

Perspectives & Roundtables

MAPIE General Meeting

2023/11/24

Thibault Cordier
PhD, Data Scientist
tcordier@quantmetry.com

Vincent Blot
PhD Student, Data Scientist
vblot@quantmetry.com

Louis Lacombe
Data Scientist
llacombe@quantmetry.com

Candice Moyet
Data Scientist
cmoyet@quantmetry.com

Nicolas Brunel
Scientific Director
nbrunel@quantmetry.com



Session 1

MAPIE

Introductory session:
Where does MAPIE stand?



15 min

Session 1: Where does MAPIE stand

Subject: present an overview of the MAPIE package (i.e., cartography of features and roadmap)

1) Presentation of MAPIE:

- MAPIE Team
- Missions of MAPIE Team
- History of MAPIE from 2021 to 2023
- Decision tree / Feature matrix of MAPIE

2) Future directions:

- Tentative Roadmap up to 2024 (main focus, priorities, etc.)
- What are your problem with MAPIE?

MAPIE Team & Contributors

MAPIE Team 2024



Thibault Cordier



Vincent Blot



Louis Lacombe



Candice Moyet



Nicolas Brunel

Emeritus Contributors of MAPIE



Vianney Taquet



Sofiane Ziane



JumpingDino



Grégoire Martinon



Thomas Morzadec



Arnaud Capitaine



dan1elherbst



remiaddon



AdirthaBorghain



alize-papp



cmougan



AndreaPi

Contributors of MAPIE

MAPIE - Model Agnostic Prediction Interval Estimator



scikit-learn-contrib / MAPIE Public



Fork

80



Starred

1k



Python library, open source and scikit-learn compatible, for estimating confidence intervals in classification and regression tasks.

- **MAPIE** is an open-source Python library hosted on scikit-learn-contrib project that allows you to:
 - 1) easily **compute conformal prediction intervals/sets** with controlled marginal coverage rate for regression, classification (binary and multi-class) and time series.
 - 2) easily **control risks** (such as coverage, recall or any other non-monotone risk) for more complex tasks (multi-label classification, semantic segmentation, ...).
 - 3) easily **wrap any model** (*scikit-learn, tensorflow, pytorch, ...*).
- **MAPIE** is designed and conceived for **academic and industrial uses**.

SPONSORS



Quantmetry
Part of Capgemini Invent



Missions of MAPIE Team

Github Project

scikit-learn-contrib / MAPIE

Code Issues 29 Pull requests 8 Discussions Actions Projects 2 Wiki Security

MAPIE Public

27 branches 24 tags Go to file Add file Code

README.rst

MAPIE - Model Agnostic Prediction Interval Estimator

Quantifying the uncertainties and controlling the risks of ML model predictions is of crucial importance for developing systems. Uncertainty quantification (UQ) involves all the stakeholders who develop and use AI models.

MAPIE is an open-source Python library hosted on scikit-learn-contrib project that allows you to:

- easily estimate conformal prediction intervals (or prediction sets) given a degree of confidence or risk for single-class problems [3-9].
- easily control risks (such as coverage, recall or any other non-monotone risk) by estimating relevant prediction sets.
- easily wrap your favorite scikit-learn-compatible model for the purposes just mentioned.

Here's a quick instantiation of MAPIE models for regression and classification problems related to uncertainty quantification:

```
# Uncertainty quantification for regression problem
from mape_regression import MapeRegressor
mape_regression = MapeRegressor(estimator=regressor, method='plus', cv=5)

# Uncertainty quantification for classification problem
from mape_classification import MapeClassifier
mape_classifier = MapeClassifier(estimator=classifier, method='score', cv=5)
```

Scientific Publications

MAPIE: an open-source library for uncertainty quantification

Vianney Daquet¹, Vincent Blot¹, Thomas Morzadec¹, Louis Lacombe¹, Armand Capitant², Nicolas Brunel^{1,2}

¹ Quantmetry, 52, rue d'Anjou, 75008, Paris
² Laboratoire de Mathématiques et de Modélisation d'Evry, EN

Abstract

Estimating uncertainties associated with the predictions (ML) models is of crucial importance to assess their reliability. In this submission, we introduce MAPIE (Model Agnostic Prediction Interval Estimator), an open-source Python library that quantifies the uncertainties and controls the risks of ML models for single-output regression and multi-class classification. MAPIE implements conformal prediction methods, also computes uncertainties with strong theoretical guarantees on and with mild assumptions on the model or on the underlying data. MAPIE is hosted on scikit-learn-contrib and is fully "scikit-learn compatible", i.e. it accepts any type of regressor or classifier compatible with the scikit-learn API.

Proceedings of Machine Learning Research 2024:1-33, 2023 Conformal and Probabilistic Prediction with Applications

Flexible and Systematic Uncertainty Estimation with Conformal Prediction via the MAPIE library

Thibault Caviller¹, Vincent Blot¹, Louis Lacombe¹, Thomas Morzadec¹, Armand Capitant², Nicolas Brunel^{1,2}

¹ Quantmetry, 52, rue d'Anjou, 75008, Paris, France
² Laboratoire de Mathématiques et de Modélisation d'Evry, ENSIE, Paris-Saclay University
³ Laboratoire interdisciplinaire des sciences du numérique, CNRS, Paris-Saclay University

Editor: Harris Papadopoulos, Khong An Nguyen, Henrik Boström and Lars Caron

Abstract

Conformal prediction (CP) is an attractive theoretical framework for estimating the uncertainty of machine learning models. In this paper, we introduce the MAPIE library, which provides a flexible and systematic way to estimate the uncertainty of machine learning models using conformal prediction.

