

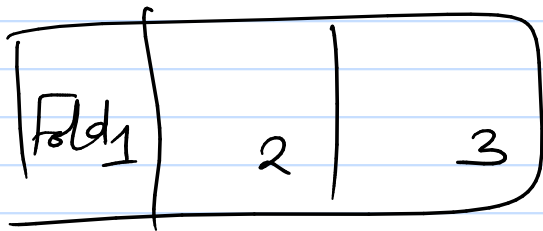
Model  
assessment

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$$X^{\text{full}} = \begin{bmatrix} X \\ X^{\text{test}} \end{bmatrix} \quad \text{dcol}$$

Selection

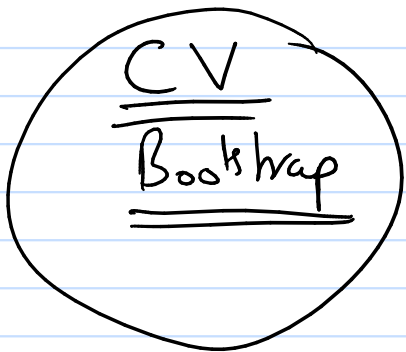
CV



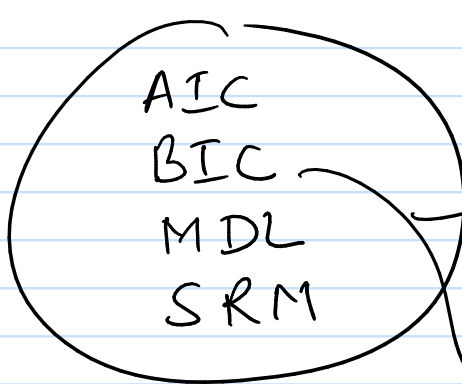
$$\lambda = 10^{-2}$$

$$\hat{\beta}_{12}$$

perf 1  
perf 3  
perf 2

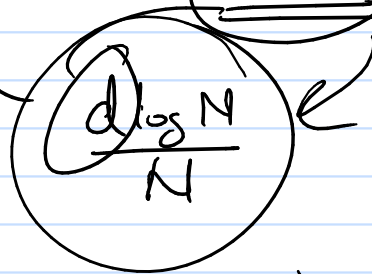


estimate Err  
directly

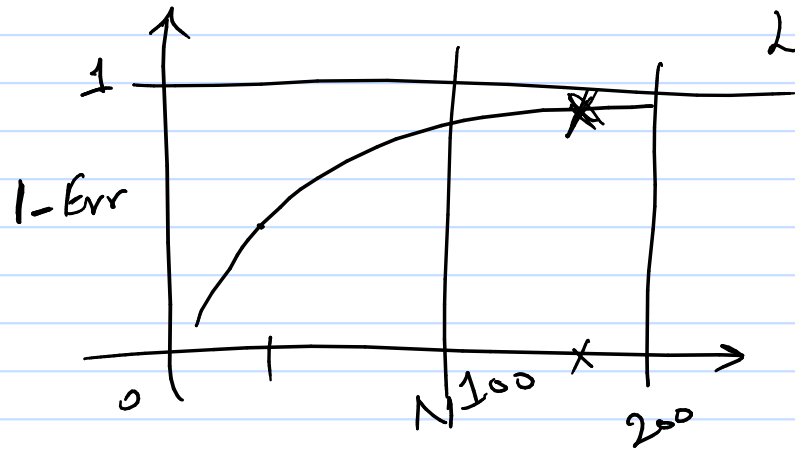


$\lambda = 10^{-2} \rightarrow \hat{\beta}$   
use training error.

+ Bonus



Learning curve.

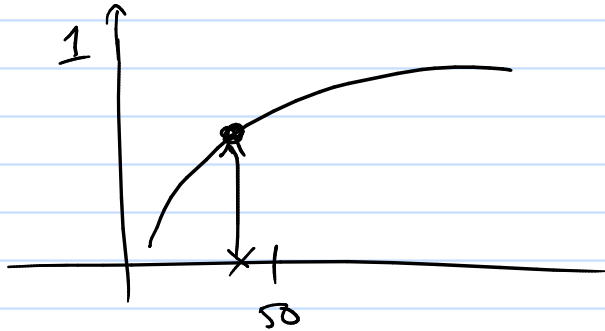


k-fold

5

40 obs in each fold.

