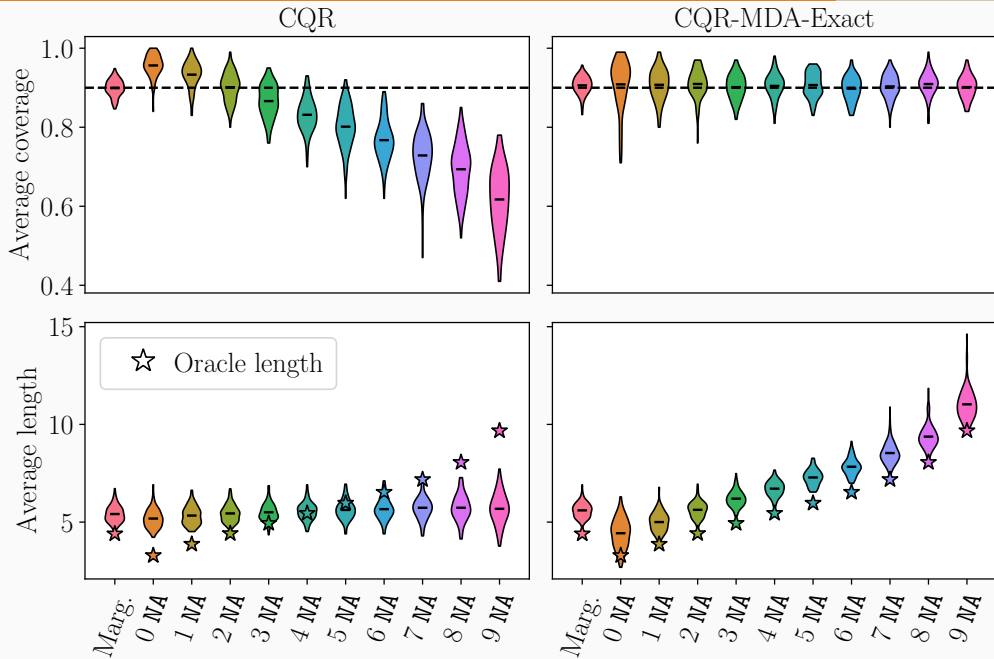
# Synthetic experiments (Gaussian linear model, $d = 10$ )



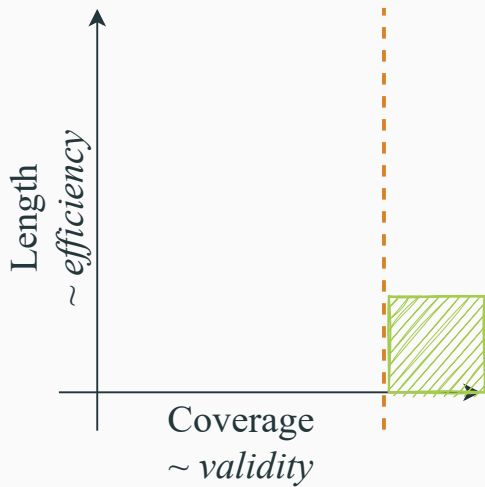
# Synthetic experiments (Gaussian linear model, $d = 10$ )



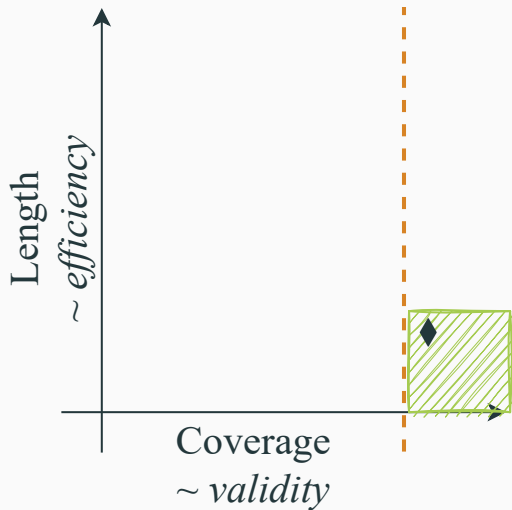
# Synthetic experiments (Gaussian linear model, $d = 10$ )



