

An introduction to conformal prediction and distribution-free inference

CIRM tutorial (part 1)

Rina Foygel Barber (University of Chicago)
CIRM December 2024

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

Introduction

Regression & prediction

Supervised learning setting:

Training data $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$

Goals:

- Inference on the regression — model distribution of Y given X
- Predictive inference — predict value of Y given X
for test points $(X_{n+1}, Y_{n+1}), (X_{n+2}, Y_{n+2}), \dots$

Regression & prediction — classical approach

- We assume a parametric model on (X, Y) or on $Y | X$
- We perform estimation & inference on the parameters....
-& then we can provide prediction intervals:

Regression & prediction — classical approach

- We assume a parametric model on (X, Y) or on $Y | X$
e.g., for linear regression, $Y = X^\top \beta + \mathcal{N}(0, \sigma^2)$
- We perform estimation & inference on the parameters....
e.g., for linear regression, distribution of $\hat{\beta}$ and $\hat{\sigma}^2$
-& then we can provide prediction intervals:
e.g., for linear regression, $Y_{n+1} \in X_{n+1}^\top \hat{\beta} \pm \dots$

Regression & prediction — nonparametric approach

- We allow a nonparametric model for (X, Y) or $Y | X$, with assumptions/constraints
- We perform estimation & inference on the model....
-& then we can provide prediction intervals:

Regression & prediction — nonparametric approach

- We allow a nonparametric model for (X, Y) or $Y | X$, with assumptions/constraints

e.g., assume $\mathbb{E}[Y | X]$ is smooth

- We perform estimation & inference on the model....

e.g., $\hat{\mu}(x) = \text{estimate of } \mathbb{E}[Y | X = x]$, via a Gaussian kernel

-& then we can provide prediction intervals:

e.g., $Y_{n+1} \in \hat{\mu}(X_{n+1}) \pm \dots\dots$

Regression & prediction — ML approach

- Train an overparametrized model for $Y | X$
- Provide predictions for new feature vectors
- Use a data-driven strategy for uncertainty quantification

Regression & prediction — ML approach

- Train an overparametrized model for $Y | X$

e.g., train a neural net on $\{(X_i, Y_i)\}$

- Provide predictions for new feature vectors

e.g., \hat{Y}_{n+i} = neural net's prediction for feature X_{n+i}

- Use a data-driven strategy for uncertainty quantification

e.g., holdout data / cross-validation / bootstrapping / etc

What can go wrong?

- For the parametric approach — our model may be wrong
- For the nonparametric approach — our assumptions (e.g., smoothness) may not hold
- For the ML approach — is data-driven inference guaranteed to give valid answers?

Our choices:

- Rely on assumptions being correct
- Or, test empirically whether our assumptions hold
- Or, use inference methods that don't rely on assumptions (or, only rely on weaker assumptions)

Setting:

- Features $X \in \mathcal{X}$, response $Y \in \mathbb{R}$ (or $Y \in \mathcal{Y}$)
- Available training data $(X_1, Y_1), \dots, (X_n, Y_n) \rightsquigarrow$ fit model $\hat{\mu}$
- Goal: given X_{n+1}, X_{n+2}, \dots , predict Y_{n+1}, Y_{n+2}, \dots

Regression & prediction — data-driven predictive inference

Setting:

- Features $X \in \mathcal{X}$, response $Y \in \mathbb{R}$ (or $Y \in \mathcal{Y}$)
- Available training data $(X_1, Y_1), \dots, (X_n, Y_n) \rightsquigarrow$ fit model $\hat{\mu}$
- Goal: given X_{n+1}, X_{n+2}, \dots , predict Y_{n+1}, Y_{n+2}, \dots

Prediction?

$$\hat{Y}_{n+i} = \hat{\mu}(X_{n+i})$$

or predictive inference?

$$Y_{n+i} \in \hat{\mu}(X_{n+i}) \pm \underbrace{(\text{margin of error})}_{\text{how to calculate?}}$$

Using the training set for inference

Using the training loss:

If fitted model $\hat{\mu}$ overfits to training data, generally

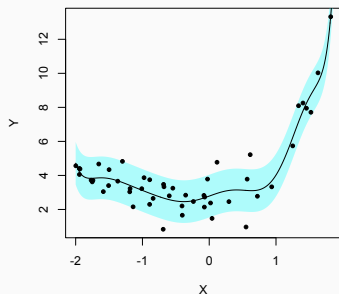
$$\underbrace{|Y_{n+i} - \hat{\mu}(X_{n+i})|}_{\text{test error}} \gg \underbrace{\frac{1}{n} \sum_{i=1}^n |Y_i - \hat{\mu}(X_i)|}_{\text{avg. training error}}$$

even if training & test data are from the same distribution

Regression & prediction — data-driven predictive inference

Simulation: suppose we construct prediction intervals as

$$\mathcal{C}(X_{n+i}) = \hat{\mu}(X_{n+i}) \pm \underbrace{\text{Quantile}_{1-\alpha}(|Y_1 - \hat{\mu}(X_1)|, \dots, |Y_n - \hat{\mu}(X_n)|)}_{\text{residuals on training data}}$$

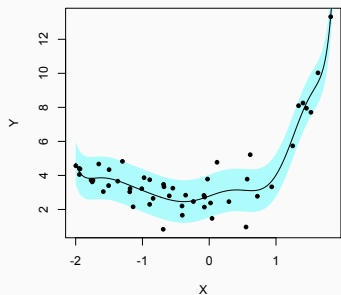


Train: 90% coverage

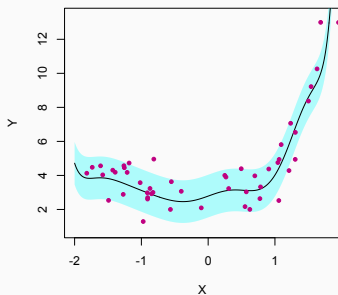
Regression & prediction — data-driven predictive inference

Simulation: suppose we construct prediction intervals as

$$\mathcal{C}(X_{n+i}) = \hat{\mu}(X_{n+i}) \pm \underbrace{\text{Quantile}_{1-\alpha}(|Y_1 - \hat{\mu}(X_1)|, \dots, |Y_n - \hat{\mu}(X_n)|)}_{\text{residuals on training data}}$$



Train: 90% coverage



Test: 78% coverage

Using a holdout set for inference

To avoid overfitting — use a holdout set (“calibration set”)

- Split the training data, $n = n_0 + n_1$
- Fit model $\hat{\mu}$ on pretraining set $\{(X_i, Y_i)\}_{1 \leq i \leq n_0}$
- Compute residuals on calibration set, $\{|Y_i - \hat{\mu}(X_i)|\}_{n_0 < i \leq n}$
- Prediction interval:

$$\mathcal{C}(X_{n+i}) = \hat{\mu}(X_{n+i}) \pm \text{Quantile}_{1-\alpha}(\{|Y_i - \hat{\mu}(X_i)|\}_{n_0 < i \leq n})$$

Using a holdout set for inference

To avoid overfitting — use a holdout set (“calibration set”)

- Split the training data, $n = n_0 + n_1$
- Fit model $\hat{\mu}$ on pretraining set $\{(X_i, Y_i)\}_{1 \leq i \leq n_0}$
- Compute residuals on calibration set, $\{|Y_i - \hat{\mu}(X_i)|\}_{n_0 < i \leq n}$
- Prediction interval:

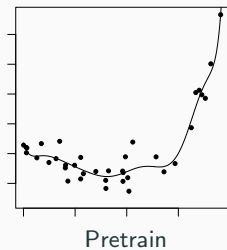
$$\mathcal{C}(X_{n+i}) = \hat{\mu}(X_{n+i}) \pm \text{Quantile}_{1-\alpha}(\{|Y_i - \hat{\mu}(X_i)|\}_{n_0 < i \leq n})$$

fitted on pretraining data

computed on calibration data

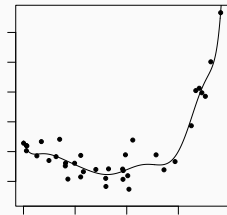
Using a holdout set for inference

Simulation:

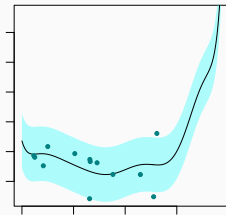


Using a holdout set for inference

Simulation:



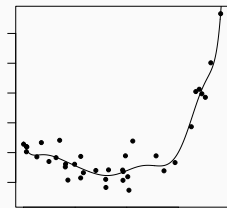
Pretrain



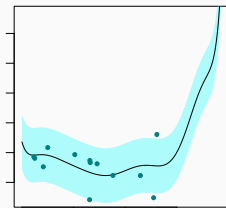
Calibration: 90% coverage

Using a holdout set for inference

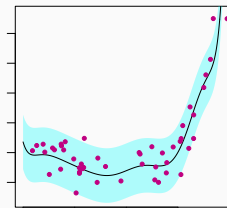
Simulation:



Pretrain



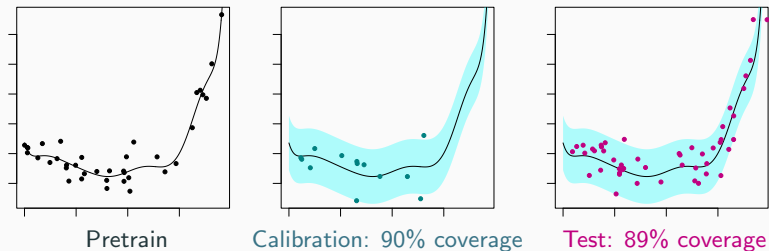
Calibration: 90% coverage



Test: 89% coverage

Using a holdout set for inference

Simulation:



Note: lower sample size $\rightsquigarrow \hat{\mu}$ is less accurate \rightsquigarrow intervals are wider

Using a holdout set for inference

- The naive method fits a more accurate $\hat{\mu}$,
but the margin of error is too small due to overfitting
- A holdout set method fits a less accurate $\hat{\mu}$,
but the margin of error is correctly calibrated
- Can we use cross-validation (CV) to get the best of both?
Will return to this!

Distribution-free prediction: aims

The goal of distribution-free inference is to provide guarantees that are valid universally over all data distributions.

For the problem of predictive inference...

- Can we construct a prediction interval $\mathcal{C}(X_{n+i}) \subseteq \mathcal{Y}$ such that

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

Distribution-free prediction: aims

The goal of distribution-free inference is to provide guarantees that are valid universally over all data distributions.

For the problem of predictive inference...

- Can we construct a prediction interval $\mathcal{C}(X_{n+i}) \subseteq \mathcal{Y}$ such that

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

- Want to avoid overly conservative solutions ($\mathcal{C}(X_{n+1}) = \mathcal{Y}$)
- Want to be able to use *any* regression method to construct \mathcal{C} (classical or ML methods)

Intro to exchangeability

Introduction to exchangeability

For the rest of this talk: let $Z_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$

The i.i.d. data setting

Assume $\underbrace{Z_1, \dots, Z_n}_{\text{training}}, \underbrace{Z_{n+1}, Z_{n+2}, \dots}_{\text{test}}$ are i.i.d. from some distrib. P

Introduction to exchangeability

For the rest of this talk: let $Z_i = (X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$

The i.i.d. data setting

Assume $\underbrace{Z_1, \dots, Z_n}_{\text{training}}, \underbrace{Z_{n+1}, Z_{n+2}, \dots}_{\text{test}}$ are i.i.d. from some distrib. P

Can we call this “distribution-free”?

- No assumptions on P (e.g., P does not need to be smooth)
- But, this does not allow for dependence across time / distribution shift / etc
- We will return to these settings later

Introduction to exchangeability

The exchangeable data setting

Assume that the data points

$$\underbrace{Z_1, \dots, Z_n}_{\text{training}}, \underbrace{Z_{n+1}, Z_{n+2}, \dots}_{\text{test}}$$

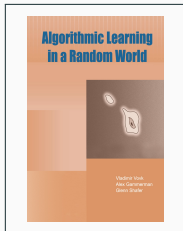
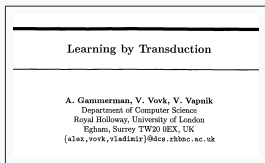
are exchangeable, i.e., $(Z_1, \dots, Z_m) \stackrel{d}{=} (Z_{\sigma(1)}, \dots, Z_{\sigma(m)})$ for every m and every permutation σ .

- The i.i.d. data setting is a special case
- Conditionally i.i.d. data is another special case
- Note: finite sequences can be exchangeable but not i.i.d. (de Finetti's theorem does not apply)

Background on conformal prediction

Conformal prediction: background

Background on the conformal prediction (CP) framework:
key idea = statistical inference via exchangeability of the data



Gammerman, Vovk, Vapnik
UAI 1998

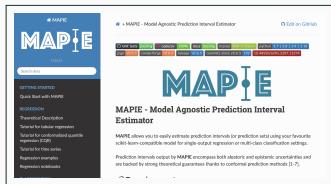
Vovk, Gammerman, Shafer
2005 — see alrw.net



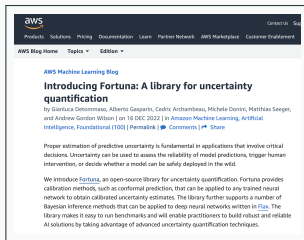
Lei, G'Sell, Rinaldo,
Tibshirani, Wasserman
JASA 2018

Conformal prediction: background

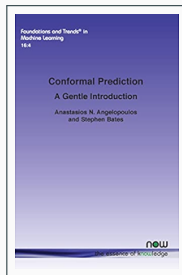
Recent developments — software packages & user-friendly tutorials



The screenshot shows the GitHub repository for MAPIE (Model Agnostic Prediction Interval Estimator). The repository name is "MAPIE" and the description is "MAPIE - Model Agnostic Prediction Interval Estimator". The repository is owned by "MAPIE" and has a "Fork" button. The repository is described as "MAPIE allows you to easily estimate prediction intervals for prediction sets using your favourite scikit-learn compatible model for single-output regression or multi-class classification settings. Prediction intervals output by MAPIE encompass both aleatoric and epistemic uncertainties and are backed by strong theoretical guarantees thanks to conformal prediction methods [1, 2]."



The screenshot shows a blog post from the AWS Machine Learning Blog titled "Introducing Fortuna: A library for uncertainty quantification". The post is by Gianluca DeTommaso, Alberto Gaspari, Cedric Archambeau, Michele Donini, Matthias Seeger, and Andrew Gordon Wilson, dated 16 DEC 2022. The post discusses the importance of predictive uncertainty in applications and introduces Fortuna, an open-source library for uncertainty quantification. Fortuna provides calibration methods, such as conformal prediction, that can be applied to any trained neural network to obtain calibrated uncertainty estimates. The library further supports a number of Bayesian inference methods that can be applied to deep neural networks written in PyTorch. The library makes it easy to run benchmarks and will enable practitioners to build robust and reliable AI solutions by taking advantage of advanced uncertainty quantification techniques.



The cover of the book "Foundations and Trends in Machine Learning 18.4: Conformal Prediction: A Gentle Introduction" by Anastasios N. Angelopoulos and Stephen Bates. The cover is purple and features the title and authors' names. The logo for "now" (now publishers) is visible at the bottom right, with the tagline "the essence of knowledge".