Code Maintenance

Filters: 23 Open 141 Closed

Labels 13 Milestones New issue

3 Bugs to be fixed (High Priority)

- Return resampled predictions even when alpha is None #252 opened by certmanco (developers) (enhancement) (good first issue)
- Trivial typo in quantile_regression.py time_series_regression.py #264 opened 3 weeks ago by Carl-McCubride-Ellis (1 linked pull request)
- Tutorials in the wrong order #348 opened on Sep 9 by LacombeLouis (bug) (development)
- Adaptive Conformal Predictions for Time Series (enhancement) #334 opened on Aug 3 by thibaultcaviller (developers) (documentation)

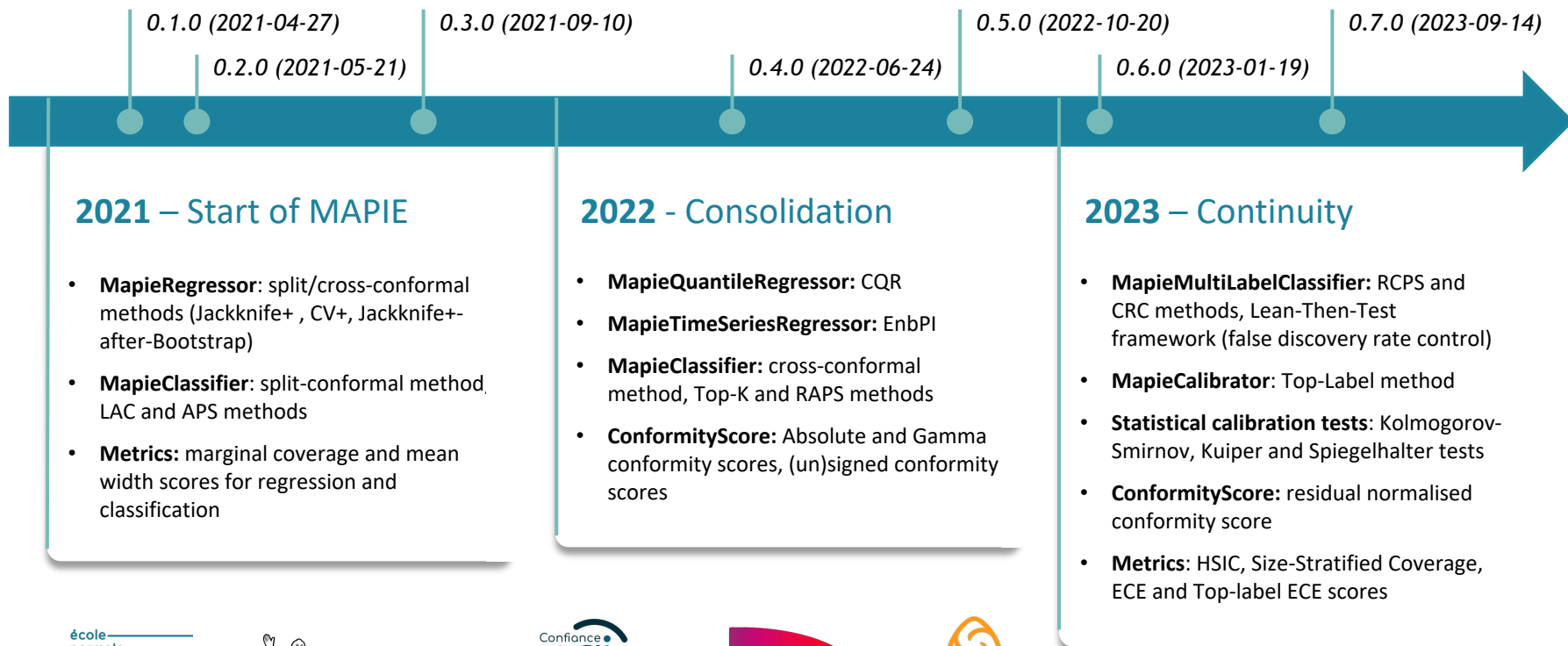
4 Bugs to be fixed (Low Priority)

- Fix _fit_params and _predict_params #252 opened by certmanco (contributors) (enhancement) (good first issue)
- Automated estimation of number of resamplings gives the size of the train data #251 opened by team-maroqu (contributors) (enhancement) (good first issue)
- Giving a fraction of samples instead of a number of samples in the subsample class #238 opened by gmartinoQM (developers) (enhancement) (good first issue)

4 Methods to be implemented (short Term)

- Adaptive Conformal Predictions for Time Series #234 opened by thibaultcaviller (developers) (enhancement) (good first issue)
- Support Binary classification #236 opened by path (contributors) (enhancement) (good first issue)
- Covariate shift conformal #272 opened by gmartinoQM (developers) (enhancement)

History of MAPIE - from 2021 to 2023



2021 – Start of MAPIE

- **MapieRegressor**: split/cross-conformal methods (Jackknife+ , CV+, Jackknife+ after-Bootstrap)
- **MapieClassifier**: split-conformal method, LAC and APS methods
- **Metrics**: marginal coverage and mean width scores for regression and classification

2022 - Consolidation

- **MapieQuantileRegressor**: CQR
- **MapieTimeSeriesRegressor**: EnbPI
- **MapieClassifier**: cross-conformal method, Top-K and RAPS methods
- **ConformityScore**: Absolute and Gamma conformity scores, (un)signed conformity scores

2023 – Continuity

- **MapieMultiLabelClassifier**: RCPS and CRC methods, Lean-Then-Test framework (false discovery rate control)
- **MapieCalibrator**: Top-Label method
- **Statistical calibration tests**: Kolmogorov-Smirnov, Kuiper and Spiegelhalter tests
- **ConformityScore**: residual normalised conformity score
- **Metrics**: HSIC, Size-Stratified Coverage, ECE and Top-label ECE scores

What can you find in the release 0.7.0 of MAPIE?