## Before more experiments, visualisation

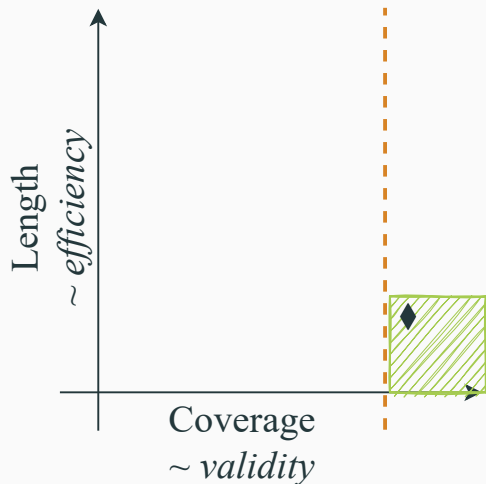


## Before more experiments, visualisation



◆ : marginal coverage, i.e.  
 $\mathbb{P}(Y \in \hat{C}_\alpha(X, M))$

## Before more experiments, visualisation



◆ : marginal coverage, i.e.

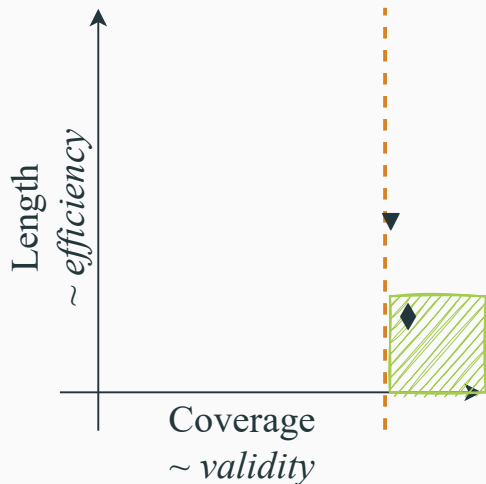
$$\mathbb{P}(Y \in \hat{C}_\alpha(X, M))$$

▼ : lowest coverage, i.e.

$$\min_{m \in \mathcal{M}} \mathbb{P}(Y \in \hat{C}_\alpha(X, m) | M = m)$$



## Before more experiments, visualisation



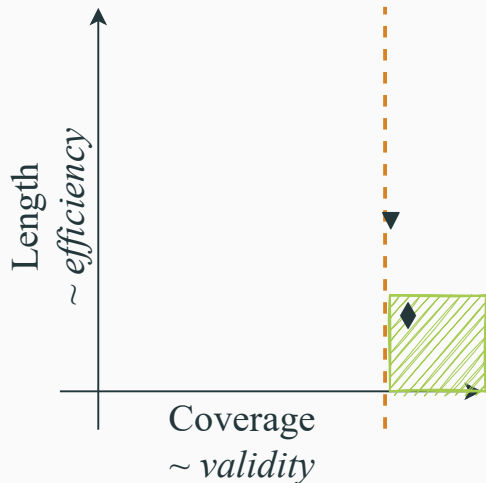
◆ : marginal coverage, i.e.

$$\mathbb{P}(Y \in \hat{C}_\alpha(X, M))$$

▼ : lowest coverage, i.e.

$$\min_{m \in \mathcal{M}} \mathbb{P}(Y \in \hat{C}_\alpha(X, m) | M = m)$$

## Before more experiments, visualisation



◆ : marginal coverage, i.e.

$$\mathbb{P}(Y \in \hat{C}_\alpha(X, M))$$

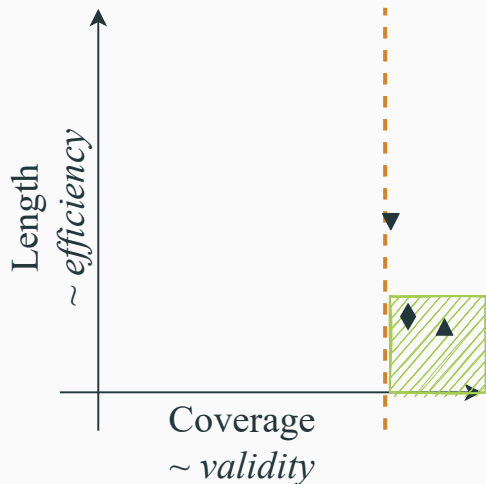
▼ : lowest coverage, i.e.

$$\min_{m \in \mathcal{M}} \mathbb{P}(Y \in \hat{C}_\alpha(X, m) | M = m)$$

▲ : highest coverage, i.e.

$$\max_{m \in \mathcal{M}} \mathbb{P}(Y \in \hat{C}_\alpha(X, m) | M = m)$$

## Before more experiments, visualisation



◆ : marginal coverage, i.e.

