

PATTERN

Page No.

Date

Intra

Feature Extraction
 Bayes decision theory
 Max likelihood estimation
 Prior Density estimation
 Nearest Neighbour
 Discrimination
 LDA
 Unsupervised

Classification Clustering

Feature
 Neural Network
 Bayes classification
 Decision tree
 LDA
 SVD
 SVA

→ Computers classify images based on features.
 Features should not be less, nor more, it should be balanced.

→ Satellite image classification

Satellite images are $36.75 \times 36.75 \text{ m}^2$. Suppose resolution is 512×512 . So image is very large. Classification of satellite images is very imp. (i) For eg, in case of Tsunami, how much area is affected, is given by satellite images. (ii) For eg, noting changes at border area.

→ eg of image classification

(i) medical purpose

x-ray, blood classification, tumor etc.

(ii) object and face recognition, voice recognition

PATTERN DOMAIN APP PATTERN IP PATTERN CLASSES

(i) Doc Image analysis OCR Doc Image char, words
 (ii) Doc classification Internet Text Doc Semantic search categories

Junk Mail Email Spam / Inbox /
 Filtering Promotion

(iii) MM Database retrieval Internet Video clips Watched games
 search

(iv) Speech recognition Telephone speech spoken words
 assistance waveform

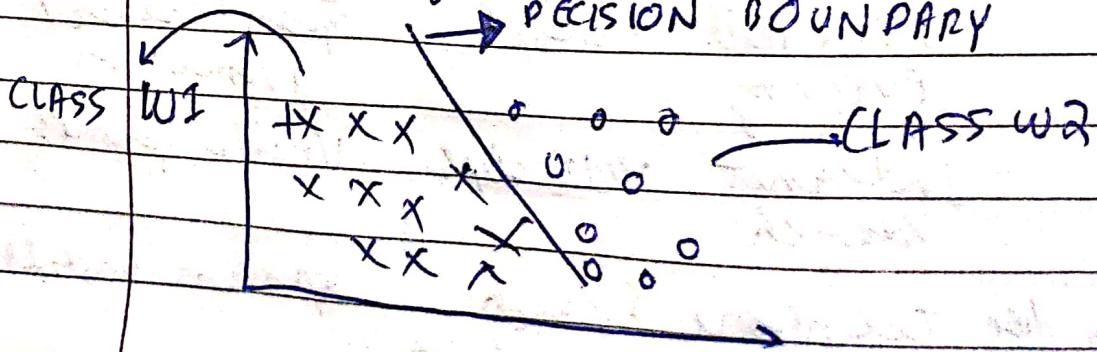
- (iv) NLP Info Extraction Sentences Parts of speech
- (vii) Biometric Personal Face, iris, Authorize user
recognition Identification fingerprint for Access control
- (vi) Medical Comp. Aided Microscopic Cancelling / diagnosis. Image Healthy cell.
- (viii) Remote Forecasting Multispectral Land use sensing crop yield Image categories
- (ix) Bioinformatics Sequence Analysis DNA Sequence Known types of genes.
- (x) Data mining Searching for Points in meaningful seq multi-dimensional space

BAYES DECISION PROBLEM

Many of the features, taken together, so, they can describe an object with some degree of accuracy. These are called FEATURE VECTOR. They have to be concatenated in a particular order, and whichever order we concatenate them throughout our problem, that is MODELLING OF PATTERN as well as RECOGNISATION OF PATTERN

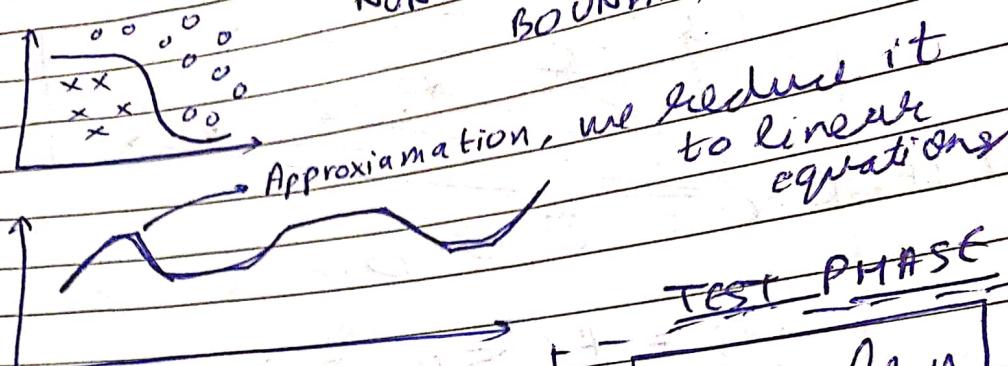
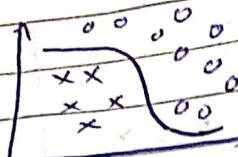
On putting those diff features in a particular order, they are called FEATURE VECTOR

→ PRECISION BOUNDARY



This is LINEAR DECISION BOUNDARY

NON-LINEAR DECISION BOUNDARY



TEST PHASE

Physical sensor

Data acquisition
sampling

Pre-processing

Feature extraction

FEATURES

Model learning /
estimation

MODEL

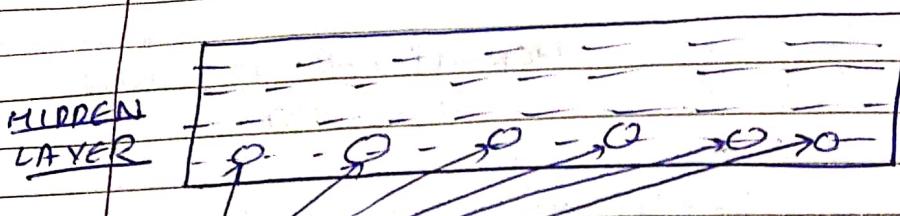
Classification
based on Model

Post-processing

decision
finding

DECISION

OIP LAYER ○○○○○○○○○○



LIP LAYER ○○○○○○○○○○

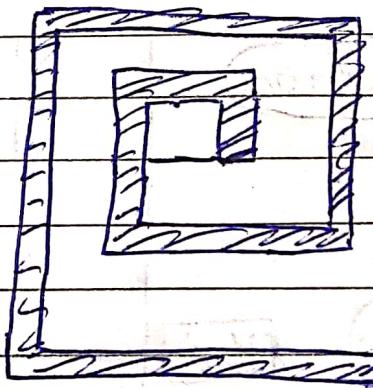
$$O_j = f(\sum w_{ij} I_j)$$

→ OVERLAPPING BOUNDARY.

eg → Person sitting of chair

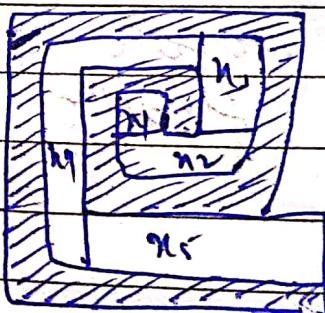
Normal = A

Shaded = B



Can not be classified
by linear/non-linear boundary

MULTI-LAYER PERCEPTRON PROBLEM



$\min A = \frac{\min v}{D}$

→ CONDITIONAL PROBABILITY
 $P(A|B) = \frac{P(A \cap B)}{P(B)} = \text{probability of } A \text{ when } B \text{ has already occurred}$

1. % of adult who are men & alcoholic
 $= 2.25$. Person given is male. Find
 no. of alcoholic.

$$P(A|M) = \frac{P(A \cap M)}{P(M)} = \frac{2.25}{100 \times \frac{1}{2}}$$

Alcoholic ← Male

$$\text{Assuming no. of male} = \frac{4.5}{100} = 0.045$$

= no. of female

EQUALY

→ BAYES THEOREM

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A|B) P(B)}{P(A)}$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

2. what is prob. of 2 girls in a couple
 when which given at least one is a girl.

ANS. $P(A|B)$

$\xrightarrow{\text{at least}} \xleftarrow{\text{first is girl}}$

$$= \frac{P(A \cap B)}{P(B)}$$

record a girl

$$P(A \cap B) = 1/4$$

P.C

$P(2G \text{ at least } 1G)$

$$P(2G) = 1/4$$

{GG, GJ}

$$P(1G) = 3/4$$

{GRB, BG, GRG}

$$P(2G|1G) = \frac{P(1G|2G) P(2G)}{P(1G)}$$

$$= \frac{1 \times 1/4}{3/4} = \frac{1}{3}$$

Always I
as if 2G
given, then
I will be
girl

3. If I randomly draw a green ball
Find prob green from first bucket

| 000 |
xx X |

B₁

| 00X |
X 00 |

B₂

X = golden
0 = red

$$A. P(G|B_1) = \frac{1}{2} \quad P(G|B_2) = \frac{1}{3}$$

$$P(B_1|G) = \frac{P(B_1 \cap G)}{P(G)} = \frac{P(G|B_1) P(B_1)}{P(G)}$$

$$= \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = \frac{3}{5}$$

$$\frac{\frac{1}{4} + \frac{1}{3}}{\cancel{4} \cancel{3}} = \frac{3+2}{12} = \frac{5}{12}$$

$$P(G) = P(G|B_1) P(B_1) + P(G|B_2) P(B_2)$$

$$= P(G \cap B_1) + P(G \cap B_2)$$

BAYES DECISION THEORY

Q-1) CLASSIFICATION BASED ON PRIOR PROBABILITY

w_1 , w_2
Accept product Reject Product
 $P(w_1)$ $P(w_2)$

$$* P(w_1) > P(w_2) \Rightarrow w_1$$

$$P(w_1) < P(w_2) \Rightarrow w_2$$

This is called a priori probability

If $P(w_1) = 0.9$ $P(w_2) = 0.1$, then product will always be accepted, and this is not ideal. So, we will add an additional feature x .

(2) CLASSIFICATION BASED ON CLASS PDF

$$* P(x|w_1) \quad P(x|w_2)$$

class/condition PDF

These can be calculated using training set.

