

## Post-Midsem L2

### Training error in classification

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i \neq \hat{y}_i)$$

$\hookrightarrow \mathbb{I}(y_i \neq \hat{y}_i) = 0$  (if classified correctly)

$= 1$  (if classified incorrectly)

### Test error rate

$$\text{Average} (\mathbb{I}(y_0 \neq \hat{y}_0))$$

→ Test error can be minimized by a simple Bayes classifier

→ The Bayes classifier assigns each observation to the most likely class, given its predictor value : simply assign a test observation with predictor vector  $x_0$  to the class ' $j$ ' for which  $P(Y=j | X=x_0)$  is maximum over all  $j$



▷ In 2 class scenario  $\begin{cases} \rightarrow 1 & P > 0.5 \\ \rightarrow 2 & P \leq 0.5 \end{cases}$

$$\text{Bayes Error rate} = 1 - \underbrace{E \left( \max_j P(y=j|X) \right)}_{\text{Correctness}}$$

KNN approximate

$$P(y=j | X=x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i = j)$$

Bayes Theorem

$$P(H|E) = \frac{P(E|H) * P(H)}{P(E)}$$

$\nearrow$  Prob(Hypotheses such that Evidence is true)  
 $\nwarrow$  Prob(Evidence such that Hypotheses is true)  
 $\nwarrow$  Prob(Hypotheses)  
 $\nwarrow$  Prob(Evidence)

$$P(\text{Positive}) = P(D \cap \text{Pos}) + P((\sim D) \cap \text{Pos})$$

$$= P(\text{Pos} | D) * P(D)$$

$$+ P(\text{Pos} | (\sim D)) * P(\sim D)$$



# Naive Bayes (an approximate Bayes classifier)

MAP ← Maximum A Posteriori Probability

$$\text{MAP}(H) = \max(P(H|E)) = \max\left(\frac{(P(E|H) * P(H))}{P(E)}\right)$$

$P(H|\text{Multiple Evidences})$

$$= \frac{P(E_1|H) * P(E_2|H) * \dots * P(E_n|H) * P(H)}{P(\text{Multiple Evidences})} \quad - F(P_1)$$

Assumption each of the predictors are independent of each other.

## Post Midsem-L3

Ex:

Predictors < outlook, temperature, humidity, wind >

↓  
< Play Badminton? > (Yes or No)

$$P(\text{No}) = \frac{5}{14}$$

$$P(\text{Yes}) = \frac{9}{14}$$

$$\textcircled{1} \quad \begin{matrix} \text{Sunny} & \text{Yes} \\ P(\text{Yes} | \text{Sunny}) & = \frac{2}{9} \end{matrix}$$

$$\begin{matrix} P(\text{No} | \text{Sunny}) & = \frac{3}{5} \\ \text{Sunny} & \text{No} \end{matrix}$$

~~②~~

$$(2) P(\text{Overcast} | Y) = 4/9$$

$$P(\text{Overcast} | N) = 0/5$$

(3)

$$P(\text{Rain} | \text{Yes}) = 3/9$$

$$P(\text{Rain} | \text{No}) = 2/5$$



$$(4) P(\text{Hot} | Y) = 2/9$$

$$P(\text{Hot} | N) = 2/5$$

(5)

$$P(\text{Mild} | Y) =$$

$$P(\text{Mild} | N) =$$



And so on

and we put in formula

-  $F(P_1)$

### Advantages

- Fast, Scalable
- Binary or Multi-class
- Multiple different types
- popular for spam email classification

### Disadvantages

- Assume that predictors are independent.



