

An introduction to conformal prediction and distribution-free inference

(Lecture 2)

Rina Foygel Barber (University of Chicago)
Columbia Statistics short course May/June 2024

<http://rinafb.github.io/>

Topics — Lecture 2

Full conformal prediction

- Motivation
- The full conformal prediction method
- Generalizing to score functions
- Computation
- Comparing conformal and classical methods
- The role of randomization

Summary

- Summary & preview

Full conformal prediction

Motivation

Can we avoid data splitting?

Recall....

- The naive method fits a more accurate $\hat{\mu}$,
but the margin of error is too small due to overfitting
- The holdout set method (i.e., split CP) fits a less accurate $\hat{\mu}$,
but the margin of error is correctly calibrated
- And, CV type methods tend to work in practice
but can fail in unstable settings / no distrib.-free theory

Can we avoid data splitting?

Why does distribution-free theory hold for split CP but not for CV?

$$\mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \hat{q} \rightsquigarrow \text{coverage if } \underbrace{|Y_{n+1} - \hat{\mu}(X_{n+1})|}_{=R_{n+1}} \leq \hat{q}$$

- For split conformal, \hat{q} is quantile of calibration residuals

$$R_i = |Y_i - \hat{\mu}(X_i)|, \quad i = n_0 + 1, \dots, n$$

and $\hat{\mu}$ is pretrained $\Rightarrow R_{n_0+1}, \dots, R_n, R_{n+1}$ are exchangeable

- For CV, \hat{q} is quantile of leave-one-out residuals

$$R_i = |Y_i - \hat{\mu}_{-i}(X_i)|, \quad i = 1, \dots, n$$

$\Rightarrow R_1, \dots, R_n, R_{n+1}$ are *not* exchangeable

Full conformal prediction

The full conformal prediction method

Full conformal prediction

Intuition—

- Split CP fits $\hat{\mu}$ to part of the data, to ensure R_i 's are exch.
- CV uses all the data for fitting $\hat{\mu}$, but R_i 's not exch.
- Full CP: use all the data for fitting $\hat{\mu}$ *and* ensure R_i 's are exch.

An additional assumption:

The symmetric algorithm assumption

For any Z_1, \dots, Z_m and any $\sigma \in S_m$,

$$\mathcal{A}(Z_1, \dots, Z_m) = \mathcal{A}(Z_{\sigma(1)}, \dots, Z_{\sigma(m)}).$$

Full conformal prediction

Full CP, oracle version: imagine we could observe Y_{n+1}

- Fit model to training+test data

$$\hat{\mu} = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1}))$$

- Compute residuals

$$R_i = |Y_i - \hat{\mu}(X_i)|, i = 1, \dots, n; R_{n+1} = |Y_{n+1} - \hat{\mu}(X_{n+1})|$$

- Check if $R_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(R_1, \dots, R_n)$



If data points are exchangeable, and \mathcal{A} is symmetric,
then R_1, \dots, R_{n+1} are exchangeable

⇒ this event has $\geq 1 - \alpha$ probability

Full conformal prediction

Running full conformal in practice:

Test value $y \in \mathcal{Y}$



- Fit model to training+test data
 $\hat{\mu}^y = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y))$
- Compute residuals
 $R_i^y = |Y_i - \hat{\mu}^y(X_i)|, i = 1, \dots, n, R_{n+1}^y = |y - \hat{\mu}^y(X_{n+1})|$
- Check if $R_{n+1}^y \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(R_1^y, \dots, R_n^y)$

$y \rightsquigarrow \{\text{Yes, No}\}$



if Yes: add y to $\mathcal{C}(X_{n+1})$



if No: discard y

Full conformal prediction

Given training data $\{(X_i, Y_i)\}_{i=1,\dots,n}$, test point X_{n+1} , alg. \mathcal{A} :

Full conformal prediction (with residual score)

$$\mathcal{C}(X_{n+1}) = \{y \in \mathcal{Y} : R_{n+1}^y \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(R_1^y, \dots, R_n^y)\}$$

where

$$R_i^y = |Y_i - \hat{\mu}^y(X_i)|, i = 1, \dots, n, \quad R_{n+1}^y = |y - \hat{\mu}^y(X_{n+1})|,$$

for fitted model $\hat{\mu}^y = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y))$.

