

Conformal Prediction with Missing Values

Aymeric Dieuleveut

joint work with Margaux Zaffran, Julie Josse, Yaniv Romano

7e journée de Statistique Mathématique

January 18, 2024





Margaux Zaffrant
École Polytechnique
Inria
Paris - France



Julie Josse
PreMeDICaL
INRIA
Montpellier - France



Yaniv Romano
Technion - Israel Institute
of Technology
Haifa - Israel

Introduction to missing values

Quantifying predictive uncertainty with missing values

Conclusion

- 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 blood platelets upon arrival at hospital, given 7 pre-hospital features.

- 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 blood platelets upon arrival at hospital, given 7 pre-hospital features.

These covariates are not always observed.

Missing values: ubiquitous in data science practice

Data: $(X^{(k)}, Y^{(k)})_{k=1}^n \in (\mathbb{R}^d \times \mathbb{R})^n$

Y	X_1	X_2	X_3	X_4	X_5	X_6
8.26	0.72	0.18	0.55	0.05	0.73	0.50
19.41	0.60	0.58	NA	NA	NA	0.40
19.75	0.54	0.43	0.96	0.77	0.06	0.66
7.32	NA	0.19	NA	0.02	0.83	0.04
13.55	0.65	0.69	0.50	0.15	NA	0.87
20.75	0.43	0.74	0.61	0.72	0.52	0.35
9.26	0.89	NA	0.84	0.01	0.73	NA
9.68	0.963	0.45	0.65	0.04	0.06	NA

Missing values: ubiquitous in data science practice

Data: $(X^{(k)}, Y^{(k)})_{k=1}^n \in (\mathbb{R}^d \times \mathbb{R})^n$

Y	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆
8.26	0.72	0.18	0.55	0.05	0.73	0.50
19.41	0.60	0.58	NA	NA	NA	0.40
19.75	0.54	0.43	0.96	0.77	0.06	0.66
7.32	NA	0.19	NA	0.02	0.83	0.04
13.55	0.65	0.69	0.50	0.15	NA	0.87
20.75	0.43	0.74	0.61	0.72	0.52	0.35
9.26	0.89	NA	0.84	0.01	0.73	NA
9.68	0.963	0.45	0.65	0.04	0.06	NA

Missing values: ubiquitous in data science practice

Data: $(X^{(k)}, Y^{(k)})_{k=1}^n \in (\mathbb{R}^d \times \mathbb{R})^n$

Y	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆
8.26	0.72	0.18	0.55	0.05	0.73	0.50
19.41	0.60	0.58	NA	NA	NA	0.40
19.75	0.54	0.43	0.96	0.77	0.06	0.66
7.32	NA	0.19	NA	0.02	0.83	0.04
13.55	0.65	0.69	0.50	0.15	NA	0.87
20.75	0.43	0.74	0.61	0.72	0.52	0.35
9.26	0.89	NA	0.84	0.01	0.73	NA
9.68	0.963	0.45	0.65	0.04	0.06	NA

If each entry has a probability 0.01 of being missing:

$$d = 6 \rightarrow \approx 94\% \text{ of rows kept}$$

$$d = 300 \rightarrow \approx 5\% \text{ of rows kept}$$

*One of the ironies of Big Data is that missing data play an ever more significant role.*¹

¹Zhu et al. (2019), *High-dimensional PCA with heterogeneous missingness*, JRSS B

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe (NA, 6, 2). Then $m = (1, 0, 0)$ and $X_{\text{obs}(m)} = (6, 2)$.

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, 2)$. Then $m = (0, 1, 0)$ and $X_{\text{obs}(m)} = (-1, 2)$.

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, \text{NA})$. Then $m = (0, 1, 1)$ and $X_{\text{obs}(m)} = (-1)$.

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, \text{NA})$. Then $m = (0, 1, 1)$ and $X_{\text{obs}(m)} = (-1)$.

There are 2^d **patterns** (statistical and computational challenges).

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, \text{NA})$. Then $m = (0, 1, 1)$ and $X_{\text{obs}(m)} = (-1)$.

There are 2^d **patterns** (statistical and computational challenges).

- Three **mechanisms**² can generate missing values.

²Rubin (1976), *Inference and missing data*, Biometrika

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, \text{NA})$. Then $m = (0, 1, 1)$ and $X_{\text{obs}(m)} = (-1)$.

There are 2^d **patterns** (statistical and computational challenges).

- Three **mechanisms**² can generate missing values.
 - ↪ **Missing Completely At Random** (MCAR): $\mathbb{P}(M = m|X) = \mathbb{P}(M = m)$
for all $m \in \{0, 1\}^d$.

²Rubin (1976), *Inference and missing data*, Biometrika

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, \text{NA})$. Then $m = (0, 1, 1)$ and $X_{\text{obs}(m)} = (-1)$.

There are 2^d **patterns** (statistical and computational challenges).

- Three **mechanisms**² can generate missing values.
 - ↪ **Missing Completely At Random** (MCAR): $\mathbb{P}(M = m|X) = \mathbb{P}(M = m)$
for all $m \in \{0, 1\}^d$. $M \perp\!\!\!\perp X$, missingness does not depend on the variables.

²Rubin (1976), *Inference and missing data*, Biometrika

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, \text{NA})$. Then $m = (0, 1, 1)$ and $X_{\text{obs}(m)} = (-1)$.

There are 2^d **patterns** (statistical and computational challenges).

- Three **mechanisms**² can generate missing values.
 - ↪ Missing Completely At Random (MCAR)
 - ↪ **Missing At Random (MAR)**

²Rubin (1976), *Inference and missing data*, Biometrika

Handling missing values depends on pattern and mechanism

- $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ random variables.
- $M \in \{0, 1\}^d$ is defined as $M_j = 1 \Leftrightarrow X_j$ is missing.
 M is called the **mask** or the **missing pattern**.

Example

We observe $(-1, \text{NA}, \text{NA})$. Then $m = (0, 1, 1)$ and $X_{\text{obs}(m)} = (-1)$.

There are 2^d **patterns** (statistical and computational challenges).

- Three **mechanisms**² can generate missing values.
 - ↪ Missing Completely At Random (MCAR)
 - ↪ Missing At Random (MAR)
 - ↪ **Missing Non At Random (MNAR)**

²Rubin (1976), *Inference and missing data*, Biometrika

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** ϕ (e.g. the mean).

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



$u^{(1)}$	-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** ϕ (e.g. the mean).
2. Train your algorithm (Random Forest, Neural Nets, etc.) on the imputed

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

Supervised learning with missing values: impute-then-regress

Impute-then-regress procedures are widely used.

1. Replace NA using an imputation function ϕ (e.g. the mean).
2. Train your algorithm (Random Forest, Neural Nets, etc.) on the imputed

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

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

Supervised learning with missing values: impute-then-regress

Impute-then-regress procedures are widely used.

1. Replace NA using an imputation function ϕ (e.g. the mean).
2. Train your algorithm (Random Forest, Neural Nets, etc.) on the imputed

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

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

- ✓ Le Morvan et al. (2021)³ show that for **any deterministic imputation** and **universal learner** this procedure is **Bayes-consistent**.

³ Le Morvan, Josse, Scornet & Varoquaux (2021), *What's a good imputation to predict with missing values?*, NeurIPS

Supervised learning with missing values: impute-then-regress

Impute-then-regress procedures are widely used.

1. Replace NA using an imputation function ϕ (e.g. the mean).
2. Train your algorithm (Random Forest, Neural Nets, etc.) on the imputed

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

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

- ✓ Le Morvan et al. (2021)³ show that for any deterministic imputation and universal learner this procedure is Bayes-consistent.
- ✗ Ayme et al. (2022)⁴ show that even for very **simple distributions** (linear model, Gaussian noise), this rate of convergence may suffer from **curse of dimensionality**.

³ Le Morvan, Josse, Scornet & Varoquaux (2021), *What's a good imputation to predict with missing values?*, NeurIPS

⁴ Ayme, Boyer, Dieuleveut & Scornet (2022), *Near-optimal rate of consistency for linear models with missing values*, ICML

Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

Impute-then-Regress+Conformalization

Missing Data Augmentation

Experimental results

Conclusion

Predictive uncertainty quantification with missing values

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

Predictive uncertainty quantification with missing values

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

1. Marginal Validity (MV)

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

Predictive uncertainty quantification with missing values

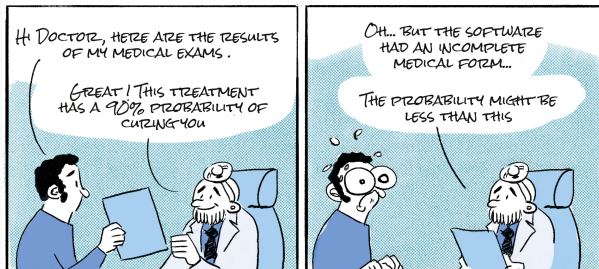
Goal: predict $Y^{(n+1)}$ with **confidence** $1 - \alpha$, i.e. build the smallest C_α such that:

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 @theo.reminger

Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

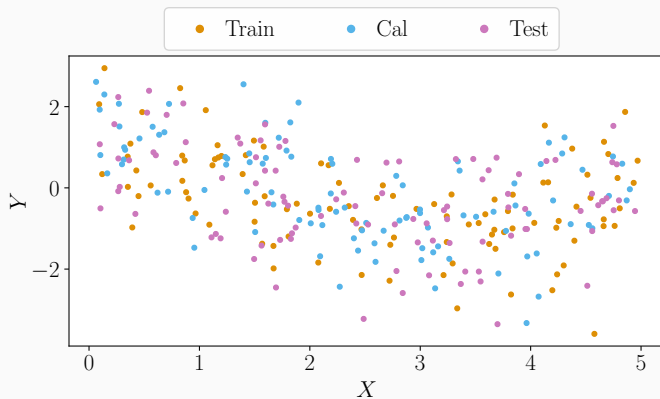
Impute-then-Regress+Conformalization

Missing Data Augmentation

Experimental results

Conclusion

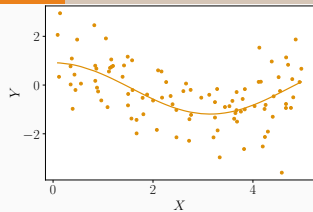
Split Conformal Prediction (Vovk et al., 2005): scheme



Randomly split the data to obtain a **proper training set** and a **calibration set**. Keep the **test set**.

Split Conformal Prediction (Vovk et al., 2005): scheme

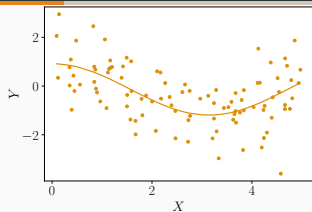
1)



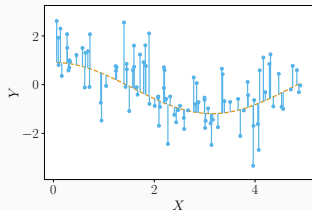
► Learn $\hat{\mu}$.

Split Conformal Prediction (Vovk et al., 2005): scheme

1)



2)



► Learn $\hat{\mu}$.

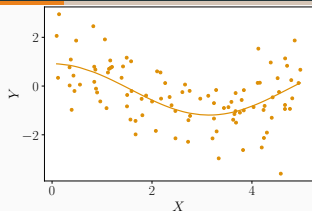
► Predict with $\hat{\mu}$.

► Get the residuals $\hat{\epsilon}_i$ and form the set of scores $\mathcal{S} = \{|\hat{\epsilon}_i|, i \in \text{Cal}\} \cup \{+\infty\}$.

► Get their $(1 - \alpha)$ empirical quantile: $Q_{1-\alpha}(\mathcal{S})$.

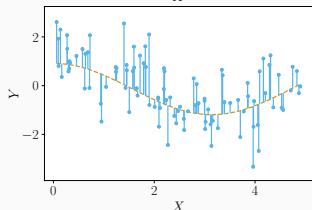
Split Conformal Prediction (Vovk et al., 2005): scheme

1)



► Learn $\hat{\mu}$.

2)

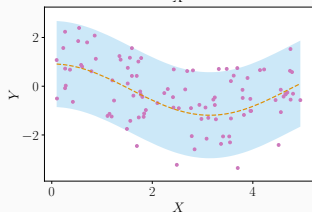


► Predict with $\hat{\mu}$.

► Get the residuals $\hat{\epsilon}_i$ and form the set of scores $\mathcal{S} = \{|\hat{\epsilon}_i|, i \in \text{Cal}\} \cup \{+\infty\}$.

► Get their $(1 - \alpha)$ empirical quantile: $Q_{1-\alpha}(\mathcal{S})$.

3)

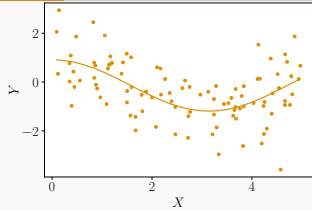


► Predict with $\hat{\mu}$.

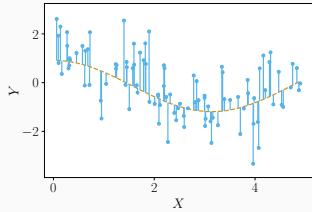
► Build $\hat{C}_\alpha(x)$: $[\hat{\mu}(x) \pm Q_{1-\alpha}(\mathcal{S})]$.

Split Conformal Prediction (Vovk et al., 2005): scheme

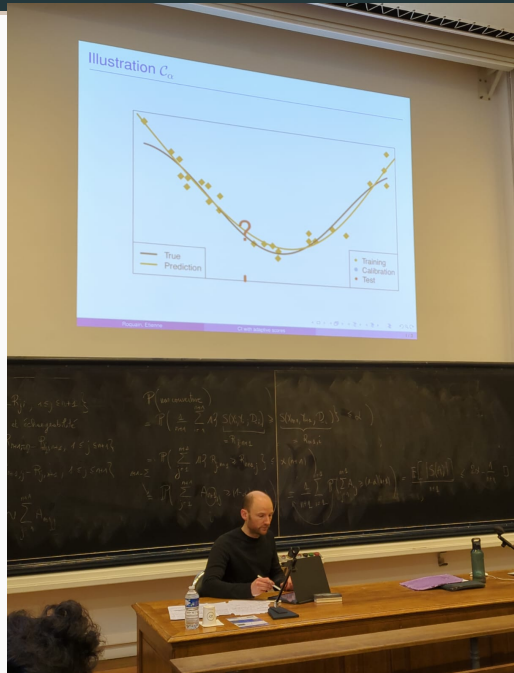
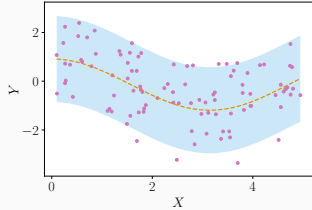
1)



2)



3)



of

Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

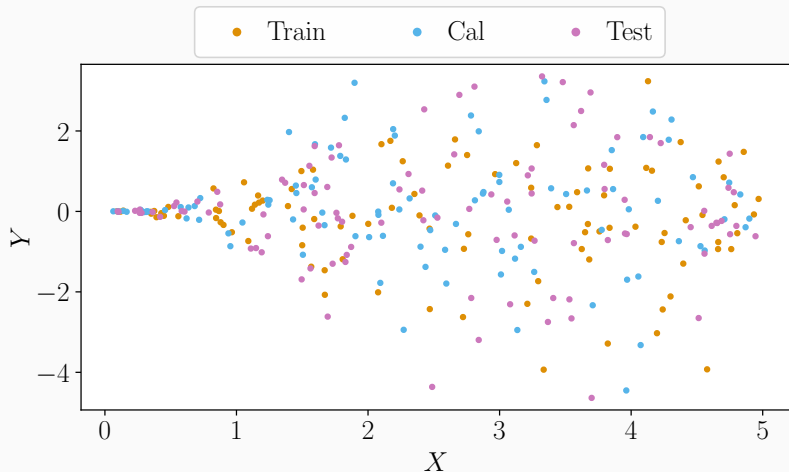
Impute-then-Regress+Conformalization

Missing Data Augmentation

Experimental results

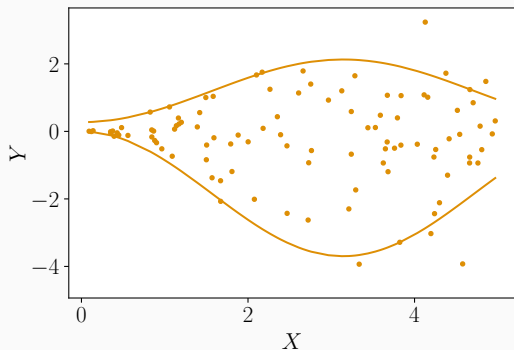
Conclusion

Conformalized Quantile Regression (CQR)⁴: toy example



⁴Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

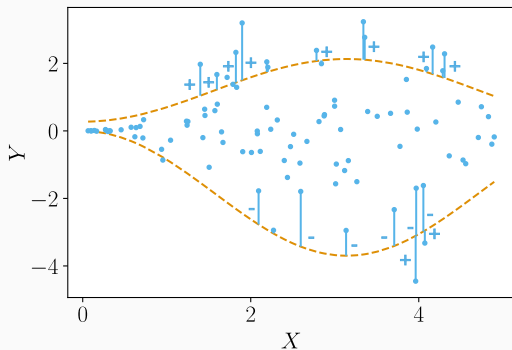
Conformalized Quantile Regression (CQR)⁴: training step



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

⁴Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

Conformalized Quantile Regression (CQR)⁴: calibration step

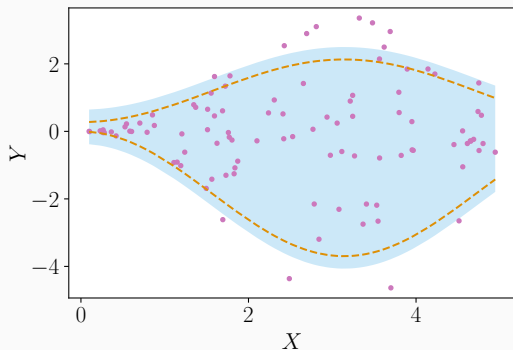


- ▶ Predict with $\widehat{QR}_{\text{lower}}$ and $\widehat{QR}_{\text{upper}}$
- ▶ Get the scores $\mathcal{S} = \{S^{(k)}\}_{\text{Cal}} \cup \{+\infty\}$
- ▶ Compute the $(1 - \alpha)$ empirical quantile of \mathcal{S} , noted $q_{1-\alpha}(\mathcal{S})$

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

⁴Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

Conformalized Quantile Regression (CQR)⁴: prediction step



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

► Build

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

⁴Romano et al. (2019), *Conformalized Quantile Regression*, NeurIPS

CQR enjoys finite sample guarantees proved in Romano et al. (2019), as a particular case of Conformal Prediction (CP).