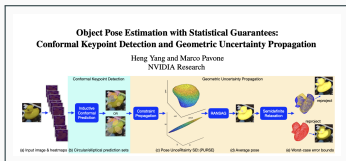
& a new theory textbook (75% on arXiv, forthcoming from CUP)

Theoretical Foundations of Conformal Prediction

Anastasios N. Angelopoulos¹, Rina Foygel Barber², Stephen Bates³

Conformal prediction: background

Recent developments — successful applications in biological sciences, machine learning, & many more domains

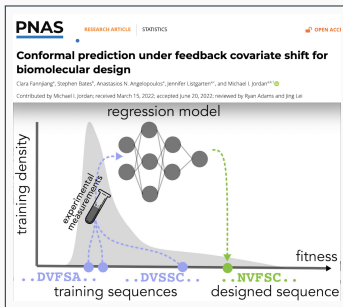


Journal of Healthcare Informatics Research (2022) 6:241–252
<https://doi.org/10.1007/s41666-021-00113-8>

REVIEW ARTICLE

Conformal Prediction in Clinical Medical Sciences

Janette Vazquez¹ · Julio C. Facelli¹



Split conformal prediction

Split conformal prediction

The split conformal prediction method

- 1 Using pretraining data Z_1, \dots, Z_{n_0} ,
construct fitted model $\hat{\mu}$ using any regression algorithm:

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

- 2 Compute quantile \hat{q} of calibration set residuals:

$$\hat{q} = \text{Quantile}_{(1-\alpha)(1+1/n_1)} \left(\{|Y_i - \hat{\mu}(X_i)|\}_{n_0 < i \leq n} \right)$$

- 3 For test point $n + 1$ return prediction interval

$$\mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \hat{q}$$

Split conformal prediction

Theorem¹

If Z_1, \dots, Z_{n+1} are exchangeable, then split conformal satisfies:

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

¹Vovk et al 2005, *Algorithmic Learning in a Random World*

Split conformal prediction

Theorem¹

If Z_1, \dots, Z_{n+1} are exchangeable, then split conformal satisfies:

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

Proof:

Define $S_i = |Y_i - \hat{\mu}(X_i)|$ for $i = \underbrace{n_0 + 1, \dots, n}_{\text{calibration}}, \underbrace{n+1}_{\text{test}}$

¹Vovk et al 2005, *Algorithmic Learning in a Random World*

Split conformal prediction

Theorem¹

If Z_1, \dots, Z_{n+1} are exchangeable, then split conformal satisfies:

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

Proof:

Define $S_i = |Y_i - \hat{\mu}(X_i)|$ for $i = \underbrace{n_0 + 1, \dots, n}_{\text{calibration}}, \underbrace{n+1}_{\text{test}}$

$$Y_{n+1} \in \mathcal{C}(X_{n+1}) \iff S_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n_1)}(S_{n_0+1}, \dots, S_n)$$

¹Vovk et al 2005, *Algorithmic Learning in a Random World*

Split conformal prediction

Theorem¹

If Z_1, \dots, Z_{n+1} are exchangeable, then split conformal satisfies:

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

Proof:

Define $S_i = |Y_i - \hat{\mu}(X_i)|$ for $i = \underbrace{n_0 + 1, \dots, n}_{\text{calibration}}, \underbrace{n+1}_{\text{test}}$

$$\begin{aligned} Y_{n+1} \in \mathcal{C}(X_{n+1}) &\iff S_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n_1)}(S_{n_0+1}, \dots, S_n) \\ &\iff S_{n+1} \leq \text{Quantile}_{1-\alpha}(S_{n_0+1}, \dots, S_n, S_{n+1}) \end{aligned}$$

¹Vovk et al 2005, *Algorithmic Learning in a Random World*

Split conformal prediction

Exchangeability for holdout set methods

If $\hat{\mu} = \mathcal{A}(Z_1, \dots, Z_{n_0})$, & the data points are exchangeable, then

$$\underbrace{|Y_{n_0+1} - \hat{\mu}(X_{n_0+1})|, \dots, |Y_n - \hat{\mu}(X_n)|}_{\text{calibration residuals}}, \underbrace{|Y_{n+1} - \hat{\mu}(X_{n+1})|}_{\text{test residual}}$$

are exchangeable.

Split conformal prediction

Exchangeability for holdout set methods

If $\hat{\mu} = \mathcal{A}(Z_1, \dots, Z_{n_0})$, & the data points are exchangeable, then

$$\underbrace{|Y_{n_0+1} - \hat{\mu}(X_{n_0+1})|, \dots, |Y_n - \hat{\mu}(X_n)|}_{\text{calibration residuals}}, \underbrace{|Y_{n+1} - \hat{\mu}(X_{n+1})|}_{\text{test residual}}$$

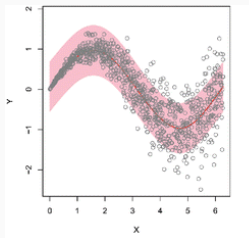
are exchangeable.

$$\implies \mathbb{P} \{S_{n+1} \leq \text{Quantile}_{1-\alpha}(S_{n_0+1}, \dots, S_n, S_{n+1})\} \geq 1 - \alpha$$

□

The conformal score

Due to the construction of the split conformal method,
 $\mathcal{C}(X_{n+1})$ has the same width regardless of the value of X_{n+1}



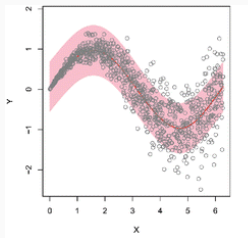
(figure from Lei et al 2018)

Why?

- $\mathcal{C}(X_{n+1}) = [\hat{\mu}(X_{n+1}) \pm \hat{q}] = \{y \in \mathbb{R} : |y - \hat{\mu}(X_{n+1})| \leq \hat{q}\}$

The conformal score

Due to the construction of the split conformal method,
 $\mathcal{C}(X_{n+1})$ has the same width regardless of the value of X_{n+1}



(figure from Lei et al 2018)

Why?

- $\mathcal{C}(X_{n+1}) = [\hat{\mu}(X_{n+1}) \pm \hat{q}] = \{y \in \mathbb{R} : |y - \hat{\mu}(X_{n+1})| \leq \hat{q}\}$
- Equivalently: we are using $|y - \hat{\mu}(X_{n+1})|$ as a *score* to determine whether y is contained in $\mathcal{C}(X_{n+1})$ or not

The split conformal prediction method² (general score)

- ① Using pretraining data Z_1, \dots, Z_{n_0} ,
construct **score function** $s : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ using any algorithm
- ② Compute quantile \hat{q} of calibration set scores:

$$\hat{q} = \text{Quantile}_{(1-\alpha)(1+1/n_1)}(S_{n_0+1}, \dots, S_n)$$

where $S_i = s(X_i, Y_i)$

- ③ For test point $n + 1$ return prediction interval

$$\mathcal{C}(X_{n+1}) = \{y \in \mathcal{Y} : s(X_{n+1}, y) \leq \hat{q}\}$$

²Vovk et al 2005, *Algorithmic Learning in a Random World*

The conformal score

The residual score:

$s(x, y) = |y - \hat{\mu}(x)|$, where $\hat{\mu}$ fitted on pretraining data

$$\implies \mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \hat{q}$$

The conformal score

An alternative score function — (will see more examples later on)

The scaled residual score:³

$$s(x, y) = \frac{|y - \hat{\mu}(x)|}{\hat{\sigma}(x)}, \text{ where } \hat{\mu}, \hat{\sigma} \text{ fitted on pretraining data}$$
$$\implies \mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \hat{q} \cdot \hat{\sigma}(X_{n+1})$$

³Lei et al 2018, *Distribution-Free Predictive Inference for Regression*

The conformal score

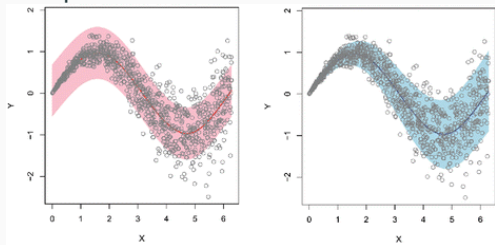
An alternative score function — (will see more examples later on)

The scaled residual score:³

$$s(x, y) = \frac{|y - \hat{\mu}(x)|}{\hat{\sigma}(x)}, \text{ where } \hat{\mu}, \hat{\sigma} \text{ fitted on pretraining data}$$

$$\implies \mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \hat{q} \cdot \hat{\sigma}(X_{n+1})$$

Compare to residual score:



(figure from Lei et al 2018)

³Lei et al 2018, *Distribution-Free Predictive Inference for Regression*

Split conformal: summary

Split CP allows us to start with any pretrained model/score, and then calibrate it to have valid predictive coverage (as long as we can assume exchangeability!)