Now we know the features, we need to classify it.

$$P(w_i \cap x) = P(w_i|x) P(x)$$

$$= P(x|w_i) P(w_i)$$

$$\Rightarrow P(w_i|x), P(x) = P(x|w_i) P(w_i)$$

$$P(w_i|x) = P(x|w_i) P(w_i)$$

$$P(x)$$

$$\Rightarrow P(x) = \sum_{i=1}^2 P(x|w_i) P(w_i)$$

POSTERIOR PROBABILITY

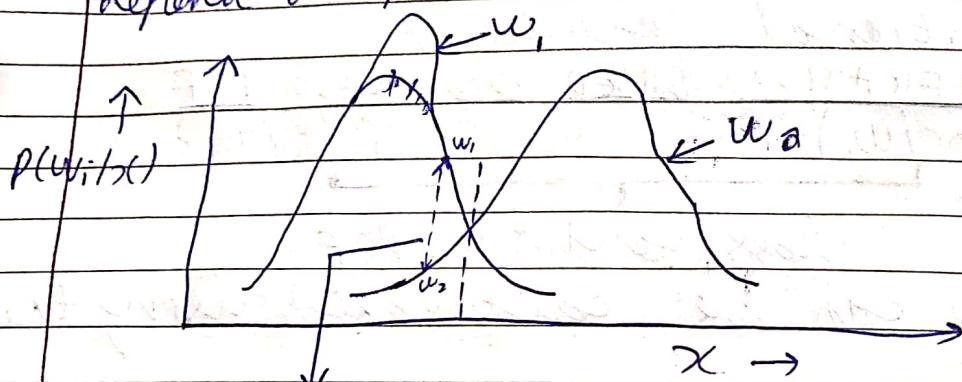
Page No. _____
Date _____

$$\Rightarrow P(w_1/x) > P(w_2/x) \Rightarrow w_1 \\ P(w_1/x) < P(w_2/x) \Rightarrow w_2$$

$$\frac{P(x/w_1) \cdot P(w_1)}{P(x)} > \frac{P(x/w_2) \cdot P(w_2)}{P(x)}$$

training set \downarrow prior probability \downarrow training set \downarrow
 $P(x)$ prior probability

* if $P(x/w_1) = P(x/w_2)$, then we only depend on priors.



if we select w_1 ,

$$P(\text{Error}/x) = P(w_2/x)$$

To minimise the error, we select w_1 if $P(w_1/x)$

and $P(w_2/x)$ if $P(w_2/x) > P(w_1/x)$ and vice versa.

e.g. x is colour Intensity. If intensity is low = accept, otherwise reject

AVERAGE ERROR

$$P(\text{Error}) = \int_{-\infty}^{\infty} P(\text{Error}/x) dx$$

$$= \int_{-\infty}^{\infty} P(\text{Error}/x) P(x) dx$$

If we minimize $P(\text{Error}/x)$ + γx , then avg error is minimised

DECISION RULE

$$P(\text{error}/x) = \min \{ P(w_1/x), P(w_2/x) \}$$

and if $P(w_1/x) > P(w_2/x) \rightarrow \text{class } w_1$
else $\rightarrow \text{class } w_2$

$$P(x|w_2)$$

$$\rightarrow P(x|w_i)$$

↓
class conditional PDF

$$P(w_i|x) = \frac{P(x|w_i) P(w_i)}{P(x)}$$

A POSTERIORI PROB

$$P(w_2|x) = \frac{P(x|w_2) P(w_2)}{P(x)}$$

$$P(x) = \sum_{i=1}^2 P(x|w_i) \cdot P(w_i) \quad (2.2)$$

BAYESIAN THEORY FOR
CONTINUOUS FEATURES

PRIOR PROBABILITIES

~~Ques Quality of course~~ good fair bad

(i) Prob Prior 0.2 0.4 0.4

~~# class people~~ good fair bad

$$P(x|w_i)$$

Interesting

$$0.8$$

$$0.5$$

$$0.1$$

Boring

$$0.2$$

$$0.5$$

$$\approx 0.9$$

~~loss function~~

~~P_l(a_i|w_i)~~ good fair bad

Taking

$$0$$

$$5$$

$$10$$

Not taking

$$20$$

$$5$$

$$0$$

GENERALIZATION

* Use more than 2 states of natural classes

* Use more than 1 feature

↳ feature vector

* Allows other action other than merely deciding states of natural (NOT ONLY ACCEPT OR REJECT, INTERMEDIATE)

(3) CLASSIFICATION BASED ON RISKS

* Introduce a loss function more general than prob. of errors

$c_i \rightarrow$ states of nature
 w_1, w_2, \dots, w_c

$a \rightarrow$ Actions

$\alpha_1, \alpha_2, \dots, \alpha_c$

Loss function \rightarrow

$\lambda(\alpha_i | w_j) \rightarrow$ Loss incurred

action \downarrow for taking action α_i ; when
 actual state true state of nature is w_j

$x \rightarrow d$ -dimensional feature vector.

$$p(w_j | x) = p(x | w_j) p(w_j)$$

ACTION α_i :

$$c \quad p(x) = \sum_{j=1}^c p(x | w_j) p(w_j)$$

$$R(\alpha_i | x) = \sum_{j=1}^c \lambda(\alpha_i | w_j) p(w_j | x)$$

risk function

conditional risk / expected loss

TWO CATEGORY CLASSIFICATION

\Rightarrow we classify it into MIN RISK

e.g.

w_1

α_1

$$\lambda(\alpha_1 | w_1) = \lambda_{11}$$

w_2

α_2

$$\begin{aligned} R_1 &> R_2 \\ R_1 - R_2 &> 0 \\ p(w_1 | x) &> p(w_2 | x) \end{aligned}$$

$$R(\alpha_1 | x) = \lambda(\alpha_1 | w_1) p(w_1 | x) + \lambda(\alpha_2 | w_2) p(w_2 | x)$$

$$= \lambda_{11} p(w_1 | x) + \lambda_{12} p(w_2 | x)$$

$$R(\alpha_2 | x) = \lambda(\alpha_2 | w_2) p(w_2 | x) + \lambda(\alpha_1 | w_1) p(w_1 | x)$$

$$\text{IF } R(\alpha_2 | x) > R(\alpha_1 | x) \quad \text{OR}$$

$$R(\alpha_2 | x) - R(\alpha_1 | x) > 0$$

$$\text{IF } (\lambda_{21} - \lambda_{11}) p(w_1 | x) > (\lambda_{22} - \lambda_{12}) p(w_2 | x)$$

ideally $\lambda_{ii} = 0$, as it is of class I,
 and we classified it in class I, so
 no errors

$$\therefore \lambda_{21} - \lambda_{11} \geq 0 \quad \text{and} \quad \lambda_{12} - \lambda_{22} \geq 0$$

$$\therefore p(w_1 | x) > p(w_2 | x) \rightarrow w_1$$

SO MIN-RISK CLASS IS CLASSIFIED

(2.3) \rightarrow MIN ERROR RATE CLASSIFICATION

$\alpha_i =$ true state of nature is w_i

$$\lambda(\alpha_i | w_j) = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases} \quad (\text{As no. errors})$$

$$R(\alpha_i | x) = \sum_{j=1}^c \lambda(\alpha_i | w_j) p(w_j | x)$$

ideally, we must try to minimize $R(\alpha_i | x)$. In order to do so, we will have to maximize $p(w_i | x)$. So, decision rule: decide w_i if $p(w_i | x) > p(w_j | x)$

$$= 1 - p(w_i | x)$$

Ex. (2.4) DISCRIMINANT FUNC (2.4)

$$x \rightarrow g_1(x) \quad ? \quad \rightarrow x \in ?$$

$$g_0(x)$$

$g(x) \rightarrow$ DISCRIMINANT FUNC

CLASSIFIER

$w_1, w_2, \dots, w_c \rightarrow$ No. of classes

$$g_i(x) : i=1, 2, \dots, c$$

$$g_i(x) > g_j(x) \quad \text{if } j \neq i$$

$$\Rightarrow x \in w_i$$

Page No. _____
Date _____

DISCRIMINANT
FUNCTION
CLASSIFIER

* \star

Discriminant func should be higher. If we classify acc to risk, then risk must be minimum - Thus we'll:

Min Risk classifier:

$$g_i(x) = -R(x; i/x)$$

(as maximum discriminant will correspond to minimum risk, we use a negation sign)

$$R(x; i/x) = \text{Error Rate} \rightarrow \sum P(w_j/x) = 1 - P(w_i/x)$$

ON BASIS OF CLASS'S CONDITIONAL PROB

$$g_i(x) = P(w_i/x)$$

(For min-risk rate classification maximum posteriori rule = maximum discriminant func)

Applying some monotonically increasing function on $P(w_i/x) = P(x/w_i) P(w_i)$ we can remove $P(x)$ as it is cancelled while comparing 2 classes $g_i(x)$ and $g_j(x)$

$$\Rightarrow P(w_i) = \sum_{j=1}^c P(x/w_j) P(w_j)$$

~~Taking log~~

$$② g_i(x) = P(x/w_i) P(w_i)$$

To taking log and as production is costly

All ① ② ③ produce same classification, but some can be

③ $\ln(g_i(x)) = \ln P(x/w_i) + \ln P(w_i)$ simpler to understand TWO-CATEGORY CLASSIFICATION and compute

* If $g_1(x) > g_2(x) \Rightarrow w_1$

$$g_1(x) = g_2(x)$$

$$g_1(x) < g_2(x) \Rightarrow w_2$$

Page No. _____
Date _____

called DICHOTOMIZER

$$g(x) = g_1(x) - g_2(x)$$

$$g(x) = P(w_1/x) - P(w_2/x)$$

$$= \ln P(x/w_1) + \ln P(x/w_2) \rightarrow \text{FROM (2)}$$

$$= P(x/w_1) \cdot P(w_1) - P(x/w_2) \cdot P(w_2)$$

DECISION RULE w_1 if $g(x) > 0$

Taking log of both sides

$$\ln g(x) = \ln (P(x/w_1) \cdot P(w_1) - P(x/w_2) \cdot P(w_2))$$

$$= \ln \left(\frac{P(x/w_1) \cdot P(w_1)}{P(x/w_2) \cdot P(w_2)} \right)$$

$$= \ln \left(\frac{P(x/w_1)}{P(x/w_2)} \cdot \frac{P(w_1)}{P(w_2)} \right) \text{ USING (3)}$$

$$\ln g(x) = \ln \left(\frac{P(x/w_1)}{P(x/w_2)} \right) + \ln \left(\frac{P(w_1)}{P(w_2)} \right) \uparrow$$

$$\ln g(x) = \ln P(x/w_1) + \ln P(w_1) - \ln P(x/w_2) - \ln P(w_2)$$

NANIE BAYES

$$= \ln \left(\frac{P(x/w_1)}{P(x/w_2)} \right) + \ln \left(\frac{P(w_1)}{P(w_2)} \right) \text{ for}$$

Q.	Sno.	color	TYPE	SPORTS CAR	ORIGIN	Stolen
1	R	Sports	Imported	Domestic	Yes	✓
2	R	Sports	Imported	"	No	✗
3	R	Sports	Imported	"	Yes	✗
4	Y	Sports	Imported	"	No	✗
5	Y	Sports	Imported	Imported	Yes	✓
6	Y	SUV	Imported	"	No	✗
7	Y	SUV	Imported	"	Yes	✓
8	Y	SUV	Domestic	Domestic	No	✗
9	R	SUV	Imported	Imported	No	✗
10	R	Sports	Imported	Imported	Yes	✓