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 enjoys finite sample guarantees proved in Romano et al. (2019), as a particular case of Conformal Prediction (CP).

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*:

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

- ✓ Distribution-free, only requires exchangeability
- ✓ Any quantile regression algorithm (neural nets, random forest...)
- ✓ Finite sample

CQR: theoretical guarantees

CQR enjoys finite sample guarantees proved in Romano et al. (2019), as a particular case of Conformal Prediction (CP).

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*:

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

- ✓ Distribution-free, only requires exchangeability
- ✓ Any quantile regression algorithm (neural nets, random forest...)
- ✓ Finite sample

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

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

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

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

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

SCP is defined by the conformity score function



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 the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#\text{Cal} + 1$ **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: $\mathbf{s}(\hat{A}(X_i), Y_i) := \max\left(\widehat{\text{QR}}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{\text{QR}}_{\text{upper}}(X_i)\right)$ 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 $\#\text{Tr}$) and a **calibration set** (size $\#\text{Cal}$)
2. Get \hat{A} by *training the algorithm \mathcal{A} on the **proper training set***
3. On the **calibration set**, obtain $\#\text{Cal} + 1$ **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: $\mathbf{s}(\hat{A}(X_i), Y_i) := \max\left(\widehat{QR}_{\text{lower}}(X_i) - Y_i, Y_i - \widehat{QR}_{\text{upper}}(X_i)\right)$ 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

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

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

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

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**

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

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

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

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

\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

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are *exchangeable*⁵. 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

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

Proof: application of the quantile lemma.

⁵Only the calibration and test data need to be exchangeable.

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

Theorem

Suppose $(X_i, Y_i)_{i=1}^{n+1}$ are *exchangeable*⁵. 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

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

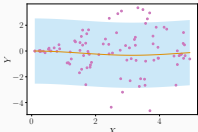
Proof: application of the quantile lemma.

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

⁵Only the calibration and test data need to be exchangeable.

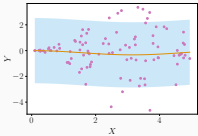
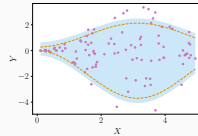
SCP: what choices for the regression scores?

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

	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		

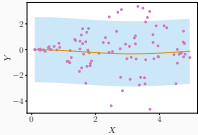
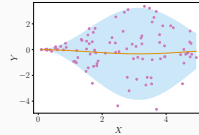
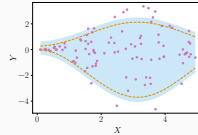
SCP: what choices for the regression scores?

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

	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

Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

Impute-then-Regress+Conformalization

Missing Data Augmentation

Experimental results

Conclusion

CP is marginally valid (MV) after imputation

To apply conformal prediction we need **exchangeable** data.

Lemma (Z. et al. (2023a))

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 function⁶ ϕ :

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

⁶Even if the imputation is not accurate, the guarantee will hold.

CP is marginally valid (MV) after imputation

To apply conformal prediction we need **exchangeable** data.

Lemma (Z. et al. (2023a))

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 function⁶ ϕ :

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

\Rightarrow CQR, and Conformal Prediction, applied on an imputed data set still enjoys marginal guarantees⁷:

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

⁶Even if the imputation is not accurate, the guarantee will hold.

⁷The upper bound also holds under continuously distributed scores.

CP is marginally valid (MV) after imputation

To apply conformal prediction we need **exchangeable** data.

Lemma (Z. et al. (2023a))

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 function⁶ ϕ :

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

\Rightarrow CQR, and Conformal Prediction, applied on an imputed data set still enjoys marginal guarantees⁷:

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

⁶Even if the imputation is not accurate, the guarantee will hold.

⁷The upper bound also holds under continuously distributed scores.

CP is marginally valid (MV) after imputation

To apply conformal prediction we need **exchangeable** data.

Lemma (Z. et al. (2023a))

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*⁶ ϕ :

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

\Rightarrow CQR, and Conformal Prediction, applied on an imputed data set still enjoys marginal guarantees⁷:

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

⁶Even if the imputation is not accurate, the guarantee will hold.

⁷The upper bound also holds under continuously distributed scores.

CQR is marginally valid on imputed data sets

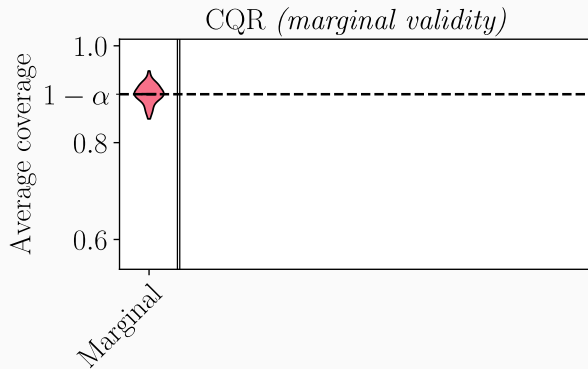
$$Y = \beta^T X + \varepsilon,$$

$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.

CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

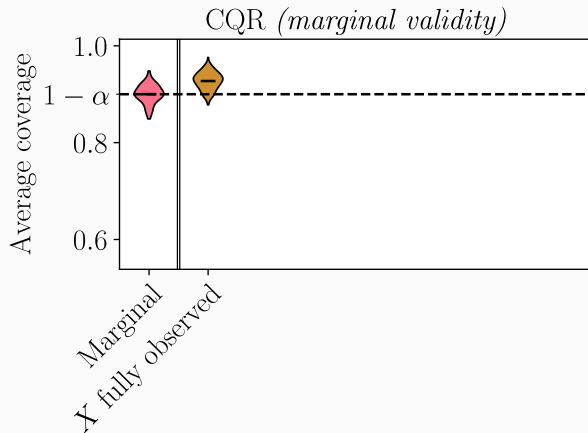
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

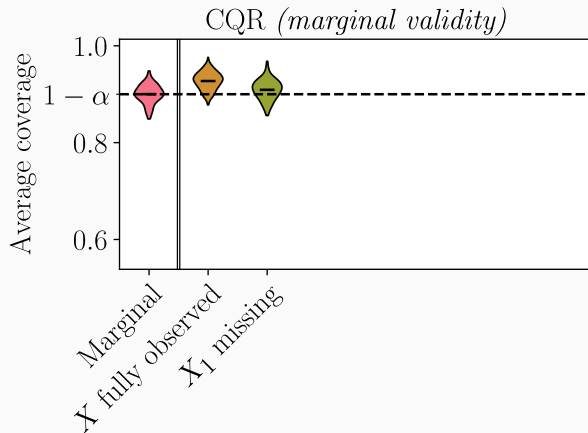
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

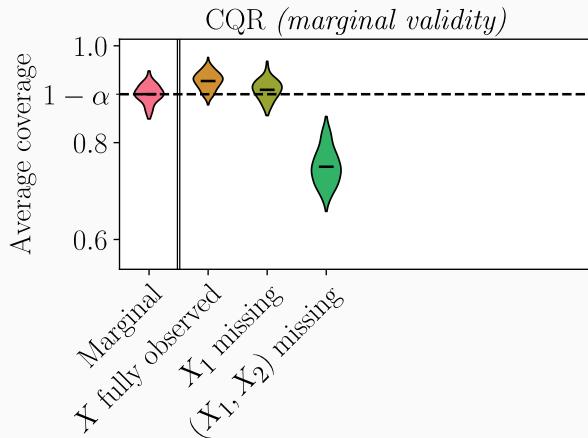
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

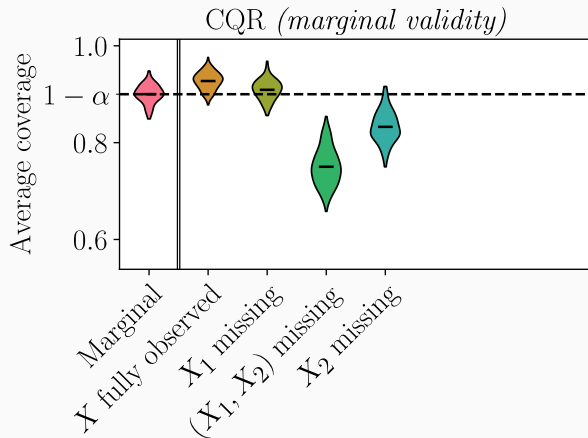
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

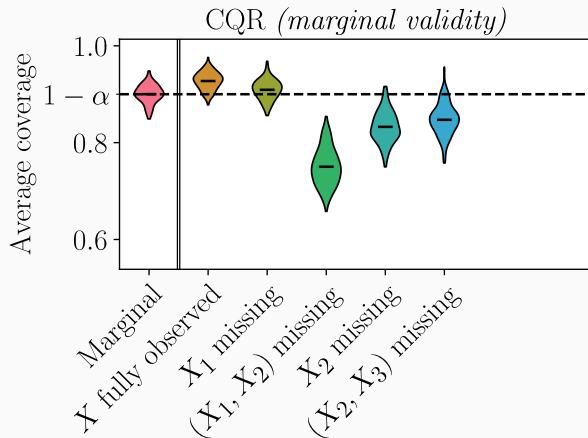
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

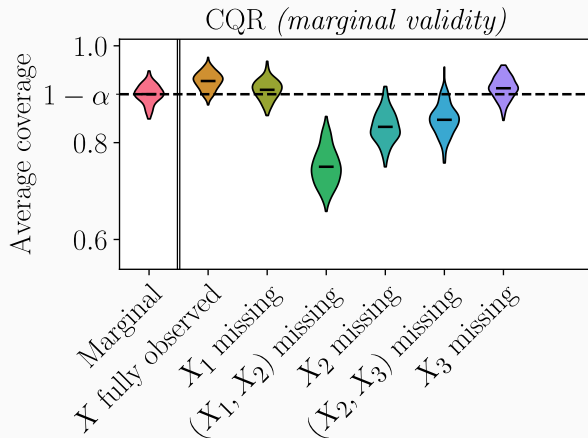
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

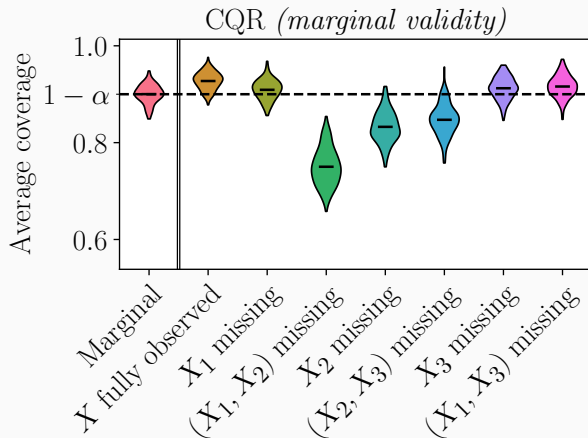
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

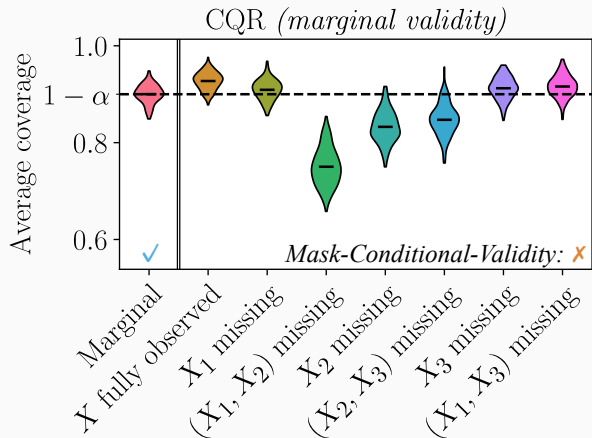
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



CQR is marginally valid on imputed data sets

$$Y = \beta^T X + \varepsilon,$$

$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



Warning: the predictive intervals cover properly **marginally**, but suffer from high **disparities depending on the missing patterns**.

Gaussian linear model

- $Y = \beta^T X + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2) \perp (X, M)$, $\beta \in \mathbb{R}^d$.
- for all $m \in \{0, 1\}^d$, there exist μ^m and Σ^m such that $X|(M = m) \sim \mathcal{N}(\mu^m, \Sigma^m)$.

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

Gaussian linear model

- $Y = \beta^T X + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2) \perp (X, M)$, $\beta \in \mathbb{R}^d$.
- for all $m \in \{0, 1\}^d$, there exist μ^m and Σ^m such that $X|(M = m) \sim \mathcal{N}(\mu^m, \Sigma^m)$.

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

Proposition (Oracle int. under Gaussian lin. mod., Z. et al. (2023a))

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

Gaussian linear model

- $Y = \beta^T X + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2) \perp (X, M)$, $\beta \in \mathbb{R}^d$.
- for all $m \in \{0, 1\}^d$, there exist μ^m and Σ^m such that $X|(M = m) \sim \mathcal{N}(\mu^m, \Sigma^m)$.

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

Proposition (Oracle int. under Gaussian lin. mod., Z. et al. (2023a))

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

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

Gaussian linear model

- $Y = \beta^T X + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2) \perp (X, M)$, $\beta \in \mathbb{R}^d$.
- for all $m \in \{0, 1\}^d$, there exist μ^m and Σ^m such that $X|(M = m) \sim \mathcal{N}(\mu^m, \Sigma^m)$.

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

Proposition (Oracle int. under Gaussian lin. mod., Z. et al. (2023a))

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

- 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)**

Goals reminder: achieve MCV!