Drawback: model $\hat{\mu}$ (or score s) less accurate due to data splitting

Full conformal prediction

Full conformal prediction

Intuition—

- Split CP fits $\hat{\mu}$ to part of the data, to ensure S_i 's are exch.
- Full CP: use all the data for fitting $\hat{\mu}$ *and* ensure S_i 's are exch.

Full conformal prediction

Intuition—

- Split CP fits $\hat{\mu}$ to part of the data, to ensure S_i 's are exch.
- Full CP: use all the data for fitting $\hat{\mu}$ *and* ensure S_i 's are exch.

An additional assumption:

The symmetric algorithm assumption

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

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

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

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

- Check if $S_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(S_1, \dots, S_n)$

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

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

- Check if $S_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(S_1, \dots, S_n)$



If data points are exchangeable, and \mathcal{A} is symmetric,
then S_1, \dots, S_{n+1} are exchangeable

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

Full conformal prediction

Running full conformal in practice:

- 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

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

- Check if $S_{n+1}^y \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(S_1^y, \dots, S_n^y)$

Full conformal prediction

Running full conformal in practice:

- 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

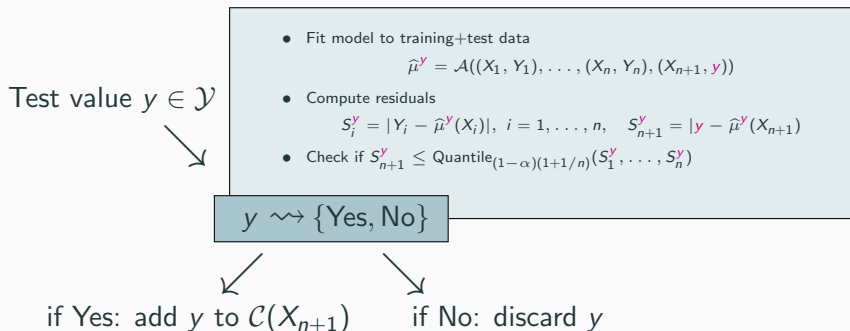
$$S_i^y = |Y_i - \hat{\mu}^y(X_i)|, \quad i = 1, \dots, n, \quad S_{n+1}^y = |y - \hat{\mu}^y(X_{n+1})|$$

- Check if $S_{n+1}^y \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(S_1^y, \dots, S_n^y)$

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

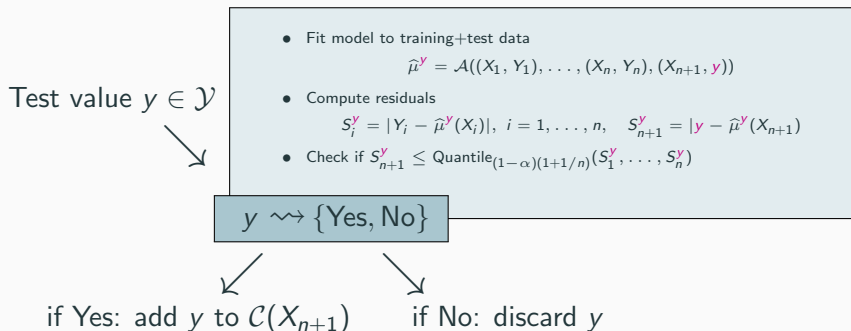
Full conformal prediction

Running full conformal in practice:



Full conformal prediction

Running full conformal in practice:



Note: split CP can be viewed as a special case of full CP:

\mathcal{A} returns a *pretrained* model $\hat{\mu}$ — doesn't depend on data

Theorem: full conformal⁴

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:

- 1 Need to verify

$$Y_{n+1} \in \mathcal{C}(X_{n+1}) \iff S_{n+1} \leq \text{Quantile}_{1-\alpha}(S_1, \dots, S_n, S_{n+1})$$

- 2 Need to verify that S_1, \dots, S_n, S_{n+1} are exchangeable

⁴Vovk et al 2005, *Algorithmic Learning in a Random World*

Full conformal prediction

$$\textcircled{1} Y_{n+1} \in \mathcal{C}(X_{n+1}) \iff S_{n+1} \leq \text{Quantile}_{1-\alpha}(S_1, \dots, S_n, S_{n+1})$$

By construction,

$$\text{If } y = Y_{n+1} \rightsquigarrow S_i^y = S_i \text{ for all } i = 1, \dots, n+1$$

train \mathcal{A} on $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y)$

train \mathcal{A} on $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1})$

Full conformal prediction

$$\textcircled{1} Y_{n+1} \in \mathcal{C}(X_{n+1}) \iff S_{n+1} \leq \text{Quantile}_{1-\alpha}(S_1, \dots, S_n, S_{n+1})$$

By construction,

$$\text{If } y = Y_{n+1} \rightsquigarrow S_i^y = S_i \text{ for all } i = 1, \dots, n+1$$

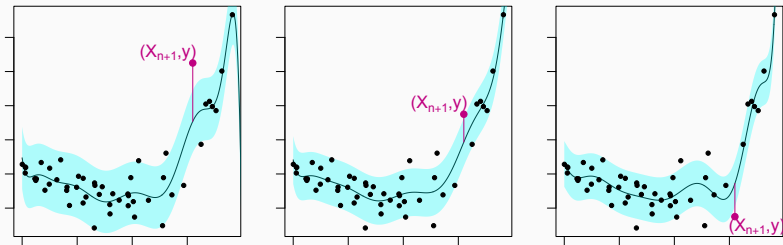
train \mathcal{A} on $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y)$

train \mathcal{A} on $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1})$

$$\begin{aligned} Y_{n+1} \in \mathcal{C}(X_{n+1}) &\iff S_{n+1}^{Y_{n+1}} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(S_1^{Y_{n+1}}, \dots, S_n^{Y_{n+1}}) \\ &\iff S_{n+1} \leq \text{Quantile}_{(1-\alpha)(1+1/n)}(S_1, \dots, S_n) \\ &\iff S_{n+1} \leq \text{Quantile}_{1-\alpha}(S_1, \dots, S_n, S_{n+1}) \end{aligned}$$

Full conformal prediction

How full conformal is run:



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

Full conformal prediction: general score

Full CP can be run with any conformal score function

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)