$$\mathbb{P}(Y \in \hat{C}_\alpha(X, M))$$

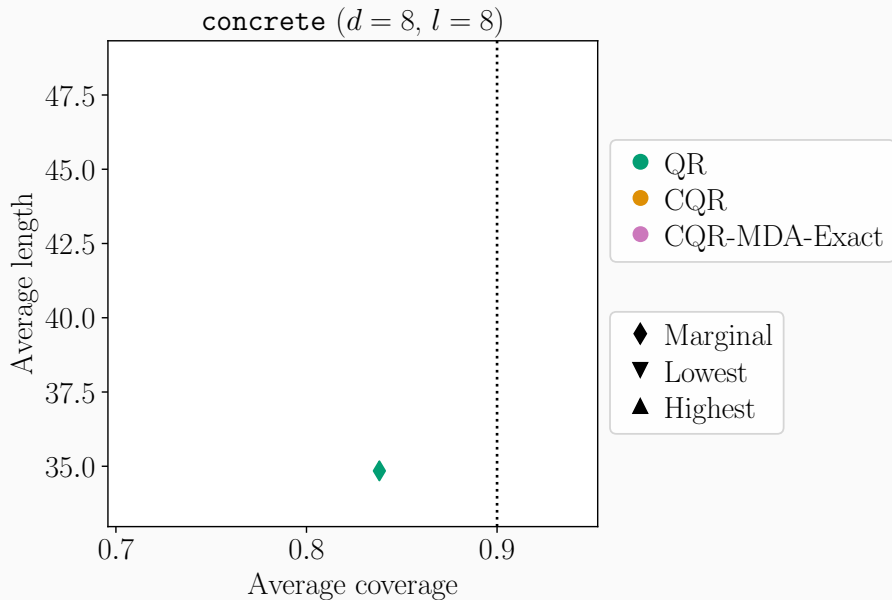
▼ : lowest coverage, i.e.

$$\min_{m \in \mathcal{M}} \mathbb{P}(Y \in \hat{C}_\alpha(X, m) | M = m)$$

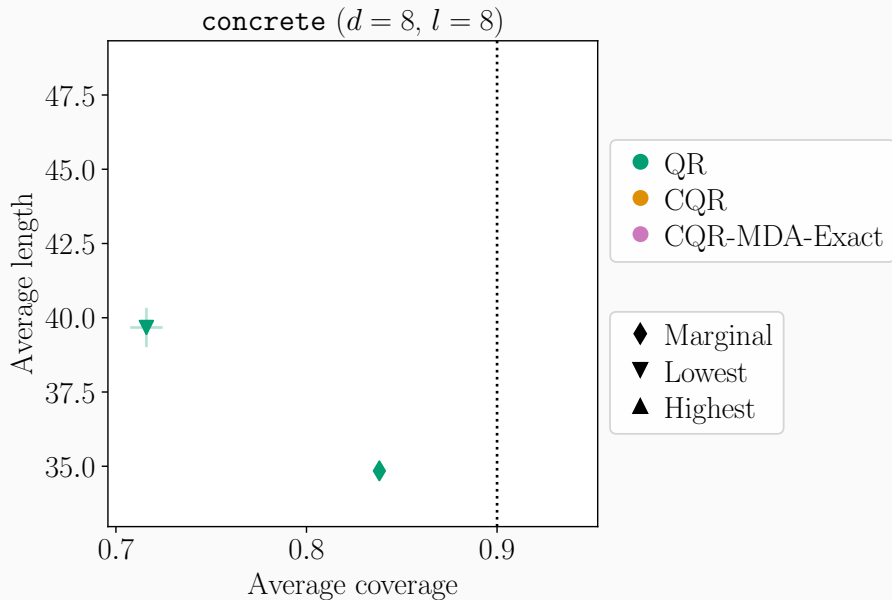
▲ : highest coverage, i.e.