160 obs  $\rightarrow \hat{\beta}$



for every choice of model.

k fold  
5

④ → Pick 4 folds

$N = 50$      $p = 5000$

$X_i \sim N(0,1)$

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

① Scatter  $Y, X_i$

$G = \mathbb{1}[Y > 0]$

pick 100

Reality: 50%.

build a classifier. → 3%

①  
k=1  
k-NN

Naive Bayes:  $W = [w_1 \ w_2 \ w_3]^T$   $p = 3$

$$P_{\theta}(Y = \text{spam} \mid \underline{W} = (w_1, w_2, w_3)) > P_{\theta}(Y = \text{not spam} \mid W = w)$$

$$P_{\theta}(Y = y \mid W = w) \propto P_{\theta}(\underline{W} = w \mid Y = y) \cdot P_{\theta}(Y = y)$$

↑  
spam

$$\star P(w_1 = w_1 \mid Y = y) \cdot P(w_2 = w_2 \mid Y = y) \cdot P(w_3 = w_3 \mid Y = y)$$

$$\rightarrow = \frac{\#(w_1 \text{ appears } \& \text{ spam})}{\# \text{ spam}}$$

$$P(Y = \text{Spam}) \approx \frac{\# \text{ Spams}}{N} \hat{\theta}_1 : \text{MLE.}$$

$$\frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

$$P_{\theta_1}(Y_1 = y_1) = \theta_1^{\mathbb{1}[y_1 = \text{spam}]} \cdot (1 - \theta_1)^{\mathbb{1}[y_1 = \text{not spam}]}$$

$$\log P(Y_1 = y_1) + \log P(Y_2 = y_2) \dots + \log P(Y_N = y_N)$$

$$\sum_{i=1}^N \left( \mathbb{1}[y_i = \text{spam}] \log \theta_1 + \mathbb{1}[y_i = \text{not spam}] \log (1 - \theta_1) \right)$$

# Spams
(N - # Spams)



$\hat{\mu}, \hat{\sigma}$



Gaussian  
mixture  
m-dels.

Bootstrap:

$\sum$

$x_1 y_1$

$x_2 y_2$

$\vdots$

$x_N y_N = \text{loss}$

$\sum_{N}^{*1}$

$\hat{f}^{*1}$

$\sum_{N}^{*2}$

$\hat{f}^{*2}$

$\sum_{N}^{*(B)}$

$\hat{f}^{*(B)}$

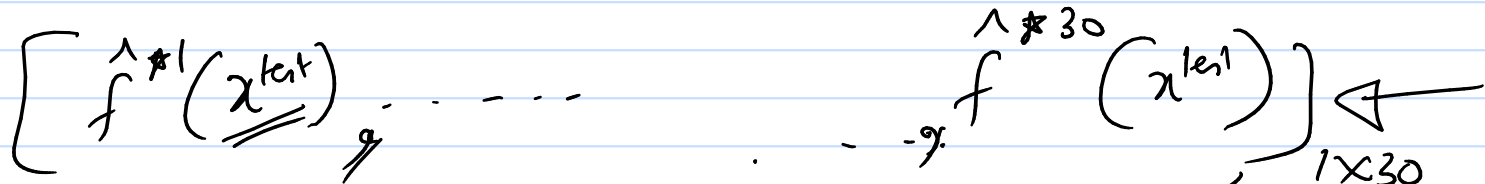
30

①  $\frac{1}{B} \sum_{b=1}^B \hat{f}^{*(b)}(x^{\text{test}})$

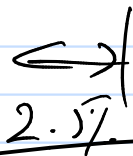
② Estimate "Err": average (held out data for each  $\hat{f}^{*(b)}$ )  
performance on



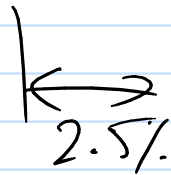
Confidence bands around estimates:  $\hat{\beta}$



