# Predictive uncertainty quantification with missing covariates

## On the hardness of distribution-free group conditional coverage

Margaux Zaffran

December 16, 2024

International Conference on Statistics and Data Science

UC Berkeley

Inria

l'X
ÉCOLE
POLYTECHNIQUE
IP PARIS

**Julie Josse**
*PreMeDICaL*
*Inria*

**Yaniv Romano**
*Technion – Israel Institute of Technology*

**Aymeric Dieuleveut**
*CMAP*
*École Polytechnique*

↪ Aymeric will present methodological results tomorrow at 9am in room Fregate!

**Distribution-free predictive uncertainty quantification**

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

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

  and $\mathcal{C}_\alpha$ should be as small as possible, in order to be informative.

▶ *Construction* of the predictive intervals should be
  - agnostic to the learning model[1]
  - agnostic to the data distribution

▶ *Validity* should be ensured
  - in finite samples
  - for all data distribution and underlying learnt model

---

[1] The underlying model can be any probabilistic model tailored for the application task at hand.

Conformal prediction (Vovk et al., 2005; Papadopoulos et al., 2002; Lei et al., 2018) builds an estimated predictive set $\widehat{C}_\alpha$ based on $n$ data points.

> **Conformal prediction achieves marginal validity (Vovk et al., 2005)**
>
> $\widehat{C}_\alpha$ outputted by conformal prediction is such that for any distribution $\mathcal{D}$ on $(\mathcal{X}, \mathcal{Y})$, it holds:
>
> $$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left( Y^{(n+1)} \in \widehat{C}_\alpha \left( X^{(n+1)} \right) \right) \geq 1 - \alpha.$$

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

$\widehat{C}_\alpha$ = estimated predictive set based on $n$ data points.

**Distribution-free $X$-conditional validity**

$\widehat{C}_\alpha$ achieves distribution-free $X$-conditional validity if for any distribution $\mathcal{D}$, it holds:

$$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}}\left(Y^{(n+1)} \in \widehat{C}_\alpha\left(X^{(n+1)}\right) | X^{(n+1)}\right) \overset{a.s.}{\geq} 1 - \alpha.$$

# Limits of distribution-free conditional predictive uncertainty quantification

**Impossibility results (Vovk, 2012; Lei and Wasserman, 2014)[2]**

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

$$\mathbb{P}_{\mathcal{D}^{\otimes(n)}} \left\{ \text{mes}\left( \widehat{C}_\alpha(x) \right) = \infty \right\} \geq 1 - \alpha.$$

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

[2]An analogous statement is also available for the classification framework.
Non exhaustive references.

$\widehat{C}_\alpha =$ estimated predictive set based on $n$ data points.

$\mathcal{G}$ a set of "groups" (i.e., define $G$ a random variable taking its values in $\mathcal{G}$).

> **Distribution-free $\mathcal{G}$-conditional validity ($\mathcal{G}$CV)**
>
> $\widehat{C}_\alpha$ achieves distribution-free $\mathcal{G}$-conditional validity if for any distribution $\mathcal{D}$ on $(\mathcal{X}, \mathcal{G}, \mathcal{Y})$, it holds that:
>
> $$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left( Y^{(n+1)} \in \widehat{C}_\alpha \left( X^{(n+1)}, G^{(n+1)} \right) | G^{(n+1)} \right) \overset{a.s.}{\geq} 1 - \alpha.$$

> **General $\mathcal{G}$CV hardness result** (Z., Josse, Romano and Dieuleveut, 2024)[3]
>
> If any $\widehat{C}_\alpha$ is distribution-free $\mathcal{G}$-conditionally valid then **for any distribution** $\mathcal{D}$, for any group $g \in \mathcal{G}$ such that $\mathcal{D}_{\mathcal{G}}(g) > 0$, it holds:
>
> $$\mathbb{P}_{\mathcal{D}^{\otimes(n+1)}} \left( \text{mes} \left( \widehat{C}_\alpha \left( X^{(n+1)}, g \right) \right) = \infty \right) \geq 1 - \alpha - \Delta_{g,n}$$
> $$\geq 1 - \alpha - \mathcal{D}_G(g)\sqrt{n+1}.$$

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

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

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

---
[3]An analogous statement is also available for the classification framework.

$G \perp\!\!\!\perp X$ **hardness result** (Z., Josse, Romano and Dieuleveut, 2024)

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

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

$Y \perp\!\!\!\perp G \,|X$ **hardness result** (Z., Josse, Romano and Dieuleveut, 2024)

If any $\widehat{C}_\alpha$ is MCV under $Y \perp\!\!\!\perp G \,|X$, then for any distribution $\mathcal{D}$ such that $Y \perp\!\!\!\perp G \,|X$, for any mask $m$ such that $\frac{1}{\sqrt{2}} \geq \mathcal{D}_G(g) > 0$, it holds:

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

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

Analogous statements are also available for the classification framework.