$$\max_{m \in \mathcal{M}} \mathbb{P}(Y \in \hat{C}_\alpha(X, m) | M = m)$$

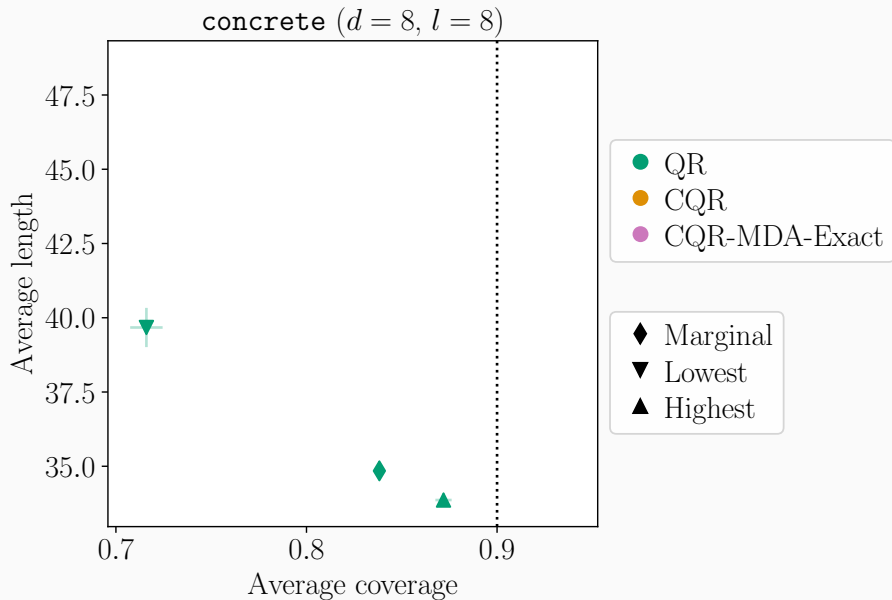
## Semi-synthetic experiments



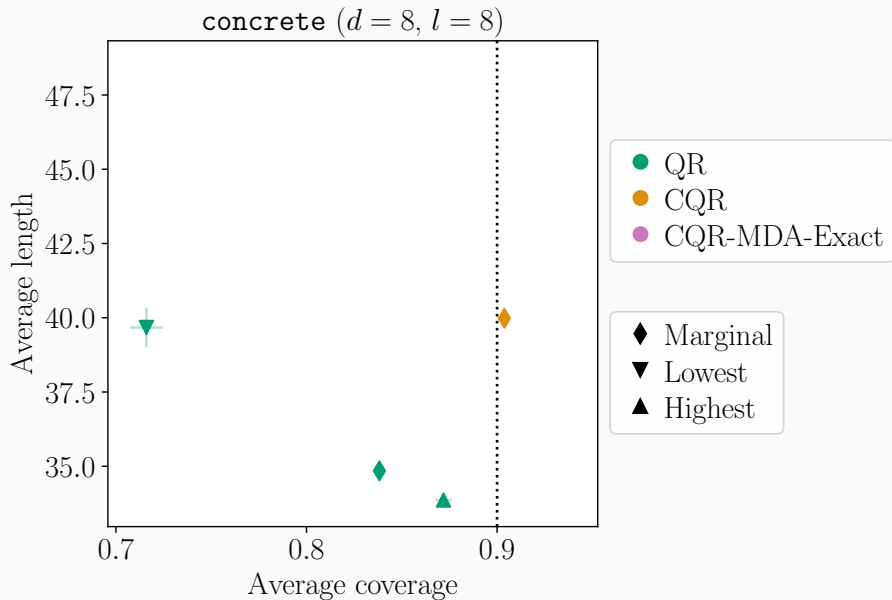
## Semi-synthetic experiments



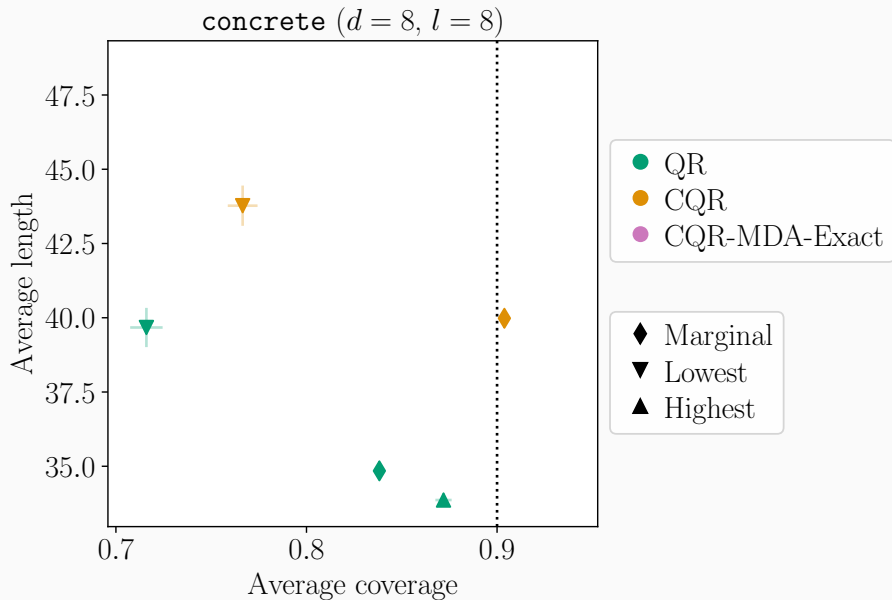