PATA ↑

Given car Red & Domestic. tell whether it is stolen or not

ANSWER

COLOR

$$\text{ANS} \rightarrow P(R/\text{No}) = 3/5$$

$$P(Y/\text{Yes}) = 2/5$$

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

$$P(Y/\text{Yes}) = 2/5$$

TYPE

$$P(\text{Sports}/\text{Yes}) = 4/5$$

$$P(\text{SUV}/\text{Yes}) = 1/5$$

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

$$P(\text{SUV}/\text{No}) = 3/5$$

Origin

$$P(\text{Domestic}/\text{Yes}) = 2/5$$

$$P(\text{Imported}/\text{Yes}) = 3/5$$

$$P(\text{Domestic}/\text{No}) = 3/5$$

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

$$\frac{P(X/\text{Yes})}{P(X/\text{Yes})} \cdot \frac{P(Y/\text{Yes})}{P(Y/\text{Yes})} = \frac{P(\text{Red}/\text{Yes})}{P(\text{Red}/\text{Yes})} \cdot \frac{P(\text{SUV}/\text{Yes})}{P(\text{SUV}/\text{Yes})} \cdot \frac{P(\text{Domestic}/\text{Yes})}{P(\text{Domestic}/\text{Yes})}$$

$$P(\text{Yes}/\text{No}) \cdot P(\text{Red}/\text{No}) = P(\text{Red}/\text{No}) \cdot P(\text{SUV}/\text{No})$$

$$P(\text{Red}/\text{No}) \cdot P(\text{Domestic}/\text{No}) = P(\text{Domestic}/\text{No}) \cdot P(\text{No})$$

$$P(\text{Yes}/\text{Yes}) = P(\text{Red}/\text{Yes}) \cdot P(\text{SUV}/\text{Yes}) - P(\text{Red}/\text{Yes}) \cdot P(\text{SUV}/\text{Yes}) \cdot P(\text{Domestic}/\text{Yes})$$

$$P(\text{Yes}/\text{X}) = P(\text{X}/\text{Yes}) \cdot P(\text{Yes}/\text{X})$$

$$P(\text{Yes}/\text{X}) = \frac{3}{5} \times \frac{1}{5} \times \frac{2}{5} = \frac{6}{125}$$

$$P(\text{X}/\text{Yes}) = P(\text{X}_1/\text{Yes}) \cdot P(\text{X}_2/\text{Yes})$$

$$P(\text{X}/\text{Yes}) = P(\text{Red}/\text{Yes}) \cdot P(\text{SUV}/\text{Yes})$$

$$P(\text{X}/\text{Yes}) = P(\text{Domestic}/\text{Yes})$$

$$P(\text{No}/\text{X}) = \frac{2}{5} \times \frac{3}{5} \times \frac{3}{5} = \frac{24}{125}$$

As $P(\text{Yes}/\text{X}) < P(\text{No}/\text{X})$, it is not stolen.
SELECT optimal course

	good	fair	bad
Prob(prior)	0.2	0.4	0.4
class condition	good	fair	bad
Pr(X/w)			
Interesting	0.8	0.5	0.1
Boring	0.2	0.5	0.9
less fun	good	fair	bad
A(a; w)			
Taking the course	0	5	10
Not taking	20	5	0

$$\text{Ans. } P_r(\text{Interesting}) = \frac{P(I/\text{good}) P(\text{good}) + P(I/\text{fair}) P(\text{fair})}{P(\text{good}) + P(I/\text{bad}) P(\text{bad})}$$

$$= 0.8 \times 0.2 + 0.5 \times 0.4 + 0.1 \times 0.4$$

$$= 0.16 + 0.20 + 0.04 = 0.4$$

$$P(\text{Boeing}) = 1 - 0.4 = 0.6$$

$$P(\text{good}/\text{Interesting}) = \frac{P(I/g)}{P(I)}$$

$$= 0.8 \times 0.2 = 0.4$$

$$P(\text{fair}/\text{Interesting}) = \frac{P(I/f)}{P(I)}$$

$$= 0.5 \times 0.4 = 0.2$$

$$P(\text{boring}/I) = 1 - (0.5 + 0.4) = 1 - 0.9 = 0.1$$

$$R(X_i/X) = 1 - P(W_i/X)$$

$$P(x_i | I) = 1 - P(\text{good} | I)$$

$$= 1 - 0.4 = 0.6$$

$$P(x_2 | I) = 1 - P(f | I)$$

$$= 1 - 0.5 = 0.5$$

w = class = Taking course, Not Taking
x = Interesting / Boring

$$R(\text{Taking} | I) = P(g | I)$$

x = feature = Interesting / Boring

w = class = good / fail / boring

$$R(\text{taking} | I) = P(g | I) \cdot P(\text{Taking} | \text{good})$$

$$+ P(f | I) \cdot P(\text{Taking} | \text{fail}) +$$

$$P(\emptyset | I) \cdot P(\text{Taking} | \text{Boring})$$

$$= 0.4 \times 0 + 0.5 \times 5 + 0.1 \times 10 = 3.5$$

$$R(\text{not taking} | I) = 0.4 \times 20 + 0.5 \times 5 + 0.1 \times 0 = 10.5$$

Similarly, R(taking $\frac{1}{2}$) + 2.5 =
and R(taking $\frac{1}{2}$)

Minimum risk = optimal solution

DISTANCE :

MINKOWSKI DISTANCE :

$$d(x_i, x_j) = \sqrt[h]{(x_{i1} - x_{j1})^h + (x_{i2} - x_{j2})^h + \dots + (x_{ir} - x_{jr})^h}$$

EUCLIDEAN DISTANCE :

$$d(x_i, x_j) = \sqrt{h} (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ir} - x_{jr})^2$$

MANHATTAN DISTANCE:

$$d(x_i, x_j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ir} - x_{jr}|$$

SQ-EUCLIDEAN

$$= (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ir} - x_{jr})^2$$

WEIGHTED EUCLIDEAN

Given weightage to each and every feature.

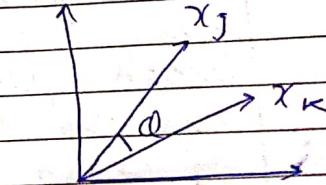
$$d(x_i, x_j) = \sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2 + \dots + w_r(x_{ir} - x_{jr})^2}$$

COSINE SIM

$$\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \|\vec{y}\|}$$

$$= \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

$$\sin(\vec{x}, \vec{y}) = \frac{|\vec{x} \times \vec{y}|}{\|\vec{x}\| \|\vec{y}\|}$$



SUPERVISED LEARNING PROBLEM

Can be classified in (i) Regression
(ii) Classification

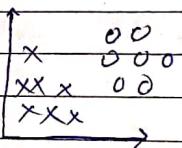
→ Regression :

e.g. → age of person.

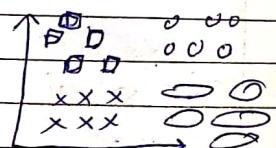
A regression problem is when output is real or continuous value such as weight of person, salary of person

→ Classification:
problem is when output variable is a category such as person is male or female or diagnosis of tumour is cancer or not. A classical Classification model attempts to draw conclusion from observed values for given one or more inputs, a classification model tries to predict value of one or more outcomes.

In classification data is categorised under different labels due to same parameters. Given an input and then labels are predicted from data. It can be demonstrated in form of if-then rules. The classification process deals with problems where data can be divided to binary or discrete values



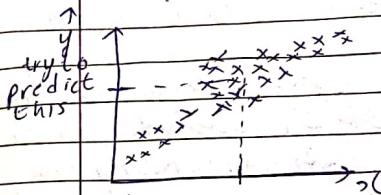
BINARY
CLASSIFICATION



MULTI-CLASS
CLASSIFICATION

A regression is a process of finding a function for distinguishing data into continuous real values instead of using classes or discrete values. It can also identify the distribution movement depending on historical data because

a regression predictive model predict a quantity. Therefore, skill of model must be exploited as a endorsed prediction.



CLASSIFICATION

- (i) Mapping func is used for mapping of values to pre-defined classes
- (ii) is used for prediction of discrete values
- (iii) nature of predicted data is unordered
- (iv) we measure accuracy of classification model to calculate performance eg - decision tree, KNN, logistic regression

REGRESSION

- (i) mapping func is used for mapping of values to continuous output.
- (ii) used to predict continuous values
- (iii) nature of predicted data is ordered
- (iv) we calculate performance by measuring mean square error
- (v) eg - random forest, linear regression

NEAREST-NEIGHBOUR CLASSIFIER

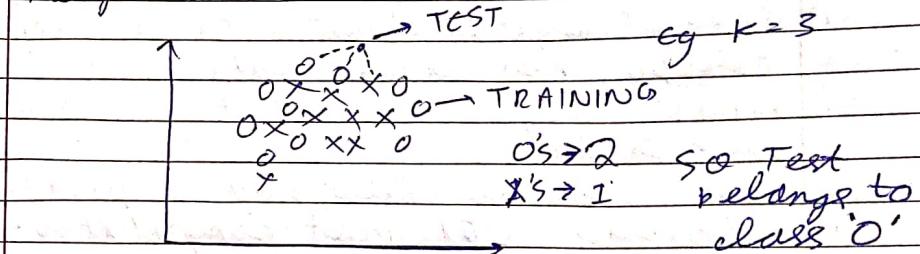
- * It is simplest classification algo in supervised learning
- * The method of classifying pattern based on class label of closest training pattern in feature space
- * The common algo used for nearest neighbours is KNN (K-Nearest Neighbour) or MKNN (Modified KNN)

- * The accuracy using nearest neighbour classifier is good. It is guaranteed to yield an error rate no more than twice of Bayes error rate. (Optimal error rate)
- * There is no training time required for this classifier. In other words, there is no learning time for training the classifier.
- * Every time a test pattern is to be classified it has to be compared with all training pattern to find closest pattern. This classification computation time can be ~~used~~^{high} if training data is large or dimensionality of data is high.

K-NEAREST NEIGHBOUR(KNN)

- * An object is classified for majority of votes of class of its neighbours.
- * The object is assigned to class most common amongst its K nearest neighbour.
- * This algo may give more correct classification than basic nearest neighbour where $K=1$.
- * The value of K has to be specified by user and it should be best choice depending on data.
- * Larger value of K reduce effect of noise on classification.

- * The K value can be chosen by using a validation set and choosing the K-value giving best accuracy on validation set.
- * Disadvantage of KNN is computation time especially when training data is large.
- * Nearest neighbour majority use Euclidean distance to find nearest neighbour.