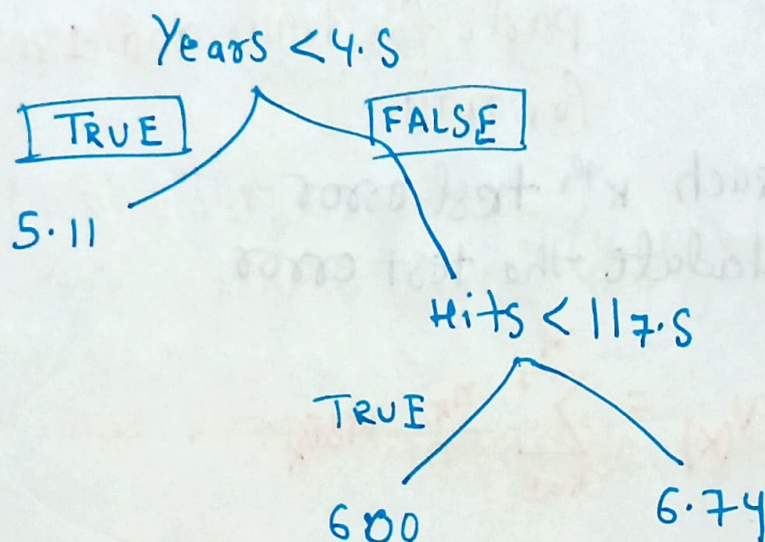
## Tree-based methods

→ Tree based algorithm decides based on multiple division of the attributes (as a tree)

Ex.

[Regression + Classification]

Decision Tree



→ Add bagging, random forests and boosting can ~~be~~ help increase the accuracy of the decision trees.

Post-Midsem L4

▸ Estimate test error on basis of availability of test data

▸ validation - set (Train/Test) approach has disadvantages.

→ overestimate

→ Not using complete train data



# K-Fold Cross-Validation

$1, 2, \dots, K$

For each part take 'k-1' parts for train and 1 part for train.

▶ using each  $k^{\text{th}}$  test error we calculate the test error

$$CV(k) = \sum_{k=1}^K \frac{n_k}{n} MSE_k$$

$$MSE_k = \sum_{i \in C_k} (y_i - \hat{y}_i)^2 / n_k$$

{LOOCV}

[Cross-Validation Error]

Leave one out cross-validation

Special case

$K=n$

total number of data points

$K=2$

Cross validation

So generally  $K=5$  to  $10$

$$CV(n) = \frac{1}{n} \sum_{i=1}^n \left( \frac{(y_i - \hat{y}_i)^2}{1 - h_i} \right)^2$$

leverage point

▶ Taking only one data point in each part



High leverage  $\rightarrow$  Have a high impact on the model

$$CV_k = \sum_{k=1}^K \frac{n_k}{n} E_k$$

$$E_k = \sum_{i \in C_k} I(y_i \neq \hat{y}_i) / n_k$$

CV Test Error for logistic regression

### Boot Strap

$\rightarrow$  Estimate the parameters

$2x$

$$\alpha = \frac{\text{Var}(X) - \text{Cov}(X, Y)}{\text{Var}(X) + \text{Var}(Y) - 2 \text{Cov}(X, Y)}$$

invest in asset return 'X'  
and  $1-\alpha$  in asset returns 'Y'

## Post Mid L5

### Pruning a tree

→ Larger the decision tree better the results but might possible to overfit the data.

~~A smaller~~ <sup>(earlier)</sup> some split seems like getting more RSS but might later on could result better RSS.

$T_0$   $\xrightarrow{\text{prune}}$  Subtree

### [Cost Complexity Pruning]

↳ weakest link pruning

$T$ : A subtree of  $T$

$$\sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

$\alpha$  Tuning parameter

$\hat{y}_{R_m}$  Predicted value of Region  $m$

Decrease with size of  $T$  increase



## Classification Tree

↳ Discrete values  
(Logistics)

$$E = 1 - \max_K (\hat{p}_{mk})$$

$k^{\text{th}}$  class

$m^{\text{th}}$  region

↑  
classification  
error rate

↳ Not much reliable

Gini index

$$\hookrightarrow G_i = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk})$$

□ Doubt: why 'E' is not reliable?