Summary table of algorithms implemented in MAPIE

Task	Feature	Algorithm	Reference
PI/PS	<div style="background-color: #f08080; padding: 2px;">MapieRegressor</div> <div style="background-color: #0056b3; color: white; padding: 2px;">MapieClassifier</div>	Jackknife/CV+	Rina Foygel Barber, Emmanuel J. Candès, Aaditya Ramdas, and Ryan J. Tibshirani. "Predictive inference with the jackknife+." <i>Ann. Statist.</i> , 49(1):486–507, (2021).
		Jackknife/CV+ ab	Kim, Byol, Chen Xu, and Rina Barber. "Predictive inference is free with the jackknife+–after-bootstrap." <i>Advances in NeurIPS</i> 33 (2020): 4138-4149.
Prediction intervals (PI)	<div style="background-color: #ff9933; padding: 2px;">AbsoluteConformityScore</div> <div style="background-color: #ff9933; padding: 2px;">GammaConformityScore</div> <div style="background-color: #ff9933; padding: 2px;">ResidualNormalizedScore</div>	Absolute Score	Vovk, Vladimir, Alexander Gammerman, and Glenn Shafer. <i>Algorithmic Learning in a Random World</i> . Springer Nature, 2005
		Gamma Score	Cordier, Thibault, Vincent Blot, Louis Lacombe, Thomas Morzadec, Arnaud Capitaine, Nicolas Brunel "Flexible and Systematic Uncertainty Estimation with Conformal Prediction via the MAPIE library", <i>COPA</i> (2023)
		Normalised Score	Papadopoulos, Harris, Proedrou, Kostas, Vovk, Volodya, and Gammerman, Alex. "Inductive confidence machines for regression". In <i>Machine Learning: ECML</i> (2002).
	<div style="background-color: #f08080; padding: 2px;">MapieTimeSeriesRegressor</div>	EnbPI	Xu, Chen, and Yao Xie. "Conformal prediction interval for dynamic time-series." <i>International Conference on Machine Learning</i> . PMLR, (2021).
	<div style="background-color: #f08080; padding: 2px;">MapieQuantileRegressor</div>	CQR	Romano, Yaniv, Evan Patterson, and Emmanuel Candès. "Conformalized quantile regression." <i>Advances in neural information processing systems</i> 32 (2019).
Prediction sets (PS)	<div style="background-color: #0056b3; color: white; padding: 2px;">MapieClassifier</div>	LAC / LABEL	Sadinle, Mauricio, Jing Lei, and Larry Wasserman. "Least ambiguous set-valued classifiers with bounded error levels." <i>Journal of the American Statistical Association</i> 114.525 (2019): 223-234.
		APS	Romano, Yaniv, Matteo Sesia, and Emmanuel Candès. "Classification with valid and adaptive coverage." <i>Advances in NeurIPS</i> 33 (2020): 3581-3591.
		Top-K	Angelopoulos, Anastasios, et al. "Uncertainty sets for image classifiers using conformal prediction." <i>International Conference on Learning Representations</i> (2021).
		RAPS	Angelopoulos, Anastasios, et al. "Uncertainty sets for image classifiers using conformal prediction." <i>International Conference on Learning Representations</i> (2021).
Control Risks (CR)	<div style="background-color: #0056b3; color: white; padding: 2px;">MapieMultiLabelClassifier</div>	RCPS	Bates, Stephen, et al. "Distribution-free, risk-controlling prediction sets." <i>Journal of the ACM (JACM)</i> 68.6 (2021): 1-34.
		CRC	Angelopoulos, Anastasios N., Stephen, Bates, Adam, Fisch, Lihua, Lei, and Tal, Schuster. "Conformal Risk Control." (2022).
		LTT	Angelopoulos, Anastasios N., Stephen, Bates, Emmanuel J. Candès, et al. "Learn Then Test: Calibrating Predictive Algorithms to Achieve Risk Control." (2022).
Calib.	<div style="background-color: #0056b3; color: white; padding: 2px;">MapieCalibrator</div>	Top-label	Gupta, Chirag, and Aaditya K. Ramdas. "Top-label calibration and multiclass-to-binary reductions." <i>arXiv preprint arXiv:2107.08353</i> (2021).

Release 0.7.0

Decision tree / Feature matrix of MAPIE

If I have a predictive model, MAPIE can give me guarantees and insights on the quality of the predictions.



1. Compute uncertainty intervals / sets:
"I want to ensure that the true labels are covered."

Task	Split	Cross	Metrics
Regression	-	Adaptive	Mc/Ma
	Time Series	Adaptive	Mc
Classification	Binary	See "3. Calibrate my model" (1)	
	Multiclass	Adaptive	Mc/Ma
	Multilabel	See "2. Control a risk" (2)	
Object Detection	Binary	Calib.: OK	Mc/Ma
	Multiclass	Class.: OK	Mc/Ma
Instance Segmentation		Class.: OK	Mc/Ma

2. Control the risk of my model: "I want to guarantee that my risk is under control with a probability guarantee."

Task	Risk	Split	Cross	Metrics
Multilabel Classification	with Precision and Recall Guarantees (2)	Adaptive	No literature	Mc
Selective Regression	with MSE Control	No literature	No literature	-
	with OOD Detection	No literature	No literature	-
Selective Classification	with Accuracy Control	No literature	No literature	-
	with OOD Detection	No literature	No literature	-
Selective Generation	with Auxiliary Control	No literature	No literature	-
Object Detection	with Coverage, and Recall Guarantees	No literature	No literature	-
Instance Segmentation	with mIOU Guarantee	No literature	No literature	-

3. Calibrate my model:
"I want to ensure that scores given by my model are probabilities"

Task	Split	Cross	Metrics
Classification	Binary (1)	Available on scikit-learn	Mc/H
	Multiclass	Adaptive	Mc
	Multilabel	No literature	-