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

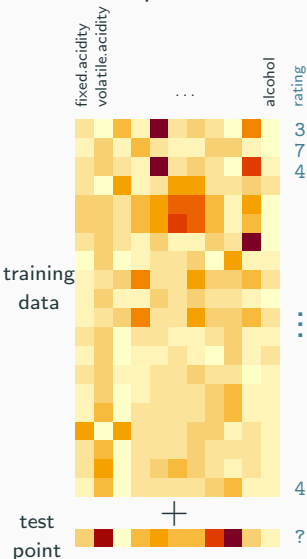
where

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

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

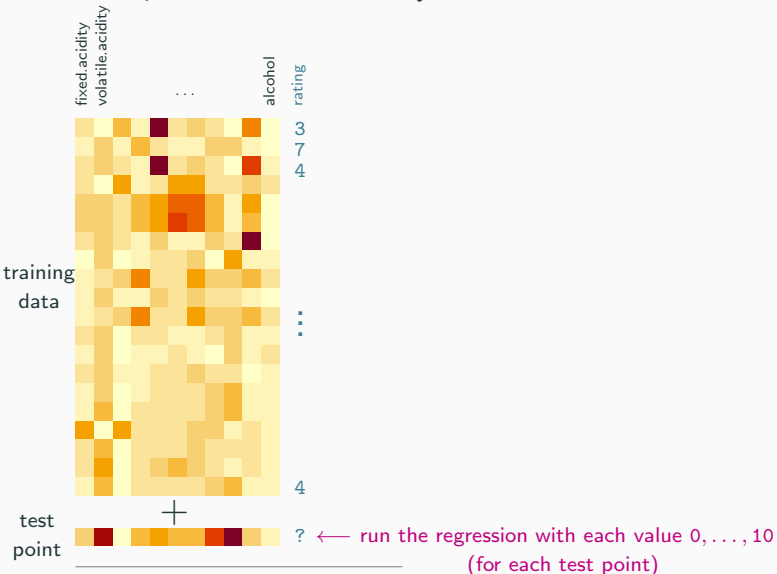
Example: the Wine Quality data set⁵



⁵Cortez et al 2009, Wine Quality data set, UCI Machine Learning Repository

Full conformal prediction: general score

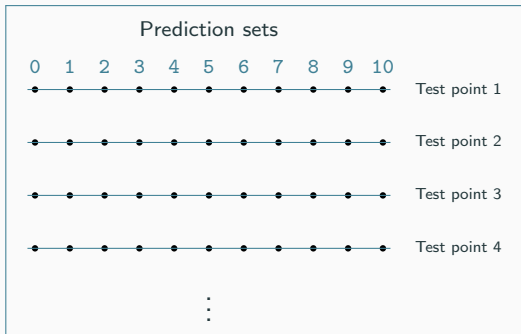
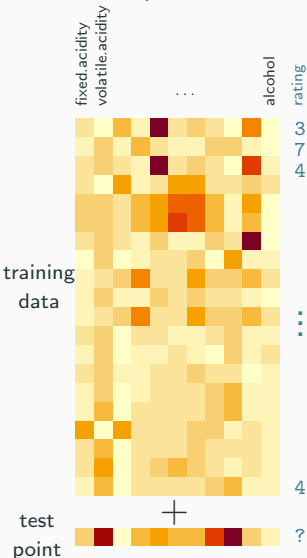
Example: the Wine Quality data set⁵



⁵Cortez et al 2009, Wine Quality data set, UCI Machine Learning Repository

Full conformal prediction: general score

Example: the Wine Quality data set⁵

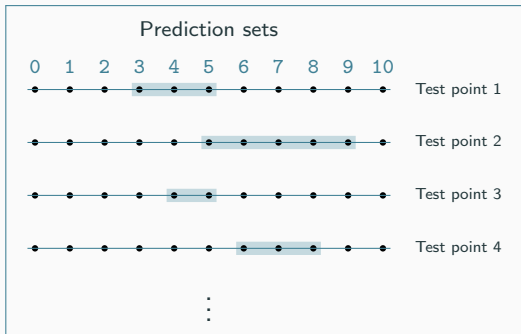
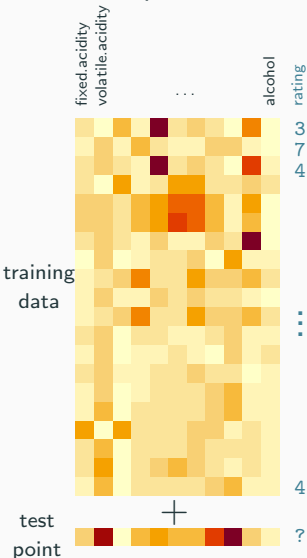


← run the regression with each value 0, ..., 10
(for each test point)

⁵Cortez et al 2009, Wine Quality data set, UCI Machine Learning Repository

Full conformal prediction: general score

Example: the Wine Quality data set⁵

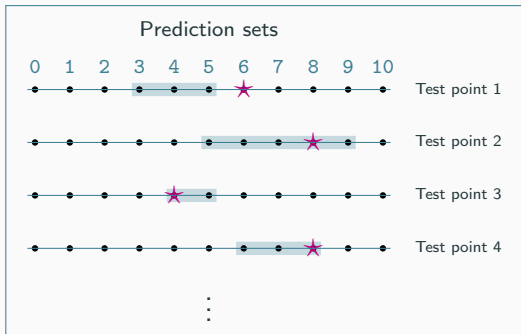
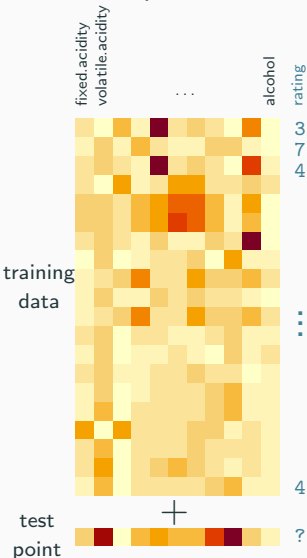


← run the regression with each value $0, \dots, 10$
(for each test point)

⁵Cortez et al 2009, Wine Quality data set, UCI Machine Learning Repository

Full conformal prediction: general score

Example: the Wine Quality data set⁵



run the regression with each value 0, ..., 10
(for each test point)

⁵Cortez et al 2009, Wine Quality data set, UCI Machine Learning Repository

Full conformal prediction: computational challenges

For a real-valued response...

- Running full CP requires refitting model for every value $y \in \mathbb{R}$

Full conformal prediction: computational challenges

For a real-valued response...

- Running full CP requires refitting model for every value $y \in \mathbb{R}$

Summary of approaches used in practice:

- Most common — restrict to a grid of y values (but no theory)
- Can use a discretized version of \mathcal{A} to restore theory⁶
- Specialized methods for specific algorithms/settings, e.g., Ridge,⁷ Lasso,⁸ stable algorithms⁹

⁶Chen, Chun, & B. 2017, *Discretized conformal prediction for efficient distribution-free inference*

⁷Burnaev & Vovk 2014, *Efficiency of conformalized ridge regression*

⁸Lei 2017, *Fast Exact Conformalization of Lasso using Piecewise Linear Homotopy*

⁹Ndiaye 2022, *Stable Conformal Prediction Sets*

Another look at the theory

Recall theoretical guarantee for split CP & full CP:

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

Limitations:

- Coverage is *marginal*, may not hold conditional on X_{n+1} — what if we undercover for certain subpopulations?
- Requires exchangeability — what if there is distribution drift?
- Full CP requires symmetric \mathcal{A}

See Part 2 for some methods to address these limitations

Conformal + CV

Using cross-validation for inference

Summarizing different methods:

- Split CP fits $\hat{\mu}$ to part of the data \rightsquigarrow distrib.-free theory
- Full CP: use all the data for $\hat{\mu}$ *and* achieves distrib.-free theory, but computationally very expensive
- Can cross-validation based methods offer a compromise?

Using cross-validation for inference

Leave-one-out CV (the “jackknife”):

$$\mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(\{|Y_i - \hat{\mu}_{-i}(X_i)|\}_{i=1,\dots,n})$$

fitted on all data



$\hat{\mu}_{-i}$ fitted on $\{(X_\ell, Y_\ell)\}_{\ell \neq i}$
(the leave-one-out model)



Using cross-validation for inference

Leave-one-out CV (the “jackknife”):

$$\mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(\{|Y_i - \hat{\mu}_{-i}(X_i)|\}_{i=1,\dots,n})$$

fitted on all data

$\hat{\mu}_{-i}$ fitted on $\{(X_\ell, Y_\ell)\}_{\ell \neq i}$
(the leave-one-out model)

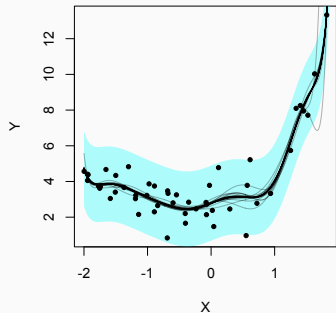
More computationally efficient: K -fold CV