## Semi-synthetic experiments



## Semi-synthetic experiments

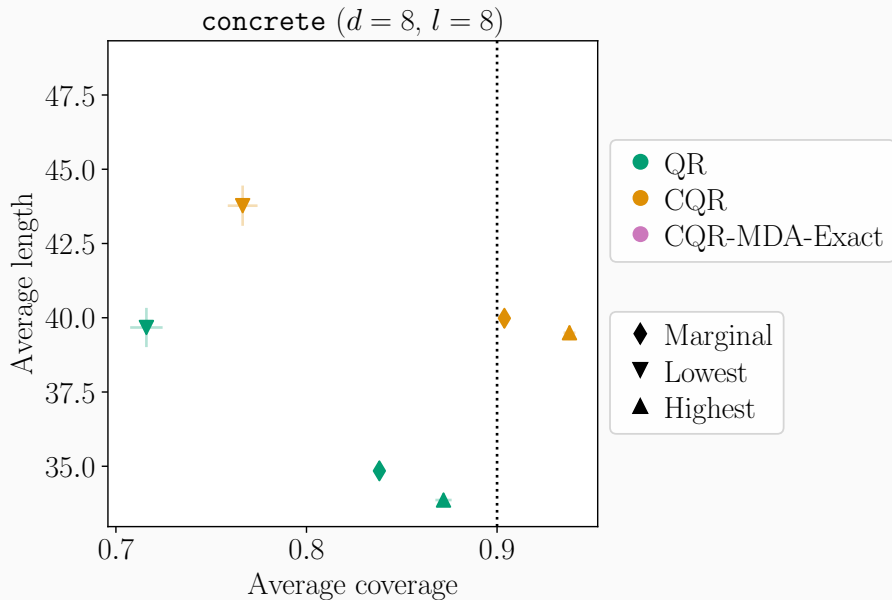


## Semi-synthetic experiments

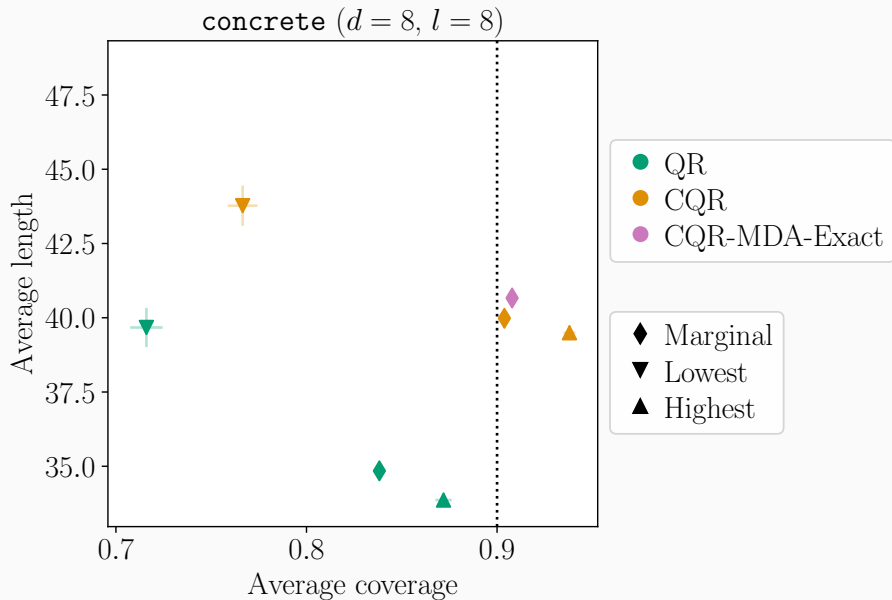




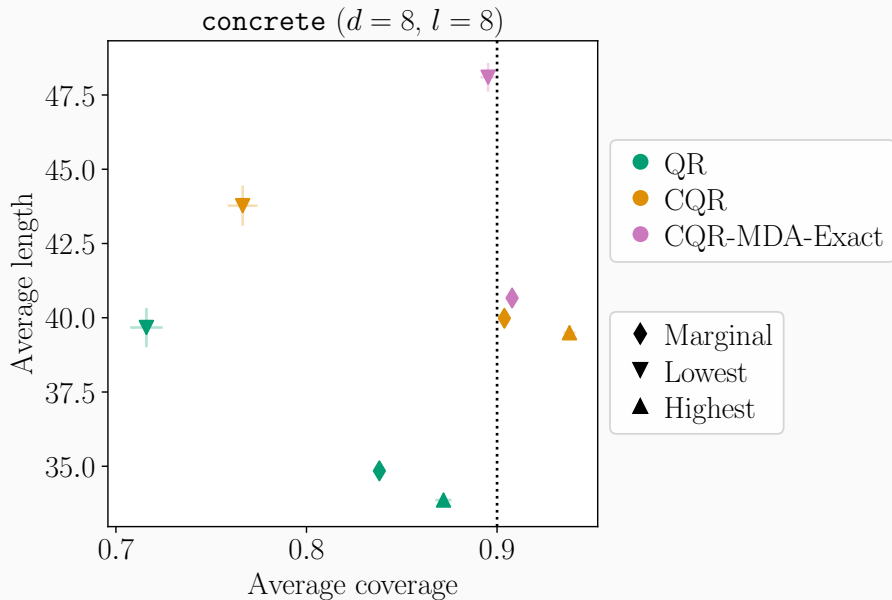