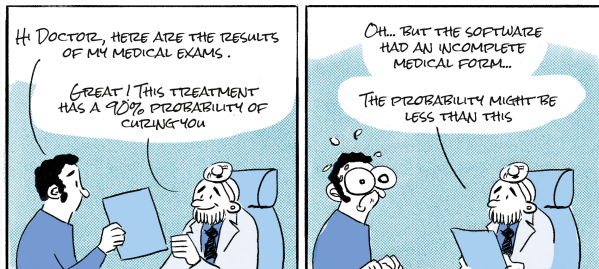
Goal: predict $Y^{(n+1)}$ with **confidence** $1 - \alpha$, i.e. build the smallest C_α such that:

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) ✗

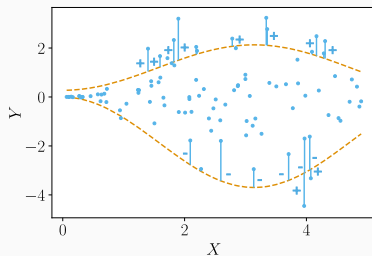
$$\forall m \in \{0, 1\}^d : \mathbb{P} \left\{ Y^{(n+1)} \in C_\alpha \left(X^{(n+1)}, m \right) \mid M^{(n+1)} = m \right\} \geq 1 - \alpha. \quad (\text{MCV})$$



Illustrations @theo.reminger

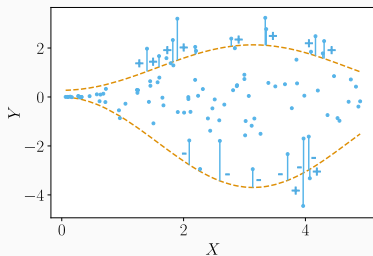
Conformalization step is independent of the important variable: the mask!

Observation: the α -correction term is computed among all the data points, regardless of their mask!



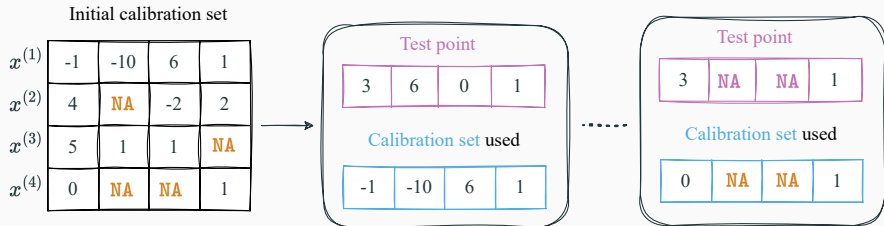
Conformalization step is independent of the important variable: the mask!

Observation: the α -correction term is computed among all the data points, regardless of their mask!



Warning: 2^d possible masks

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



Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

Impute-then-Regress+Conformalization

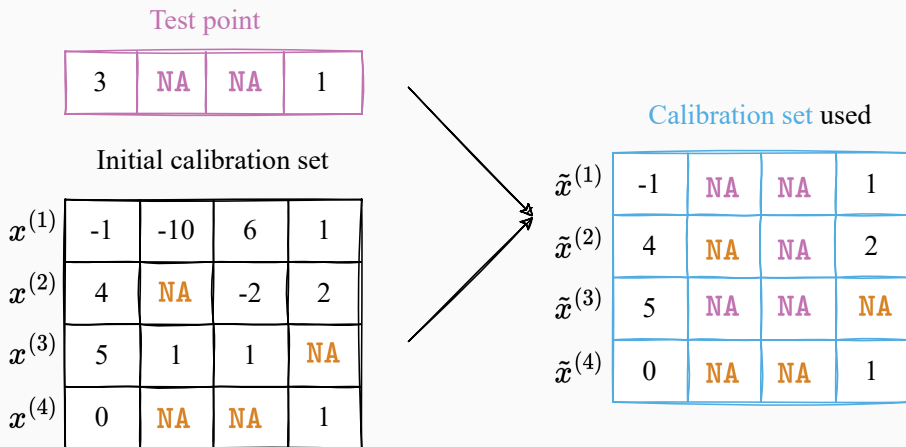
Missing Data Augmentation

Experimental results

Conclusion

Missing Data Augmentation (MDA) of the calibration set

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



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

Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

Impute-then-Regress+Conformalization

Missing Data Augmentation

MDA with Exact masking

MDA with Nested masking

Experimental results

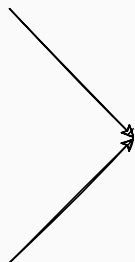
Conclusion

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
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

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



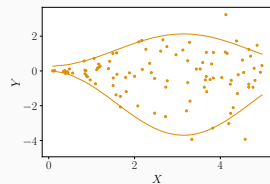
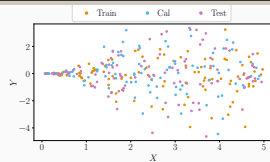
Calibration set used

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

#Cal^{M(test)} observations

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, Z. et al. (2023a))

If: i) the data is exchangeable, ii) $M \perp\!\!\!\perp X$, iii) $(Y \perp\!\!\!\perp M)|X$, 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:

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

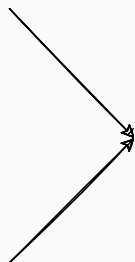
What if we kept all observations?

Test point

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

Initial calibration set

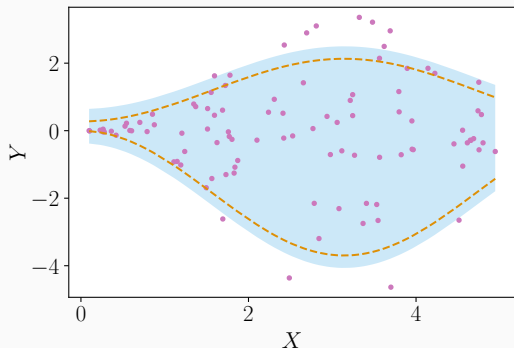
$x^{(1)}$	-1	-10	6	1
$x^{(2)}$	4	NA	-2	2
$x^{(3)}$	5	1	1	NA
$x^{(4)}$	0	NA	NA	1



Calibration set used

$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

What if we kept all observations?



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

► Build

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

Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

Impute-then-Regress+Conformalization

Missing Data Augmentation

MDA with Exact masking

MDA with Nested masking

Experimental results

Conclusion

Idea: modify the test point accordingly

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
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

Temporary test points

3	NA	NA	1
3	NA	NA	1
3	NA	NA	NA
3	NA	NA	1

and

↪ similar motivation than Barber et al. (2021)⁸ and Gupta et al. (2022)⁹.

⁸Predictive inference with the jackknife+, The Annals of Statistics

⁹Nested conformal prediction and quantile out-of-bag ensemble methods, Pattern Recognition

5. For a test point $(X^{(n+1)}, M^{(n+1)})$:

5.1 Set $\tilde{M}^{(k)} = \max(M^{(k)}, M^{(n+1)})$ for k
in the **calibration set**

	3	NA	NA	1
$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

5. For a test point $(X^{(n+1)}, M^{(n+1)})$:

- 5.1 Set $\tilde{M}^{(k)} = \max(M^{(k)}, M^{(n+1)})$ for k
in the **calibration set**
- 5.2 Impute the new **calibration set**
- 5.3 For each augmented **calibration point** k :

	3	NA	NA	1
$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

5. For a test point $(X^{(n+1)}, M^{(n+1)})$:

- 5.1 Set $\tilde{M}^{(k)} = \max(M^{(k)}, M^{(n+1)})$ for k in the calibration set
- 5.2 Impute the new calibration set
- 5.3 For each augmented calibration point k :
 - 5.3.1 Get its score $S^{(k)}$

	3	NA	NA	1
$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

5. For a test point $(X^{(n+1)}, M^{(n+1)})$:

5.1 Set $\tilde{M}^{(k)} = \max(M^{(k)}, M^{(n+1)})$ for k
in the **calibration set**

5.2 Impute the new **calibration set**

5.3 For each augmented **calibration point** k :

5.3.1 Get its score $S^{(k)}$

5.3.2 **Impute-then-predict** on the **augmented test point**
 $(X^{(n+1)}, \tilde{M}^{(k)})$, giving: $\widehat{QR}_{\alpha/2}(\tilde{X}^{(n+1),k})$ and
 $\widehat{QR}_{1-\alpha/2}(\tilde{X}^{(n+1),k})$

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

$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

3	NA	NA	1
3	NA	NA	1
3	NA	NA	NA
3	NA	NA	1

5. For a test point $(X^{(n+1)}, M^{(n+1)})$:

5.1 Set $\tilde{M}^{(k)} = \max(M^{(k)}, M^{(n+1)})$ for k
in the **calibration set**

5.2 Impute the new **calibration set**

5.3 For each augmented **calibration point** k :

5.3.1 Get its score $S^{(k)}$

5.3.2 **Impute-then-predict** on the **augmented test point**
 $(X^{(n+1)}, \tilde{M}^{(k)})$, giving: $\widehat{QR}_{\alpha/2}(\tilde{X}^{(n+1),k})$ and
 $\widehat{QR}_{1-\alpha/2}(\tilde{X}^{(n+1),k})$

5.3.3 Compute the corrected prediction interval:

$$[\widehat{QR}_{\alpha/2}(\tilde{X}^{(n+1),k}) - S^{(k)}; \widehat{QR}_{1-\alpha/2}(\tilde{X}^{(n+1),k}) + S^{(k)}] := [Z_{\text{lower}}^{(k)}; Z_{\text{upper}}^{(k)}]$$

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

$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

3	NA	NA	1
3	NA	NA	1
3	NA	NA	NA
3	NA	NA	1

5. For a test point $(X^{(n+1)}, M^{(n+1)})$:

5.1 Set $\tilde{M}^{(k)} = \max(M^{(k)}, M^{(n+1)})$ for k
in the **calibration set**

5.2 Impute the new **calibration set**

5.3 For each augmented **calibration point** k :

5.3.1 Get its score $S^{(k)}$

5.3.2 **Impute-then-predict** on the **augmented test point**
 $(X^{(n+1)}, \tilde{M}^{(k)})$, giving: $\widehat{QR}_{\alpha/2}(\tilde{X}^{(n+1),k})$ and
 $\widehat{QR}_{1-\alpha/2}(\tilde{X}^{(n+1),k})$

5.3.3 Compute the corrected prediction interval:

$$[\widehat{QR}_{\alpha/2}(\tilde{X}^{(n+1),k}) - S^{(k)}; \widehat{QR}_{1-\alpha/2}(\tilde{X}^{(n+1),k}) + S^{(k)}] := [Z_{\text{lower}}^{(k)}; Z_{\text{upper}}^{(k)}]$$

5.4 Compute the quantiles $q_{\alpha}(\{Z_{\text{lower}}^{(k)}\}_{k \in \text{Cal}})$ and $q_{1-\alpha}(\{Z_{\text{upper}}^{(k)}\}_{k \in \text{Cal}})$

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

$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

3	NA	NA	1
3	NA	NA	1
3	NA	NA	NA
3	NA	NA	1

5. For a test point $(X^{(n+1)}, M^{(n+1)})$:

5.1 Set $\tilde{M}^{(k)} = \max(M^{(k)}, M^{(n+1)})$ for k
in the **calibration set**

5.2 Impute the new **calibration set**

5.3 For each augmented **calibration point** k :

5.3.1 Get its score $S^{(k)}$

5.3.2 **Impute-then-predict** on the **augmented test point**
 $(X^{(n+1)}, \tilde{M}^{(k)})$, giving: $\widehat{QR}_{\alpha/2}(\tilde{X}^{(n+1),k})$ and
 $\widehat{QR}_{1-\alpha/2}(\tilde{X}^{(n+1),k})$

5.3.3 Compute the corrected prediction interval:

$$[\widehat{QR}_{\alpha/2}(\tilde{X}^{(n+1),k}) - S^{(k)}; \widehat{QR}_{1-\alpha/2}(\tilde{X}^{(n+1),k}) + S^{(k)}] := [Z_{\text{lower}}^{(k)}; Z_{\text{upper}}^{(k)}]$$

5.4 Compute the quantiles $q_{\alpha}(\{Z_{\text{lower}}^{(k)}\}_{k \in \text{Cal}})$ and $q_{1-\alpha}(\{Z_{\text{upper}}^{(k)}\}_{k \in \text{Cal}})$

5.5 Predict $[q_{\alpha}(\{Z_{\text{lower}}^{(k)}\}_{k \in \text{Cal}}); q_{1-\alpha}(\{Z_{\text{upper}}^{(k)}\}_{k \in \text{Cal}})]$

	3	NA	NA	1
$\tilde{x}^{(1)}$	-1	NA	NA	1
$\tilde{x}^{(2)}$	4	NA	NA	2
$\tilde{x}^{(3)}$	5	NA	NA	NA
$\tilde{x}^{(4)}$	0	NA	NA	1

3	NA	NA	1
3	NA	NA	1
3	NA	NA	NA
3	NA	NA	1

Theorem (CP-MDA-Nested marginal validity, Z. et al. (2023b))

If the data is exchangeable, then for almost all imputation function CP-MDA-Nested is such that:

$$\mathbb{P} \left(Y \in \hat{C}_\alpha (X, M) \right) \geq 1 - 2\alpha.$$

Theorem (CP-MDA-Nested marginal validity, Z. et al. (2023b))

If the data is exchangeable, then for almost all imputation function CP-MDA-Nested is such that:

$$\mathbb{P} \left(Y \in \hat{C}_\alpha(X, M) \right) \geq 1 - 2\alpha.$$

- ✓ Any missing mechanism (no need to assume $M \perp\!\!\!\perp X$)
- ✓ Does not require $(Y \perp\!\!\!\perp M) | X$
- ✗ Marginal guarantee

Theorem (CP-MDA-Nested marginal validity, Z. et al. (2023b))

If the data is exchangeable, then for almost all imputation function CP-MDA-Nested is such that:

$$\mathbb{P} \left(Y \in \hat{C}_\alpha(X, M) \right) \geq 1 - 2\alpha.$$

- ✓ Any missing mechanism (no need to assume $M \perp\!\!\!\perp X$)
- ✓ Does not require $(Y \perp\!\!\!\perp M) | X$
- ✗ Marginal guarantee

Proof element: based on Jackknife+ ideas (Barber et al., 2021).

Leaving-out the k -th data point to predict on the l -th data point

\leftrightarrow

Apply the mask of the k -th data point to the l -th data point on which you predict

Stochastic domination of the quantiles (SDQ)

Let $(\check{m}, \check{m}) \in (\{0, 1\}^d)^2$. If $\check{m} \subset \check{m}$ then for any $\delta \in [0, 0.5]$:

$$q_{1-\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})} \leq q_{1-\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})}, \text{ and } q_{\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})} \geq q_{\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})}.$$

↪ predictive uncertainty increases with bigger masks.

MDA-Nested (nearly) achieves Mask-Conditional-Validity (MCV)

Stochastic domination of the quantiles (SDQ)

Let $(\check{m}, \check{m}) \in (\{0, 1\}^d)^2$. If $\check{m} \subset \check{m}$ then for any $\delta \in [0, 0.5]$:

$$q_{1-\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})} \leq q_{1-\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})}, \text{ and } q_{\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})} \geq q_{\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})}.$$

\rightsquigarrow predictive uncertainty increases with bigger masks.

Theorem (CP-MDA-Nested (nearly) achieves MCV, Z. et al. (2023a))

If i) the data is exchangeable, ii) $M \perp\!\!\!\perp X$, iii) $(Y \perp\!\!\!\perp M)|X$, iv) SDQ holds, then for almost all imputation function "CP-MDA-Nested" is s.t. for any $m \in \{0, 1\}^d$:

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

MDA-Nested (nearly) achieves Mask-Conditional-Validity (MCV)

Stochastic domination of the quantiles (SDQ)

Let $(\check{m}, \check{m}) \in (\{0, 1\}^d)^2$. If $\check{m} \subset \check{m}$ then for any $\delta \in [0, 0.5]$:

$$q_{1-\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})} \leq q_{1-\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})}, \text{ and } q_{\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})} \geq q_{\delta/2}^{Y|(X_{\text{obs}(\check{m})}, M=\check{m})}.$$

\rightsquigarrow predictive uncertainty increases with bigger masks.