ALGORITHM:

1. locate the test data
2. choose the value of K
3. For each point in test data
 - find the euclidean dist to all training data points
 - store the Euclidean dist in a list & sort it

Almond The conditional probability of given each class value is the same Bayes Model from start as the probability of each class given diff input values calculated.

class prob.
the conditional class prob.

temp	Humid	windy	play
Hot	High	False	No
Hot	High	True	No
Hot	Normal	F	Yes
Mild	High	F	Yes
Mild	Normal	F	Yes
Cool	Normal	T	No
Cool	Normal	F	Yes
Mild	High	F	Nice
Cool	Normal	F	Yes
Mild	Normal	F	Yes
Mild	Normal	T	Yes
Mild	High	T	Yes
Hot	Normal	F	Yes
Mild	High	T	No

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

$$P(\text{Yes}) = \frac{2}{14} = \frac{2}{9}$$

$$P(\text{overcast/Yes}) = \frac{4}{9}$$

$$P(\text{Sunny/Yes}) = \frac{3}{9} \quad P(\text{Sunny/No}) = \frac{2}{5}$$

$$P(\text{Rainy/No}) = \frac{3}{5} \quad P(\text{overcast/No}) = 0$$

$$P(\text{Hot/Yes}) = \frac{2}{9} \quad P(\text{Hot/No}) = \frac{2}{5}$$

$$P(\text{Mild/Yes}) = \frac{4}{9} \quad P(\text{Mild/No}) = \frac{2}{5}$$

$$P(\text{Cool/Yes}) = \frac{3}{9} \quad P(\text{Cool/No}) = \frac{1}{5}$$

$$P(\text{Migh/Yes}) = \frac{3}{9} \quad P(\text{High/No}) = \frac{4}{5}$$

$$P(\text{Normal/Yes}) = \frac{6}{9} \quad P(\text{Normal/No}) = \frac{1}{5}$$

$$P(\text{windy False/Yes}) = \frac{6}{9} \quad P(\text{False/No}) = \frac{2}{5}$$

$$P(\text{True/Yes}) = \frac{3}{9} \quad P(\text{True/No}) = \frac{3}{5}$$

Ques. Given $x = (\text{Sunny}, \text{cool}, \text{Humid}, \text{High}, \text{Windy}, \text{True})$

find play

$$\begin{aligned} P(\text{Yes}/x) &= P(\text{Sunny} \cap \text{cool} \cap \text{Humid} \cap \text{High} \cap \text{Windy} \cap \text{True}) / P(x) \\ &= \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14} = \frac{3^6}{9^4 \times 14 \times P(x)} \end{aligned}$$

P(x)

$$P(\text{No}/x) = \frac{2 \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} \times \frac{5}{14}}{P(x)}$$

$$= \frac{24}{14 \times 5^3 \times P(x)}$$

- Choose first K-points
- Assign a class to test points based on major majority of classes present in chosen point

PROS AND CONS OF KNN :

PROS :

- (i) No assumption about data is useful.
For eg, non-linear data
- (ii) Simple algo to understand and interpret and explain.
- (iii) High accuracy (Relatively) : pretty high accuracy but not competitive in comparison to many of supervised learning model
- (iv) mainly used in classification

CONS:

- (i) Computationally expensive because every new test point, also need to calculate Euclidean distance to all other points.
- (ii) High storage because algo stores all training data and Euclidean distance of each point with every other point.
- (iii) Prediction stage might be slow
- (iv) KNN is also sensitive to irrelevant features and scale of data

NAIVE BAYES CLASSIFICATION ALGO FOR BINARY AND MULTI-CLASS CLASSIFICATION PROBLEMS

- The technique is easiest to understand when describing using binary or categorical input value
- It is called naive Bayes or idiot Bayes because calculation of probabilities for each hypothesis are simplified to make their calculation tractable rather than rather than attempting to calculate value of each attribute value ($P(d_1, d_2, d_3 | h)$). They are assumed to be conditionally independent given target value and calculated as
$$= P(d_1 | h) \times P(d_2 | h) \dots$$
- This is very strong assumption that is most unlikely in real data that means attributes do not interact
- Nevertheless, the approach performs surprisingly well on data which is not independent
- $\# \text{ USE }$

NAIVE BAYES REPRESENTATION MODEL

- It includes
 - (i) class probabilities: prob of each class in training data set

(ii) conditional prob: The conditional prob for each value given each class value

→ Learning a Naive Bayes Model from training data is fast because only the probability of each class and prob of each class given diff input values needs to be calculated.

→ Calculate class prob.

→ Calculate the conditional class prob.

Ques	outlook	temp	Humid	windy	play
1	Rainy	Hot	High	False	No
2	Rainy	Hot	High	True	No
3	overcast	Hot	High	F	Yes
4	Sunny	Mild	High	F	Yes
5	Sunny	Cool	Normal	F	Yes
6	Sunny	Cool	Normal	T	No
7	overcast	Cool	Normal	T	Yes
8	Rainy	Cool	High	F	No
9	Rainy	Cool	Normal	F	Yes
10	Sunny	Mild	Normal	F	Yes
11	Rainy	Mild	Normal	T	Yes
12	overcast	Mild	High	T	Yes
13	overcast	Hot	Normal	F	Yes
14	Sunny	Mild	High	T	No

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

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

$$P(Rainy/Yes) = \frac{2}{14} = \frac{2}{9}$$

$$P(overcast/Yes) = \frac{4}{14} = \frac{2}{7}$$

$$P(Sunny/Yes) = \frac{3}{9} \quad P(Sunny/No) = \frac{2}{5}$$

$$P(Rainy/No) = \frac{3}{5} \quad P(overcast/No) = 0$$

$$P(Hot/Yes) = \frac{2}{9} \quad P(Hot/No) = \frac{2}{5}$$

$$P(Mild/Yes) = \frac{4}{9} \quad P(Mild/No) = \frac{2}{5}$$

$$P(Cool/Yes) = \frac{3}{9} \quad P(Cool/No) = \frac{1}{5}$$

$$P(High/Yes) = \frac{3}{9} \quad P(High/No) = \frac{4}{5}$$

$$P(Normal/Yes) = \frac{6}{9} \quad P(Normal/No) = \frac{1}{5}$$

$$P(Windy/False/Yes) = \frac{6}{9} \quad P(False/No) = \frac{2}{5}$$

$$P(True/Yes) = \frac{3}{9} \quad P(True/No) = \frac{3}{5}$$

Ques. Given $x = (\text{Sunny}, \text{cool}, \text{Humid}, \text{High}, \text{Windy})$

find play

$$\begin{aligned} Ans. P(Yes/x) &= \frac{P(\text{Sunny} \cap \text{cool} \cap \text{Humid} \cap \text{High} \cap \text{Windy} \cap \text{True})}{P(x)} \\ &= \frac{3 \times 3 \times 3 \times 3 \times 9}{9 \times 9 \times 9 \times 9 \times 14} = \frac{3^6}{9^4 \times 14 \times P(x)} \end{aligned}$$

$$P(x)$$

$$P(No/x) = \frac{2 \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} \times \frac{5}{14}}{P(x)}$$

$$= 24$$

$$14 \times 5^3 \times P(x)$$

$$\begin{aligned}
 &= \frac{P(\text{Yes}/x)}{\frac{3}{7}} - \frac{P(\text{No}/x)}{\frac{4}{7}} \\
 &= \frac{1}{\frac{1}{14} \times 9 \times P(x)} - \frac{1}{\frac{1}{14} \times 12 \times P(x)} = \frac{1}{P(x)} \left(\frac{1}{9 \times 14} - \frac{1}{12 \times 14} \right) \\
 &= \frac{1}{P(x)} \left(\frac{15}{12} - \frac{12}{12} \right) = \frac{1}{P(x)} \times \frac{3}{12} = \frac{1}{P(x)} \times \frac{1}{4} > 0
 \end{aligned}$$

$\therefore P(\text{Yes}/x) > P(\text{No}/x)$
 $\rightarrow P(\text{Yes}/x) < P(\text{No}/x)$
 $\therefore \text{ANS} = \text{NO}$

If we find $P(x)$

$$\begin{aligned}
 P(x) &= P(\text{Sunny} \cap \text{cool} \cap \text{High} \cap \text{True}) \\
 &= \frac{5}{14} \times P(\text{Sunny}) P(\text{cool}) P(\text{High}) P(\text{True}) \\
 &= \frac{5}{14} \times \frac{4}{14} \times \frac{7}{14} \times \frac{6}{14} \\
 &= 0.02186
 \end{aligned}$$

$$P(\text{Yes}/x) = \frac{0.0079}{0.02186} = 0.36139$$

$$P(\text{No}/x) = \frac{0.0137}{0.02186} = 0.62671$$

\therefore values are now NORMALISED
we find $P(x)$ for normalisation of answer

TEXT CLASSIFICATION

- It is a process of assigning tags or categories to text according to its content
- It is one of the fundamental task of NLP with broad applications such as sentiment

- Text classification is the task of assigning a set of predefined categories to free text.
- Text classifier can be used to organise structures & categories pretty much anything for eg:

- Spam Detection:
- organise the topics according to their category like sports, movie, defence, politics.

Working on Text (Classification)

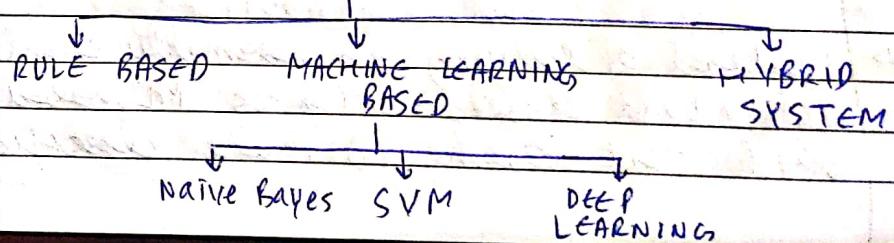
We can do it in 2 ways:

- Manual
- Automatic

→ MANUAL CLASSIFICATION usually provide better quality result but it is very expensive and time consuming

→ AUTOMATIC CLASSIFICATION require machine learning NLP (natural language processing) and other technique to automatically classify the text in more effective and faster way.

TEXT - CLASSIFICATION



- # RULE BASED APPROACH is classified into:

 - organised group by using set of hand-crafted linguistic rules
 - Tell rules the system to use semantically relevant elements of the text to identify relevant category based on its context.
 - Each rule consists of pattern and a predicted category.
eg: You want to classify news into two groups, sports and politics.
 - They are a combination of manual and automatic classification
 - They are human comprehensible and implose every time but, it also has some disadvantages for starters, these system require deep knowledge of domain
 - They are also time consuming since, the generating rule for a complex system can be quite challenging and usually requires a lot of analysis and testing which is performed by human experts
 - Rule based system are also difficult to maintain and do not scale well given that adding new rules can affect the results of pre-existing rules.