## Semi-synthetic experiments



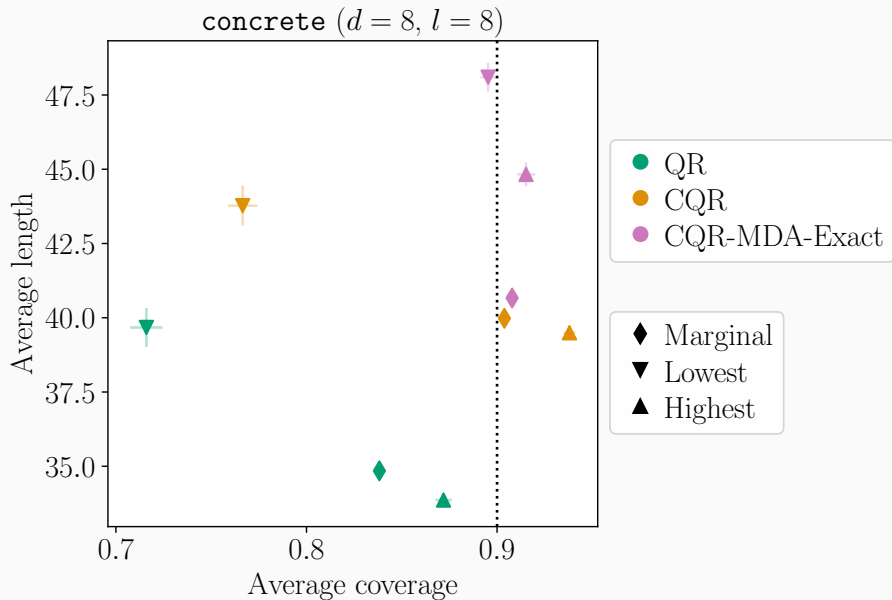
## Semi-synthetic experiments



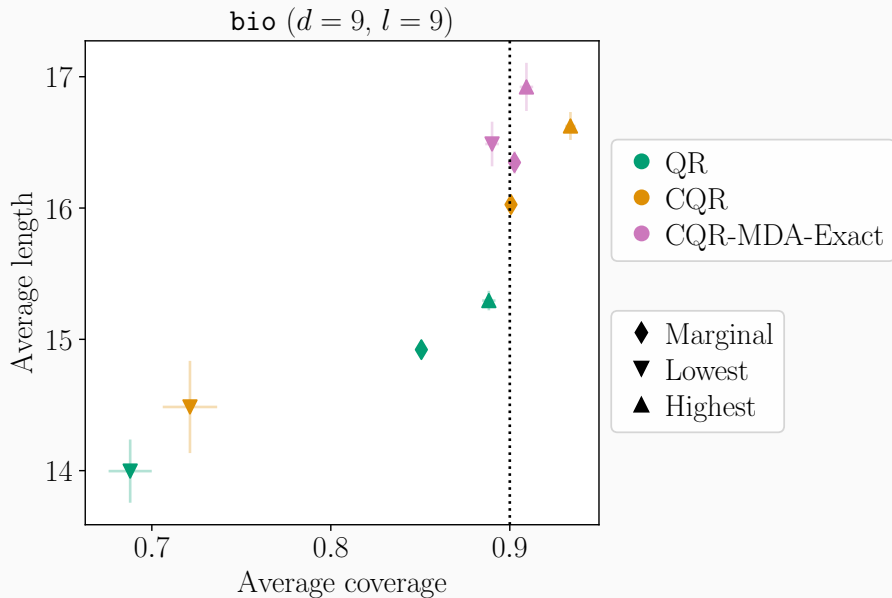
## Semi-synthetic experiments



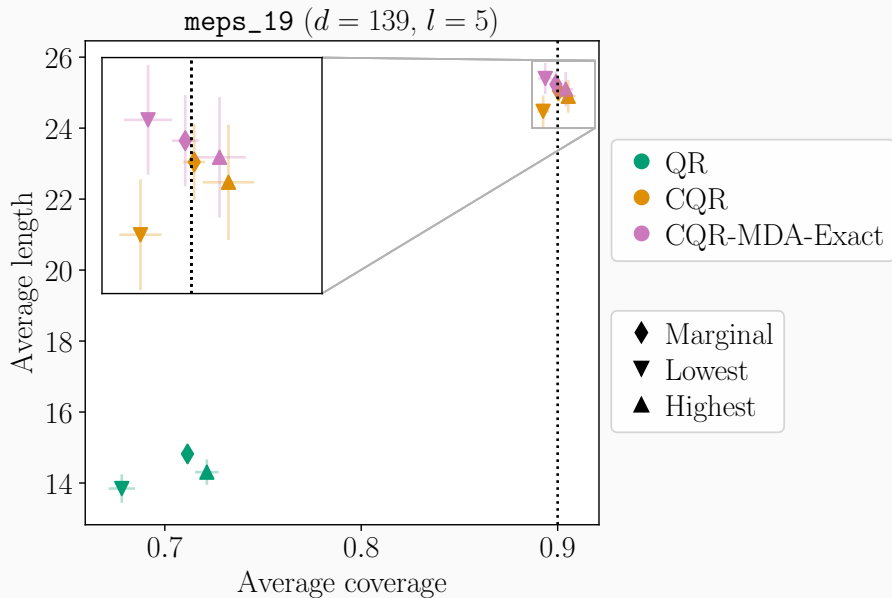
## Semi-synthetic experiments



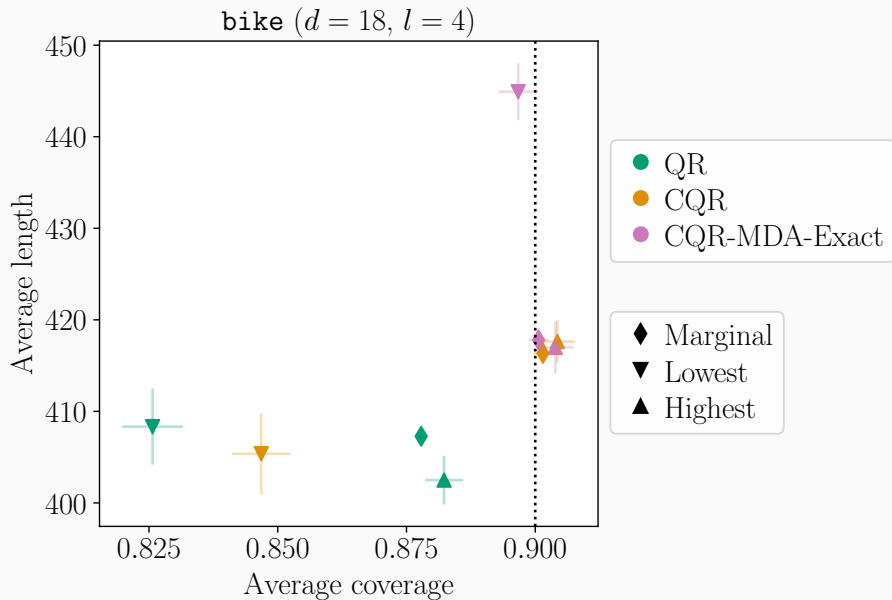
## Semi-synthetic experiments



## Semi-synthetic experiments



## Semi-synthetic experiments



- 30 hospitals
- More than 30 000 trauma patients
- 4 000 new patients per year
- 250 continuous and categorical variables  
↳ Many useful statistical tasks