Theorem (CP-MDA-Nested (nearly) achieves MCV, Z. et al. (2023a))

If i) the data is exchangeable, ii) $M \perp\!\!\!\perp X$, iii) $(Y \perp\!\!\!\perp M)|X$, iv) SDQ holds, then for almost all imputation function "CP-MDA-Nested" is s.t. for any $m \in \{0, 1\}^d$:

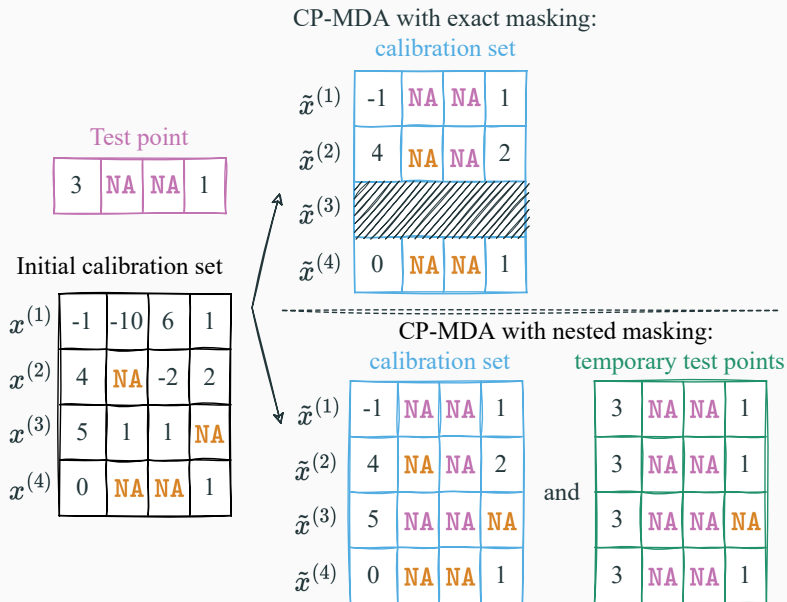
$$\mathbb{P}\left(Y \in \widehat{C}_\alpha(X, m) \mid M = m\right) \geq 1 - \alpha.$$

Change on MDA-Nested: outputs any

$[q_\alpha(\{Z_{\text{lower}}^{(k)}\}_{k \in \text{Cal}^{\check{m}}}); q_{1-\alpha}(\{Z_{\text{upper}}^{(k)}\}_{k \in \text{Cal}^{\check{m}}})]$, where \check{m} is randomly¹⁰ selected such that $m \subset \check{m}$.

¹⁰The randomness may depend on $\#\text{Cal}^{\check{m}}$.

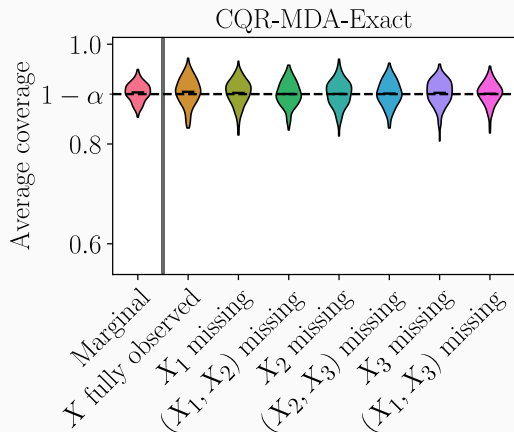
Summary of CP-MDA



MDA achieves Mask-Conditional-Validity (MCV)

$$Y = \beta^T X + \varepsilon,$$

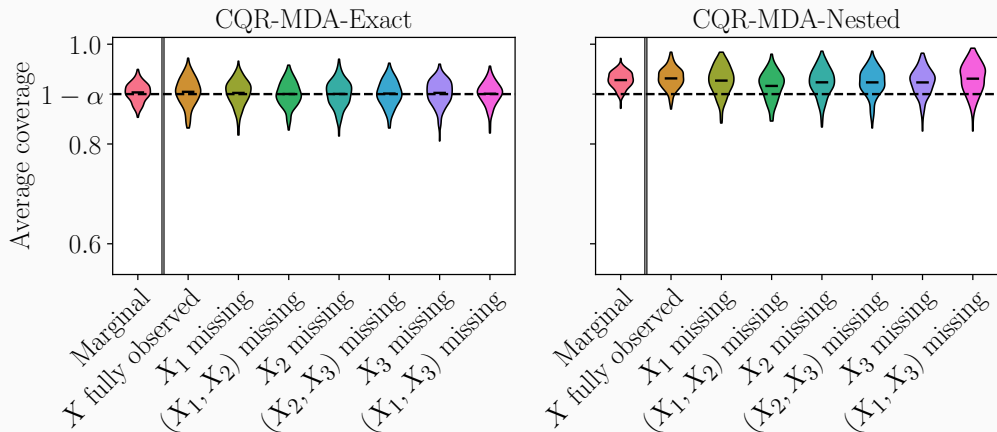
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



MDA achieves Mask-Conditional-Validity (MCV)

$$Y = \beta^T X + \varepsilon,$$

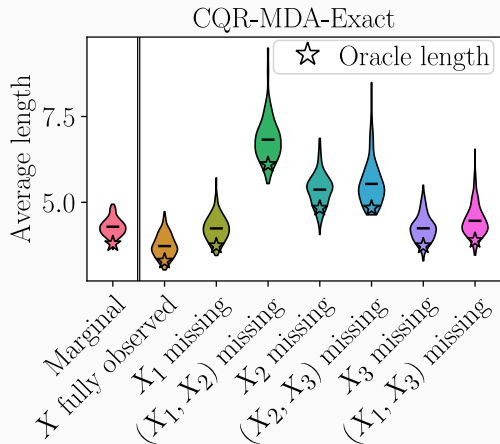
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



MDA achieves (MCV) in an informative way

$$Y = \beta^T X + \varepsilon,$$

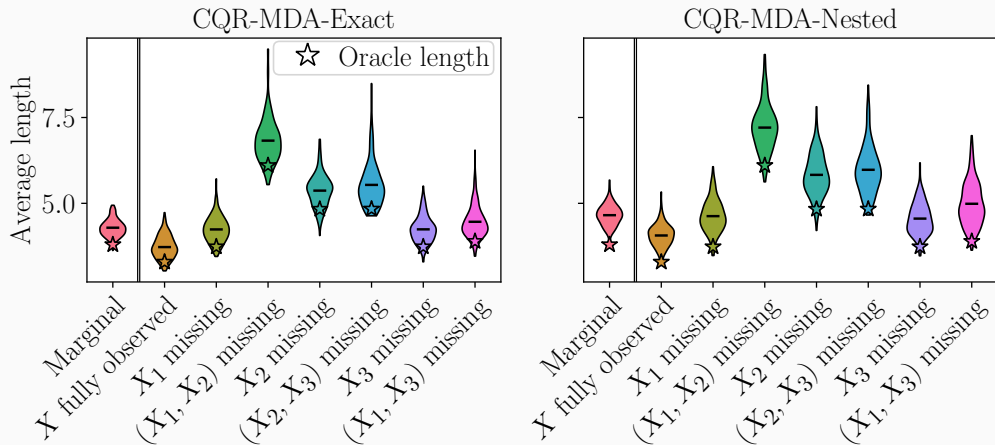
$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



MDA achieves (MCV) in an informative way

$$Y = \beta^T X + \varepsilon,$$

$\beta = (1, 2, -1)^T$, $\varepsilon \perp X$, X and ε Gaussian, 20% uniform MCAR missing values.



Introduction to missing values

Quantifying predictive uncertainty with missing values

Split Conformal Prediction

Conformalized Quantile Regression

Impute-then-Regress+Conformalization

Missing Data Augmentation

Experimental results

Conclusion

- Imputation by iterative ridge (\sim conditional expectation)

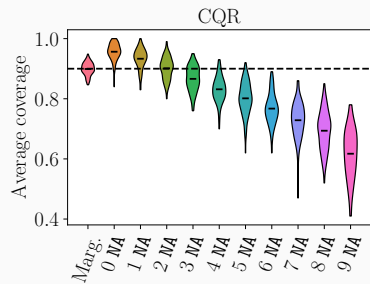
- Imputation by iterative ridge (\sim conditional expectation)
- **Concatenate the mask in the features**

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

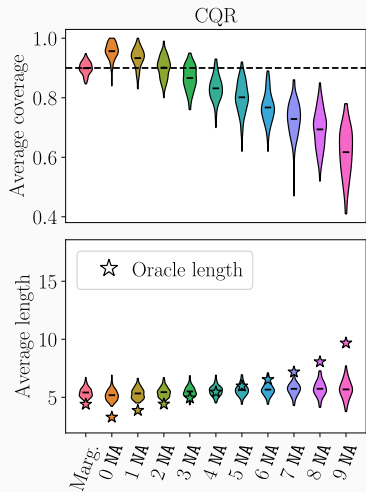
- Imputation by iterative ridge (\sim conditional expectation)
- **Concatenate the mask in the features**
- Neural network, fitted to minimize the pinball loss
- (Semi)-synthetic experiments:

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

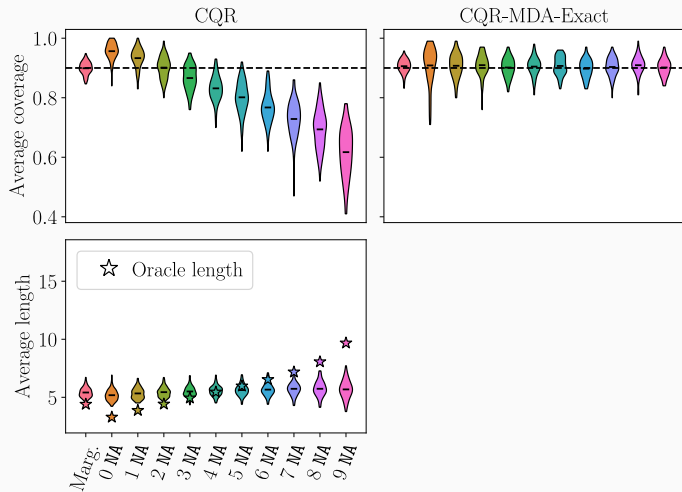
Synthetic experiments (Gaussian linear model, $d = 10$)



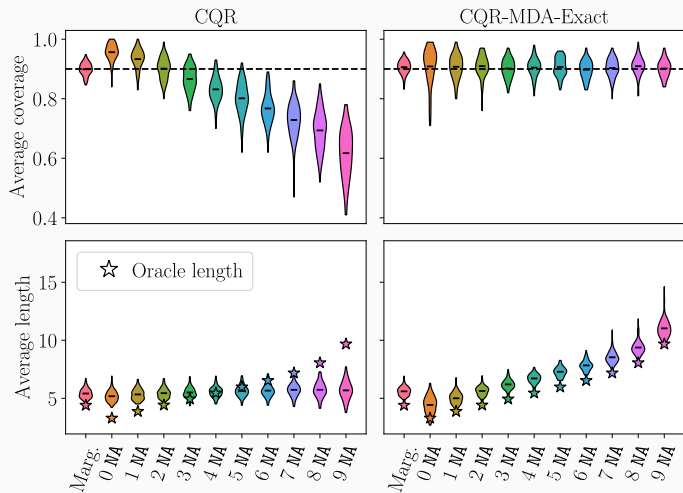
Synthetic experiments (Gaussian linear model, $d = 10$)



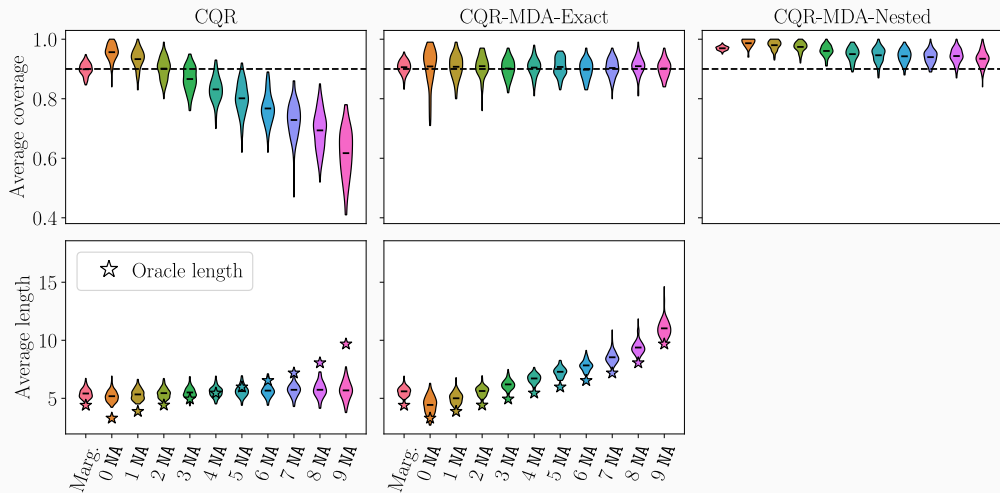
Synthetic experiments (Gaussian linear model, $d = 10$)



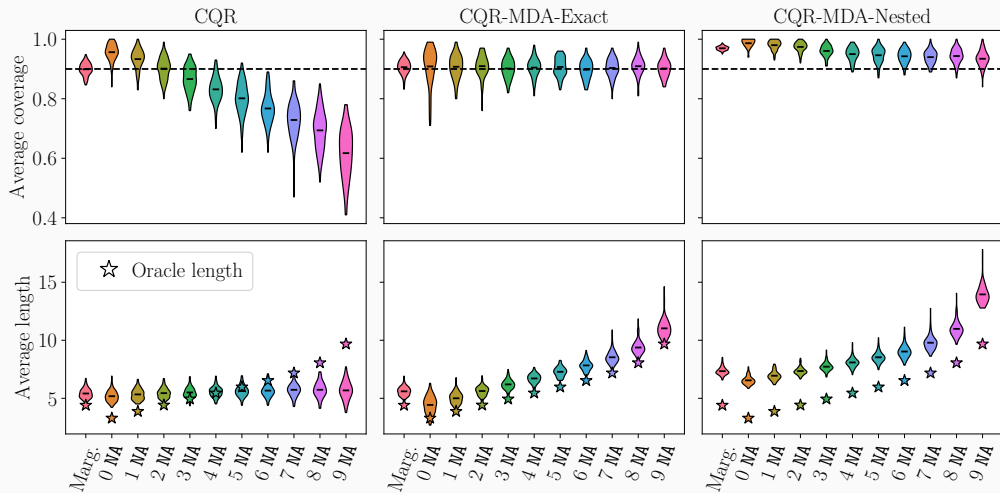
Synthetic experiments (Gaussian linear model, $d = 10$)



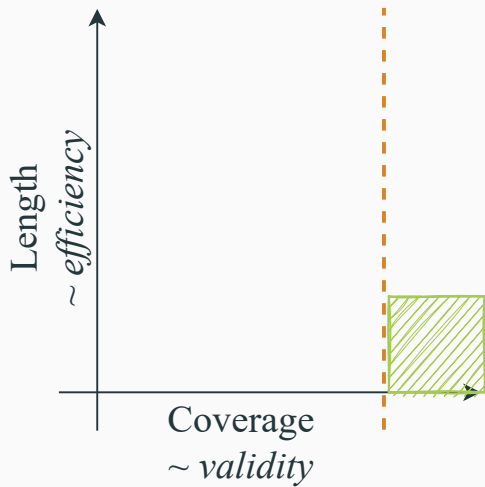
Synthetic experiments (Gaussian linear model, $d = 10$)



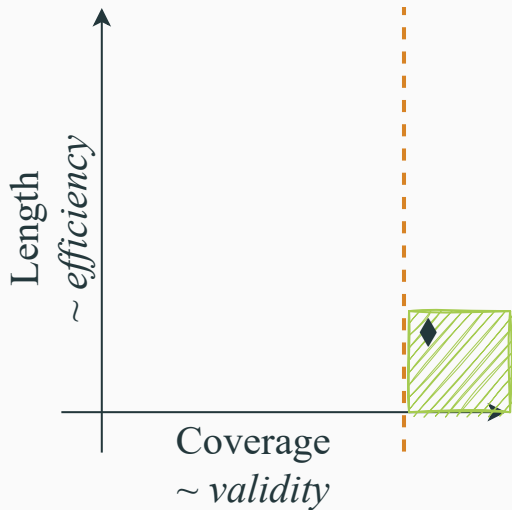
Synthetic experiments (Gaussian linear model, $d = 10$)



Before more experiments, visualisation

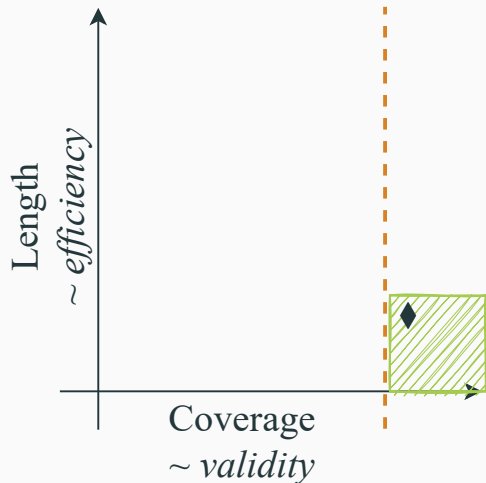


Before more experiments, visualisation



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

Before more experiments, visualisation



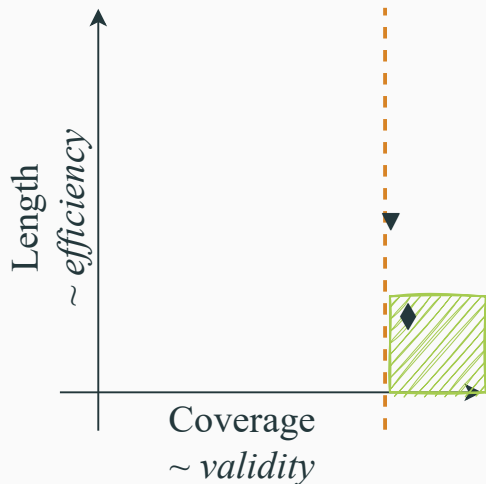
◆ : marginal coverage, i.e.