MACHINE-LEARNING

- Instead of relying on normally crafted rules, text classification is ML based in which Text learns to make classifications

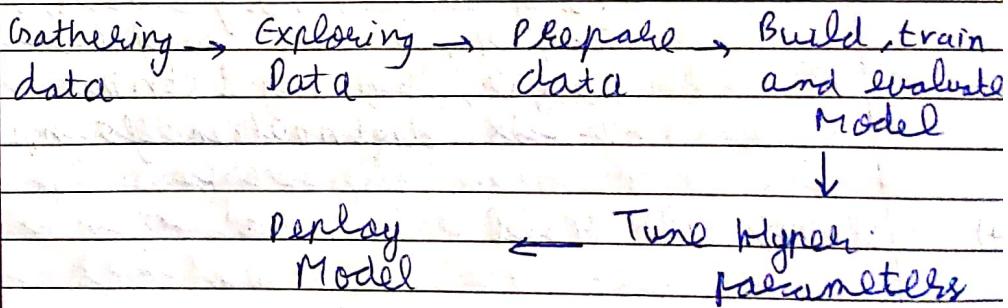
based on past observation
It uses pre-labeled (e.g. training) data.
A machine learning algo can learn
different association b/w pieces of text and
that a particular text / class.

The first step of training a classifier
using ML is feature extraction (method
to represent features into numerical data
in form of a vector). The most common
approach is "bag of words".

HYBRID SYSTEM

- It combines a Bayes classifier trained with ML and a rule based system which is used to further improve the result.
 - These hybrid system can be easily find by adding specific rules for these conflicting text that are not being correctly modelled by Bayes classifier.

TEXT-CLASSIFICATION WORK-FLOW



GATHERING DATA

- If you are using a public API to collect the data then try to understand the limitation of API before using it.
eg: Some API set limitation on the rate with which you can make queries.
- Make sure that the no. of samples for every class and features is not imbalanced, i.e. you should have comparable no. of samples per class.
- Make sure, your sample covers spaces of possible inputs not only the common classes.

EXPLORE DATA

- Building and training data is only one part of workflow. Understanding the characteristics of your data beforehand helps you to build better model.
- It will help in achieving better efficiency and higher accuracy.
- less data for computing and training

Check the data

Collect key matrices

↳ no. of samples, no. of classes, no. of samples per class, no. of words per sample, frequency distribution of words, distribution of sample length

Choose the Model.

PREPARE DATA

- Before our data can be fed to our model, it needs to be transformed to a format that model can understand.
- First, data sample we have gathered may be in specific order. So, we always shuffle the data before doing anything else because we do not want any info. associated with ordering.
- Second, machine learning algos take numbers as input so, we need to convert text to numerical vectors by 2 steps :

(i) TOKENISATION:

Divide text into words or smaller sub-text which will enable good generalisation relationship b/w text and numbers

(ii) VECTORIZATION:

Define a numerical measure to characterise these texts

BUILD, TRAIN AND EVALUATE MODEL

- Building ML model is all about assembling layers together, data-processing building blocks. These layers allow us to specify sequence of transformation we want to perform on our input. As our learning algo takes a single text input and

output a classification, we can create linear stacks of layers using sequential model.

TUNE HYPER PARAMETERS

- we have to choose numbers of hyper parameters for defining and training model.
- We rely on intuitions, examples and best practice recommendation.
- our first choice of hyper-parameters value however may not yield best result. But it can give us a good starting point for training.
- Every problem is different and tuning the hyper-parameters will help refine our model to better represent the particular problem
- Eg. some of hyper parameters are which can be tuned
 - (1) No. of layers in model
 - (2) No. of units per layer
 - (3) Drop-out rate
 - (4) Learning Rate

DEPLOY YOUR MODEL

- Following key things should be kept in mind while deploying our model:
 - (i) Make sure your production data follows same distributional as your

Page No. 3
Date 8/8

training and evaluation data
 (ii) regularly re-evaluate by collecting more training data
 (iii) if your data distribution changes.
 Re-train your model.

ques.	chills	running nose	headache	fever	flu
	Y	N	Mild	Y	N
	Y	Y	No	N	Y
	Y	N	Strong	Y	Y
	N	Y	Mild	Y	X
	N	N	No	N	N
	N	Y	Strong	Y	Y
	N	Y	Strong	N	N
	Y	Y	Mild	Y	Y

$$x = [Y \ N \ Mild \ N] \text{ flu?}$$

$$P(\text{Yes} | \text{No}) = \frac{1}{3}$$

$$\text{CHILLS } P(\text{Yes} | \text{Yes}) = \frac{3}{5} \quad P(\text{Yes} | \text{No}) = \frac{2}{3}$$

$$\text{RUNNING NOSE } P(\text{No} | \text{Yes}) = \frac{1}{5} \quad P(\text{Yes} | \text{No}) = \frac{2}{3}$$

$$\text{HEADACHE } P(\text{Mild} | \text{Yes}) = \frac{2}{5} \quad P(\text{Mild} | \text{No}) = \frac{1}{3}$$

$$\text{FEVER } P(\text{No} | \text{Yes}) = \frac{1}{5} \quad P(\text{No} | \text{No}) = \frac{2}{3}$$

$$P(\text{Yes} | x) = \frac{3}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{8} = \frac{6}{125 \times 8}$$

$$P(\text{No} | x) = \frac{1}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{3}{8} = \frac{4}{27 \times 8}$$

$$P(\text{Yes}/x) = \frac{e^3}{125 \times 84}$$

$$P(\text{No}/x) = \frac{e^2}{27 \times 84}$$

$$P(\text{Yes}/x) = P(\text{No}/x)$$

$$= \frac{1}{4} \left(\frac{3}{125} - \frac{2}{27} \right) < 0$$

$$\therefore P(\text{No}/x) > P(\text{Yes}/x)$$

\therefore NO FLU

NAIVE BAYES \rightarrow TEXT CLASSIFICATION

- Present each document by vector of words
 \rightarrow one attribute per position in document
 (doc)
- Learning, use training eg. to estimate
 $P(+), P(-), P(\text{doc}+), P(\text{doc}-)$

Naive bayes,

$$P(\text{doc} | v_j) = \prod_{i=1}^{\text{length}(\text{doc})} P(a_i = w_k | v_j)$$

$$= P(a_i = w_k | v_j)$$

Probability $P(a_i = w_k | v_j)$ is that word in
 collision $\overset{\uparrow}{i}$ is w_k given v_j .

Ques.	DOC	Text	Class
1	I	I loved the movie	+
2	F	F hated the movie	-
3	A	A great movie, good movie	+
4	n	poor acting	-
5	g	great acting, a good movie	+

Test data = "I hated the poor acting"

ANS

- convert text to features where no of distinct words in features and their frequency in value in feature vector

Class	I	loved	the	movie	hated	a	great	good	poor	acting
+	1	1	1	1	0	0	0	0	0	0
-	1	0	1	1	1	0	0	0	0	0
+	0	0	0	2	0	1	1	1	0	0
-	0	0	0	0	0	0	0	0	1	1
+	0	0	0	1	0	1	1	1	0	1

$$(i) P(+)=\frac{3}{5} \quad P(-)=\frac{2}{5} \quad \text{PRIORI PROB}$$

$$(ii) P(+|+)$$

(CLASS CONDITIONAL PROBABILITY)

$$P(w_k | x_i) = \frac{n_k + 1}{|\text{vocab}| + 1}$$

word \hookrightarrow class

where n_k = number of times w_k occurs in class x_i

n = total words in class $w_1 + w_2$
 $|w_{cl}|$ = number of distinct words
 in content

$$P(I|+) = \frac{1+1}{10+14} = \frac{2}{24} \quad P(I|-) = \frac{1+1}{10+6} = \frac{2}{16}$$

$$P(\text{loved}|+) = \frac{2}{24} \quad P(\text{loved}|-) = \frac{1}{16}$$

$$P(\text{the}|+) = \frac{2}{24} \quad P(\text{the}|-) = \frac{2}{16}$$

$$P(\text{mail}|+) = \frac{5}{24} \quad P(\text{mail}|-) = \frac{2}{16}$$

$$P(\text{hated}|+) = \frac{1}{24} \quad P(\text{hated}|-) = \frac{2}{16}$$

$$P(a|+) = \frac{3}{24} \quad P(a|-) = \frac{1}{16}$$

$$P(\text{great}|+) = \frac{3}{24} \quad P(\text{great}|-) = \frac{1}{16}$$

$$P(\text{good}|+) = \frac{3}{24} \quad P(\text{good}|-) = \frac{1}{16}$$

$$P(\text{poor}|+) = \frac{1}{24} \quad P(\text{poor}|-) = \frac{2}{16}$$

$$P(\text{acting}|+) = \frac{2}{24} \quad P(\text{acting}|-) = \frac{2}{16}$$

(iv) CALCULATION OF POSTERIORI BASED ON NAIVE BAYES ASSUMPTION

$$P(+|\text{text}) = P(\text{text}|+) P(+)$$

$$P(\text{text})$$

Ignoring

$$= P(+|I|) P(+|I)$$

$$= P(I|+) P(\text{hated}(+)) P(\text{the}(+)) P(\text{poor}(+))$$

$$P(\text{acting}(+)) P(+)$$

$$= \frac{2}{24} \times \frac{1}{24} \times \frac{2}{24} \times \frac{1}{24} \times \frac{2}{24} \times \frac{3}{5}$$

$$= \frac{24}{24^5 \times 5} = \frac{1}{5 \times 24^4}$$

$$P(-|\text{text}) = \frac{2}{16} \times \frac{2}{16} \times \frac{2}{16} \times \frac{2}{16} \times \frac{2}{16} \times \frac{2}{5}$$

$$= \frac{1}{5 \times 16^3 \times 4}$$

$$P(-|\text{text}) > P(+|\text{text})$$

∴ Text belongs to negative class

SYLLABUS

- Basic structure of pattern
- Naive Bayes ✓
- Bayesian classification ✓
- Intro to Bayesian decision theory for 2 category classification, continuous case
- Estimation of post-probability results
- Conditional Risk (R) ✓
- Min-Risk classification ✓
- Min-Error classification ✓
- Text classification → ✓
- Text classification through Naive Bayes ✓
- Maximum likelihood estimation
- KNN ✓ → Distance measure ✓
- PDF ✓

Ques 1. TEXT

A great game
The election was over
very clean match
A clean but forgettable game
It was a close election

TAG

Sports
Not sports
Sports
Sports
Not sports

TEST → IT IS A VERY CLOSE GAME

Ques 2.

TEXT

Chinese Beijing Chinese
Chinese Chinese Shanghai
Chinese Macau
Tokyo Japan Chinese

CLASS

C

C

C

J

TEST TEXTS

Chinese Chinese Chinese Tokyo Japan

Ans 1.