Full conformal prediction

Theorem: full conformal¹

If Z_1, \dots, Z_{n+1} are exchangeable, and \mathcal{A} is symmetric,
then full conformal prediction satisfies

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

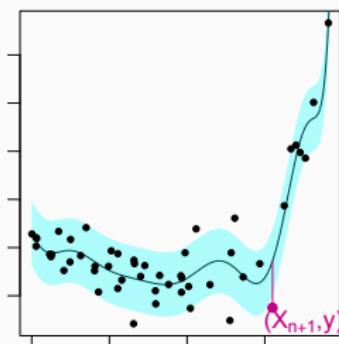
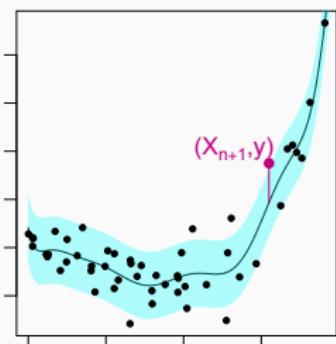
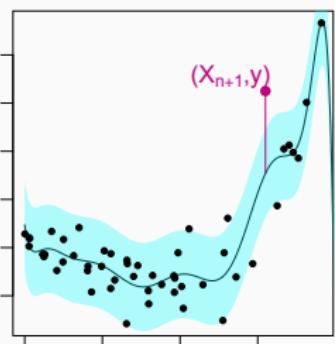
Proof:

$$\begin{aligned} Y_{n+1} \in \mathcal{C}(X_{n+1}) &\iff R_{n+1}^{Y_{n+1}} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(R_1^{Y_{n+1}}, \dots, R_n^{Y_{n+1}}) \\ &\iff R_{n+1}^{Y_{n+1}} \leq \text{Quantile}_{1-\alpha}(R_1^{Y_{n+1}}, \dots, R_n^{Y_{n+1}}, R_{n+1}^{Y_{n+1}}) \end{aligned}$$

And, this holds with prob. $\geq 1 - \alpha$ (as we proved for oracle version)

Full conformal prediction

How full conformal is run:



- $\hat{\mu}$ needs to be refitted for each X_{n+1} & each possible y

Full CP versus split CP

Split CP can be viewed as a special case of full CP:

- Let $\hat{\mu}$ be the pretrained model:

$$\hat{\mu} = \mathcal{A}(Z_1, \dots, Z_{n_0})$$

- Define a non-data-dependent algorithm:

\mathcal{A}' maps *any* data set to the pretrained model $\hat{\mu}$

- [Full CP with $\{(X_i, Y_i)\}_{n_0 < i \leq n}$ & \mathcal{A}'] = [Split CP with \mathcal{A}]

Terminology: $\begin{cases} \text{holdout method / split conformal / inductive conformal} \\ \text{conformal / full conformal / transductive conformal} \end{cases}$

Full conformal prediction

Generalizing to score functions

Full conformal prediction: general score

Full CP can be run with any score function (e.g., categorical data).

New definition of an algorithm:

\mathcal{A} maps a data set to a *score function* $s : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$.

Full conformal prediction (general score)

$$\mathcal{C}(X_{n+1}) = \{y \in \mathcal{Y} : S_{n+1}^y \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(S_1^y, \dots, S_n^y)\}$$

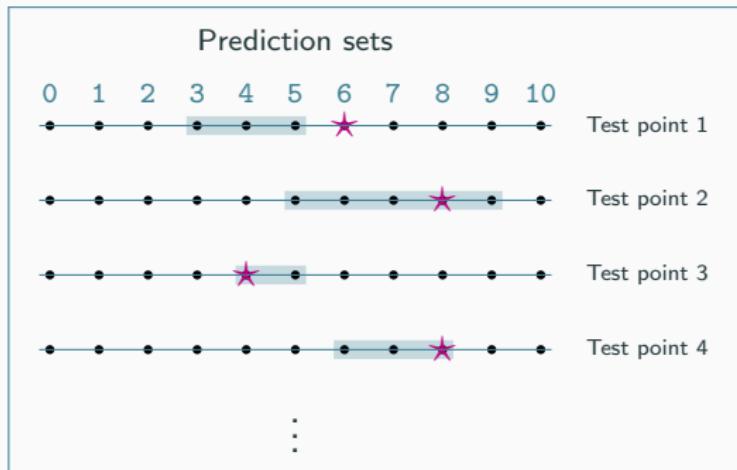
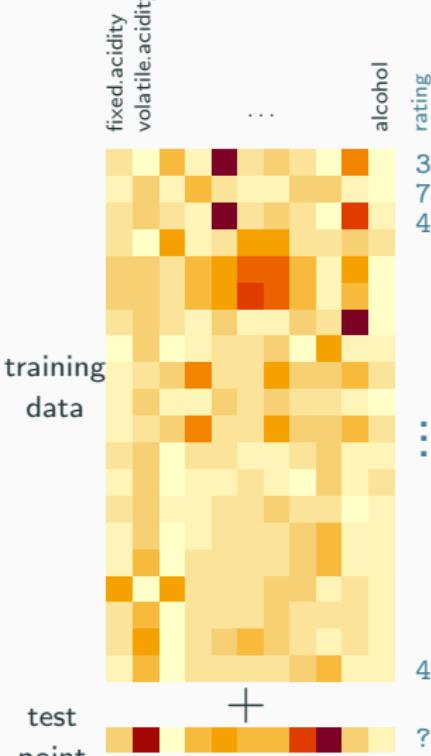
where