- Partition $\{1, \dots, n\} = A_1 \cup \dots \cup A_K$, with $|A_k| = n/K$
- Fit models $\hat{\mu}_{-A_k}$ to data $\{(X_i, Y_i)\}_{i \notin A_k}$
- Compute the margin of error using residuals

$$\{|Y_i - \hat{\mu}_{-A_k}(X_i)|\}_{k=1,\dots,K; i \in A_k}$$

Using cross-validation for inference

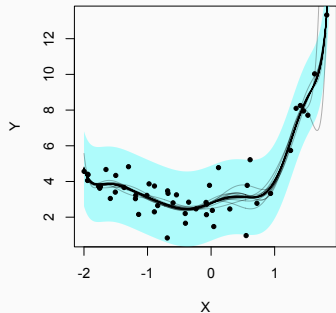
Leave-one-out CV: simulation



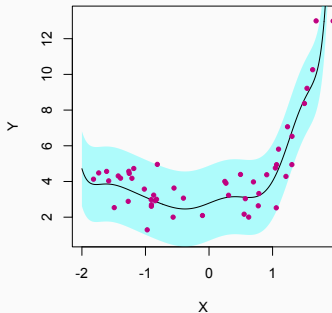
Train: 90% coverage

Using cross-validation for inference

Leave-one-out CV: simulation



Train: 90% coverage

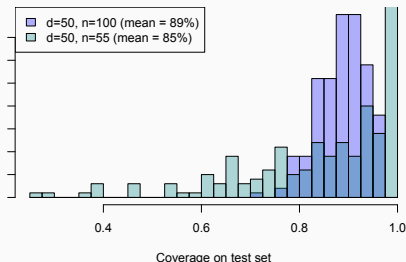


Test: 92% coverage

Using cross-validation for inference

However, no assumption-free theory for CV...

Example: least squares regression + jackknife



- Theoretical guarantees under asymptotic settings
- In practice, generally we see $\approx 1 - \alpha$ coverage, but unstable models may lead to undercoverage

Challenges for cross-validation

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

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

Challenges for cross-validation

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

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

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

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

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

Challenges for cross-validation

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

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

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

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

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

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

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

$\Rightarrow S_1, \dots, S_n, S_{n+1}$ are *not* exchangeable

$$\text{Jackknife: } \mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(S_i)$$

Jackknife can equivalently be defined as:

$$\mathcal{C}(X_{n+1}) = \left[\text{Quantile}_{\alpha}(\hat{\mu}(X_{n+1}) - S_i), \text{Quantile}_{1-\alpha}(\hat{\mu}(X_{n+1}) + S_i) \right]$$

$$\text{Jackknife: } \mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(S_i)$$

Jackknife can equivalently be defined as:

$$\mathcal{C}(X_{n+1}) = \left[\text{Quantile}_{\alpha}(\hat{\mu}(X_{n+1}) - S_i), \text{Quantile}_{1-\alpha}(\hat{\mu}(X_{n+1}) + S_i) \right] \\ - \text{Quantile}_{1-\alpha}(-\hat{\mu}(X_{n+1}) + S_i)$$

$$\text{Jackknife: } \mathcal{C}(X_{n+1}) = \widehat{\mu}(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(S_i)$$

Jackknife can equivalently be defined as:

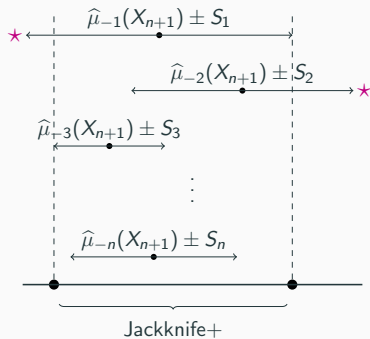
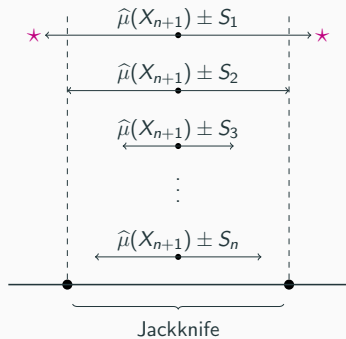
$$\mathcal{C}(X_{n+1}) = \left[\text{Quantile}_{\alpha}(\widehat{\mu}(X_{n+1}) - S_i), \text{Quantile}_{1-\alpha}(\widehat{\mu}(X_{n+1}) + S_i) \right] \\ - \text{Quantile}_{1-\alpha}(-\widehat{\mu}(X_{n+1}) + S_i)$$

A modified version of the method: the jackknife+.¹⁰

$$\mathcal{C}(X_{n+1}) = \left[-\text{Quantile}_{(1-\alpha)(1+1/n)}(-\widehat{\mu}_{-i}(X_{n+1}) + S_i), \right. \\ \left. \text{Quantile}_{(1-\alpha)(1+1/n)}(\widehat{\mu}_{-i}(X_{n+1}) + S_i) \right]$$

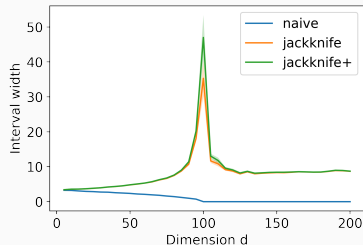
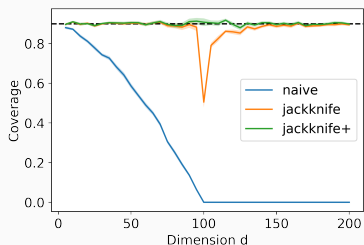
¹⁰B., Candès, Ramdas, Tibshirani 2019, *Predictive inference with the jackknife+*

Jackknife & jackknife+



Jackknife & jackknife+

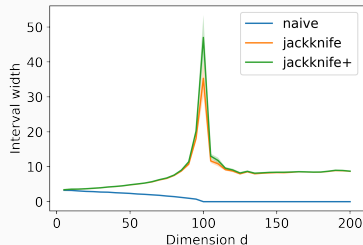
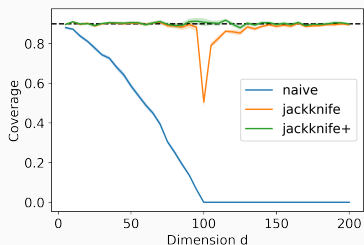
Empirical comparison (linear regression with $n = 100$):



- “Ridgeless” regression — minimum- ℓ_2 -norm solution, if $d > n$

Jackknife & jackknife+

Empirical comparison (linear regression with $n = 100$):



- “Ridgeless” regression — minimum- ℓ_2 -norm solution, if $d > n$
- Note: ridgeless regression is stable except the $d \approx n$ regime¹¹

¹¹Hastie et al 2022, *Surprises in High-Dimensional Ridgeless Least Squares Interpolation*

Theorem: coverage for jackknife+¹²

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

$$\mathbb{P}\{Y_{n+1} \in \mathcal{C}(X_{n+1})\} \geq 1 - 2\alpha$$

(In contrast, jackknife may have zero coverage, in the worst case)

¹²B., Candès, Ramdas, Tibshirani 2019, *Predictive inference with the jackknife+*

From leave-one-out to K -fold

To avoid computational cost of leave-one-out CV —
 K -fold CV (e.g., $K = 5$ or $K = 10$)

From leave-one-out to K -fold

To avoid computational cost of leave-one-out CV —
 K -fold CV (e.g., $K = 5$ or $K = 10$)

- Partition $\{1, \dots, n\}$ into K folds $A_1 \cup \dots \cup A_K$
- Fit model $\hat{\mu}_{-A_k} = \mathcal{A}\left(\{(X_i, Y_i) : i \in \{1, \dots, n\} \setminus A_k\}\right)$
- For $i \in A_k$ define $S_i = |Y_i - \hat{\mu}_{-A_k}(X_i)|$

$$\mathcal{C}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(S_1, \dots, S_n)$$

