

Cross Validation

K

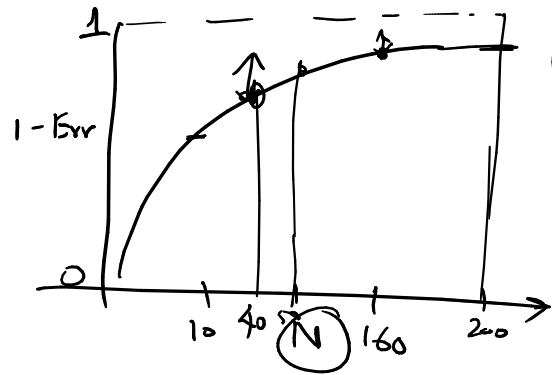
① K may impact \hat{Err}

A. $K=5$
 $N=200$

B. $K=5$
 $N=50$

$$\lambda \in [10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 100]$$

$$K=N. \quad \{x_i, y_i\}_{i=1}^N = D$$



$P_{X,Y}$
models

2). Wrong way: Preprocessing the data using both x_i s and y_i s. (training)

Right way: $K=5$

for each choice:

for each fold:

preprocessing here. on the remaining folds.]

~~*~~
N=50
P=5000

X_i are Gaussians
 $X_i \perp Y$

$$Y = 0 \cdot X_1 + 0X_2 + \dots + 0X_{5000} + \epsilon$$

Bootstrap

① Procedure

② Understanding

$$\left\| \frac{\hat{\beta}_j}{\sqrt{V_j}} \sim \underline{\underline{N(0,1)}} \right.$$

①

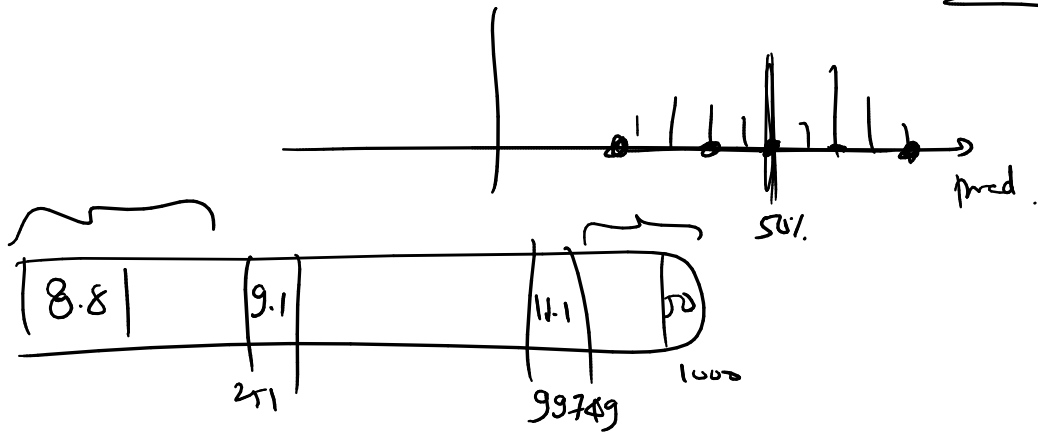
$$\{Z_i\}_{i=1}^N = Z$$

Create bootstrap samples.

5 times 1st row was sampled.
↓
Z^{*1}, ... Z^{*B}.
↓ N obs each.
↓
f^{*1}

↓ 10,000

$\hat{f}^{\#1}(x^{\text{test}})$, $\hat{f}^{\#2}(x^{\text{test}})$, \dots \quad B \quad \underline{\underline{10000}}



Estimate Err

$$E_Z E_{P_{XY}} [(\hat{f}_Z(x) - Y)^2] \quad Z \sim (P_{XY})^N$$

=

for each obs i

pick the ones it is not part of.