A great game	The election was over	very clean match	
1 1 1 0 0 0 0 0 0	0 0 1 1 1 0 0 0 0	0 0 0 0 0 1 1 1 1	0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0
0 1 0 0 1 0 0 0 0	0 1 0 0 0 0 0 0 0	0 1 0 0 0 0 0 0 0	0 1 0 0 0 0 0 0 0
1 0 0 0 1 1 0 0 0	0 0 1 1 0 0 0 0 0	1 0 0 0 0 0 0 0 0	0 0 1 1 0 0 0 0 0

but forgettable	IT	close	close
0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0
0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0
0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0
1 0 0 0	1 0 0 0	0 0 0 0	0 0 0 0
0 0 0 0	0 0 0 0	1 1 1 1	1 1 1 1

$$P(A|S) = \frac{2+1}{14+11} = \frac{3}{25}$$

$$P(A|NS) = \frac{1+1}{14+11} = \frac{2}{23}$$

$$P(\text{rely}|S) = \frac{1+7}{14+11} = \frac{2}{25}$$

$$P(\text{rely}|NS) = \frac{1}{23}$$

$$P(\text{close}|S) = \frac{0+7}{14+11} = \frac{1}{25}$$

$$P(\text{close}|NS) = \frac{2}{23}$$

$$P(\text{game}|S) = \frac{2+1}{14+11} = \frac{3}{25}$$

$$P(\text{game}|NS) = \frac{1}{23}$$

23

$$P(S) = \frac{3}{5} \quad P(NS) = \frac{2}{5}$$

$$P(\text{Text}|S) = \frac{3}{5} \times \frac{3}{25} \times \frac{2}{25} \times \frac{1}{25} \times \frac{3}{25} \\ = \frac{54}{5 \times 25^4}$$

$$P(\text{Text}|NS) = \frac{2}{5} \times \frac{2}{23} \times \frac{1}{23} \times \frac{2}{23} \times \frac{1}{23} \\ = \frac{8}{5 \times 23^4}$$

∴ Sports Text

↳ as $P(\text{Text}|S) > P(\text{Text}|NS)$

$$P(S|x) = \frac{P(x|S)P(S)}{P(x)}$$

INS. 2

Chinese	Beijing	Sanghai	Macao	Tokyo	Japan	class
2	1	0	0	0	0	C
2	0	1	0	0	0	C
1	0	0	1	0	0	C
1	0	0	0	0	1	j

$$P(\text{Chinese} | C) = \frac{5+1}{8+6} = \frac{6}{14}$$

$$P(\text{Tokyo} | C) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(\text{Japan} | C) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(\text{Chinese} | j) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{Tokyo} | j) = \frac{1+1}{9} = \frac{2}{9}$$

$$P(\text{Japan} | j) = \frac{2}{9}$$

$$P(\text{Text} | C) = \frac{3}{4} \times \left(\frac{6}{14}\right)^2 \times \frac{1}{14} \times \frac{1}{14} = \frac{18}{4 \times 14^5} = \frac{0.0016}{0.00030112}$$

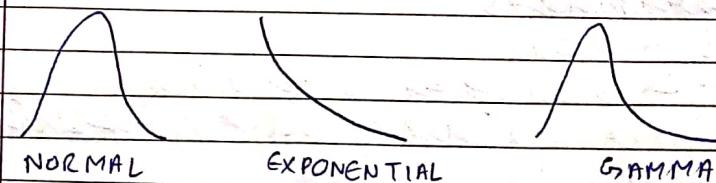
$$P(\text{Text} | j) = \frac{1}{4} \times \left(\frac{2}{9}\right)^2 = \frac{82}{4 \times 9^5} = \frac{0.0027}{0.000133}$$

∴ Chinese \vdash Japanese
 \vdash Chinese.

MAXIMUM LIKELIHOOD ESTIMATION

→ Is a method of statistical model for given observation. The method obtains the parameter estimates by finding parameter value that maximize the likelihood func.

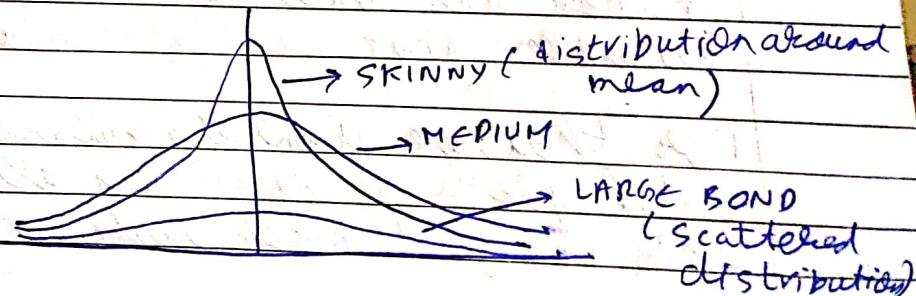
→ It is method that determines values for parameters of a model. The parameter value are found s.t. they maximize the likelihood that measures closeness of the model produce the data that were actually obtained.



NORMALLY DISTRIBUTED

(i) We expect most of the measurements to lie close to mean.

(ii) We expect measurement to be relatively symmetric about mean.



so far any given observed data, first we decide which model we think best describes process of generating data. In this example we assume that a generation process can be adequately described by gaussian (Normal) distribution.

¹ (Fig) Visual inspection of figure suggest that gaussian distribution is most suitable because most of 70 points are clustered in middle with few points scattered in left and right.

We know gaussian distribution has 2 parameters

(i) $\mu \rightarrow$ mean

(ii) $\sigma \rightarrow$ Standard deviation

Different values of μ, σ result in diff. curves.

Now, we want to calculate which curve is most reasonable for generating data points that we are observing.

Maximum likelihood estimation is a method that will find mean and standard deviation that best fit data.

For that, we want to calculate total prob of observing all data i.e. joint prob distribution of all

points which requires calculation of conditional prob. So we make assumption :-

(i) Each data point is generated independently of others for easy computation

The prob density of observing a single data point x that is generated from gaussian distribution is given by :-

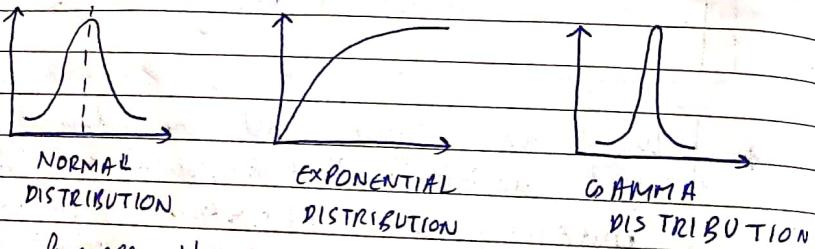
$$P(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

We can figure out the value of mean & std. deviation that results in giving the maximum value of above expression by using calculus (maxima of func).

$$\frac{\partial}{\partial \mu} \left(\frac{(x-\mu)^2}{2\sigma^2} \right) = 0$$

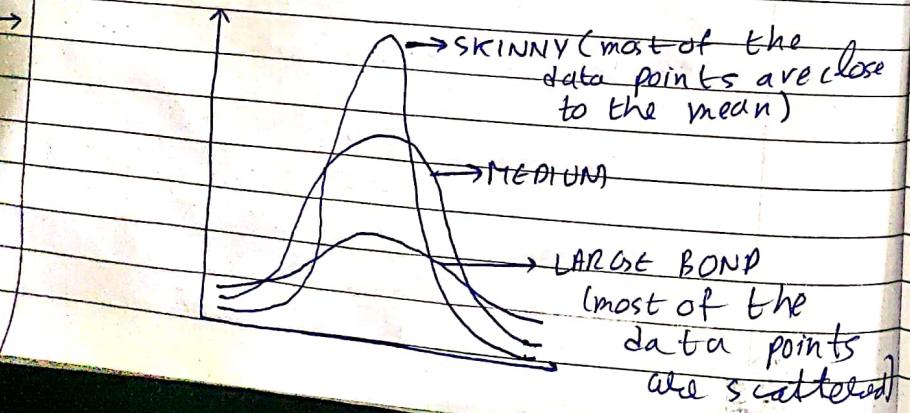
MAXIMUM LIKELIHOOD ESTIMATION

- Maximum likelihood estimation is a process of finding the curve and the parameters of the curve that best matches the data set given to us.
- We try to find the parameters of the curve in such a way that it maximises the maximum likelihood function.
- The curve we choose is based on data points available to us.



→ We choose the gaussian distribution when

- data set is symmetric about mean
- data is close to the mean



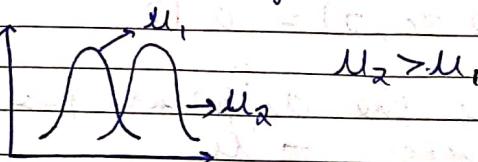
→ MAXIMUM LIKELIHOOD ESTIMATION FOR GAUSSIAN DISTRIBUTION

- (a) Gaussian distribution has 2 parameters
 $\mu \rightarrow$ mean
 $\sigma \rightarrow$ standard deviation.

Based on μ, σ , shape of curve changes

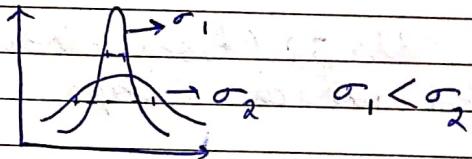
EFFECT OF μ

$\mu \uparrow$, graph shifts right
 $\mu \downarrow$, graph shifts left.



EFFECT OF σ

$\sigma \uparrow$, graph is broader
 $\sigma \downarrow$, graph is narrower about mean



- (b) Maximum likelihood involves calculation of μ, σ , that maximises the probability density function,

$$P(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

For ex:

(i) Let us assume data point $x = 32$

(ii) Fixing $\sigma = 2$, and finding μ for which slope = 0, or $\frac{dP(x)}{dx} = 0$

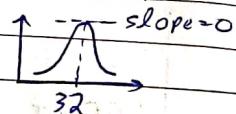
$$P(32, 28, 2) = \frac{1}{\sqrt{2\pi(4)}} \exp\left(-\frac{(32-28)^2}{2 \cdot 4}\right)$$

$$= 0.03$$

$$P(32, 30, 2) = 0.12$$

(iii) Plot the different points with x -axis $\rightarrow \mu$

$$y\text{-axis} \Rightarrow P(32, \mu, 2)$$



The point with slope = 0 is the actual μ required.