$$S_i^y = s^y(X_i, Y_i), i = 1, \dots, n, \quad S_{n+1}^y = s^y(X_{n+1}, y),$$

for fitted score function $s^y = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y))$

Full conformal prediction: general score

Example: the Wine Quality data set²



²Cortez et al 2009, Wine Quality data set, UCI Machine Learning Repository

Full conformal prediction: general score

An alternative formulation for full CP....

The algorithm \mathcal{A} maps a data set to a vector of scores:

$$(Z_1, \dots, Z_{n+1}) \mapsto (S_1, \dots, S_{n+1})$$

The three versions:

1. \mathcal{A} returns a fitted model:

$$\mathcal{A}(Z_1, \dots, Z_{n+1}) = \hat{\mu} \rightsquigarrow \text{scores } S_i = R_i = |Y_i - \hat{\mu}(X_i)|$$

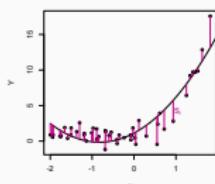
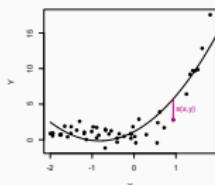
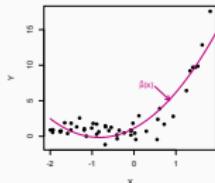
2. \mathcal{A} returns a fitted score function:

$$\mathcal{A}(Z_1, \dots, Z_{n+1}) = s \rightsquigarrow \text{scores } S_i = s(X_i, Y_i)$$

3. \mathcal{A} returns a vector of scores:

$$\mathcal{A}(Z_1, \dots, Z_{n+1}) = (S_1, \dots, S_{n+1})$$

Full conformal prediction: general score



\mathcal{A} maps a data set $\{(X_i, Y_i)\}$ to...

1. a fitted model $\hat{\mu}(x)$
2. a score function $s(x, y)$
3. a vector of scores S_i

Full conformal prediction: general score

For distribution-free coverage to hold,
need (1) exchangeable data, (2) symmetry of \mathcal{A}

What does symmetry of \mathcal{A} mean in each setting?

1. Fitted model: \mathcal{A} is invariant to permutations,

$$\hat{\mu} = \mathcal{A}(Z_1, \dots, Z_{n+1}) = \mathcal{A}(Z_{\sigma(1)}, \dots, Z_{\sigma(n+1)})$$

2. Fitted score function: \mathcal{A} is invariant to permutations,

$$s = \mathcal{A}(Z_1, \dots, Z_{n+1}) = \mathcal{A}(Z_{\sigma(1)}, \dots, Z_{\sigma(n+1)})$$

3. Vector of scores: \mathcal{A} commutes with permutations,

$$\mathcal{A}(Z_1, \dots, Z_{n+1}) = (S_1, \dots, S_{n+1}) \rightsquigarrow$$

$$\mathcal{A}(Z_{\sigma(1)}, \dots, Z_{\sigma(n+1)}) = (S_{\sigma(1)}, \dots, S_{\sigma(n+1)})$$

Full conformal prediction

Computation

Full conformal prediction: computational challenges

Full conformal prediction requires that the algorithm \mathcal{A} is re-run:

- For each test value X_{n+1} of interest
- For every possible value of Y_{n+1} (e.g, all $y \in \mathbb{R}$)

