

Lecture 14

LVII posted

CP recap

Missing data

403 Course Overview

- EDF \hat{F} and consistency of means
 - MC integration \leftarrow Application
 - How to sample - F^{-1}
 - importance s.
 - rejection s.
 - * (CV)
 - Bootstrap $\leftarrow \text{Var}(\hat{\theta}) \approx ?$
 - *
 - Conformal prediction ✓ \leftarrow
 - *
 - Imputation (fill missing data) \leftarrow
 - *
 - MCMC (Markov Chain MC)
- ↑
Resampling

Conformal prediction

new x

AFTER
TRAINING

- Conformal prediction: CI for a single prediction $\hat{y} = f(x)$

Given data set $\mathcal{D} = \{(x_i, y_i), i = 1 : N\}$

Training algorithm \mathcal{A} , so that $\mathcal{A}(\mathcal{D}) = f$ the predictor

Want CI for $\hat{y} = f(x)$ where x is a new data point

so that the CI is NOT dependent on \mathcal{A} being statistically correct (e.g. \mathcal{A} overfits, ...)

GOOD !!

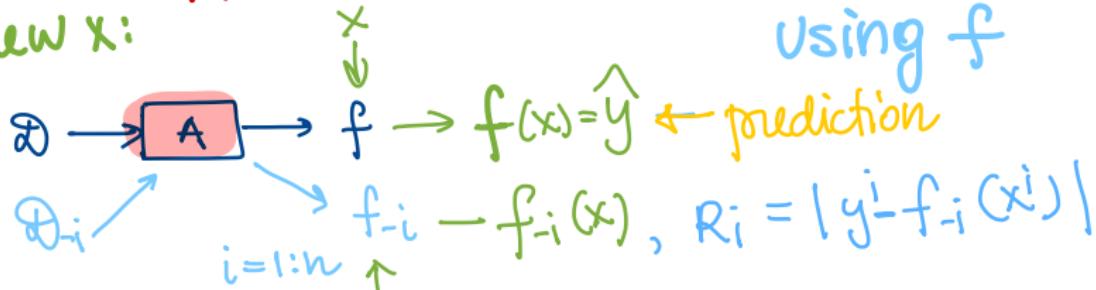
- jackknife+ is a simple algorithm for CP
- More advanced algorithms exist. This is an active area of research in statistics.



!!!! Do NOT use CP to "crossvalidate" your algorithm!

jackknife+ $C_{1\alpha} = [a, b]$ [Training, CV, Model Selection, ...]
 MUST contain ALL

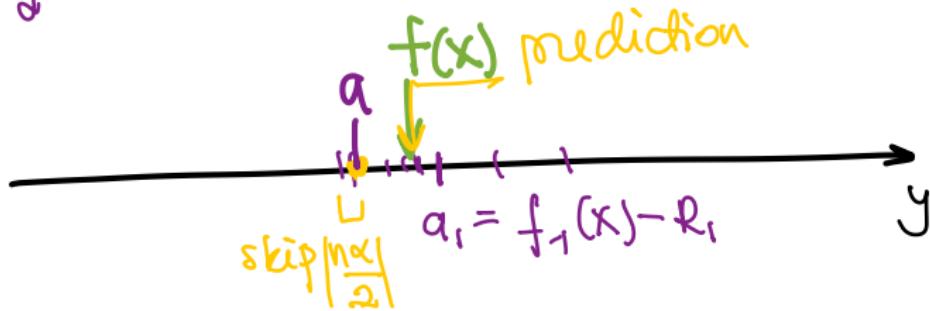
new X :



$$\alpha = 20\%$$

$$\frac{\alpha}{2} = 10\%$$

$q = \left[n \frac{\alpha}{2} \right]$ quantile $\rightarrow \{a_i, i=1:n\}$



Lecture Notes VII – Missing Data and Imputation

"
fill in

Marina Meila
mmp@stat.washington.edu

Department of Statistics
University of Washington

April 2025

The missing data problem ↪

Missing data

data $y_{1,2,\dots,N}$ some missing values

$x_{1:n}$ $y_{1:n}$ observed

$x_{n+1:N}$ $y_{\underline{n+1:N}}$ unobserved

other information, known

if want \hat{y} , param of interest

$y_{n+1:n+m} =$ votes to be removed

$x_{n+1:n+m}$ \uparrow impute $=$ addresses

$x_{n+1:n+m}$ \downarrow impute $=$ gender

Missing data: the problem

- Data sampled i.i.d. $\sim F$, but some (parts of each sample) are missing
- Here we denote data point = (X^i, Y^i) where X^i is always observed, Y^i may be missing or not.

complete

$$R_i = \begin{cases} 1 & Y^i \text{ observed} \\ 0 & \text{otherwise} \end{cases}$$

- A new random variable R , $R^i = 0$ when Y^i not observed, $R^i = 1$ when observed

- What is the distribution of R ?

MCAR Missing Completely At Random $R \perp X, Y$

MAR Missing At Random $R \perp Y | X$

dependence of R
on X, Y

MCAR

$$Y_{n+1:m} = ?$$

MAR

$$R \perp Y | X$$

and R depends on X , but
NOT on Y

A, B events, r.v.'s
 $A \perp B$

A independent of B

$$P(A, B) = P(A)P(B)$$

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

? ↑
know-A

What is the distribution of R ? *how does R depend on X and Y ?*

- ▶ Let data $\mathcal{D} = \{(x^1, y^1), \dots, (x^n, y^n)\} \cup \{x_{n+1}, \dots, x_{n+m}\}$
 - ▶ complete data $\{(x^1, y^1), \dots, (x^n, y^n)\}$ with $R^{1:n} = 1$
 - ▶ data with missing values $\{x_{n+1}, \dots, x_{n+m}\}$ with $R^{n+1:n+m} = 0$

MCAR Missing Completely At Random $R \perp X, Y$

- ▶ $\{(x^1, y^1), \dots, (x^n, y^n)\}$ has no information about $\{y^{n+1}, \dots, y^{n+m}\}$

MAR Missing At Random $R \perp Y | X$

- ▶ $\{(x^1, y^1, R^1), \dots, (x^n, y^n, R^n)\}$ together with $\{(x^{n+1}, R^{n+1}), \dots, (x^{n+m}, R^{n+m})\}$ has information about $\{y^{n+1}, \dots, y^{n+m}\}$
- ▶ Probability distribution of missing values $Y^{n+1:n+m}$ is

$$Y \sim p(Y | X, R = 0) \quad \text{and} \quad p(Y | X, R = 0) = p(Y | X, R = 1) \quad (1)$$

use it to predict Y

$$R \perp Y | X \Leftrightarrow P(Y | X, R) = P(Y | X)$$

MAR $R \perp Y | X$

and R depends on X , but
NOT on Y