# Application to learning with missing covariates

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

| $Y$ | $X_1$ | $X_2$ | $X_3$ |
|---|---|---|---|
| 22 | 5 | 6 | 3 |
| 19 | 6 | 8 | NA |
| 19 | 5 | 3 | 6 |
| 7 | NA | 9 | NA |
| 13 | 4 | 9 | 0 |
| 20 | NA | NA | 1 |
| 9 | 8 | NA | 4 |

Mask $M =$

| $(M_1$ | $M_2$ | $M_3)$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 0 | 0 |
| 1 | 1 | 0 |
| 0 | 1 | 0 |

$\hookrightarrow 2^d$ potential masks.

$\hookrightarrow M$ can depend on $X$ or $Y$ (depending on the missing mechanism[4]).

$\Rightarrow$ Statistical and computational challenges.

---

[4]Three mechanisms connecting $X$ and $M$ from Rubin (1976), *Inference and missing data*, Biometrika

Impute-then-predict 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

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

$\hookrightarrow$ we consider an impute-then-predict pipeline in this work.
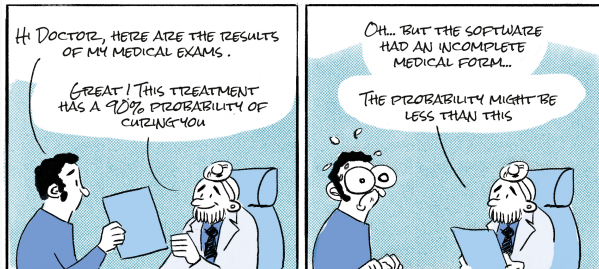
**Goal:** predict $Y^{(n+1)}$ with **confidence** $1 - \alpha$, i.e. build the smallest $\mathcal{C}_\alpha$ such that:

**1. Marginal Validity (MV)**

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

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

$$\mathbb{P}\left\{ Y^{(n+1)} \in \mathcal{C}_\alpha\left( X^{(n+1)}, M^{(n+1)} \right) \mid M^{(n+1)} \right\} \overset{a.s.}{\geq} 1 - \alpha. \tag{MCV}$$



Illustrations @theoremlinger

# Validities of 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)**

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

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

$$\mathbb{P}\left\{ Y^{(n+1)} \in \mathcal{C}_\alpha\left( X^{(n+1)}, M^{(n+1)} \right) \middle| M^{(n+1)} \right\} \overset{a.s.}{\geq} 1 - \alpha. \qquad \text{(MCV)}$$

|        | Existing approaches | New approach (Z., Josse, Romano and Dieuleveut, 2024) |
|--------|---------------------|-------------------------------------------------------|
| (MV)   | ✓ (Z., Dieuleveut, Josse, and Romano, 2023) | ✓ |
| (MCV)  | ✗ | ✓ under $M \perp\!\!\!\perp (X, Y)$ |

Thanks for listening and feel free to reach out to us!



Tomorrow at 9am in room Fregate: methodological results by Aymeric!

Chernozhukov, V., Wüthrich, K., and Zhu, Y. (2021). Distributional conformal prediction. *Proceedings of the National Academy of Sciences*, 118(48).

Gibbs, I., Cherian, J. J., and Candès, E. J. (2023). Conformal prediction with conditional guarantees. arXiv: 2305.12616.

Guan, L. (2022). Localized conformal prediction: a generalized inference framework for conformal prediction. *Biometrika*, 110(1).

Izbicki, R., Shimizu, G., and Stern, R. B. (2022). CD-split and HPD-split: Efficient conformal regions in high dimensions. *Journal of Machine Learning Research*, 23(87).

Jung, C., Noarov, G., Ramalingam, R., and Roth, A. (2023). Batch multivalid conformal prediction. In *International Conference on Learning Representations*.

## References ii

Kivaranovic, D., Johnson, K. D., and Leeb, H. (2020). Adaptive, Distribution-Free Prediction Intervals for Deep Networks. In *International Conference on Artificial Intelligence and Statistics*. PMLR.

Lei, J., G'Sell, M., Rinaldo, A., Tibshirani, R. J., and Wasserman, L. (2018). Distribution-Free Predictive Inference for Regression. *Journal of the American Statistical Association*.

Lei, J. and Wasserman, L. (2014). Distribution-free prediction bands for non-parametric regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1).

Papadopoulos, H., Proedrou, K., Vovk, V., and Gammerman, A. (2002). Inductive Confidence Machines for Regression. In *Machine Learning: ECML*. Springer.

Romano, Y., Barber, R. F., Sabatti, C., and Candès, E. (2020). With Malice Toward None: Assessing Uncertainty via Equalized Coverage. *Harvard Data Science Review*, 2(2).

Romano, Y., Patterson, E., and Candès, E. (2019). Conformalized Quantile Regression. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.

Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3).

Sesia, M. and Romano, Y. (2021). Conformal prediction using conditional histograms. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc.

Vovk, V. (2012). Conditional Validity of Inductive Conformal Predictors. In *Asian Conference on Machine Learning*. PMLR.

Vovk, V., Gammerman, A., and Shafer, G. (2005). *Algorithmic Learning in a Random World*. Springer US.

Zaffran, M., Dieuleveut, A., Josse, J., and Romano, Y. (2023). Conformal prediction with missing values. In *Proceedings of the 40th International Conference on Machine Learning*. PMLR.

Zaffran, M., Josse, J., Romano, Y., and Dieuleveut, A. (2024). Predictive uncertainty quantification with missing covariates. *arXiv:2405.15641*.