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

▼ : lowest coverage, i.e.

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

Before more experiments, visualisation



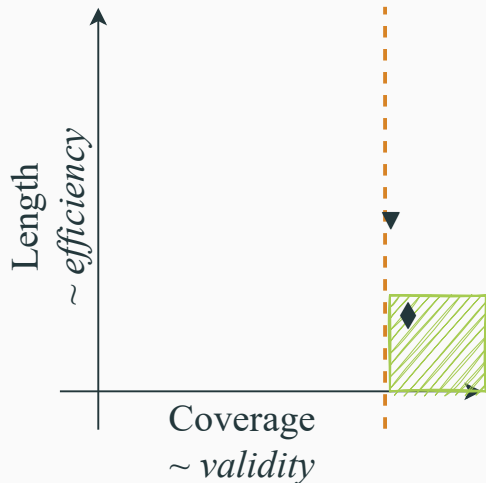
◆ : marginal coverage, i.e.

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

▼ : 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.

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

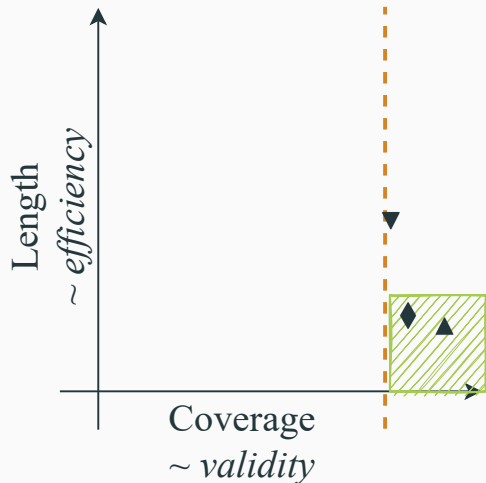
▼ : lowest coverage, i.e.

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

▲ : 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.

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

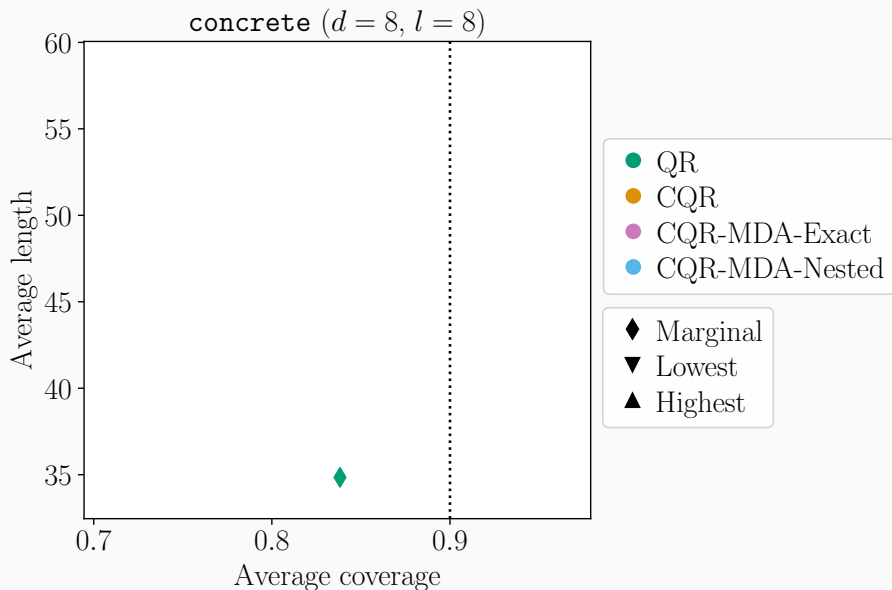
▼ : lowest coverage, i.e.

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

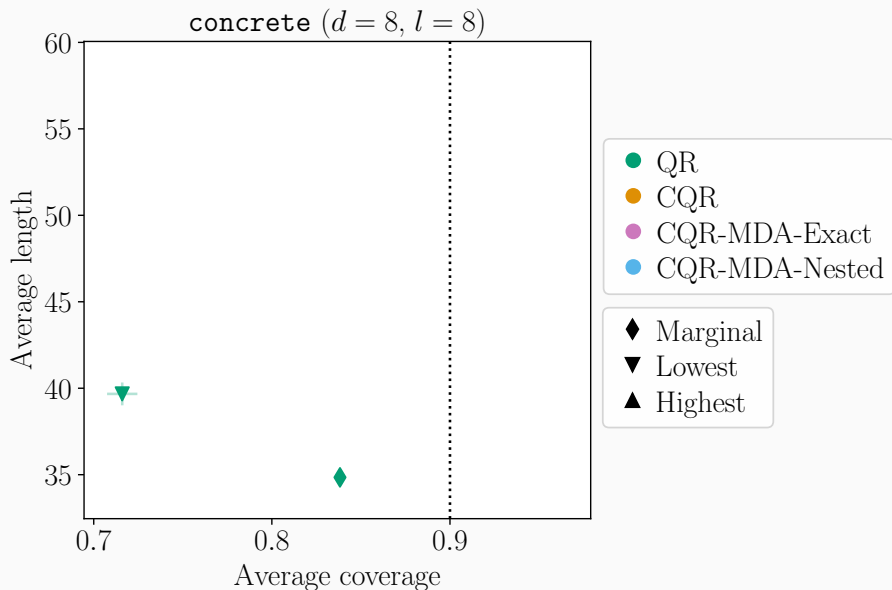
▲ : highest coverage, i.e.