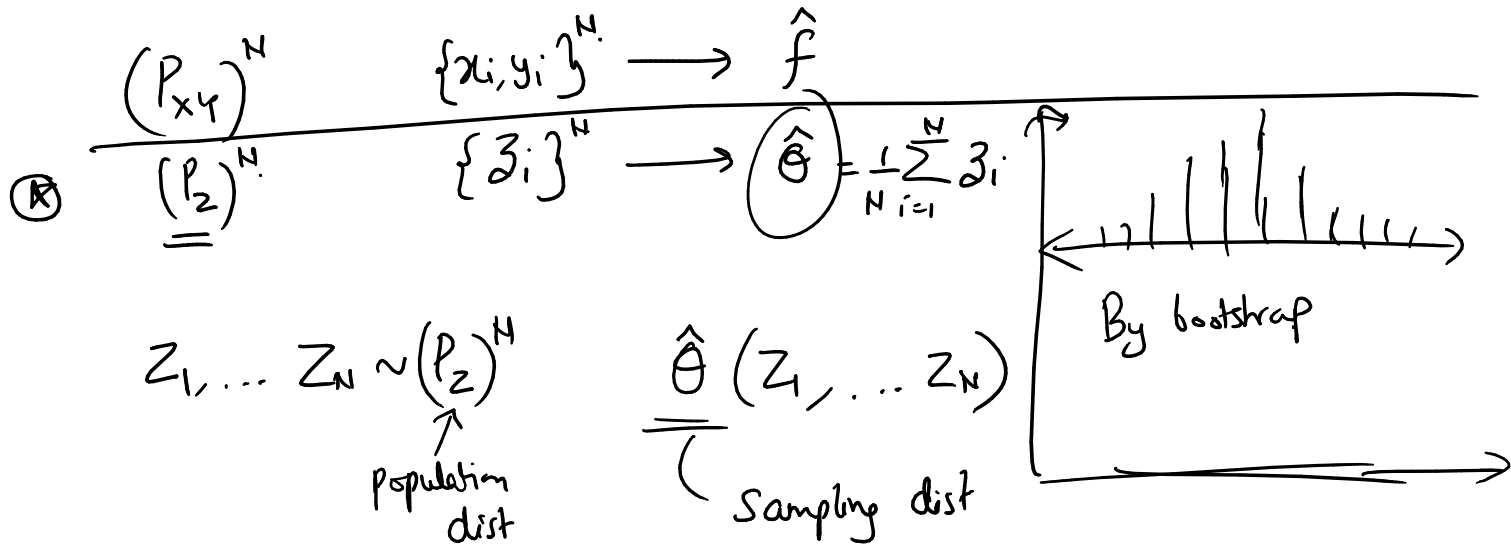
$$\{x_i, y_i\}_{i=1}^N \rightarrow \hat{f}$$

$$C_{-i} = \{1, \dots, 100\}$$

$$\underline{\underline{E_{P_{XY}} [(\hat{f}(x) - Y)^2]}}$$

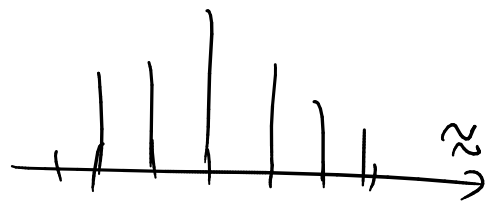
$$Err \approx \frac{1}{N} \sum_{i=1}^N \underbrace{\frac{1}{|C_{-i}|} \sum_{b \in C_{-i}} (f^{\hat{*}b}(x_i) - y_i)^2}$$

Understanding:



$$(z_1, \dots, z_N)^{*1} \sim (P_Z)^N \rightarrow \hat{\theta}_1$$

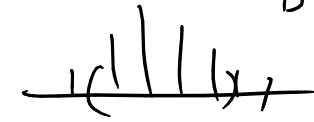
$$(z_1, \dots, z_N)^{*1} \sim (P_Z)^N \rightarrow \hat{\theta}_2$$



Sampling
dist of $\hat{\theta}(z_1, \dots, z_N)$

$$\{z_1, \dots, z_N\} \sim (P_Z)^N$$

w/R. B times



$\hat{\theta}_1$
⋮
 $\hat{\theta}_B$

MLE.

$$z_1, \dots, z_N \sim \underline{\underline{N(\mu, 1)}}$$

likelihood of data for a given parameter.