(iv) Fix $\mu = 32$ (from above) and finding σ for which slope = 0

(v) If more than one data points are given, say x_1, x_2 is given

$$P(x_1, \mu, \sigma) = L_1$$

$$P(x_2, \mu, \sigma) = L_2$$

$$L = L_1 \times L_2$$

Maximum likelihood estimation assumes that measurements of x_1 and x_2 are independent.

$$L(x_1=32, x_2=34; 28; 2)$$

$$= \frac{1}{\sqrt{2\pi(4)}} \exp\left(-\frac{(32-28)^2}{2(4)}\right) \frac{1}{\sqrt{2\pi(4)}} \exp\left(-\frac{(34-28)^2}{2(4)}\right)$$

(3) GENERAL CASE OF n points in data set

$$L = L_1 \times L_2 \times \dots \times L_n$$

PROBLEM: we want to maximize derivative of L which is computationally expensive

SOLUTION: take log of above log is a monotonically ↑ function, so we can take log of above

$$\log L = \log L_1 + \log L_2 + \dots + \log L_n$$

FINDING $\log L_n$

$$\log L_n = \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_n - \mu)^2}{2\sigma^2}\right)\right)$$

$$= \log \frac{1}{\sqrt{2\pi}} + \log \frac{1}{\sigma^2} + \log e^{-\frac{(x_n - \mu)^2}{2\sigma^2}}$$

log(L_n)

$$= \frac{-1}{2} \log 2\pi - n \log \sigma - \frac{(x_n - \mu)^2}{2\sigma^2} \text{ logee}$$

$$\log(L_n) = \frac{-1}{2} \log(2\pi) - \log \sigma - \frac{(x_n - \mu)^2}{2\sigma^2}$$

USING THIS :

$$\# \log L = -n \log 2\pi - n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

(a) Fixing $\sigma = \text{const}$ and maximising L , finding μ for which slope = 0

Differentiating wrt μ

$$\frac{d(\log L)}{d\mu} = 0 + 0 + \frac{-1}{2\sigma^2} \sum (2(x_i - \mu)(-1))$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)$$

Equating to 0

$$\rightarrow \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$x_1 + x_2 + \dots + x_n - n\mu = 0$$

$$\mu = \frac{x_1 + x_2 + \dots + x_n}{n}$$

(b) Fixing $\mu = \frac{x_1 + x_2 + \dots + x_n}{n}$, and finding σ for which slope = 0

Differentiating wrt σ

$$\rightarrow 0 - \frac{n}{\sigma} - \sum_{i=1}^n (x_i - \mu)^2 \cdot \frac{(-2) \times 1}{\sigma^3} \cdot \frac{1}{2}$$

$$= \frac{-n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2$$

Equating to 0

$$\frac{n}{\sigma} = \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2$$

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

$$\sigma = \sqrt{\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots + (x_n - \mu)^2}{n}}$$

→ This is used in text classification
with $m = \text{vocab}$ $\pi_{\text{tail}} = \frac{1}{m}$

→ This is done so that
 $\text{prob}_{\text{tail}} = 0$

(1) M-ESTIMATE

$$\pi_{\text{tail}} = \frac{N_{\text{tail}} + m \pi_{\text{tail}}}{N_{\text{all}} + m}$$

↑ occurrences of tail ↑ prior expectation
total outcomes confidence

- m-estimate helps us to improve probability analysis for less samples
- If $N_{\text{all}} = N_{\text{tail}} = 0 \rightarrow$ degenerates to prior expectation

Toss No.	1	2	3	4	5
Outcome	tail	tail	head	tail	head

$$\pi_{\text{tail}} = 0.5 \quad (i) m = 2$$

$$(ii) m = 100$$

$$(iii) m = 1$$

NS. (i) m-estimate

Toss No.	1	2	3	4	5
----------	---	---	---	---	---

N_{tail}	1.0	2.0	2	3	3
-------------------	-----	-----	---	---	---

$$\begin{array}{cccccc} \text{m-estimate} & 1+1 & 2+1 & 2+1 & 3+1 & 3+1 \\ & 1+2 & 2+2 & 3+2 & 4+2 & 5+2 \\ & = \frac{2}{3} & = \frac{3}{4} & = \frac{3}{5} & = \frac{4}{6} & = \frac{4}{7} \end{array}$$

Toss No.	1	2	3	4	5
----------	---	---	---	---	---

N_{tail}	1	2	2	3	3
-------------------	---	---	---	---	---

$$\begin{array}{cccccc} \text{m-estimate} & \frac{51}{101} & \frac{58}{102} & \frac{52}{103} & \frac{53}{104} & \frac{53}{105} \end{array}$$

Page No.		
Date		

Toss No.	1	2	3	4	5
N_{tail}	1	2	2	3	3
m-estimate	3	5	5	7	7
	4	6	8	10	12

(2) PDF (PROBABILITY DISTRIBUTION FUNCTION)

→ Standard Gaussian PDF
(FOR CONTINUOUS FEATURES)

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

→ In case of multiple features

$$P(x) = \frac{1}{(2\pi)^n \sigma^n} \sum_{i=1}^n \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma^2}\right)$$

	Attr1	Attr2	Attr3	Class	$\sigma = 1$
ex 1	3.1	2.1	2.3	+	
ex 2	4.2	6.2	7.6	+	
ex 3	7.8	1.3	0.5	+	
ex 4	2.3	5.2	2.4	-	
ex 5	6.4	3.2	4.3	-	
ex 6	1.3	5.8	3.3	-	

$$P(x_1 | +) = \frac{1}{(2\pi)^{3/2}} \left[\exp\left(-\frac{(x_1 - 3.1)^2}{2}\right) + \exp\left(-\frac{(x_1 - 4.2)^2}{2}\right) \right]$$

$$+ \exp\left(-\frac{(x_1 - 7.8)^2}{2}\right)$$

$$P(x_1 | -) = \frac{1}{(2\pi)^{3/2}} \left[\exp\left(-\frac{(x_1 - 2.3)^2}{2}\right) + \exp\left(-\frac{(x_1 - 6.4)^2}{2}\right) \right]$$

$$+ \exp\left(-\frac{(x_1 - 1.3)^2}{2}\right)$$

$$P(x_1|+) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-2-1)^2}{2}\right) + \exp\left(-\frac{(x-6-2)^2}{2}\right) + \exp\left(-\frac{(x-1-3)^2}{2}\right) \right]$$

$$P(x_1|-) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-5-2)^2}{2}\right) + \exp\left(-\frac{(x-3-2)^2}{2}\right) + \exp\left(-\frac{(x-5-8)^2}{2}\right) \right]$$

$$P(x_3|+) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-2-3)^2}{2}\right) + \exp\left(-\frac{(x-7-6)^2}{2}\right) + \exp\left(-\frac{(x-6-5)^2}{2}\right) \right]$$

$$P(x_3|-) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-2-4)^2}{2}\right) + \exp\left(-\frac{(x-4-3)^2}{2}\right) + \exp\left(-\frac{(x-3-3)^2}{2}\right) \right]$$

→ TESTING VARIABLE : $x(9, 2.5, 3.2)$

$$P(+|x) = P(x|+) P(+)$$

$$= P(x_1|+) P(x_2|+) P(x_3|+) P(+) \\ = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(9-3.1)^2}{2}\right) + \exp\left(-\frac{(9-4.2)^2}{2}\right) + \exp\left(-\frac{(9-7.8)^2}{2}\right) \right] x$$

$$\frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(2.5-2.1)^2}{2}\right) + \exp\left(-\frac{(2.5-6.2)^2}{2}\right) + \exp\left(-\frac{(2.5-1.3)^2}{2}\right) \right] x$$

... x 3

$P(-|x) =$ similarly

→ BASIC-STRUCTURE OF PATTERN RECOGNITION

(a) PATTERN

- * an object, process or event that is given a name
- * an image is not a pattern till it is given a name

(b) PATTERN RECOGNITION

- * The process by which machines can
 - (i) observe the environment
 - (ii) identify the patterns of interest
 - (iii) classify the pattern

(c) CLASSES

Set of pattern having common attributes and common source origin in most cases

(d) FEATURE

The attributes / values that are common in objects of same class, and different in pattern of different classes is called feature

→ DATA - ACQUISITION AND SENSING

The process of acquiring and extracting the relevant and useful data from the environment

e.g : extracting address from envelop

Page No. _____
Date _____

Page No. _____
Date _____

$$P(x_1|+) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-2-1)^2}{2}\right) + \exp\left(-\frac{(x-6-2)^2}{2}\right) + \exp\left(-\frac{(x-1-3)^2}{2}\right) \right]$$

$$P(x_1|-) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-5-2)^2}{2}\right) + \exp\left(-\frac{(x-3-2)^2}{2}\right) + \exp\left(-\frac{(x-5-8)^2}{2}\right) \right]$$

$$P(x_2|+) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-2-3)^2}{2}\right) + \exp\left(-\frac{(x-7-6)^2}{2}\right) + \exp\left(-\frac{(x-6-5)^2}{2}\right) \right]$$

$$P(x_2|-) = \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(x-2-4)^2}{2}\right) + \exp\left(-\frac{(x-4-3)^2}{2}\right) + \exp\left(-\frac{(x-3-3)^2}{2}\right) \right]$$

→ TESTING VARIABLE : $x(9, 2.5, 3.2)$

$$P(+|x) = P(x|+) P(+)$$

$$= P(x_1|+) P(x_2|+) P(x_3|+) P(+)$$

$$= \frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(9-3.1)^2}{2}\right) + \exp\left(-\frac{(9-4.2)^2}{2}\right) + \exp\left(-\frac{(9-7.8)^2}{2}\right) \right] x$$

$$\frac{1}{2\pi^{3/2}} \left[\exp\left(-\frac{(2.5-2.1)^2}{2}\right) + \exp\left(-\frac{(2.5-6.2)^2}{2}\right) + \exp\left(-\frac{(2.5-1.3)^2}{2}\right) \right] x$$

$$\dots x \frac{3}{6}$$

$P(-|x) =$ similarly

→ BASIC STRUCTURE OF PATTERN RECOGNITION

(a) PATTERN

- * an object, process or event that is given a name
- * an image is not a pattern till it is given a name

(b) PATTERN RECOGNITION

- * The process by which machines can
 - (i) observe the environment
 - (ii) identify the patterns of interest
 - (iii) classify the pattern

(c) CLASSES

Set of pattern having common attributes and common source origin in most cases

(d) FEATURE

The attributes / values that are common in objects of same class, and different in pattern of different classes is called feature

→ DATA ACQUISITION AND SENSING

The process of acquiring and extracting the relevant and useful data from the environment

e.g: extracting address from envelop

→ FEATURE EXTRACTION

The process of extracting relevant features about the data and preparation of feature vector containing numerical values based upon which we can classify the patterns.

→ PRE-PROCESSING

This step deals with removal of noise from the pattern. Noise can be defined as property of sensed pattern that is due to randomness of world and not due to model.

→ POST-PROCESSING

Involves calculation of accuracy, performance, evaluation of confidence etc.

* NORMALISATION OF FEATURE POINTS IN KNN

(a) To bring each feature in 0-1,

x_{ij} is replaced by $\frac{x_{ij} - \min(\text{column})}{\max(\text{col } j) - \min(\text{col })}$

(b) eq. $x_1 \quad x_2 \quad x_3 \quad \text{Class}$

10 15 20 1

25 30 1 2

2 5 3 1

Changes to

x_1	x_2	x_3	Class
$\frac{10-2}{25-2}$	$\frac{15-5}{30-5}$	$\frac{20-1}{20-1}$	1
:	:	:	2
:	:	1	

(b) Feature points: YY Y N N N N

Find $P(y)$ that maximises the likelihood

ANS. $P(y) = \begin{cases} x & y = Y \\ 1-x & y = N \end{cases}$

For given feature points

$$P(x) = x^3 (1-x)^5$$

$$\log P(x) = 3 \log x + 5 \log(1-x)$$

Diff wrt x

$$\frac{d \log P(x)}{dx} = 0 = \frac{3}{x} + \frac{5}{1-x}$$

$$\frac{3}{x} = \frac{5}{1-x}$$

$$3 - 3x = 5x$$

$$8x = 3 \quad x = \frac{3}{8}$$

$$\therefore P(y = \text{Yes}) = 3/8$$

$$P(y = \text{No}) = 5/8$$

WEIGHTED-KNN ALGORITHM

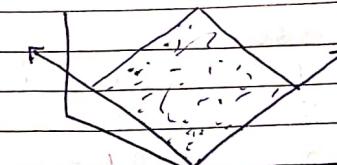
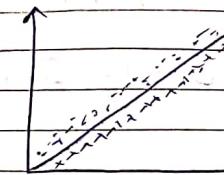
→ In weighted-KNN algorithm, higher preference or weights are assigned to closer neighbours and less weights are assigned to far away neighbours.