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

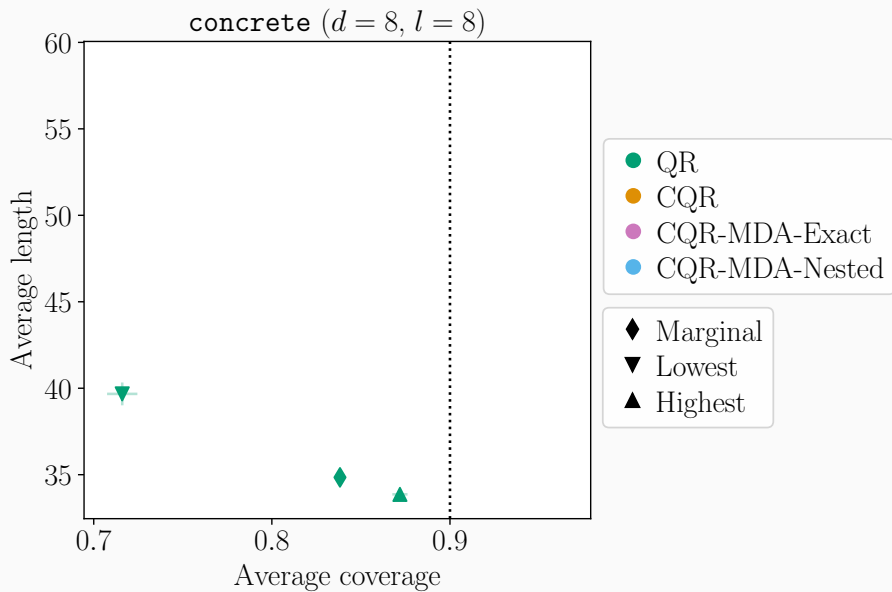
Semi-synthetic experiments



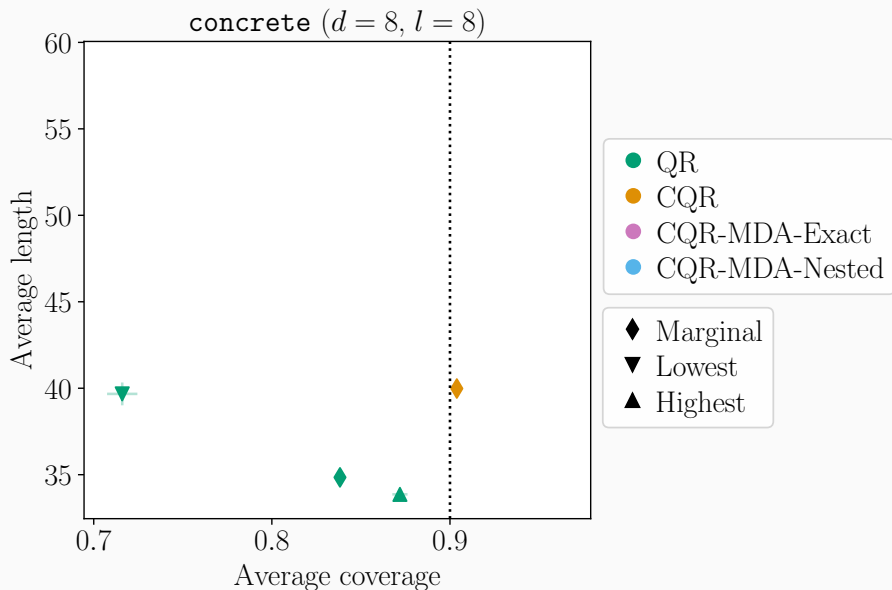
Semi-synthetic experiments



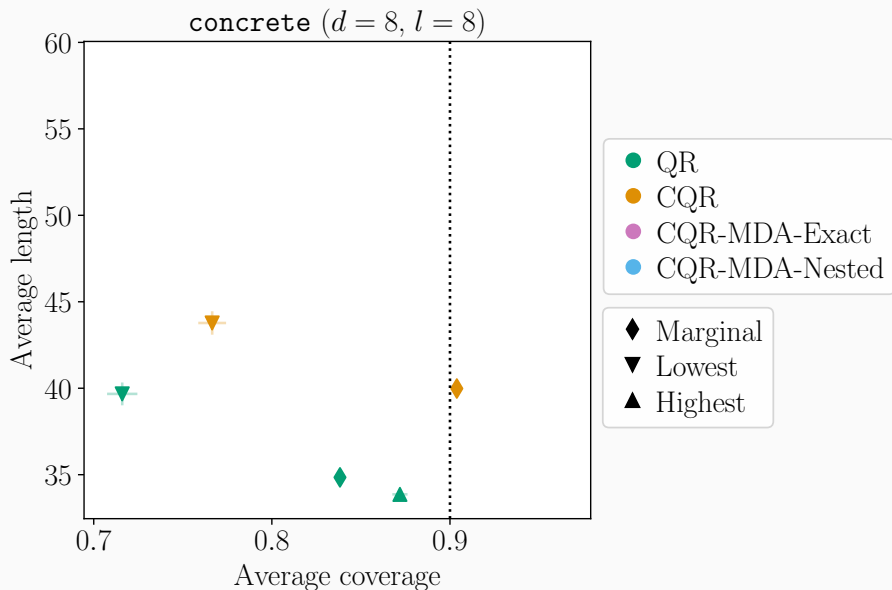
Semi-synthetic experiments



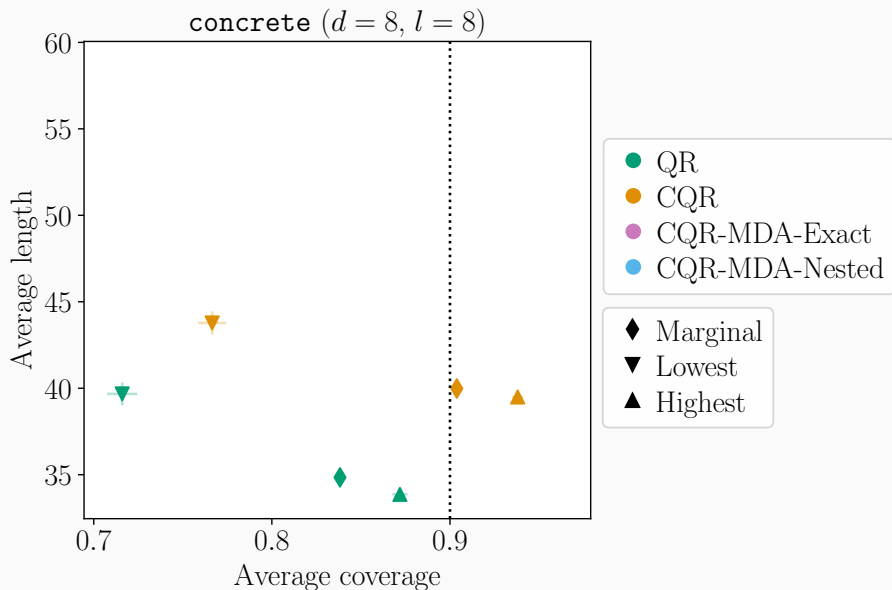
Semi-synthetic experiments



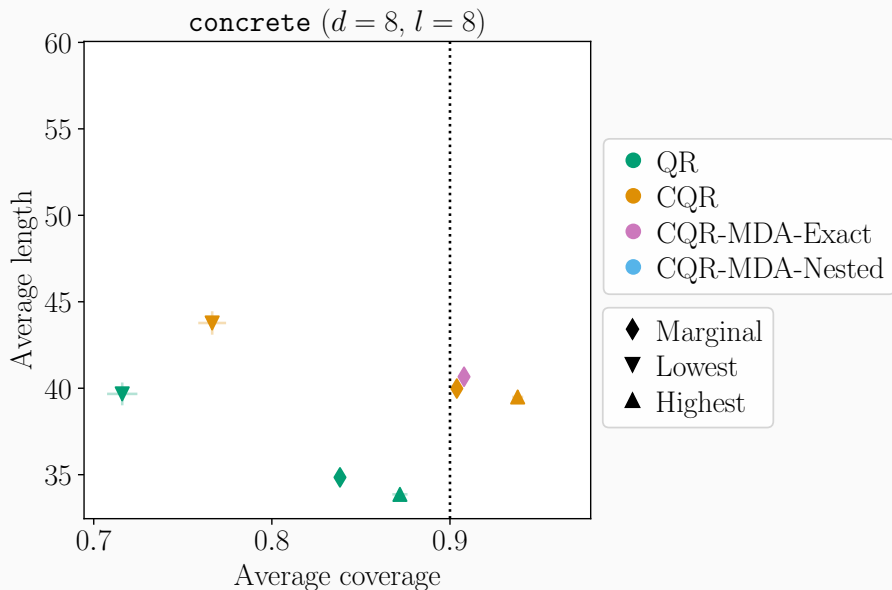
Semi-synthetic experiments



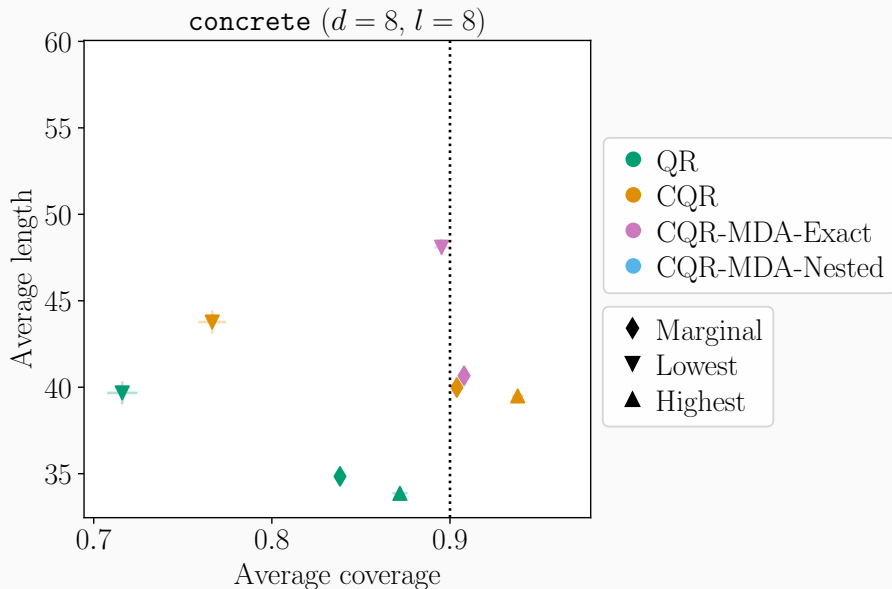
Semi-synthetic experiments



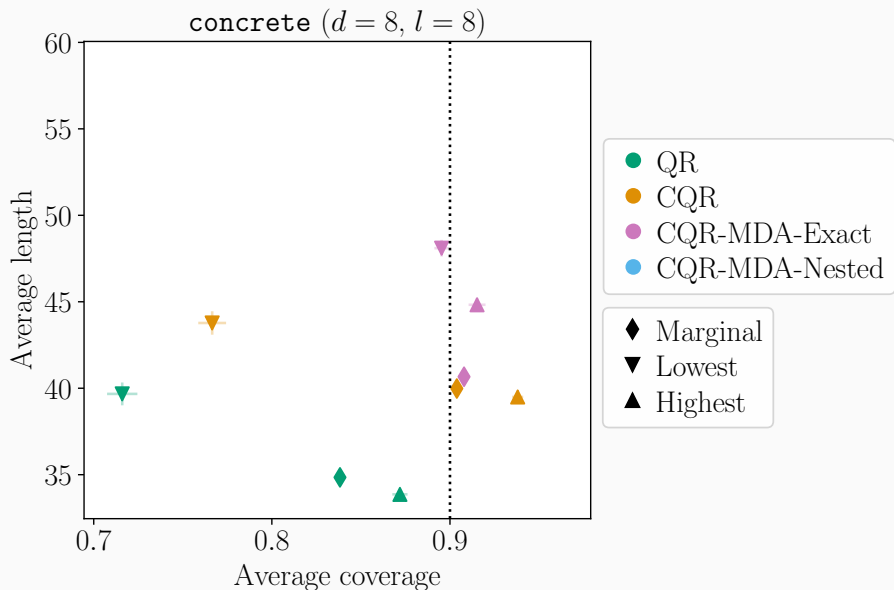
Semi-synthetic experiments



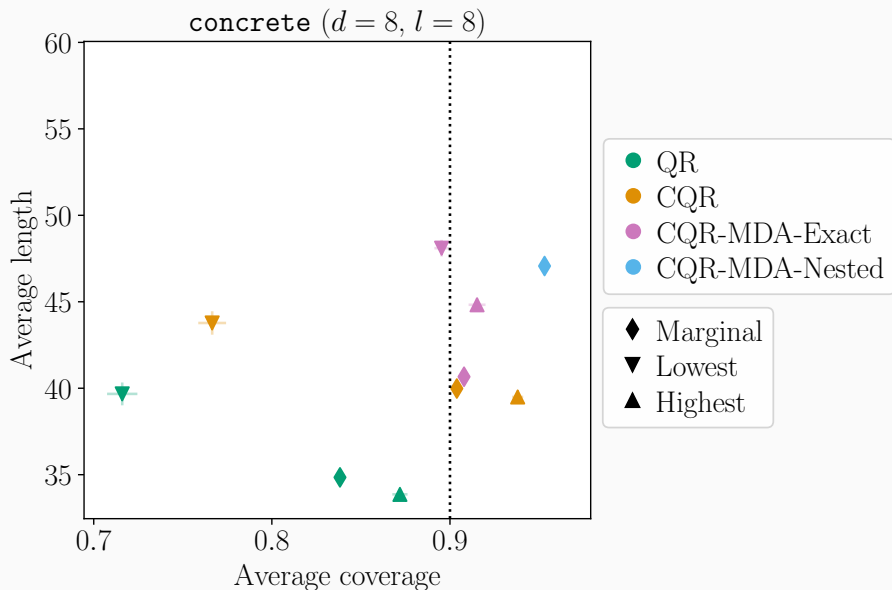
Semi-synthetic experiments



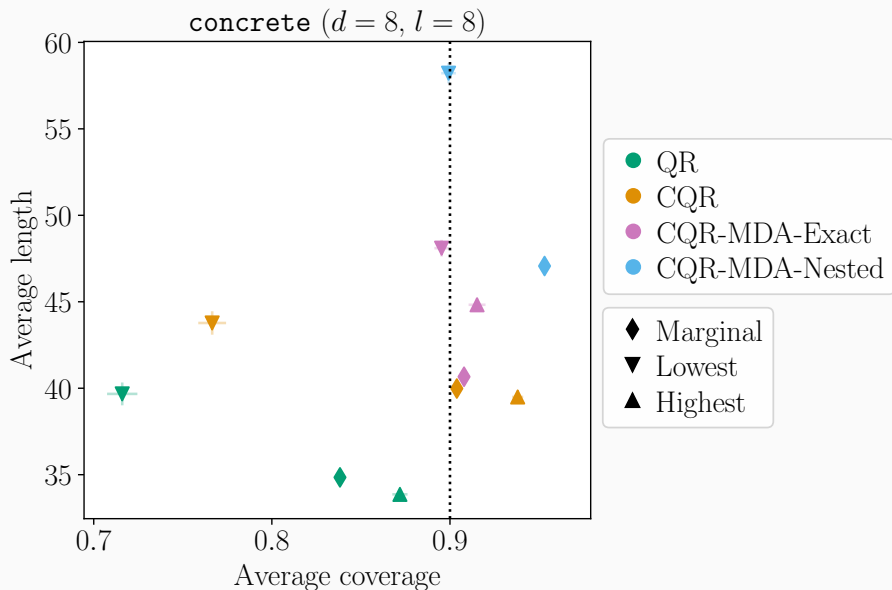
Semi-synthetic experiments



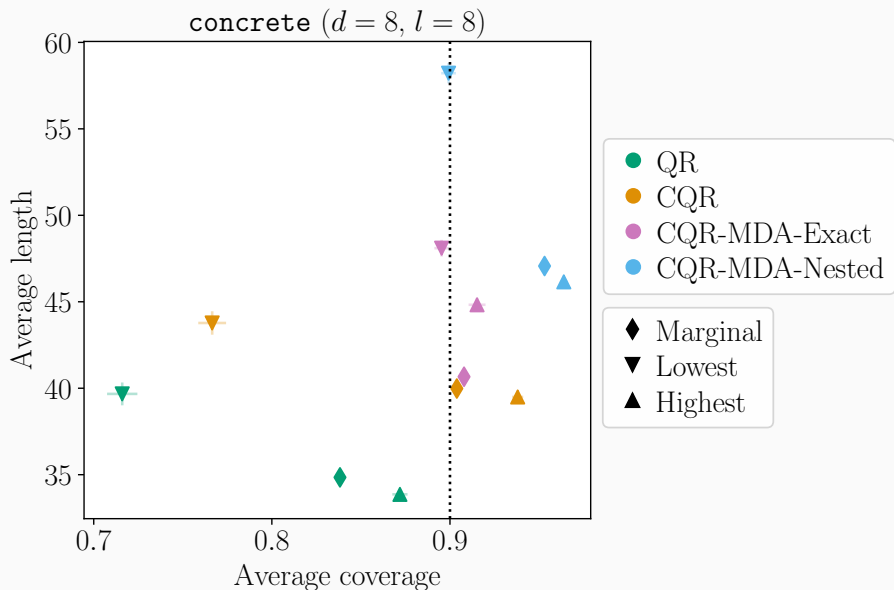
Semi-synthetic experiments

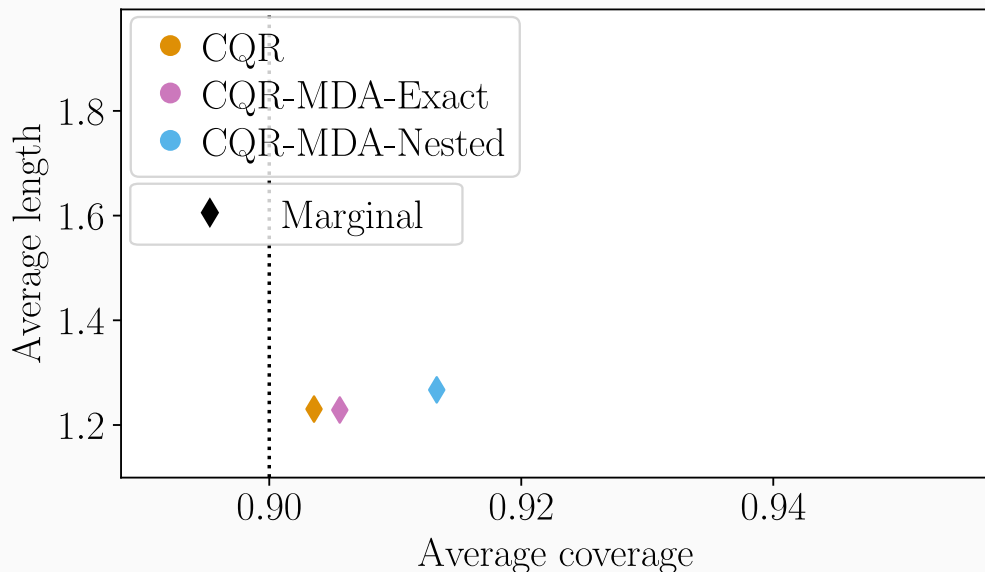


Semi-synthetic experiments

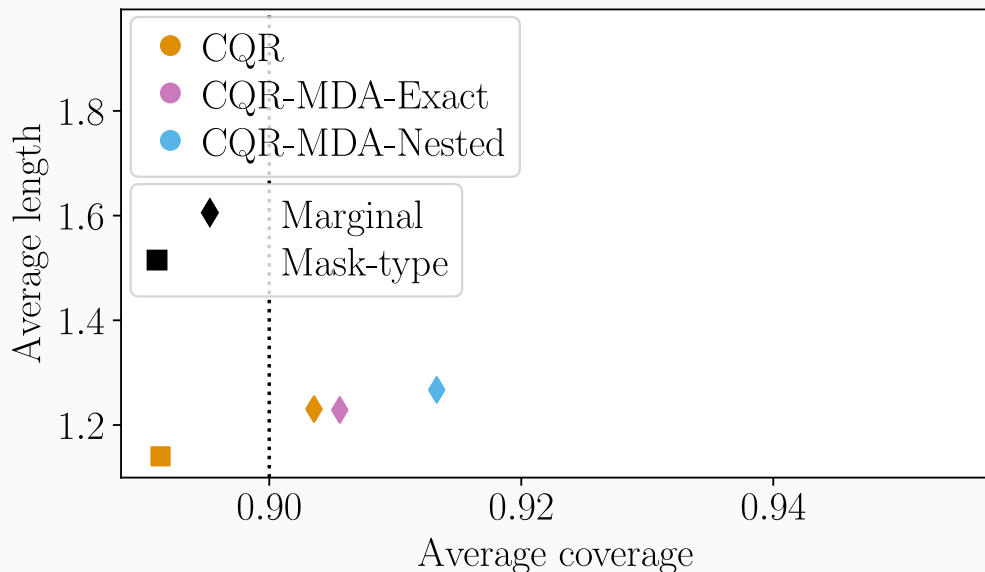


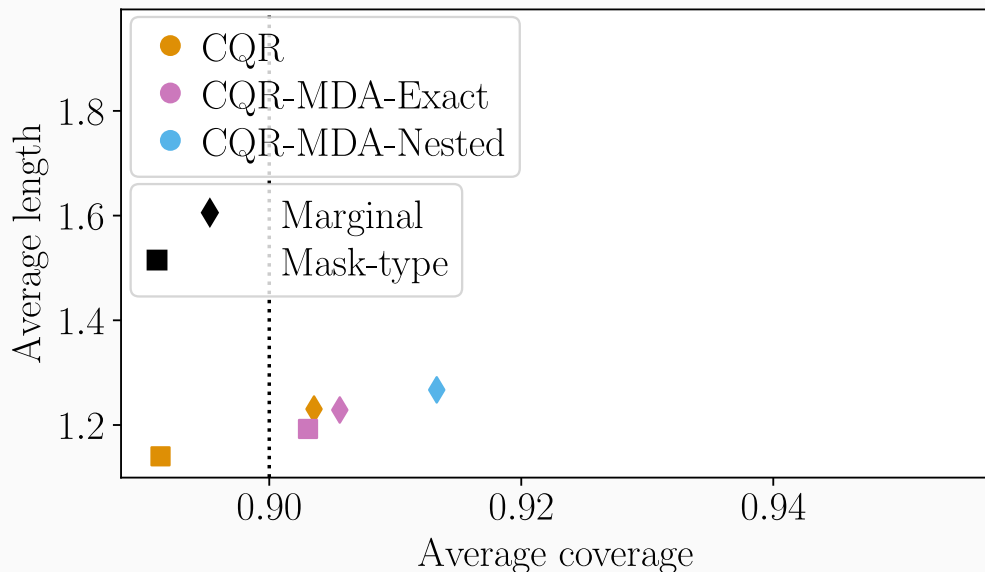
Semi-synthetic experiments



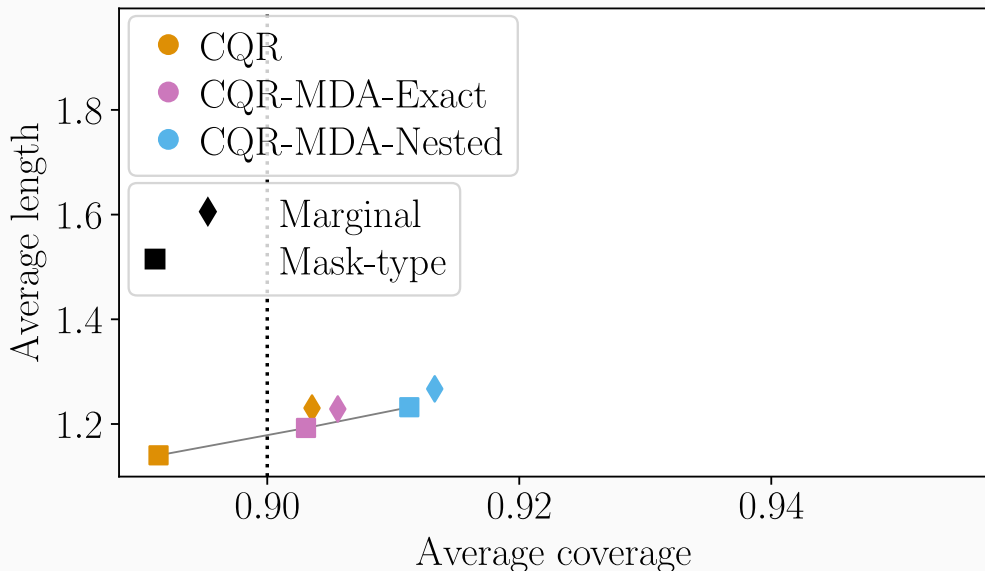


Real data experiment: TraumaBase[®], critical care medicine

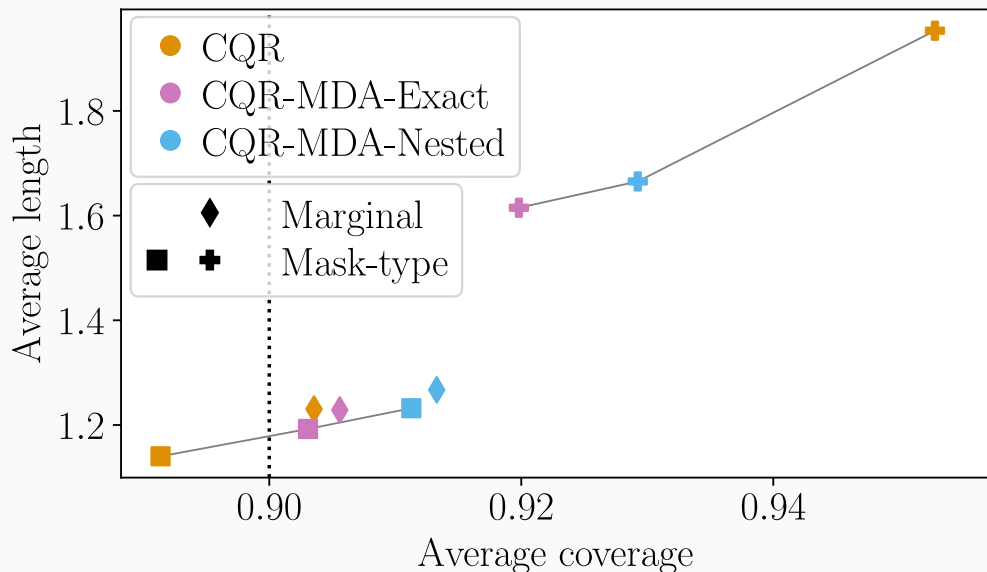




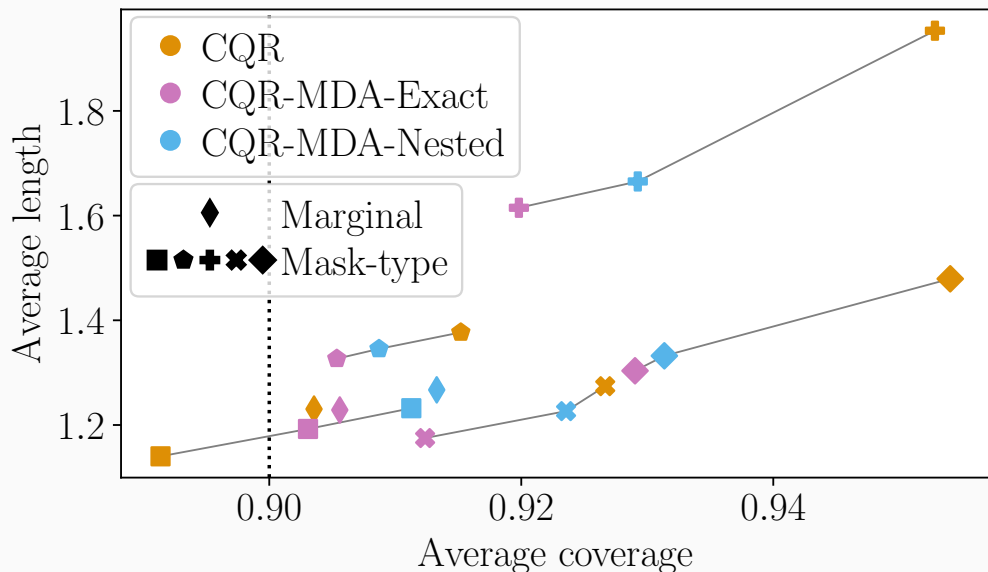
Real data experiment: TraumaBase[®], critical care medicine



Real data experiment: TraumaBase[®], critical care medicine



Real data experiment: TraumaBase[®], critical care medicine



Introduction to missing values

Quantifying predictive uncertainty with missing values

Conclusion

Take-home-messages

- CP marginal guarantees hold on the imputed data set.

Take-home-messages

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

Take-home-messages

- 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.

Take-home-messages

- 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.**

Take-home-messages

- 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).

- 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).
- Extension: consistency of universal quantile learner when chained with almost any imputation function.



- Investigate alternative methods relying on trade-offs between MDA-Exact and MDA-Nested

- Investigate alternative methods relying on trade-offs between MDA-Exact and MDA-Nested
- Relationship with Gibbs et al. (2023)¹¹
 - ✓ Beyond MCAR
 - ✗ Upper bound in $\frac{2^d}{(n+1)\mathbb{P}_M(m)}$: high value for less probable masks
 - ↪ MCV are non-overlapping groups: boils down to splitting the calibration set!

¹¹ *Conformal Prediction With Conditional Guarantees*

- Investigate alternative methods relying on trade-offs between MDA-Exact and MDA-Nested
- Relationship with Gibbs et al. (2023)¹¹
 - ✓ Beyond MCAR
 - ✗ Upper bound in $\frac{2^d}{(n+1)\mathbb{P}_M(m)}$: high value for less probable masks
 - ↪ MCV are non-overlapping groups: boils down to splitting the calibration set!

¹¹ *Conformal Prediction With Conditional Guarantees*

- Investigate alternative methods relying on trade-offs between MDA-Exact and MDA-Nested
- Relationship with Gibbs et al. (2023)¹¹
 - ✓ Beyond MCAR
 - ✗ Upper bound in $\frac{2^d}{(n+1)\mathbb{P}_M(m)}$: high value for less probable masks
 - ↪ MCV are non-overlapping groups: boils down to splitting the calibration set!
- Quantify the impact of the imputation's choice on Quantile Regression quality in finite sample

¹¹ *Conformal Prediction With Conditional Guarantees*

Thank you! Questions? :)

- Ayme, A., Boyer, C., Dieuleveut, A., and Scornet, E. (2022). Near-optimal rate of consistency for linear models with missing values. *ICML*.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. (2021). Predictive inference with the jackknife+. *The Annals of Statistics*, 49(1).
- Gibbs, I., Cherian, J. J., and Candès, E. J. (2023). Conformal prediction with conditional guarantees.
- Gupta, C., Kuchibhotla, A. K., and Ramdas, A. (2022). Nested conformal prediction and quantile out-of-bag ensemble methods. *Pattern Recognition*, 127.
- Le Morvan, M., Josse, J., Scornet, E., and Varoquaux, G. (2021). What's a good imputation to predict with missing values? *NeurIPS*.
- 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*.

- Romano, Y., Patterson, E., and Candès, E. (2019). Conformalized Quantile Regression. *NeurIPS*.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3).
- Vovk, V., Gammerman, A., and Shafer, G. (2005). *Algorithmic Learning in a Random World*. Springer US.
- Z. , M., Dieuleveut, A., Josse, J., and Romano, Y. (2023a). Conformal prediction with missing values. *ICML*.
- Z. , M., Dieuleveut, A., Josse, J., and Romano, Y. (2023b). Predictive uncertainty quantification with missing values. To be submitted.
- Zhu, Z., Wang, T., and Samworth, R. J. (2019). High-dimensional principal component analysis with heterogeneous missingness. *arXiv*.

Appendix

Towards asymptotic individualized coverage

Consistency of a universal quantile learner after imputation

Let Φ be an imputation function chosen by the user.

Denote $g_{\beta, \Phi}^* \in \operatorname{argmin}_{g: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E}[\rho_{\beta}(Y - g \circ \Phi(X, M))] := \mathcal{R}_{\beta, \Phi}(g)$.

Consistency of a universal quantile learner after imputation

Let Φ be an imputation function chosen by the user.