95%



2.5%



2.5%

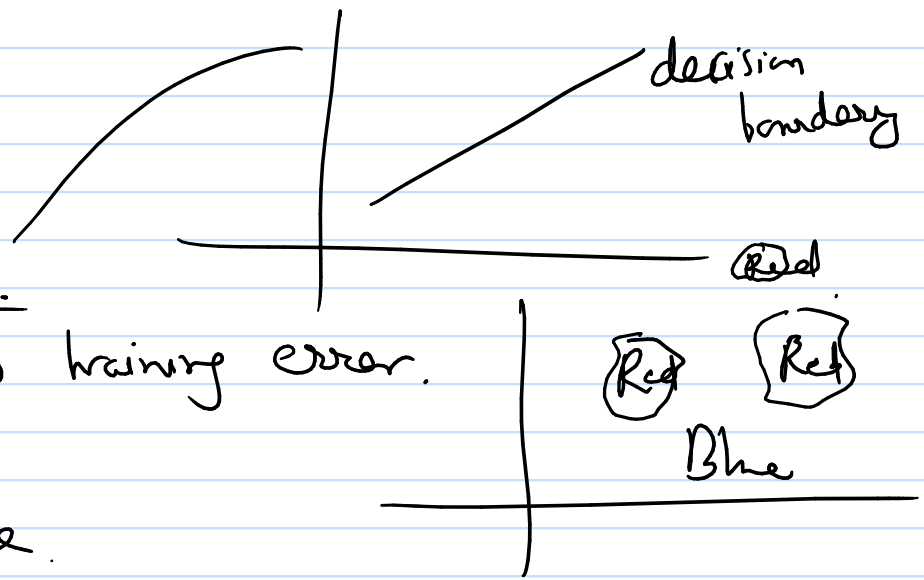
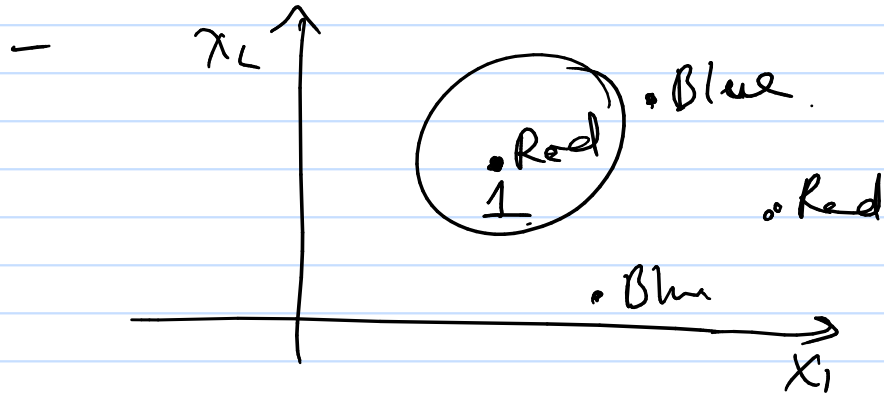
$$Y = \beta^T X + \varepsilon \leftarrow N(0, \sigma^2)$$

① PCA & Ridge regression.

② Why  $\hat{\beta} = (X^T X)^{-1} X^T Y$

L1: k-NN & Linear model.

- k-Choice k=1 gives 0 training error.