4. Test a hypothesis: "I want to make sure that assumptions on my data are satisfied or detect deviation of the model."

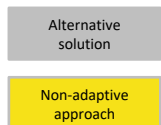
Test	Split	Cross	Metrics
Calibration test of binary classifier	Adaptive	Available in MAPIE	H
Exchangeability test (distribution drift)	No literature	No literature	-
Performance stability test (perf. drift)	No literature	No literature	-
Anomaly detection	No literature	No literature	-

Colour code:

Not available	Available in MAPIE
---------------	--------------------

Legend

Colour code:



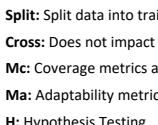
Alternative solution



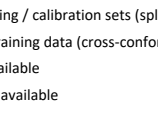
Non-adaptive approach



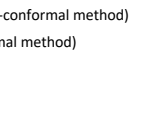
Adaptive approach



No literature



Not available in MAPIE



Available on scikit-learn

Split: Split data into training / calibration sets (split-conformal method)

Cross: Does not impact training data (cross-conformal method)

Mc: Coverage metrics available

Ma: Adaptability metrics available

H: Hypothesis Testing

Tentative Roadmap up to 2024

0.8.0 (2024-x-x)

Main Focus - P1

- **[Session 2]: Mondrian Conformal Prediction:** validity within categories thanks to Mondrian approach for both regression and classification
- **[Session 2]: Adaptive Conformal Prediction:** approaches based on non-conformity scores distribution estimation
- **[Session 3]: Hypothesis Testing:** toolbox of hypothesis tests for exchangeability and performance drift
- **Risk control:** selective regression and classification with LTT

Priorities - P2

- **Interoperability:** support natively pytorch, tensorflow, transformers, etc.
- **Large Scale Deployment:** accelerate code, be robust to memory problems...
- **Use Cases:** instance segmentation, object detection, generation (tabular data, text), etc.

To be completed...

- Other propositions, with your contributions and your comments!

Schedule

Session 1: Overview of MAPIE

 15 min



Session 2: Adaptive Uncertainty Quantification

 15 min + 10 min of Q&A




Session 3: Hypothesis Testing

 10 min + 10 min of Q&A



Session 4: Round Table

 15 min

Teams Chat:
Q&A at the end
of each session



Session 2

MAPIE

Session 2:
Adaptive Uncertainty Quantification



15 min

Session 2: Adaptive Uncertainty Quantification

1) Presentation:

- What is Adaptive Uncertainty Quantification? Why?
- What already exists in MAPIE

2) Future directions:

- Methods for conditional coverage
- Methods based on estimated scores distribution
- Metrics for adaptive uncertainty

3) Opening: round table

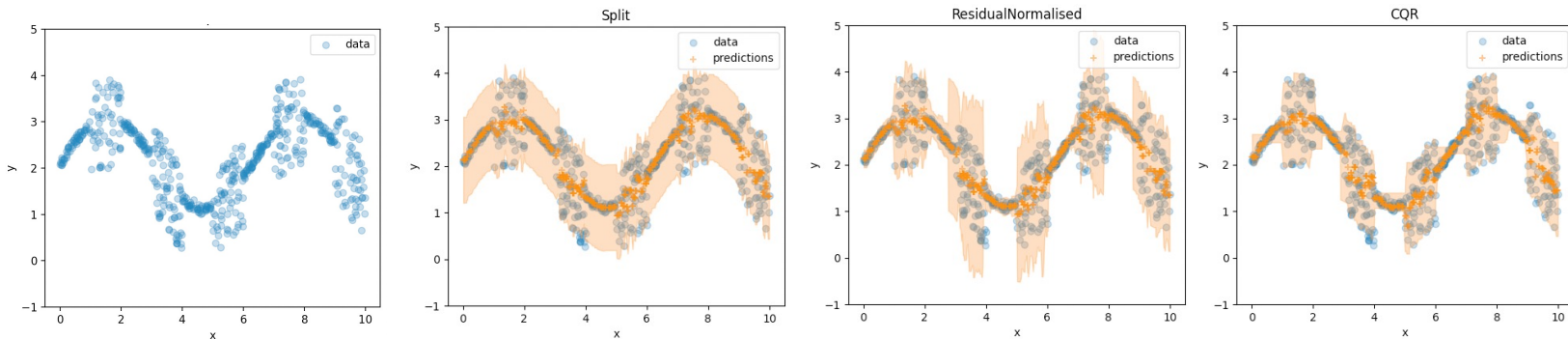
What is Adaptive Uncertainty Quantification? Why?

The ideal framework for uncertainty quantification (UQ):

- 1) Be distribution-free ✓
- 2) Be valid in finite samples ✓
- 3) Satisfy the conditional coverage guarantee: ?
 - $\mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1}) | X_{n+1} = x) = 1 - \alpha$

What is a good adaptive method for UQ?

- 1 Predictions sets should be of **various sizes**.
- 2 **Size of the prediction sets** \Leftrightarrow **Uncertainty of the model**.
- 3 **Prediction coverage should be guaranteed locally** (for any prediction), and not just globally (on average).



Issues and limitations:

- **Conditional coverage** is proven to be **impossible to satisfy** ([1-3]) without assumptions on the distribution or algorithm.
- **Marginal coverage** is easily obtained with conformal prediction but without local interpretation: $\mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1})) = 1 - \alpha$

[1] Foygel Barber, R., Candes, E. J., Ramdas, A., & Tibshirani, R. J. (2021). The limits of distribution-free conditional predictive inference. *Information and Inference: A Journal of the IMA*, 10(2), 455-482.

[2] Vladimir Vovk. (2012) Conditional validity of inductive conformal predictors. In *Asian conference on machine learning*, pages 475-490.