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N z_i$$

$$P_{\mu}(z_1, \dots, z_N) = \prod_{i=1}^N \underbrace{P(z_i)}_{\propto e^{-\frac{(z_i - \mu)^2}{2}}}$$

\swarrow

$$LL = \sum_{i=1}^N \left(-\frac{(z_i - \mu)^2}{2} \right) + C$$

least squares : MLE interpretation.

$$\frac{P(Y|X=x)}{\beta, \sigma} \sim N(\beta^T x, \sigma^2)$$
$$\{x_i, y_i\}_{i=1}^N$$

$$LL \propto \sum_{i=1}^N \log \left(\exp \left(- \frac{(\beta^T x_i - y_i)^2}{2\sigma^2} \right) \right) + C(\sigma)$$

$$\max_{\beta} LL \quad \equiv \quad \min_{\beta} \sum_{i=1}^N (y_i - \beta^T x_i)^2$$

$$P_{XY} \left[\begin{array}{l} X \text{ is not random.} \\ Y = \beta^T X + \varepsilon \\ \varepsilon \sim N(0, \sigma^2) \end{array} \right]$$

Ridge regression.

$$\min_{\beta} \left(\sum_{i=1}^N (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_2^2 \right)$$

MAP estimation

$$\max_{\beta} \underline{P(\beta | \text{data})} \propto P(\text{data} | \beta) \cdot P(\beta)$$

$$\Rightarrow \max_{\beta} \log P(\beta | \text{data}) \propto \log P(\text{data} | \beta) + \log P(\beta)$$

$$\Rightarrow \text{MAP est} \equiv \min_{\beta} \sum_{i=1}^N (y_i - x_i^T \beta)^2 + \frac{\lambda}{\tau^2} \beta^T \beta$$

$$\left. \begin{aligned} P(\beta | \text{data}) \\ \propto \underline{P(\text{data} | \beta)} \cdot \overline{P(\beta)} \\ P(\beta) \sim \mathcal{N}(0, \tau^2 \mathbf{I}) \end{aligned} \right\} \underline{\underline{- \beta^T \Sigma^{-1} \beta}}$$

$$P(\beta) \sim N(0, z^2 I)$$

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad z \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\frac{1}{\sqrt{2\pi} \det(z^2 I)} \cdot \exp\left(-\frac{1}{2} (\beta - 0)^T \underbrace{\Sigma^{-1}}_{(z^2 I)^T} (\beta - 0)\right)$$

$$\frac{1}{z^2} I$$

$$-\frac{1}{2} \beta^T \left(\frac{1}{z^2} I\right) \beta = -\frac{1}{2z^2} \beta^T \beta$$

① LASSO: Belief on β each coordinate is Laplace distributed with 0 mean.

$$\lambda \|\beta\|_1$$
$$\equiv \lambda \sum_{j=1}^p |\beta_j|$$

$$P(\beta_j) \propto \exp\left(-\frac{|\beta_j|}{c}\right)$$

Naive Bayes. for classification: MLE

$$P(x_i | g_i) = \prod_{j=1}^p P(x_{ij} | g_i)$$

$$G = \left\{ \begin{array}{cc} \text{Span} & \text{No-Span} \\ 1 & 2 \end{array} \right\}$$

$$x^{\text{test}}: \underset{k \in \{1, 2\}}{\text{argmax}} \underbrace{P(G=k | x^{\text{test}})}$$

$$x_1 \dots x_p$$

↑

$$\underbrace{P(x^{\text{test}} | G=k)} \cdot \underbrace{P(G=k)}$$

$$\hat{P}(G=1) = \frac{\# \text{ class 1}}{N}$$

$$\hat{P}(x_{j=1} | G=k) = ? \quad \text{for } j=1, \dots, p, k=1, 2.$$