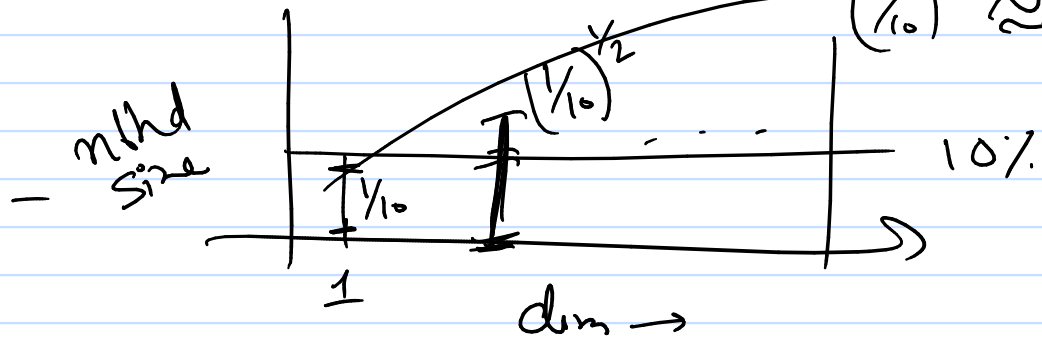


L2: Regression function  $\underline{\underline{E[Y|X=x]}} \approx \underline{\underline{k-NN}}$  : in theory

Curse of dimensionality.

$$\approx \beta^T x$$

$$\left(\frac{1}{10}\right)^{\frac{1}{100}} \approx .95$$



- Bayes classifier:  $P_{GX}$



$$RSS(\beta) = \frac{1}{N} \underbrace{(Y - X\beta)^T}_{1 \times N} \underbrace{(Y - X\beta)}_{N \times 1} = \frac{1}{N} \sum_{i=1}^N (y_i - x_i^T \beta)^2.$$

L3: Bias-Variance tradeoff.  $k \cdot nn \propto \frac{1}{k}$

Sampling properties of  $\hat{\beta}$

$$Y = \beta^T X + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2)$$

Variance.

known

independence

not random

Normal.

$$Z \propto \hat{\beta}_i - \beta_i$$

Ways to add bias: subset, Ridge & Lasso.

↑  
CV ✓

L4: Ridge & Lasso

geometric

likelihood view.

$$\frac{1}{N} \sum_{i=1}^N (y_i - x_i^T \beta)^2 + \lambda \sum_j \beta_j^2$$



lambda large.

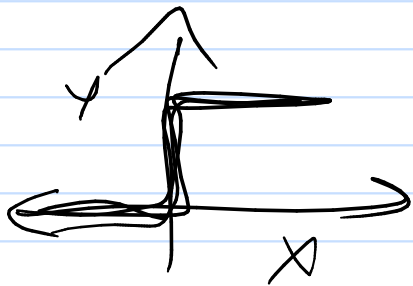
$$\hat{\beta}_{\text{ridge}} = \left( X^T X + \lambda I \right)^{-1} X^T Y$$

$$\rightarrow \text{RSS}(\beta) + \lambda \|\beta\|_2^2$$

$$(\underline{Y} - \underline{X}\beta)^T (\underline{Y} - \underline{X}\beta) + \lambda \beta^T \mathbf{I} \beta$$

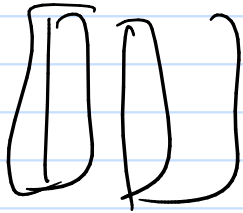
$$(\underline{Y}^T - \beta^T \underline{X}^T) (\underline{Y} - \underline{X}\beta) + \lambda \beta^T \mathbf{I} \beta$$

$$\underline{Y}^T \underline{Y} - \underline{Y}^T \underline{X} \beta - \beta^T \underline{X}^T \underline{Y} + \underbrace{\beta^T \underline{X}^T \underline{X} \beta}_{\text{RSS}} + \lambda \underbrace{\beta^T \mathbf{I} \beta}_{\text{penalty}}$$



Reg for classification.

$$\hat{\beta}_{p \times 3} = \frac{\underbrace{(X^T X)^{-1}}_{p \times p} \underbrace{X^T Y}_{p \times 3}}{A_{p \times M}}$$



LDA vs LReg.

Normality

$$y_i = \begin{bmatrix} 0 & \dots & 0 & 1 & 0 & \dots & 0 \end{bmatrix}$$

"7" ↑

1x10

$$\beta_l = l = 1, \dots, K-1$$

Cross entropy loss (i) =  $\sum_{l=1}^K y_{ij} \cdot \frac{e^{\beta_l^T x}}{1 + \sum_{k=1}^K e^{\beta_k^T x}}$

L6:

CV, Bootstrap: Model selection, assessment.