Denote $g_{\beta, \Phi}^* \in \operatorname{argmin}_{g: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} [\rho_{\beta}(Y - g \circ \Phi(X, M))] := \mathcal{R}_{\beta, \Phi}(g)$.

Comparison with: $\operatorname{argmin}_f \mathbb{E} [\rho_{\beta}(Y - f(X, M))] \text{ (informal)}$.

Consistency of a universal quantile learner after imputation

Let Φ be an imputation function chosen by the user.

Denote $g_{\beta, \Phi}^* \in \operatorname{argmin}_{g: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} [\rho_{\beta}(Y - g \circ \Phi(X, M))] := \mathcal{R}_{\beta, \Phi}(g)$.

Comparison with: $\operatorname{argmin}_f \mathbb{E} [\rho_{\beta}(Y - f(X, M))]$ (*informal*).

Proposition (Pinball-consistency of an universal learner)

For almost all \mathcal{C}^{∞} imputation function Φ , the function $g_{\beta, \Phi}^* \circ \Phi$ is Bayes optimal for the pinball-risk of level β .

Consistency of a universal quantile learner after imputation

Let Φ be an imputation function chosen by the user.

Denote $g_{\beta, \Phi}^* \in \operatorname{argmin}_{g: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} [\rho_{\beta}(Y - g \circ \Phi(X, M))] := \mathcal{R}_{\beta, \Phi}(g)$.

Comparison with: $\operatorname{argmin}_f \mathbb{E} [\rho_{\beta}(Y - f(X, M))]$ (*informal*).

Proposition (Pinball-consistency of an universal learner)

For almost all \mathcal{C}^{∞} imputation function Φ , the function $g_{\beta, \Phi}^* \circ \Phi$ is Bayes optimal for the pinball-risk of level β .

\hookrightarrow any universally consistent algorithm for **quantile regression** trained on the data imputed by Φ is pinball-**Bayes-consistent**.

Consistency of a universal quantile learner after imputation

Let Φ be an imputation function chosen by the user.

Denote $g_{\beta, \Phi}^* \in \operatorname{argmin}_{g: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E} [\rho_{\beta}(Y - g \circ \Phi(X, M))] := \mathcal{R}_{\beta, \Phi}(g)$.

Comparison with: $\operatorname{argmin}_f \mathbb{E} [\rho_{\beta}(Y - f(X, M))] \text{ (informal)}$.

Proposition (Pinball-consistency of an universal learner)

For almost all \mathcal{C}^{∞} imputation function Φ , the function $g_{\beta, \Phi}^* \circ \Phi$ is Bayes optimal for the pinball-risk of level β .

\Leftrightarrow any universally consistent algorithm for **quantile regression** trained on the data imputed by Φ is pinball-**Bayes-consistent**.

This is an extension of the result of Le Morvan et al. (2021).

Corollary

For any missing mechanism, for almost all C^∞ imputation function Φ , if $F_{Y|(X_{\text{obs}(M)}, M)}$ is continuous, a universally consistent quantile regressor trained on the imputed data set yields asymptotic conditional coverage.

Corollary

For any missing mechanism, for almost all C^∞ imputation function Φ , if $F_{Y|(X_{\text{obs}(M)}, M)}$ is continuous, a universally consistent quantile regressor trained on the imputed data set yields asymptotic conditional coverage.

$\Leftrightarrow \mathbb{P}(Y \in \widehat{C}_\alpha(x) | X = x, M = m) \geq 1 - \alpha$ for any $m \in \mathcal{M}$ and any $x \in \mathbb{R}^d$, asymptotically with a super quantile learner.

$$d = 3$$

Data generation

$$(X, Y) \in \mathbb{R}^3 \times \mathbb{R}.$$

$$Y = \beta^T X + \varepsilon$$

with $\varepsilon \sim \mathcal{N}(0, 1)$, $\beta = (1, 2, -1)^T$ and

$$(X_1, X_2, X_3) \sim \mathcal{N} \left(\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.8 & 0.8 \\ 0.8 & 1 & 0.8 \\ 0.8 & 0.8 & 1 \end{pmatrix} \right).$$

Data generation

$$(X, Y) \in \mathbb{R}^3 \times \mathbb{R}.$$

$$Y = \beta^T X + \varepsilon$$

with $\varepsilon \sim \mathcal{N}(0, 1)$, $\beta = (1, 2, -1)^T$ and

$$(X_1, X_2, X_3) \sim \mathcal{N} \left(\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.8 & 0.8 \\ 0.8 & 1 & 0.8 \\ 0.8 & 0.8 & 1 \end{pmatrix} \right).$$

All components of X each have a probability 0.2 of being missing, Completely At Random.

Simulation settings

- Method: CQR
- Basemodel: neural network
- 200 repetitions
 - train size of 250 points
 - calibration size of 250 points
 - test size of 2000 points

$d = 10$, with missing data augmentation

Data generation

$$(X, Y) \in \mathbb{R}^{10} \times \mathbb{R}.$$

$$Y = \beta^T X + \varepsilon$$

with $\varepsilon \sim \mathcal{N}(0, 1)$, $\beta = (1, 2, -1, 3, -0.5, -1, 0.3, 1.7, 0.4, -0.3)^T$ and

$$(X_1, \dots, X_{10}) \sim \mathcal{N} \left(\begin{pmatrix} 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.8 & \cdots & 0.8 \\ 0.8 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0.8 \\ 0.8 & \cdots & 0.8 & 1 \end{pmatrix} \right).$$

Data generation

$$(X, Y) \in \mathbb{R}^{10} \times \mathbb{R}.$$

$$Y = \beta^T X + \varepsilon$$

with $\varepsilon \sim \mathcal{N}(0, 1)$, $\beta = (1, 2, -1, 3, -0.5, -1, 0.3, 1.7, 0.4, -0.3)^T$ and

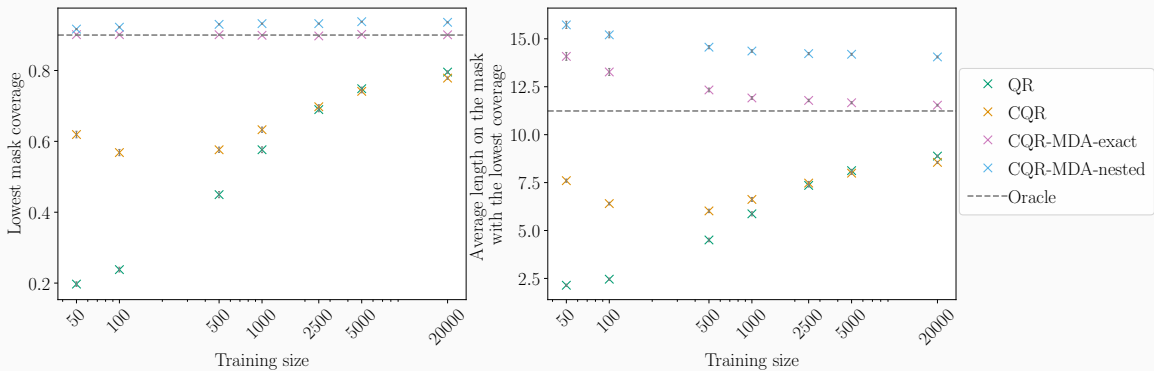
$$(X_1, \dots, X_{10}) \sim \mathcal{N} \left(\begin{pmatrix} 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.8 & \dots & 0.8 \\ 0.8 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0.8 \\ 0.8 & \dots & 0.8 & 1 \end{pmatrix} \right).$$

All components of X each have a probability 0.2 of being missing, Completely At Random.