[3] Vovk, V., Gammerman, A., & Shafer, G. (2005). *Algorithmic learning in a random world* (Vol. 29). New York: Springer.

What already exists in MAPIE for adaptive uncertainty quantification

Our MAPIE methods are mainly based on **split-CP** with **calibration dataset** $D_n^{Cal} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ and NCS $\{s_i = s(x_i, y_i)\}_{1 \leq i \leq n}$. They compute the **quantiles of the NCS** $Q_{n,\alpha}^-/Q_{n,\alpha}^+$ to estimate the **bounds of the prediction intervals** $\hat{C}_{n,\alpha}(X_{n+1})$ given a test data X_{n+1} .

Regression Methods in MAPIE

- **Conformalized Quantile Regression (CQR)**, *Romano et al. 2019*

$$\hat{C}_{n,\alpha}(X_{n+1}) = [\hat{q}_{\alpha_{low}}(X_{n+1}) - Q_{n,\alpha}^-, \hat{q}_{\alpha_{up}}(X_{n+1}) + Q_{n,\alpha}^+]$$

$$s(x, y) = \max(y - \hat{q}_{\alpha_{low}}(x), \hat{q}_{\alpha_{up}}(x) - y)$$

- **Residual Normalized Score**, *Lei et al. 2016*

$$\hat{C}_{n,\alpha}(X_{n+1}) = [\hat{\mu}(X_{n+1}) - Q_{n,\alpha}^- * \hat{\sigma}(X_{n+1}), \hat{\mu}(X_{n+1}) + Q_{n,\alpha}^+ * \hat{\sigma}(X_{n+1})]$$

$$s(x, y) = \frac{|y - \hat{\mu}(x)|}{|\hat{\sigma}(x)|}$$

- **Gamma Score**, *Cordier et al. 2023*

$$\hat{C}_{n,\alpha}(X_{n+1}) = [\hat{\mu}(X_{n+1}) * (\mathbf{1} - Q_{n,\alpha}^-), \hat{\mu}(X_{n+1}) * (\mathbf{1} + Q_{n,\alpha}^+)]$$

$$s(x, y) = \frac{|y - \hat{\mu}(x)|}{|\hat{\mu}(x)|}$$

Classification Methods in MAPIE

- **Adaptive Prediction Sets (APS)**, *Romano et al. 2020*

$$\hat{C}_{n,\alpha}(X_{n+1}) = \{\pi_1, \dots, \pi_k\} \quad \text{where } k = \inf\{s(x_n, k) \geq Q_{n,\alpha}\}$$

$$s(x_i, k) = \sum_{j=1}^k \hat{\mu}(x_i)_{\pi_j} \quad \text{where } \forall j \geq k \quad \hat{\mu}(x_i)_{\pi_j} \geq \hat{\mu}(x_i)_{\pi_k}$$

- **Regularized APS (RAPs)**, *Angelopoulos et al. 2020*

$$\hat{C}_{n,\alpha}(X_{n+1}) = \{\pi_1, \dots, \pi_k\} \quad \text{where } k = \inf\{s(x_n, k) \geq Q_{n,\alpha}\}$$

$$s(x_i, k) = \sum_{j=1}^k \hat{\mu}(x_i)_{\pi_j} + \lambda(k - k_{reg})^+ \quad \text{where } \forall j \geq k \quad \hat{\mu}(x_i)_{\pi_j} \geq \hat{\mu}(x_i)_{\pi_k}$$

Issues and limitations:

- These methods requires **an auxiliar model** (calculation costs, depend on their performance, must be trained externally, on other calibration data, ...).
- They are based only **on marginal coverage guarantee** in their design.

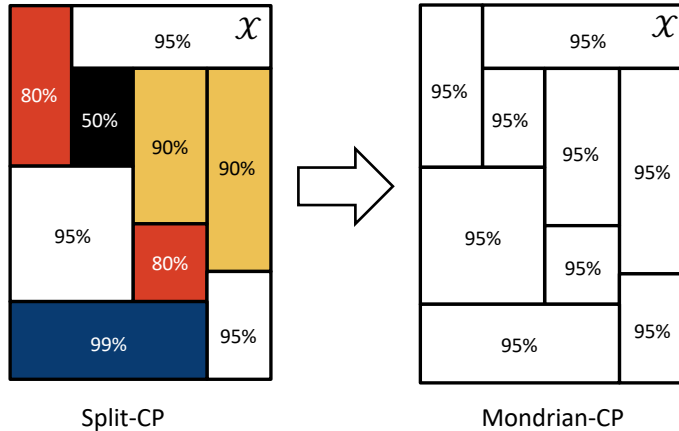
Methods for conditional coverage

Goal: Group-conditional coverage or $\mathbb{P}(Y_{n+1} \in \widehat{C}(X_{n+1}) | Z_{n+1} = z) = 1 - \alpha$ when Z_{n+1} is a categorical variable.

Mondrian Conformal Predictors

[Algorithmic Learning in a Random World, Vovk et al. 2005\(v1\), 2022\(v2\) \[1\]](#)

- Partition of \mathcal{X} \rightarrow “conformalization” on each part of the partition



Ideas / Openings

- **Mondrian-CP Process:** $\forall \mathcal{P}_i \in \mathcal{P}(\mathcal{X})$,
 1. Define α_i a risk and $n_i = |\mathcal{P}_i|$
 2. Estimate \mathcal{S}_i the NCS of the observations in \mathcal{P}_i
 3. Estimate Q_{n_i, α_i} the quantile $(1 - \alpha_i)$ of \mathcal{S}_i
- Equivalent to applying several Split-CPs to each part.
- **Generalizing the conditioning criterion:**
 - Based on an exogenous criterion
 - Based on a partition of the \mathcal{X} space
 - Based on a partition of the \mathcal{Y} space
 - Based on a partition of the $\mathcal{X} \times \mathcal{Y}$ space
- **Integration perspectives:**
 - Applicable to regression and classification tasks
 - Needs to be adapted for flexible integration and use

[1] Vovk, V., Gammerman, A., & Shafer, G. (2005/2022). Algorithmic learning in a random world (Vol. 29). New York: Springer.