- 30 hospitals
- More than 30 000 trauma patients
- 4 000 new patients per year
- 250 continuous and categorical variables  
↪ Many useful statistical tasks

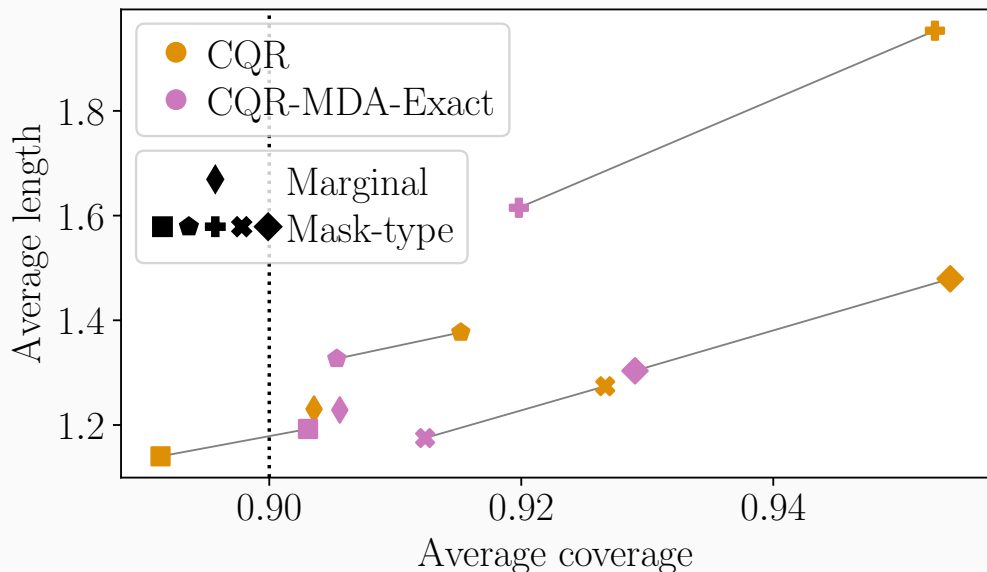
Predict the level of platelets upon arrival at hospital, given 7 covariates chosen by medical doctors.

- 30 hospitals
- More than 30 000 trauma patients
- 4 000 new patients per year
- 250 continuous and categorical variables  
↳ Many useful statistical tasks

Predict the level of platelets upon arrival at hospital, given 7 covariates chosen by medical doctors.

These covariates are not always observed: from 0% to 24% of missing values by features, with a total average of 7%.

# Real data experiment: TraumaBase<sup>®</sup>, critical care medicine



What about splitting the data?

**Predictive uncertainty quantification with missing values**

Learning with Missing Data

Conformal Prediction with Missing Values

Missing Data Augmentation

Experimental Results

**Conclusions**

- Consistency of universal quantile learner when chained with almost any imputation function.
- CP-MDA-Nested, an algorithm which does not discard any calibration point.

Paper →

Poster →

Code →



- CP marginal guarantees hold on the imputed data set.

- CP marginal guarantees hold on the imputed data set.
- Missingness introduces additional heteroskedasticity, creating a need for quantile regression based methods.

- CP marginal guarantees hold on the imputed data set.
- Missingness introduces additional heteroskedasticity, creating a need for quantile regression based methods.
- CQR fails to attain coverage conditional on the missing pattern.



- CP marginal guarantees hold on the imputed data set.
- Missingness introduces additional heteroskedasticity, creating a need for quantile regression based methods.
- CQR fails to attain coverage conditional on the missing pattern.
- **Missing data augmentation is the first method to output predictive intervals with missing values.**

- CP marginal guarantees hold on the imputed data set.
- Missingness introduces additional heteroskedasticity, creating a need for quantile regression based methods.
- CQR fails to attain coverage conditional on the missing pattern.
- **Missing data augmentation is the first method to output predictive intervals with missing values.**
- Missing data augmentation attains conditional coverage with respect to the missing pattern (in MCAR setting).

Thanks for listening! Any question? :)

## References

- 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*.
- Papadopoulos, H., Proedrou, K., Vovk, V., and Gammerman, A. (2002). Inductive Confidence Machines for Regression. In *Machine Learning: ECML 2002*. Springer.
- Romano, Y., Patterson, E., and Candès, E. (2019). Conformalized Quantile Regression. *NeurIPS*.
- Vovk, V., Gammerman, A., and Shafer, G. (2005). *Algorithmic Learning in a Random World*. Springer US.