

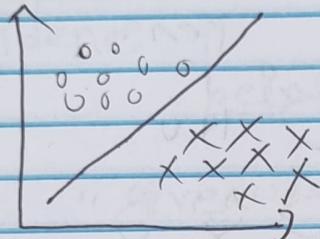
rec. of 24  
02-07-24

## # Discriminant Analysis:

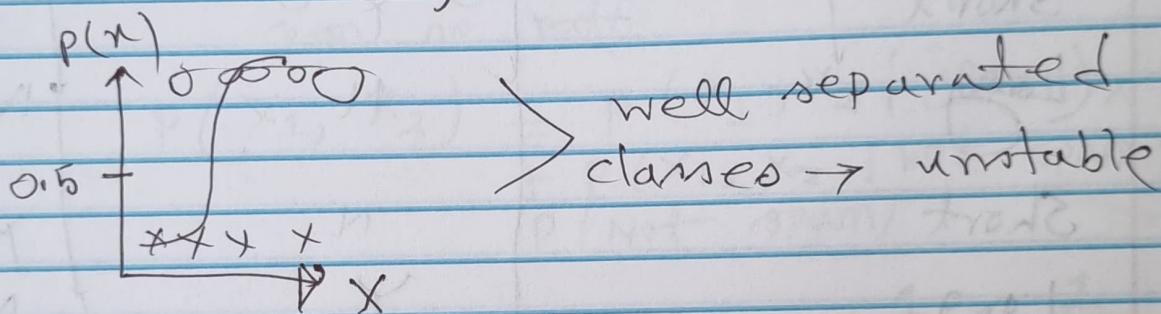
More than 2 dataset

↓  
logistic regression is not good

2 classes → <sup>not</sup> stable → well separated data



But not stable in estimating



## # Generative Models:

→ Generate the function

↓  
from the data.

→ parametric model

↓  
assumption from data

and then use  
Bayes theorem to  
flip probabilities

Bayes theorem → finding probability

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

e.g.:

Bayes theorem:

Email	$x_1$ length	$x_2$ contains "giveaway"	spam
Short	yes	Yes	Y
Long		No	NO
Short		Y	Y
Long		N	N
Short		N	Y
Long		Y	N
Short		N	Y
Long		Y	Y

$$P(\text{spam} | x_1, x_2) = \frac{P(x_1, x_2 | \text{spam}) \cdot P(\text{spam})}{P(x_1, x_2)}$$

i) calculate probability of  $P(\text{spam}) = \frac{4}{8}$

$$= \frac{1}{2}$$

$$\textcircled{3} \text{ Likelihood of } p(x_1, x_2 | \text{spam}) = \frac{2}{4}$$

$$\begin{cases} x_1 = \text{short} \\ x_2 = \text{Yes} \end{cases} = \frac{1}{2}$$

$$p(x_1, x_2 | \text{not spam})$$

\textcircled{3} Evidence : The probability of email being short (op) contain "giveaway" irrespective of being spam or not

$$p(x_1, x_2) = p(x_1, x_2 | \text{spam}) \times p(\text{spam}) + p(x_1, x_2 | \text{not spam}) \cdot p(\text{not spam})$$

$$= \frac{1}{2} \times \frac{1}{2} + \frac{1}{2} \times \frac{1}{2}$$

short yes yes

$$= \frac{1}{4} + \frac{1}{4} = 0.375$$

$$+ \frac{1}{4} = \frac{1}{2}$$

Chapter 4  
Bayes theorem

$$= \frac{1}{4} = 0.25$$

Naive Bayes

$$P(\text{spam} | x_1, x_2) = \frac{P(x_1, x_2 | \text{spam})}{P(\text{spam})} \times P(x_1, x_2)$$

$$= \frac{\frac{1}{2} \times \frac{1}{2}}{\frac{1}{4}} = 1$$

# LDA :  $\rightarrow$  parametric

$\rightarrow$  The data is of a uniform distribution ( $\mathcal{N}$ )

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

$\Rightarrow$  normal Distribution function

Assume  $f_x$   
then put in  
naive bayes

but in LDA  $\rightarrow$  putting this  
 $f_x$  to the naive bayes

Bayes theorem  
then putting normal distrib  
inside this  
LDA (contd)

$$\text{probability of } P_k(x) = \frac{\pi_k \times \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2} \left(\frac{x-\mu_k}{\sigma}\right)^2}}{\sum_{l=1}^K \pi_l \times \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2} \left(\frac{x-\mu_l}{\sigma}\right)^2}}$$

# QDA →

→ parametric

→ Assumes the data has a quadratic form

→ Assumes all classes / Labels have their own covariance matrix.