Full conformal prediction: computational challenges

A preliminary observation:³

for $\mathcal{Y} = \mathbb{R}$, we can truncate CP to the empirical range:

$$\mathcal{C}(X_{n+1}) = \left\{ y \in [\min_{1 \leq i \leq n} Y_i, \max_{1 \leq i \leq n} Y_i] : s(X_{n+1}, y) \leq \hat{q} \right\}$$

split CP version

Calculate coverage:

$$\begin{aligned} & \mathbb{P}\{Y_{n+1} \in \mathcal{C}(X_{n+1})\} \\ &= \mathbb{P}\left\{s(X_{n+1}, Y_{n+1}) \leq \hat{q} \text{ and } Y_{n+1} \in [\min_{1 \leq i \leq n} Y_i, \max_{1 \leq i \leq n} Y_i]\right\} \\ &\geq \mathbb{P}\{s(X_{n+1}, Y_{n+1}) \leq \hat{q}\} - \mathbb{P}\left\{Y_{n+1} \notin [\min_{1 \leq i \leq n} Y_i, \max_{1 \leq i \leq n} Y_i]\right\} \\ &\geq 1 - \alpha - \frac{2}{n+1} \end{aligned}$$

³Chen, Wang, Ha, & B. 2016, *Trimmed conformal prediction for high-dimensional models*.

Full conformal prediction: computational challenges

Conformal prediction, in theory:



Conformal prediction, in practice:



Problems to solve:

- Theory to guarantee coverage rate $\approx 1 - \alpha$?
- Avoid wider intervals due to discretized grid?

Full conformal prediction: computational challenges

Why do we lose the coverage guarantee?

- If fit $\hat{\mu}$ on $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y)$ for y in a grid...
 - Equivalent to: fit $\hat{\mu}$ on $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, [Y_{n+1}])$
- \nearrow
- Y_{n+1} rounded to grid
- Essentially, using an algorithm that does not treat data symmetrically

Full conformal prediction: computational challenges

Goal: theoretical validity + computational efficiency.

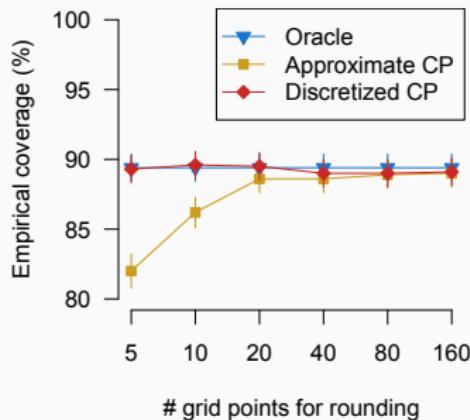
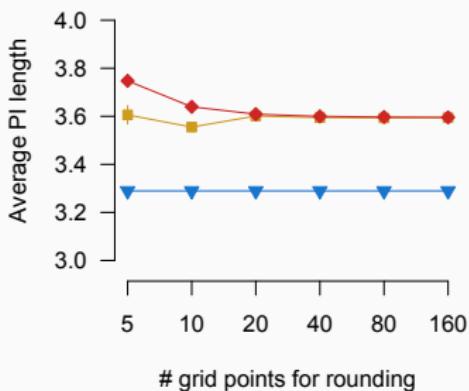
Discretized conformal prediction:⁴

- Define a new algorithm $\tilde{\mathcal{A}} = \mathcal{A} \circ [\cdot]$
 - Then run full CP with $\tilde{\mathcal{A}}$ as the regression algorithm
- $$\tilde{\mathcal{A}}((X_1, Y_1), \dots, (X_{n+1}, Y_{n+1})) = \mathcal{A}((X_1, [Y_1]), \dots, (X_{n+1}, [Y_{n+1}]))$$

⁴Chen, Chun, & B. 2017, *Discretized conformal prediction for efficient distribution-free inference*

Full conformal prediction: computational challenges

Discretized conformal prediction:⁵



⁵Chen, Chun, & B. 2017, *Discretized conformal prediction for efficient distribution-free inference*

Full conformal prediction: computational challenges

A special case: methods that are linear in Y

The algorithm \mathcal{A} takes the following form:

For $\hat{\mu}$ fitted to data points $(X_1, Y_1), \dots, (X_{n+1}, Y_{n+1})$,

$$\begin{pmatrix} \hat{\mu}(X_1) \\ \vdots \\ \hat{\mu}(X_{n+1}) \end{pmatrix} = \mathbf{H} \cdot \begin{pmatrix} Y_1 \\ \vdots \\ Y_{n+1} \end{pmatrix}$$

where $\mathbf{H} = \mathbf{H}(X_1, \dots, X_{n+1})$ depends only on the features.