Methods based on estimated scores distribution (without covariate shift)

Goal: Generalizing group-conditional coverage with a proxy / a relaxation of conditional coverage.

Proxy by estimating conditioned NCS distribution

Adaptive Conformal Prediction by Reweighting Nonconformity Scores, Amoukou et al. 2023 [1]

- The prediction set is built as follows:

$$C_{n,\alpha}(X_{n+1}) = \left\{ \mathbf{y} \mid s(X_{n+1}, \mathbf{y}) \leq \widehat{Q} \left(\mathbf{1} - \tilde{\alpha}; \widehat{F}_S(\cdot \mid X_{n+1} = x) \right) \right\}$$

where $\widehat{F}_S(s \mid X_{n+1} = x) = \sum_{i=1}^n w(x_i, x) \mathbf{1}\{s_i \leq s\}$ is the conditional f.d.r. estimated by a Random Forest, and $\tilde{\alpha}$ is selected for reaching target marginal coverage $1 - \alpha$.

- Conditional-training can be obtained with additional correction, and asymptotic conditional coverage.
- Extension of **Localized conformal prediction: a generalized inference framework for conformal prediction**, Guan 2022 [2] based on Nadaraya-Watson estimators.

Interpolation between marginal and conditional coverage overlapping groups

Conformal Prediction with Conditional Guarantees, Gibbs et al. 2023 [3]

- Generalization of Mondrian to overlapping groups by relaxing conditional coverage:

$$\mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1}) \mid X_{n+1} = x) = 1 - \alpha, \quad \text{for all } x$$
$$\iff$$

$$\mathbb{E}[f(X_{n+1})(\mathbf{1}\{Y_{n+1} \in \hat{C}(X_{n+1})\} - (1 - \alpha))] = 0, \quad \text{for all measurable } f.$$

- The prediction set is built as follows:

$$C_{n,\alpha}(X_{n+1}) = \left\{ \mathbf{y} \mid s(X_{n+1}, \mathbf{y}) \leq \widehat{g}_S(X_{n+1}, \mathbf{y})(X_{n+1}) \right\}$$

by replacing the constant quantile of prediction scores $\{s_i = s(x_i, y_i)\}$ by a quantile regression $g(\cdot) \in \mathcal{F}$ (a RKHS) on X_{n+1} (with pinball loss):

- Possibility of control for finite (e.g. partition of sets gives back Mondrian) and infinite dimensional function space \mathcal{F} .

[1] Amoukou, S. I., & Brunel, N. J. (2023). Adaptive Conformal Prediction by Reweighting Nonconformity Score. arXiv preprint arXiv:2303.12695.

[2] Guan, L. (2023). Localized conformal prediction: A generalized inference framework for conformal prediction. Biometrika, 110(1), 33-50.

[3] Gibbs, I., Cheriai, J. J., & Candès, E. J. (2023). Conformal Prediction With Conditional Guarantees. arXiv preprint arXiv:2305.12616.

What metrics already exist in MAPIE for adaptive uncertainty

Why: To find out whether the model is **uniformly good at being adaptive** to model uncertainty

How: Global measure expressing local coverage (without depending on X or other parameters)

Existing in MAPIE

- **Size-Stratified Coverage (SSC)**, Angelopoulos et al. 2021 [1]

$$\text{SSC}(X, Y; \mathcal{G}) = \min_{g \in \mathcal{G}} \frac{1}{|J_g|} \sum_{i \in J_g} \mathbb{1}_{\{Y_i \in C(X_i)\}}$$

- **Hilbert-Schmidt Independence Criterion (HSIC)**, usage proposed by Feldman et al. 2021 [2]

Correlation measure between the coverage (X) and the interval size ($\ell(X)$). By considering two separable RKHS on X and $\ell(X)$, HSIC is defined as the Hilbert Schmidt norm of the cross-covariance operator.

$$\text{HSIC}(X, \ell(X); \mathcal{F}, \mathcal{G}) = \|\mathbb{C}_{X\ell(X)}\|_{HS}^2$$

- **[SOON] Coverage Width-based Criterion (CWC)**, usage proposed by Jensen et al. 2022 [3]

Trade-off between the prediction interval normalized average width (PINAW) and the prediction interval coverage probability (PICP).

$$\text{CWC}(\eta) = (1 - \text{PINAW})e^{-\eta(\text{PICP} - (1-\alpha))^2}$$

Metrics related to Mondrian-CP

- **Group-Conditional Coverage:**

Given a partition of the observations with respect to a categorical criterion (sub-groups) defined by an exogenous rule given by user, compute the **coverage within the categories**:

$$\forall z \in \mathcal{Z}, \\ \text{cov}(z) = \mathbb{E}_{X,Y} [\mathbb{1}_{\{Y \in C(X)\}} \mid Z = z]$$

Issues and limitations:

- These measures are **not easy to use and interpret**; they require arbitrary binarization, a mixing parameter or depend on substitution variables.

[1] Angelopoulos, A. N., & Bates, S. (2023). Conformal prediction: A gentle introduction. Foundations and Trends® in Machine Learning, 16(4), 494-591.

[2] Feldman, S., Bates, S., & Romano, Y. (2021). Improving conditional coverage via orthogonal quantile regression. Advances in neural information processing systems, 34, 2060-2071.

[3] Jensen, V., Bianchi, F. M., & Anfinson, S. N. (2022). Ensemble conformalized quantile regression for probabilistic time series forecasting. IEEE Transactions on Neural Networks and Learning Systems.