From leave-one-out to K -fold

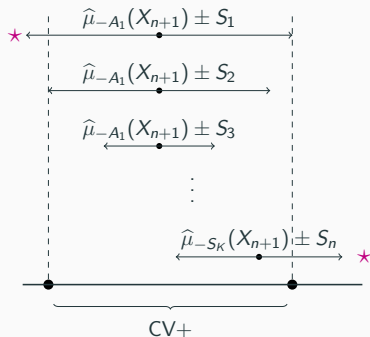
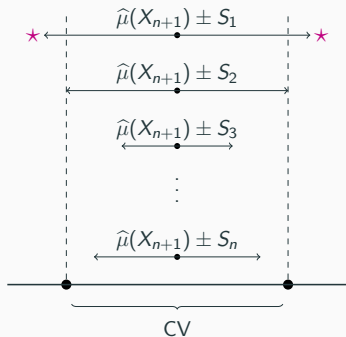
Generalize jackknife+ to the K -fold setting \rightsquigarrow CV+

K-fold CV+

- Partition $\{1, \dots, n\}$ into K folds $A_1 \cup \dots \cup A_K$
- Fit model $\hat{\mu}_{-A_k} = \mathcal{A}\left(\{(X_i, Y_i) : i \in \{1, \dots, n\} \setminus A_k\}\right)$
- For $i \in A_k$ define $S_i = |Y_i - \hat{\mu}_{-A_k}(X_i)|$
- Prediction set

$$\mathcal{C}(X_{n+1}) = \left[\begin{array}{l} -\text{Quantile}_{(1-\alpha)(1+1/n)}\left(\{-\hat{\mu}_{-A_k}(X_{n+1}) + S_i\}\right), \\ \text{Quantile}_{(1-\alpha)(1+1/n)}\left(\{\hat{\mu}_{-A_k}(X_{n+1}) + S_i\}\right) \end{array} \right]$$

From leave-one-out to K -fold



Cross-conformal prediction

CV+ is related to a more general method:

Cross-conformal prediction^{13,14}

- Partition $\{1, \dots, n\}$ into K folds $A_1 \cup \dots \cup A_K$
- Fit score function $s^{(k)} = \mathcal{A}\left(\{(X_i, Y_i) : i \in \{1, \dots, n\} \setminus A_k\}\right)$
- For $i \in A_k$ define $S_i = s^{(k)}(X_i, Y_i)$
- Prediction set

$$\mathcal{C}(X_{n+1}) = \left\{ y \in \mathcal{Y} : \sum_{k=1}^K \sum_{i \in A_k} \mathbb{1}\{S_i \geq s^{(k)}(X_{n+1}, y)\} \geq \alpha(n+1) \right\}$$

¹³Vovk 2015, *Cross-conformal predictors*

¹⁴Vovk et al 2018, *Cross-conformal predictive distributions*

Relating cross-conformal & CV+

For the residual score function $s(x, y) = |y - \hat{\mu}(x)|$,

$$\mathcal{C}_{\text{cross-conf.}}(\mathcal{X}_{n+1}) \subseteq \mathcal{C}_{\text{CV+}}(\mathcal{X}_{n+1})$$

Comparison:

- Cross-conformal is more flexible (can use any score function)
- CV+ always returns an interval (by construction)

Theorem: coverage for CV+ and cross-conformal

If Z_1, \dots, Z_{n+1} are i.i.d., and \mathcal{A} is symmetric,
then K -fold cross-conformal satisfies

$$\mathbb{P}\{Y_{n+1} \in \mathcal{C}(X_{n+1})\} \geq \begin{cases} 1 - 2\alpha - 2/K &^{15} \\ 1 - 2\alpha - 2K/n &^{16} \end{cases}$$

As a special case, the same is true for K -fold CV+.

\implies For any K ,

$$\mathbb{P}\{Y_{n+1} \in \mathcal{C}(X_{n+1})\} \geq 1 - 2\alpha - \frac{2}{\sqrt{n}}.$$

¹⁵B., Candès, Ramdas, Tibshirani 2019, *Predictive inference with the jackknife+*

¹⁶Vovk et al 2018, *Cross-conformal predictive distributions*

Conformal methods vs model-based methods

Conformal or classical?

In a practical application....

Should we use a model?

- Model is probably a good approximation
- Obtain more precise answers
- But, may lose coverage if assumptions don't hold

Should we use conformal prediction?

- Coverage doesn't depend on assumptions
- But, coverage guarantee is only marginal
- Would we get wider intervals (less informative)?

Conformal or classical?

Answer: use both, & get the best of both worlds!

Conformal prediction is a *family* of methods

- Choosing a score specifies a particular method
- Can incorporate models / assumptions into the score

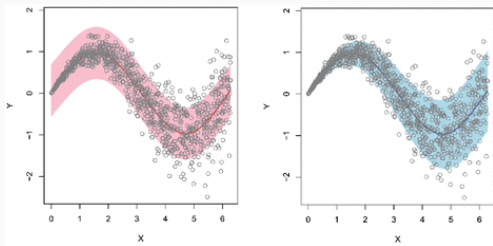
The conformal score: examples

Recall....

① The residual score: $s(x, y) = |y - \hat{\mu}(x)|$

② The scaled residual score:¹⁷

$$s(x, y) = \frac{|y - \hat{\mu}(x)|}{\hat{\sigma}(x)}, \text{ where } \hat{\mu}, \hat{\sigma} \text{ fitted on pretraining data}$$



(figure from Lei et al 2018)

¹⁷Lei et al 2018, *Distribution-Free Predictive Inference for Regression*

The conformal score: examples

③ Conformalized quantile regression:¹⁸

$$s(x, y) = \max \{y - \hat{\gamma}_{hi}(x), \hat{\gamma}_{lo}(x) - y\}$$

where $\hat{\gamma}_{lo}, \hat{\gamma}_{hi}$ fitted on pretraining data

estimated quantiles of $Y|X$

$$\implies \mathcal{C}(X_{n+1}) = [\hat{\gamma}_{lo}(X_{n+1}) - \hat{q}, \hat{\gamma}_{hi}(X_{n+1}) + \hat{q}]$$

¹⁸Romano et al 2019, *Conformalized quantile regression*

The conformal score: examples

3 Conformalized quantile regression:¹⁸

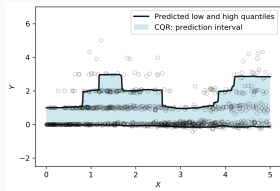
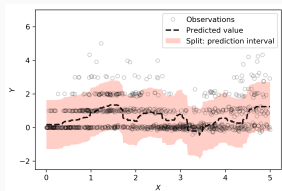
$$s(x, y) = \max \{y - \hat{\gamma}_{hi}(x), \hat{\gamma}_{lo}(x) - y\}$$

where $\hat{\gamma}_{lo}, \hat{\gamma}_{hi}$ fitted on pretraining data

estimated quantiles of $Y|X$

$$\implies \mathcal{C}(X_{n+1}) = [\hat{\gamma}_{lo}(X_{n+1}) - \hat{q}, \hat{\gamma}_{hi}(X_{n+1}) + \hat{q}]$$

Compare to residual score:



(figure from Romano et al 2019)

¹⁸Romano et al 2019, *Conformalized quantile regression*

The conformal score: examples

④ Distributional conformal prediction:¹⁹

$s(x, y) = |\widehat{F}(y|x) - 0.5|$ where $\widehat{F}(\cdot|x)$ is fitted on pretraining data
estimated conditional CDF
of Y given $X = x$

$$\implies \mathcal{C}(X_{n+1}) = \left[\widehat{F}^{-1}(0.5 - \widehat{q} | X_{n+1}), \widehat{F}^{-1}(0.5 + \widehat{q} | X_{n+1}) \right]$$

¹⁹Chernozhukov et al 2019, *Distributional conformal prediction*

The conformal score: examples

- ⑤ The high-density score:²⁰

$s(x, y) = -\hat{f}(y|x)$ where $\hat{f}(\cdot|x)$ is fitted on pretraining data
estimated conditional density
of Y given $X = x$

$$\implies \mathcal{C}(X_{n+1}) = \left\{ y \in \mathcal{Y} : \hat{f}(y|X_{n+1}) \geq -\hat{q} \right\}$$

²⁰Izbicki et al 2020, *Flexible distribution-free conditional predictive bands using density estimators*