Examples:

- Linear regression & ridge regression
- Regression with polynomial terms / interaction terms / etc
- Nearest neighbor methods
- Kernel smoothing type methods

Full conformal prediction: computational challenges

How to run full conformal prediction in this special case:⁶

1. Compute $\mathbf{H} = \mathbf{H}(X_1, \dots, X_{n+1})$
2. Then for $\hat{\mu}^y = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y))$,

$$\hat{\mu}^y(X_i) = \mathbf{H}_{i,*} \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \\ y \end{pmatrix} = \mathbf{H}_{i,*} \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \\ 0 \end{pmatrix} + y \cdot \mathbf{H}_{i,*} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} = a_i + b_i y$$

3. Prediction interval can now be computed in finite time:

$$\mathcal{C}(X_{n+1}) = \{y : |y - (a_{n+1} + b_{n+1}y)| \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(|Y_i - (a_i + b_i y)|)\}$$

⁶Burnaev & Vovk 2014, *Efficiency of conformalized ridge regression*

Full conformal prediction: computational challenges

A special case: stable conformal prediction⁷

Assume stability of \mathcal{A} ,

$$|\hat{\mu}^y(X_i) - \hat{\mu}^{y_0}(X_i)| \leq \tau_i, \quad i = 1, \dots, n+1$$

- Stability assumption \Rightarrow scores are bounded above & below:

$$\underbrace{|Y_i - \hat{\mu}(X_i)|}_{=R_i} \in \underbrace{|Y_i - \hat{\mu}^{y_0}(X_i)|}_{=R_i^{y_0}} \pm \tau_i$$

- Bound prediction set above & below:

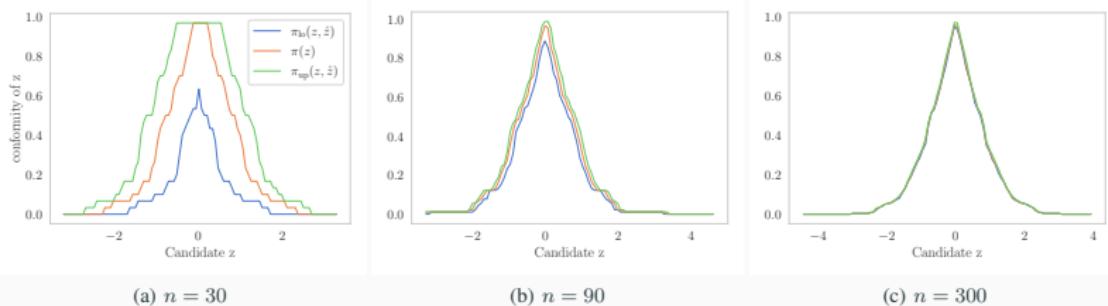
$$\mathcal{C}(X_{n+1}) \subseteq \left\{ y : R_{n+1}^{y_0} - \tau_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(R_i^{y_0} + \tau_i) \right\}$$

$$\mathcal{C}(X_{n+1}) \supseteq \left\{ y : R_{n+1}^{y_0} + \tau_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(R_i^{y_0} - \tau_i) \right\}$$

⁷ Ndiaye 2022, *Stable Conformal Prediction Sets*

Full conformal prediction: computational challenges

In general, large $n \rightsquigarrow \tau_i \approx 0$, so the bounds converge:



(figure from Ndiaye 2022)

But, the method requires knowing the stability properties of \mathcal{A} .

Full conformal prediction: computational challenges

Summary of approaches:

- In practice — restrict to a grid of y values (but no theory)
- Specialized methods for specific algorithms/settings,
e.g., Ridge⁸, Lasso⁹, stable algorithms¹⁰
- Discretized CP — use a discretized version of \mathcal{A}
to restore theoretical guarantees¹¹

⁸Burnaev & Vovk 2014, *Efficiency of conformalized ridge regression*

⁹Lei 2017, *Fast Exact Conformalization of Lasso using Piecewise Linear Homotopy*

¹⁰Ndiaye 2022, *Stable Conformal Prediction Sets*

¹¹Chen, Chun, & B. 2017, *Discretized conformal prediction for efficient distribution-free inference*