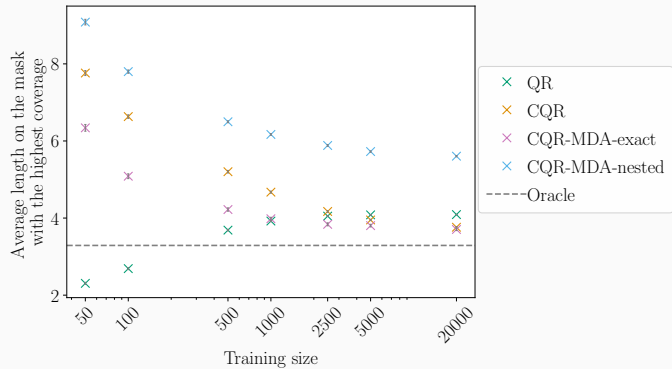
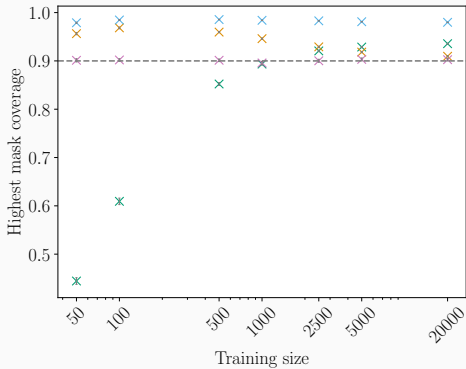
Simulation settings: varying training size

- Method: CQR
- Basemodel: neural network
- Imputation: iterative (\approx conditional expectation)
- Mask as features: yes
- 100 repetitions
 - train size varies
 - calibration size of 1000 points
 - test size of 2000 points

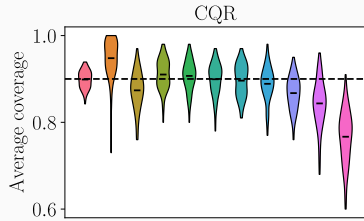
Results on the worst group



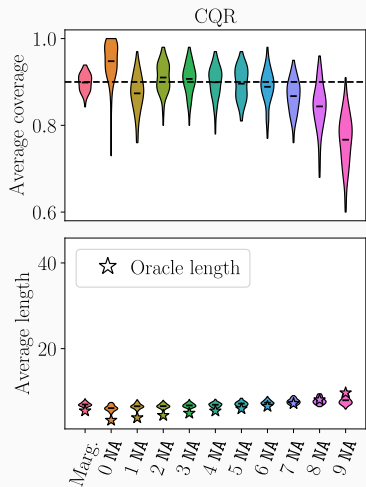
Results on the best group



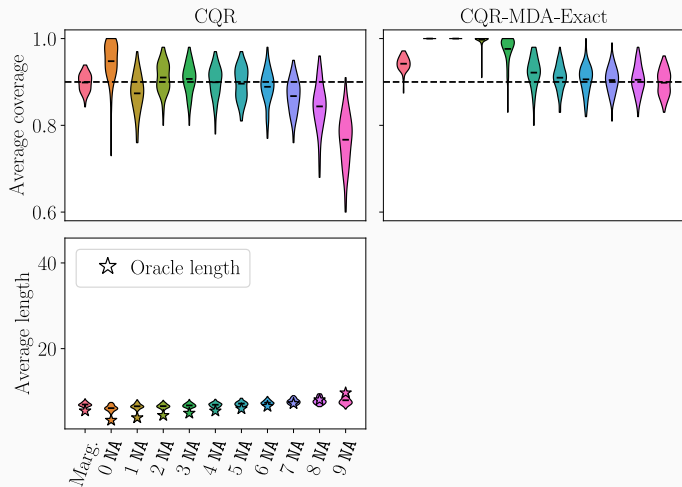
Synthetic experiments, 40% of missing values (Gaussian linear model, $d = 10$)



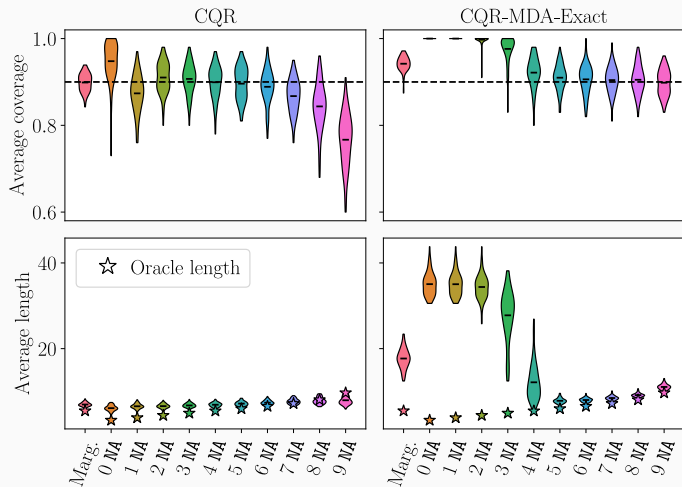
Synthetic experiments, 40% of missing values (Gaussian linear model, $d = 10$)



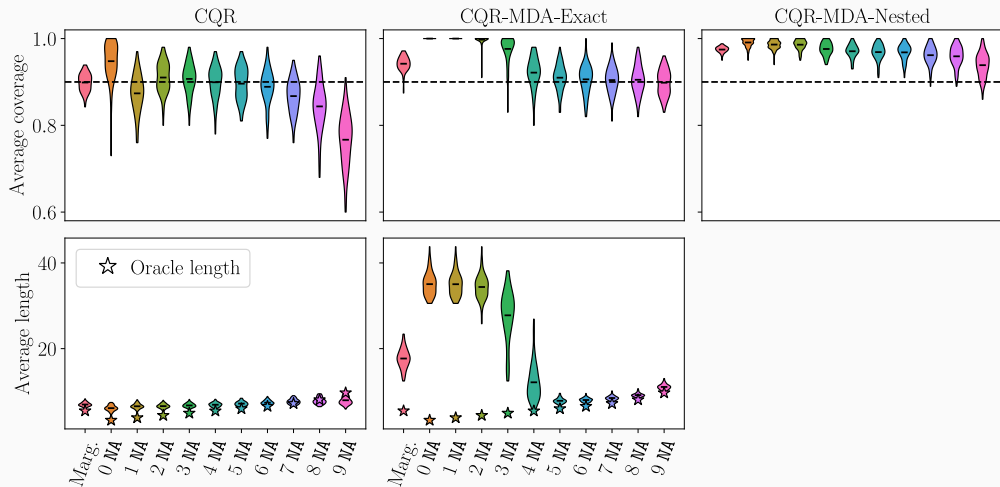
Synthetic experiments, 40% of missing values (Gaussian linear model, $d = 10$)



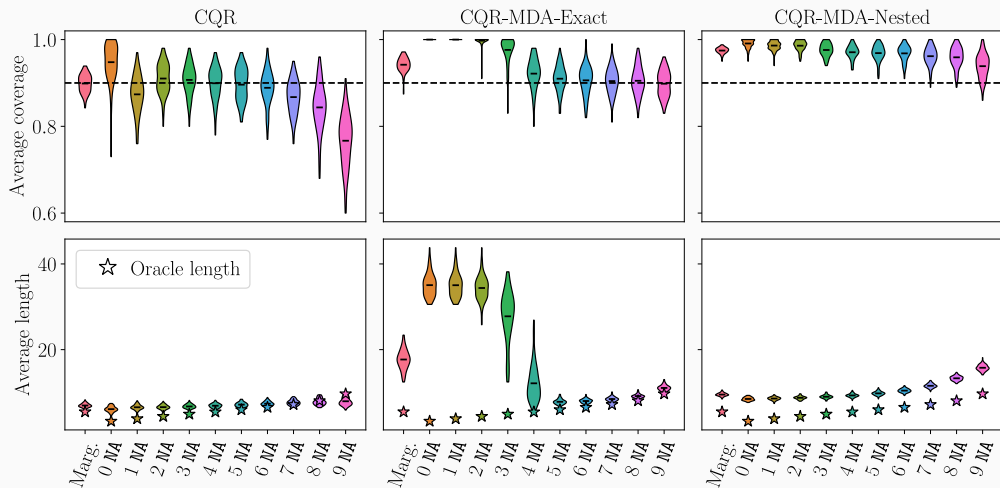
Synthetic experiments, 40% of missing values (Gaussian linear model, $d = 10$)



Synthetic experiments, 40% of missing values (Gaussian linear model, $d = 10$)



Synthetic experiments, 40% of missing values (Gaussian linear model, $d = 10$)

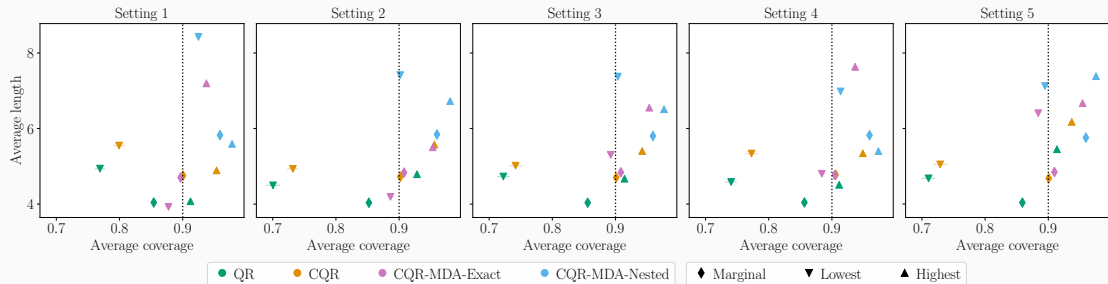


Simulation settings: beyond MCAR

- 6 variables (denote this set X_{missing}) out of 10 can be missing (the 4 others form the set X_{observed})
 - $X_{\text{missing}} = \{X_1, X_2, X_3, X_5, X_8, X_9\}$;
- Proportion of missing entries fixed to be 20%.

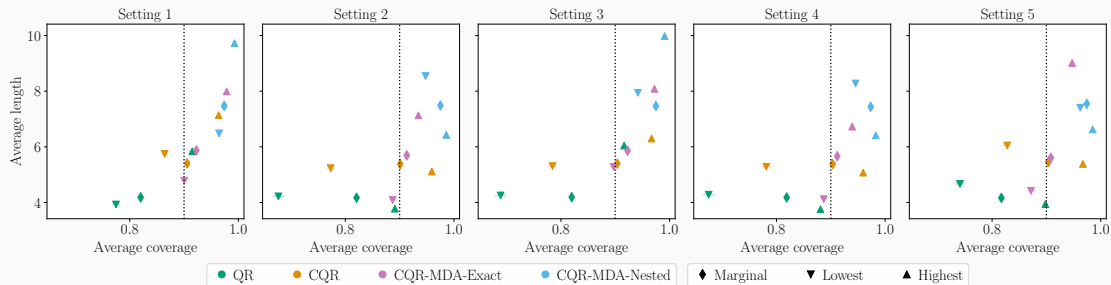
MAR missingness

- Probability of the variables in X_{missing} to be missing given by a logistic model of arguments X_{observed} .
- This setting is declined 5 times, with different weights for the logistic model.



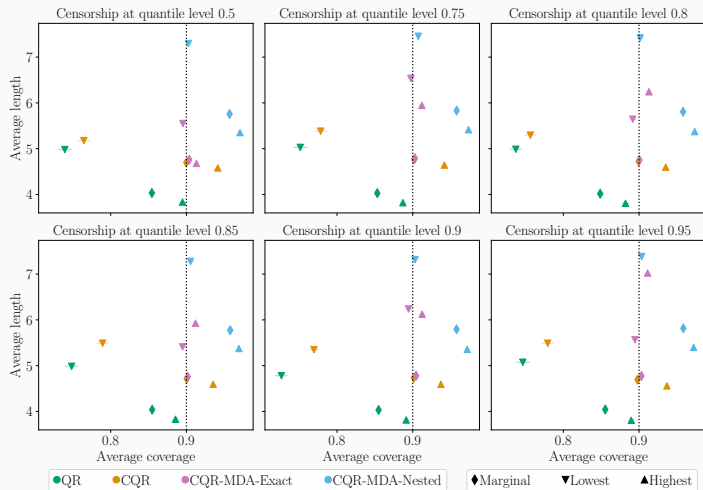
MNAR self masked missingness

- Probability of each variable in X_{missing} to be missing given by a logistic model of argument the same variable of X_{missing} .
- This setting is declined 5 times, with different weights for the logistic model.

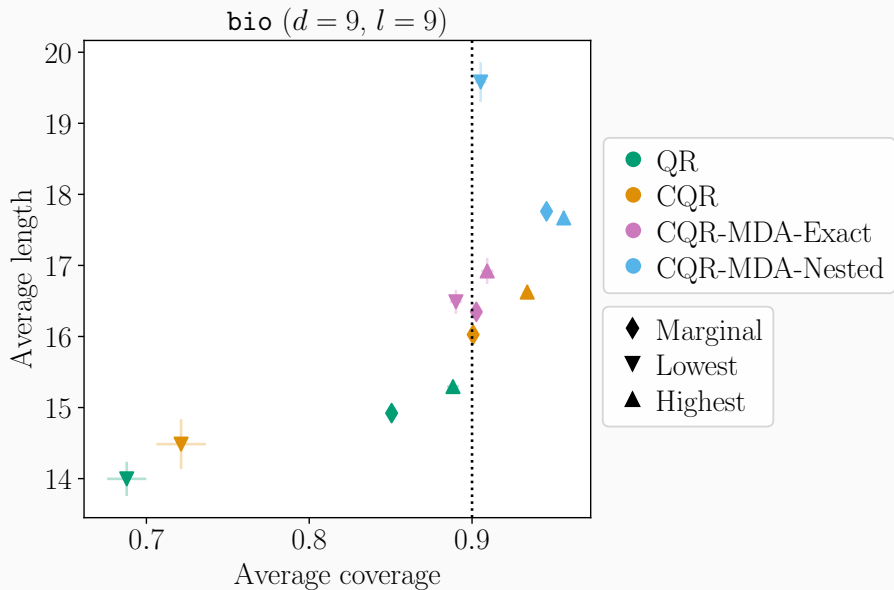


MNAR quantile censorship missingness

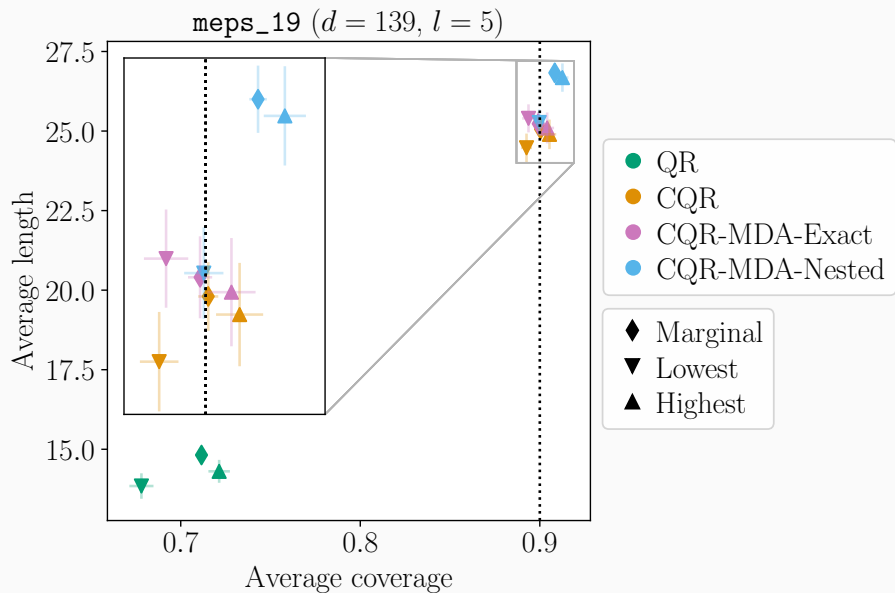
- Missing values are introduced at random in each q -quantile of the variables in X_{missing} .
- 6 different settings: q varies between 0.5, 0.75, 0.8, 0.85, 0.9 and 0.95.



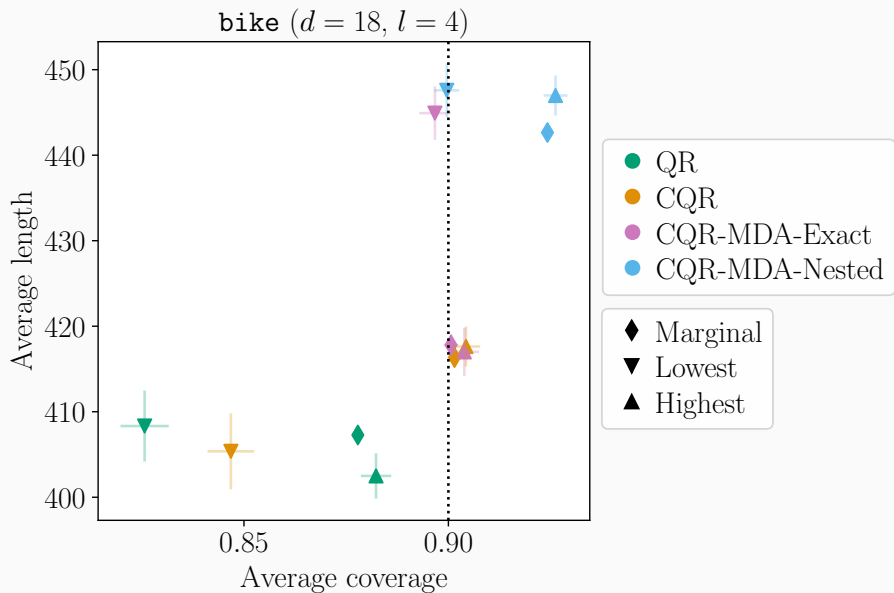
Semi-synthetic experiments



Meps_19 data set



Bike data set



TraumaBase®

Data set description i

- Age: the age of the patient (no missing values);
- Lactate: the conjugate base of lactic acid, upon arrival at the hospital (17.66% missing values);
- Delta_hemo: the difference between the hemoglobin upon arrival at hospital and the one in the ambulance (23.82% missing values);
- VE: binary variable indicating if a Volume Expander was applied in the ambulance. A volume expander is a type of intravenous therapy that has the function of providing volume for the circulatory system (2.46% missing values);
- RBC: a binary index which indicates whether the transfusion of Red Blood Cells Concentrates is performed (0.37% missing values);

- SI: the shock index. It indicates the level of occult shock based on heart rate (HR) and systolic blood pressure (SBP), that is $SI = \frac{HR}{SBP}$, upon arrival at hospital (2.09% missing values);
- HR: the heart rate measured upon arrival of hospital (1.62% missing values).