Metrics for adaptive uncertainty

Why: To find out whether the model is **uniformly good at being adaptive** to model uncertainty

How: Global measure expressing local coverage (without depending on X or other parameters)

Ideas / Openings

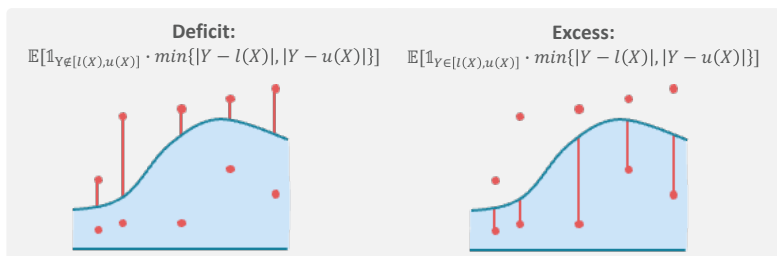
- **Deficit and Excess, Seedat et al. 2023 [1]**

Deficit: interval shortfall, when the true value y lies outside the predicted interval.

- measures the **under-coverage** of the prediction sets.

Excess: additional width included not needed to capture the true value y

- measures the **over-coverage** of the prediction sets.



- **Correlation between interval size and residuals**

Beyond measuring adaptability

Statistical inference for fairness auditing, Cherian et al. 2023 [2]

- Construct a statistical certificate for controlling the disparity of a performance metric $L(\hat{\mu}(X), Y)$ between a group G and the global target of the model:

$$\underbrace{\epsilon(G)}_{\text{disparity}} := \underbrace{\mathbb{E}_P[L(f(X), Y) \mid (X, Y) \in G]}_{\text{group-specific}} - \underbrace{\theta_P}_{\text{target}}.$$

- Bootstrap is used for computing lower bound ϵ_{lb} such that:

$$\lim_{n \rightarrow \infty} \mathbb{P}(\epsilon_{lb} \leq \epsilon(G) \forall G) \geq 1 - \alpha$$

- and we can test $H_0(G): \epsilon(G) \leq \epsilon$, with adapted threshold and FWER gives:

$$\lim_{n \rightarrow \infty} \mathbb{P}(\exists G \text{ falsely certified}) \leq \alpha$$

[1] Seedat, N., Jeffares, A., Imrie, F., & van der Schaar, M. (2023, April). Improving adaptive conformal prediction using self-supervised learning. In International Conference on Artificial Intelligence and Statistics (pp. 10160-10177). PMLR.

[2] Cherian, J. J., & Candès, E. J. (2023). Statistical Inference for Fairness Auditing. arXiv preprint arXiv:2305.03712.

Round table



10 min

Session 3

MAPIE

Session 3: Hypothesis Testing



10 min

Session 3: Hypothesis Testing

1) Presentation:

- Calibration tests in MAPIE
- Why hypothesis testing?
- Two families of hypothesis testing

2) Future directions:

- Hypothesis testing for detecting significant performance drift

3) Opening: round table

- Which hypothesis tests do you need?

What already exists in MAPIE for hypothesis testing?

Even if the calibration of binary classifier is not implemented in MAPIE as it already exists in scikit-learn, one would like to test if the model is calibrated or not:

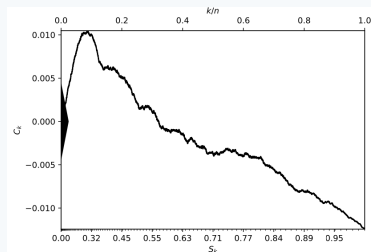
[H0] My model is well calibrated vs. [H1] My model is not calibrated

Cumulative differences

$$C_k = \frac{1}{N} \sum_{i=1}^N (y_i - s_i)$$

Characteristic length

$$\sigma = \frac{1}{N} \sqrt{\sum_{i=1}^N s_i(1 - s_i)}$$



Kolmogoriv-Smirnov test [1, 2, 3]

- Test statistic:

$$G = \max_{1 \leq k \leq N} |C_k|$$

G/σ converges in distribution to the maximum absolute value of 1D Brownian motion

p-value: [2, 3] propose closed-form formulas for the cumulative distribution function (CDF) of G/σ is given by :

$$p = 1 - CDF(G/\sigma)$$

Kuiper test [1, 2, 3]

- Test statistic:

$$H = \max_{1 \leq k \leq N} |C_k| - \min_{1 \leq k \leq N} |C_k|$$

H/σ converges in distribution to the maximum absolute value of 1D Brownian motion

p-value: [2, 3] propose closed-form formulas for the cumulative distribution function (CDF) of H/σ is given by :

$$p = 1 - CDF(H/\sigma)$$

Spiegelhalter test [4]

- Test statistic:

$$B = \frac{1}{N} \sum_{i=1}^N (y_i - s_i)(1 - 2s_i) + \frac{1}{N} \sum_{i=1}^N s_i(1 - s_i)$$

$$Z = \frac{B - \mathbb{E}[B]}{\sqrt{\text{Var}[B]}} = \frac{\sum_{i=1}^N (y_i - s_i)(1 - 2s_i)}{\sqrt{\sum_{i=1}^N (1 - 2s_i)^2 s_i(1 - s_i)}}$$

p-value: This statistic follows a normal distribution of cumulative distribution CDF, so that we state the associated p-value:

$$p = 1 - CDF(Z)$$

[1] Arrieta-Ibarra I, Gujral P, Tannen J, Tygert M, Xu C. Metrics of calibration for probabilistic predictions. The Journal of Machine Learning Research. 2022 Jan 1;23(1):15886-940.

[2] Tygert M. Calibration of P-values for calibration and for deviation of a subpopulation from the full population. arXiv preprint arXiv:2202.00100. 2022 Jan 31.