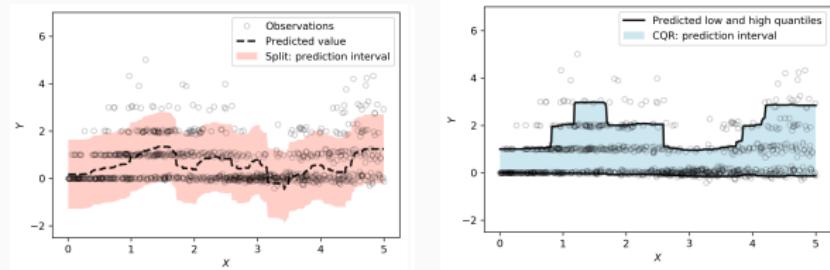
Full conformal prediction

Comparing conformal and classical methods

Comparing conformal & classical methods

For a particular data set, how do we choose between CP versus an existing model/algorithm? What if:

- CP provides a distribution-free coverage guarantee
- But, the existing model/algorithm is a better fit to the data distribution, & provides a narrower interval



(figure from Romano et al 2019)

By choosing a good score function, can have the best of both

Comparing conformal & classical methods

CP is a “wrapper” method that pairs with any model/algorithm, so we can:

- Choose the model that we believe fits the data
- Use it to guide our definition of the score function
- Apply CP $\rightsquigarrow \mathcal{C}(X_{n+1})$ is an adjusted version of the original model’s answer

Drawbacks:

computational cost for full CP / reduced sample size for split CP

Comparing conformal & classical methods

An example strategy...

Construct model-based prediction regions at each confidence level:

$$\dots \subseteq \mathcal{C}_{0.8}^{\text{mod}}(x) \subseteq \mathcal{C}_{0.9}^{\text{mod}}(x) \subseteq \mathcal{C}_{0.95}^{\text{mod}}(x) \subseteq \dots$$

Define a score $s(x, y) = \inf \{t : y \in \mathcal{C}_t^{\text{mod}}(x)\}$

Then split CP returns $\mathcal{C}(X_{n+1}) = \mathcal{C}_{\hat{q}}^{\text{mod}}(X_{n+1})$



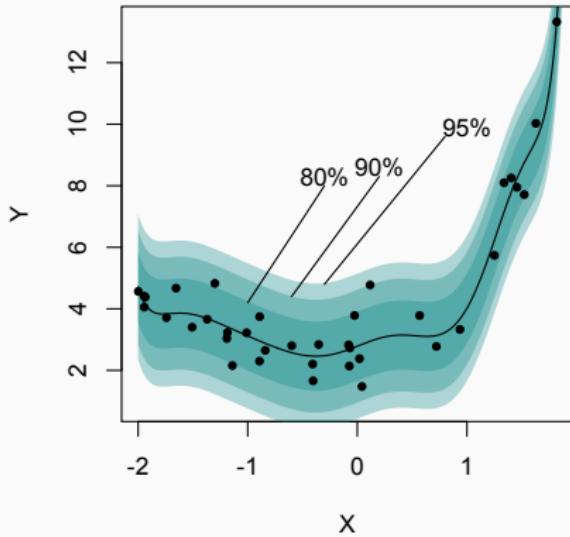
If model is correct, $\hat{q} \approx 1 - \alpha$

If model is misspecified, $\hat{q} > 1 - \alpha$

Comparing conformal & classical methods

Interpretation via nested sets:¹²

$$\dots \subseteq \mathcal{C}_{0.8}^{\text{mod}}(x) \subseteq \mathcal{C}_{0.9}^{\text{mod}}(x) \subseteq \mathcal{C}_{0.95}^{\text{mod}}(x) \subseteq \dots$$



¹²Gupta et al 2022, *Nested conformal prediction and quantile out-of-bag ensemble methods*

Comparing conformal & classical methods

Theoretical guarantees that CP converges to the “oracle” answer:

Setting 1: estimating the conditional mean¹³

For the residual score, $s(x, y) = |y - \hat{\mu}(x)|$,
if $\hat{\mu}(x) \rightarrow \mathbb{E}[Y | X = x]$ & noise is symmetric,

$$\mathcal{C}(X_{n+1}) \rightarrow \mathcal{C}^*(X_{n+1})$$

Setting 2: estimating the conditional density¹⁴

For the density score, $s(x, y) = -\hat{f}(y|x)$, if $\hat{f}(y|x) \rightarrow f(y|x)$,

$$\mathcal{C}(X_{n+1}) \rightarrow \mathcal{C}^*(X_{n+1})$$

¹³Lei et al 2018, *Distribution-Free Predictive Inference for Regression*

¹⁴Izbicki et al 2020, *Flexible distribution-free conditional predictive bands using density estimators*

Full conformal prediction

The role of randomization

Coverage or overcoverage?