→ why → not logistics? → naturally built for binary mat  
generative model

↓ generate fn out of data

↓ data look like this → predicting on

↓ Then LDA → assumes ND-fn

variant of LDA then flipping this into the bayes theorem.  
QDA (assumes quadratic fn form)

lec: 08  
02-12-29

- Logistic Regression
- LDA
- KNN
- Naive Bayes classifier

# Bayes Theorem:

$$P(Y|X) = \frac{P(X|Y) \times P(Y)}{P(X)}$$

Naive Bayes classifier:

i) All predictors are independent

conditional independence:

The effect of the value of a predictor ( $X$ ) on a given class ( $Y$ ) is independent of other predictors.

$P(Y)$  → prior probability

$P(X|Y)$  → Likelihood

$P(X)$  → evidence

$$P(Y|X) = \frac{P(X_1|Y) \times P(X_2|Y) \times \dots \times P(X_p|Y)}{P(X)}$$

$x_1, x_2, \dots, x_p$

prior                      likelihood

#

free	money	class
1	0	spam
1	1	spam
0	1	not spam
0	0	" "
1	0	spam
0	1	not "
1	1	spam
1	0	spam

	spam	not spam
free	5	0
money	2	2
	7	2

word free  
appeared  
5 times  
in  
spam

① calculate prior

$$P(S) = \frac{5}{8} \quad P(NS) = \frac{3}{8}$$

② calculate likelihood

$$p(\text{Free}|s) = \frac{\text{Number of free occurrences}}{\text{Total spam emails}}$$

$$= \frac{5}{5} = 1$$

$$p(\text{Money}|s) = \frac{\text{Number of money "}}{\text{Total spam emails}}$$

$$= \frac{2}{5}$$

$$p(\text{Free}|ns) = \frac{\text{Number of free occ'c'}}{\text{Total not spam emails}}$$

$$= \frac{0}{3} = 0$$

$$p(\text{Money}|ns) = \frac{\text{Num of Money}}{\text{Total}}$$

$$= \frac{2}{3}$$

~~p(Y|X)~~

~~Test :~~

"free money click the  
link below"

$$P(S \text{ free, money}) = P(\text{free}|S) \times P(\text{money}|S) \times P(S)$$

$$= 1 \times \frac{3}{8} \times \frac{5}{8}$$
$$= \frac{15}{64} = 0.234 \text{ probability}$$

$$P(NS \text{ free, Money}) = P(\text{free}|NS) \times P(\text{Money}|NS) \times P(NS)$$

$$= 0.7 \times \frac{2}{3} = 0$$

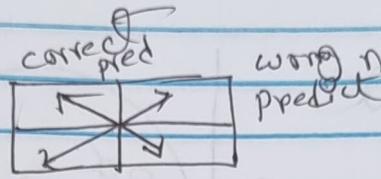
of free money  
as being the  
spam.

free credit card, free -

## # confusion Matrix:

evalution of classification Model:

True positive  $\rightarrow$  predictor + outcome =



How can we evaluate accuracy?  
correct predict out of all  
dataset

$$\frac{\text{correct pred}}{\text{Total}}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Precision :

$$\frac{TP}{\text{Total positive}} = \frac{1}{2} \cdot \frac{TP}{TP+FP}$$

Actual  
Positive      Negative

Recall :

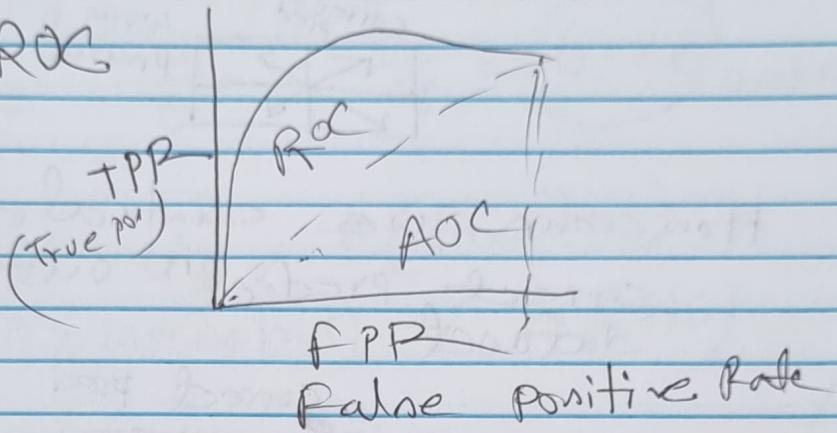
$\downarrow$

$$\frac{TP}{TP + FN}$$

TP	2	3
FN	4	2

#  $F_1$ -score  $\rightarrow$  Mean / Avg for  
 ↓ the precision and Recall  
 closer to the accuracy

# AUC-ROC



FPR  $\rightarrow$  should be 0

TPR  $\rightarrow$   $\uparrow\uparrow$

High TPR  $\rightarrow$  then  
 FPR doesn't matter

$\downarrow$   
 Area should  
 be bigger