[3] D. A. Darling, A. J. F. Siegert. The First Passage Problem for a Continuous Markov Process. Ann. Math. Statist. 24 (4) 624 - 639, December, 1953.

[4] Spiegelhalter DJ. Probabilistic prediction in patient management and clinical trials. Statistics in medicine. 1986 Sep;5(5):421-33.

Why hypothesis testing?

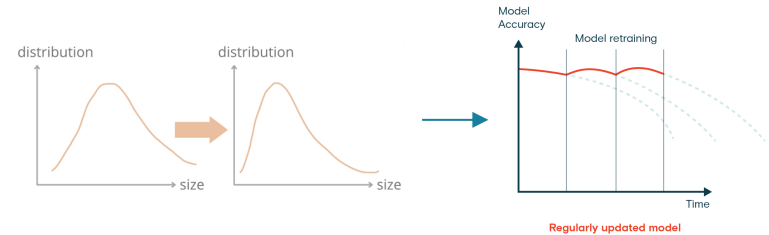
Preventive and on-line ways of monitoring the strength of evidence against the assumption of exchangeability.

Theoretical use for CP (Before using MAPIE)



- Can I use the conformal prediction framework?
- **Hypothesis testing (upstream)**
 - [H0] exchangeability vs [H1] non exchangeability

Business application (When the model is deployed)



- How do I detect changes in the data?
- Do I need to retain my model?
- **Performance testing (downstream)**
 - [H0] stable performance vs [H1] derived performance

Two families of hypothesis testing

Distribution drift / exchangeability testing

[Algorithmic Learning in a Random World, Vovk et al. 2022 \[1\]](#)

• Part III – Testing Randomness

- “*Conformal testing is a way of testing the IID assumption based on conformal prediction.*” [2]
- “*Valid testing procedures are equated with test martingales*” [2]
- Usual Testing (batch) vs. Conformal Testing (online)

• Overview of uses suggested by Vovk

- Testing Exchangeability
- Testing for Concept and Label Shift
- When to retrain: CUSUM, Shiryayev-Roberts, Variable & Fixed Training Schedules



Testing randomness by Vladimir Vovk 2020 [2]

Testing exchangeability: fork-convexity, supermartingales, and e-processes, by Aaditya Ramdas et al. 2021 [3]

• High-potential subject to be explored

- Call for proposals for tests relevant to MAPIE users

Performance drift testing

Bounded metrics:

Regression: coverage

Classification: coverage, precision & recall



Tracking the risk of a deployed model and detecting harmful distribution shifts by Aleksandr Podkopaev, Aaditya Ramdas 2021 [4]

Unbounded metrics:

Regression: MSE, SSR



Tips for moving from an unbounded metric to a bounded metric or to other directions.

[1] Vovk, V., Gammerman, A., & Shafer, G. (2005, 2022). Algorithmic learning in a random world (Vol. 29). New York: Springer.

[2] Vovk, V. (2021). Testing randomness online. *Statistical Science*, 36(4), 595-611.

[3] Ramdas, A., Ruf, J., Larsson, M., & Koolen, W. M. (2022). Testing exchangeability: Fork-convexity, supermartingales and e-processes. *International Journal of Approximate Reasoning*, 141, 83-109..

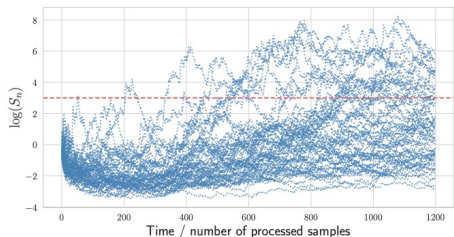
[4] Podkopaev, A., & Ramdas, A. (2021). Tracking the risk of a deployed model and detecting harmful distribution shifts. *arXiv preprint arXiv:2110.06177*.

Hypothesis testing for detecting significant performance drift

Motivation: Use a metric to control for **harmful** distribution drifts that have a **business** impact to **re-train models** at a given time.

Issue

In a classification setting, the harmful shift only occurs once the argmax changes value. In this scenario, we see that even though the marginal probability class of class 1 goes from 0.1 to 0.45 in increments of 0.05, test martingales would raise an error multiple times.



Issues and limitations:

- The data needs to be **IID or independent**, not exchangeable.
- **Risk needs to be upper bounded (or lower)**, hence, for classification: precision, recall and for regression: coverage → not only business metrics.

Hypothesis Testing in a business setting

Tracking the risk of a deployed model and detecting harmful distribution shifts, Aleksandr Podkopaev and Aaditya Ramdas, 2021

- $l(\cdot, \cdot)$: the loss function, chosen to be monitored
- $f: X \rightarrow Y$: the predictors
- $R(f) = \mathbb{E}[l(f(X), Y)]$: expected loss; called the risk of f
- $\hat{U}_S(f)$: upper confidence bound on the source risk
- $\hat{L}_T^{(t)}(f)$: lower confidence bound on the target risk continuously updated for the target risk as new data points are observed at time t
- $H_0 = R_t(f) \leq R_S(f) + \varepsilon_{tol}$

Algorithm 1 Sequential testing for an absolute increase in the risk.

Input: Predictor f , loss ℓ , tolerance level ε_{tol} , sample from the source $\{(X_i, Y_i)\}_{i=1}^{n_S}$.

- 1: **procedure**
- 2: Compute the upper confidence bound on the source risk $\hat{U}_S(f)$;
- 3: **for** $t = 1, 2, \dots$ **do**
- 4: Compute the lower confidence bound on the target risk $\hat{L}_T^{(t)}(f)$;
- 5: **if** $\hat{L}_T^{(t)}(f) > \hat{U}_S(f) + \varepsilon_{tol}$ **then**
- 6: Reject H_0 (equation 1) and fire off a warning.

Round table



10 min

Session 4

MAP|E

Session 4:
Round Table / Retex



15 min





Thank you for your attention.