Return to the theory:

Theorem: full conformal¹⁵

If Z_1, \dots, Z_{n+1} are exchangeable, and \mathcal{A} is symmetric,
then full conformal prediction satisfies

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

Proof idea:

$$S_1, \dots, S_{n+1} \text{ exch.} \Rightarrow \mathbb{P}\{S_{n+1} \leq \text{Quantile}_{1-\alpha}(S_1, \dots, S_{n+1})\} \geq 1 - \alpha$$



May be $> 1 - \alpha$ due to:

- (1) ties between S_i 's
- (2) $(n+1)(1-\alpha)$ noninteger

¹⁵Vovk et al 2005, *Algorithmic Learning in a Random World*

The conformal p-value

A reinterpretation of conformal prediction—

Full conformal prediction (p-value version)¹⁶

For each $y \in \mathcal{Y}$ define a conformal p-value:

$$p^y = \frac{1 + \sum_{i=1}^n \mathbb{1} \{ S_i^y \geq S_{n+1}^y \}}{n + 1}$$

where

$$S_i^y = s^y(X_i, Y_i), i = 1, \dots, n, \quad S_{n+1}^y = s^y(X_{n+1}, y),$$

for fitted score function $s^y = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y))$

Then define prediction interval: $\mathcal{C}(X_{n+1}) = \{y \in \mathcal{Y} : p^y > \alpha\}$.

¹⁶Vovk et al 2005, *Algorithmic Learning in a Random World*

The conformal p-value

Interpretation: p^y is a p-value testing the hypothesis $Y_{n+1} = y$,
and $\mathcal{C}(X_{n+1})$ is the set of all y 's that are *not* rejected:

$$\mathcal{C}(X_{n+1}) = \{y \in \mathcal{Y} : p^y > \alpha\}$$

$S_1^{Y_{n+1}}, \dots, S_{n+1}^{Y_{n+1}}$ exchangeable \iff $p^{Y_{n+1}}$ is superuniform
i.e., $\mathbb{P}\{p^{Y_{n+1}} \leq \alpha\} \leq \alpha$

The conformal p-value

Use randomization (also called “smoothing”)
to make the p-value exactly $\text{Unif}[0, 1]$:

$$p^y = \frac{\sum_{i=1}^n \mathbb{1}\{S_i^y > S_{n+1}^y\} + \textcolor{red}{U} \cdot [1 + \sum_{i=1}^n \mathbb{1}\{S_i^y = S_{n+1}^y\}]}{n+1}$$

where $U \sim \text{Unif}[0, 1]$

- Can also be applied to split CP as a special case
- Note: typically use the *same* draw of U for *all* values $y \in \mathcal{Y}$

Randomized algorithms

Another source of randomness: the algorithm \mathcal{A} may use...

- Stochastic gradient descent (SGD)
- Random initialization
- Resampling / ensembling / CV
- etc.

random seed $\sim \text{Unif}[0, 1]$

New notation: score function $s = \mathcal{A}(Z_1, \dots, Z_m; \xi)$

The symmetric algorithm assumption (randomized version)

For any Z_1, \dots, Z_m and any $\sigma \in \mathcal{S}_m$,

$$\mathcal{A}(Z_1, \dots, Z_m; \xi) = \mathcal{A}(Z_{\sigma(1)}, \dots, Z_{\sigma(m)}; \xi')$$

for some coupling of $\xi \sim \text{Unif}[0, 1]$ and $\xi' \sim \text{Unif}[0, 1]$.

Summary

Summary & preview

Summary: lecture 2

Lecture 2 topics:

- Full conformal prediction: method & theory
- Computational strategies for full CP
- Perspectives on score functions, symmetry, & randomization

Full CP allows us to start with any algorithm/model,
and then calibrate it to have valid predictive coverage
without data splitting (but at a **high computational cost**)

Preview: lecture 3

Cross-validation based methods offer a compromise between split CP (computational efficiency) & full CP (statistical effic.)

- Cross-conformal prediction
- Jackknife+ and CV+
- Connections to algorithmic stability