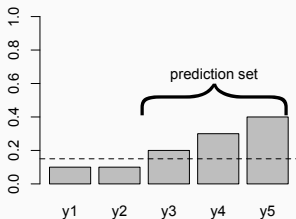
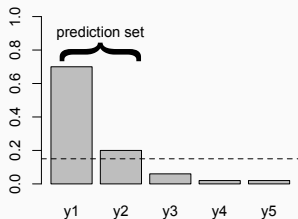
The conformal score: examples

If the response Y is categorical, with values $\mathcal{Y} = \{y_1, \dots, y_K\}$ —

⑥ The high-probability score:

$s(x, y_k) = -\hat{p}_k(x)$ where $\hat{p}_k(x)$ is fitted on pretraining data
estimate of $\mathbb{P}\{Y = y_k \mid X = x\}$

$$\implies \mathcal{C}(X_{n+1}) = \{y_k : \hat{p}_k(X_{n+1}) \geq -\hat{q}\}$$



Model-based theory for conformal

A general recipe:

- Suppose that if our model is correct, the “oracle” answer can be written in the form

$$\mathcal{C}^*(X_{n+1}) = \{y : s^*(X_{n+1}, y) \leq q^*\}$$

Model-based theory for conformal

A general recipe:

- Suppose that if our model is correct, the “oracle” answer can be written in the form

$$\mathcal{C}^*(X_{n+1}) = \{y : s^*(X_{n+1}, y) \leq q^*\}$$

- Compare to the split conformal prediction set:

$$\mathcal{C}(X_{n+1}) = \{y : s(X_{n+1}, y) \leq \hat{q}\}$$

\rightsquigarrow for conformal to approximate the oracle, need

$$s \approx s^* \text{ and } \hat{q} \approx q^*$$

Model-based theory for conformal

A general recipe:

- Suppose that if our model is correct, the “oracle” answer can be written in the form

$$\mathcal{C}^*(X_{n+1}) = \{y : s^*(X_{n+1}, y) \leq q^*\}$$

- Compare to the split conformal prediction set:

$$\mathcal{C}(X_{n+1}) = \{y : s(X_{n+1}, y) \leq \hat{q}\}$$

↪ for conformal to approximate the oracle, need

$$s \approx s^* \text{ and } \hat{q} \approx q^*$$

↑
relies on fitting a good model
using the pretraining data

↑
relies on concentration of quantiles
for calibration set scores

Model-based theory for conformal: examples

- 1 If the true model is $Y = \mu(X) + \epsilon$ where $\epsilon \perp\!\!\!\perp X$ is symmetric & unimodal noise,

$$\begin{aligned} \mathcal{C}^*(X_{n+1}) &= \mu(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(|\epsilon|) \\ &= \{y : s^*(X_{n+1}, y) \leq \text{Quantile}_{1-\alpha}(|\epsilon|)\} \end{aligned}$$

for $s^*(x, y) = |y - \mu(x)|$

²¹Lei et al 2018, *Distribution-Free Predictive Inference for Regression*

Model-based theory for conformal: examples

- ① If the true model is $Y = \mu(X) + \epsilon$ where $\epsilon \perp\!\!\!\perp X$ is symmetric & unimodal noise,

$$\begin{aligned} \mathcal{C}^*(X_{n+1}) &= \mu(X_{n+1}) \pm \text{Quantile}_{1-\alpha}(|\epsilon|) \\ &= \{y : s^*(X_{n+1}, y) \leq \text{Quantile}_{1-\alpha}(|\epsilon|)\} \end{aligned}$$

for $s^*(x, y) = |y - \mu(x)|$

\implies if we use the residual score, and if $\hat{\mu} \rightarrow \mu$, then²¹

$$\mathcal{C}(X_{n+1}) \approx \mathcal{C}^*(X_{n+1}) \text{ as } n \rightarrow \infty$$

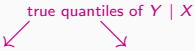
²¹Lei et al 2018, *Distribution-Free Predictive Inference for Regression*

Model-based theory for conformal: examples

- 2 In the regression setting, suppose we would like an equal-tailed, conditional coverage guarantee.

Then the optimal set is

true quantiles of $Y \mid X$


$$\mathcal{C}^*(X_{n+1}) = [\gamma_{\alpha/2}(X_{n+1}), \gamma_{1-\alpha/2}(X_{n+1})]$$

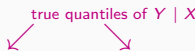
²²Romano et al 2019, *Conformalized quantile regression*; Sesia & Candès 2020, *A comparison of some conformal quantile regression methods*

Model-based theory for conformal: examples

- ② In the regression setting, suppose we would like an equal-tailed, conditional coverage guarantee.

Then the optimal set is

true quantiles of $Y \mid X$


$$C^*(X_{n+1}) = [\gamma_{\alpha/2}(X_{n+1}), \gamma_{1-\alpha/2}(X_{n+1})]$$

Can rewrite as

$$C^*(X_{n+1}) = \{y : s^*(X_{n+1}, y) \leq 0\}$$

$$\text{where } s^*(x, y) = \max \{y - \gamma_{1-\alpha/2}(x), \gamma_{\alpha/2}(x) - y\}$$

²²Romano et al 2019, *Conformalized quantile regression*; Sesia & Candès 2020, *A comparison of some conformal quantile regression methods*

Model-based theory for conformal: examples

- ② In the regression setting, suppose we would like an equal-tailed, conditional coverage guarantee.

Then the optimal set is

$$C^*(X_{n+1}) = [\gamma_{\alpha/2}(X_{n+1}), \gamma_{1-\alpha/2}(X_{n+1})]$$

true quantiles of $Y \mid X$

Can rewrite as

$$C^*(X_{n+1}) = \{y : s^*(X_{n+1}, y) \leq 0\}$$

where $s^*(x, y) = \max \{y - \gamma_{1-\alpha/2}(x), \gamma_{\alpha/2}(x) - y\}$

If $\hat{\gamma}_{lo} \rightarrow \gamma_{\alpha/2}$ and $\hat{\gamma}_{hi} \rightarrow \gamma_{1-\alpha/2}$,²²

$$C(X_{n+1}) \approx C^*(X_{n+1}) \text{ as } n \rightarrow \infty$$

²²Romano et al 2019, *Conformalized quantile regression*; Sesia & Candès 2020, *A comparison of some conformal quantile regression methods*

Model-based theory for conformal: examples

- ③ In the categorical setting, with conditional PMF $p(y | x)$ — smallest possible prediction set with *marginal* coverage is

$$\begin{aligned}\mathcal{C}^*(X_{n+1}) &= \{y : p(y | x) \geq t\} \\ &= \{y : s^*(x, y) \leq -t\} \text{ where } s^*(x, y) = -p(y | x)\end{aligned}$$

\implies if we use the high-probability score, & $\hat{p} \rightarrow p$, then²³

$$\mathcal{C}(X_{n+1}) \approx \mathcal{C}^*(X_{n+1}) \text{ as } n \rightarrow \infty$$

²³Sadinle et al 2019, *Least ambiguous set-valued classifiers with bounded error levels*

Theorem (informal): asymptotic results for split CP²⁴

Assume Z_1, Z_2, \dots are i.i.d., & the data split satisfies $n_0, n_1 \rightarrow \infty$.

If $s_n \rightarrow s^*$, then

$$|\mathcal{C}(X_{n+1}) \Delta \mathcal{C}^*(X_{n+1})| \rightarrow 0.$$

²⁴Duchi et al 2024, *Predictive inference in multi-environment scenarios*; Angelopoulos, B., Bates 2024, *Theoretical Foundations of Conformal Prediction*

Summary

Summary: part 1

- Conformal allows us to start with any algorithm,
& calibrate it to achieve (marginal) predictive coverage
- Tradeoff between statistical & computational efficiency:
Split CP, full CP, and CV-based versions
- Conformal + model-based methods \rightsquigarrow “best of both worlds”

In part 2, we will ask if conformal can be extended to handle:

- The streaming-data setting
- Distribution shift & distribution drift
- Conditional coverage rather than only marginal coverage
- & other extensions