→ $w_i = \frac{1}{d(x_i, y)}$ where y = testing point
 x_i = training point

ALGORITHM:

- (i) Find dist b/w test point and all training points
 - (ii) Find the k-nearest points.
 - (iii) $w_i = \frac{1}{d(x_i, y)}$ weight assigned to each point such that closer points have a higher weight. This w_i can be varied based on the application
 - (iv) Compute m_j for each class
- $m_j = \sum_{i=1}^k \text{class}(x_i, j) \times w_i$
- where $n = \text{no. of features}$
 where $\delta(x_i, j) = 1 \text{ if } \text{class}(x_i) = j$
 0 else
- (v) class with max m_j is the required class of test point.

DIMENSIONALITY REDUCTION



Data is in 3D, but it is about a line only. so instead of representing in 3D wrt (0,0) origin, we can take eqn of line and represent in 2D

Data is in 3D. Can be represented in 2D

RANK OF MATRIX

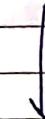
number of linearly independent columns of matrix is rank of matrix

M T W Th Fr Sat Sun

	M	T	W	Th	Fr	Sat	Sun
A	1	1	1	0	0		
B	2	2	2	0	0		
C	1	1	1	0	0		
D	5	5	5	0	0		
E	0	0	0	2	2		
F	0	0	0	3	3		
G	0	0	0	1	1		

$$\text{eg2 } A = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$$

$$\text{Row 3} = \text{Row 1} - \text{Row 2}$$



Basic axis $[0 \ 0 \ 0]$

change it to $[1 \ 2 \ 1]$
 $[-2 \ -3 \ 1]$

$$[1 \ 2 \ 1] + 0[-2 \ -3 \ 1] \rightarrow \text{first row } [1 \ 2 \ 1]$$

$$0[1 \ 2 \ 1] + 1[-2 \ -3 \ 1] \rightarrow \text{second row } [0 \ 1 \ 1]$$

$$\text{third row } [1 \ 1 \ 1]$$

actually $[1 \ -1 \ 1]$

for eg 1, base axis
 called vector
 Date _____
 [1 1 1 0 0]
 [0 0 0 1 1]

from
 1 [1 0] = [1 1 1 0 0] + 0x[0 0 0 1 1]
 2 [2 0]
 3 [3 1 0]
 4 [4 5 0] → called POINTS
 5 [0 2]
 6 [0 3]
 7 [0 1]

WHY DIMENSIONALITY REDUCTION?

- (i) Discovers hidden correlation
- (ii) Reduces noisy and redundant features
- (iii) Interpretation and visualization
- (iv) Easier storage and processing of data

SVD - SINGULAR VALUE DECOMPOSITION

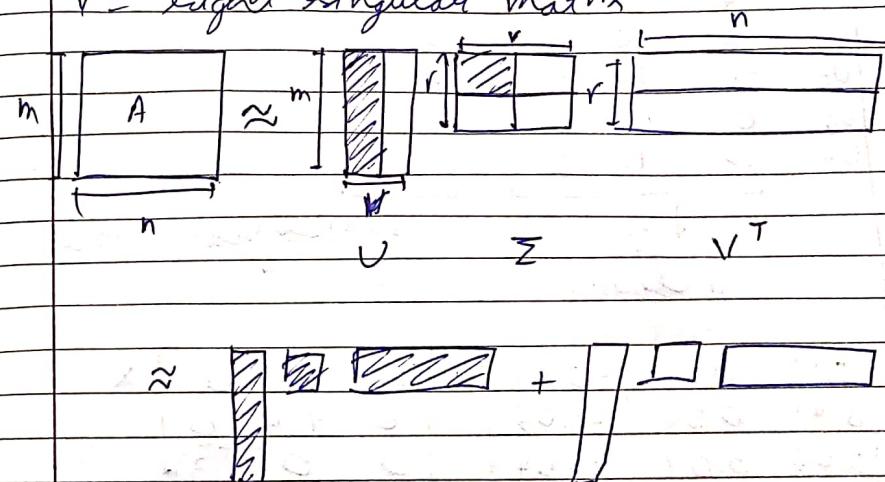
$$A_{m \times n} = U_{m \times r} \Sigma_{r \times r} (V_{n \times r})^T$$

It is always possible to decompose a real matrix A into $U \Sigma V^T$ where

- (i) U, Σ , and V are unique
- (ii) U and V are column orthonormal meaning $\rightarrow U^T U = I$
- (iii) Σ is diagonal matrix where each $V^T V = I$

entry is positive and sorted in decreasing order

A = input data matrix ($m \times n$)
 U = left singular matrix
 Σ = singular matrix \rightarrow determinant $\neq 0$
 V = right singular matrix



As Σ is decreasing matrix, most of the values at the end are 0. So if we know which values to neglect, our task will be reduced

$$\begin{aligned} A &= U \Sigma V^T \\ &= U \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \bar{\sigma}_3 & 0 \\ 0 & 0 & 0 & \sigma_4 \end{bmatrix} V^T \end{aligned}$$

→ as σ_3 and σ_4 are close to 0, we can neglect this part

$$\|M_p\| = \sqrt{\sum_{ij} (A_{ij} - B_{ij})}$$

eg → Following is data of movies

U_1	1	1	1	0	0
U_2	3	3	3	0	0
U_3	4	4	4	0	0
U_4	5	5	5	0	0
U_5	0	2	0	4	4
U_6	0	0	0	5	5
U_7	0	1	0	2	2

sci-fic
movies

↳ Romantic movies

$$\begin{bmatrix}
 0.13 & 0.02 & -0.01 \\
 0.41 & 0.02 & -0.03 \\
 0.55 & 0.09 & -0.04 \\
 0.68 & 0.11 & -0.05 \\
 0.15 & -0.59 & 0.65 \\
 0.07 & -0.73 & 0.67 \\
 0.07 & -0.29 & 0.32
 \end{bmatrix} \times
 \begin{bmatrix}
 12.4 & 0 & 0 \\
 0 & 9.5 & 0 \\
 0 & 0 & 1.3
 \end{bmatrix} \times$$

This is the min value
amongst all 3. So as

this is closest to 0

We neglect it. As
we neglect last value, we also
neglect last column of U and last
row of V^T

$$\begin{bmatrix}
 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\
 0.12 & -0.102 & 0.12 & -0.69 & 0.69 \\
 0.40 & -0.80 & 0.40 & 0.09 & 0.09
 \end{bmatrix}$$

V- PRODUCT

It gives the coordinates of point in
coordinate axis.

$$U \Sigma V^T = \begin{bmatrix} 4.61 & 0.19 & -0.01 \\ 5.08 & 0.66 & -0.03 \\ 6.82 & 0.85 & -0.05 \\ 8.43 & 1.04 & -0.06 \\ 1.86 & -5.60 & -0.84 \\ 0.86 & -6.93 & - \\ 0.86 & -2.75 & - \end{bmatrix}$$

eg → Test $q = \begin{bmatrix} 5 & 0 & 0 & 0 & 0 \\ 4 & 5 & 0 & 0 & 0 \\ 0 & 4 & 5 & 0 & 0 \end{bmatrix}$

$$\begin{bmatrix} 0 & 4 & 5 & 0 & 0 \end{bmatrix} \times V$$

amount movie
belongs to romantic

$$\begin{bmatrix} 0 & 4 & 5 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 0.56 & 0.12 \\ 0.59 & -0.12 \\ 0.56 & 0.12 \\ 0.09 & -0.69 \\ 0.09 & 0.69 \end{bmatrix} = \begin{bmatrix} 5.2 & 0.4 \end{bmatrix}$$

Amount the
movie
belongs to
sci-fi

PRINCIPLE COMPONENT ANALYSIS

- Mean
 - Std deviation] 1D
 - Variance]
- Co-variance $\Sigma - 20$

$$\text{Var} = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x}) \quad \text{ALWAYS}$$

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1} \quad \begin{matrix} \uparrow \\ (\text{use } n \text{ in MEAN}) \end{matrix}$$

eg → Hours Marks Find covariance

9 39

15

56

$$\text{ANS. } \bar{x}_x = \frac{\sum x_i}{n-1}$$

25

93

$$= \frac{-39.92}{12} = 13.92$$

14

61

$$\bar{x}_y = \frac{\sum y_i}{n-1}$$

10

50

$$= \frac{749}{12} = 62.42$$

18

75

$$\text{Cov}(x, y)$$

0

32

$$= (-39.92)(-50.42)$$

16

85

≠

$$= (-13.92)(-30.42)$$

$$M_i - \bar{M}$$

$$M_i - \bar{M}$$

$$M_i - \bar{H}$$

$$M_i - \bar{H}$$

$$-4.92$$

$$-23.42$$

$$2.08$$

$$22.58$$

$$1.08$$

$$-6.42$$

$$-8.92$$

$$-20.42$$

$$11.08$$

$$30.58$$

$$9.08$$

$$7.58$$

$$0.08$$

$$-1.42$$

$$2.08$$

$$3.58$$

$$-3.92$$

$$-12.42$$

$$6.08$$

$$7.58$$

$$4.08$$

$$12.58$$

$$-13.92$$

$$-30.42$$

$$(H_i - \bar{H})(M_i - \bar{M})$$

115.23

-6.93

338.83

-0.11

48.69

51.33

423.45

46.97

182.15

38.15

7.45

106.89

$$\text{cov}(x, y) = 104.59$$

$$C = \text{MATRIX} = \begin{bmatrix} \text{cov}(xx) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(yy) \end{bmatrix}$$

$$\bar{M}_x = \frac{18.1}{10} = 1.81 \quad \bar{M}_y = 1.91$$

$$\text{cov}(x, y) = \frac{5.539}{10} = 0.5539$$

$$\begin{aligned} \text{Var}(x) &= 0.4761 + 1.7161 + 0.1521 + 0.0081 \\ &\quad + 1.6641 + 0.2401 + 0.0361 + 0.6561 \\ &\quad + 0.0961 + 0.717 + 0.5041 \\ &= \frac{5.549}{10} = 0.5549 \end{aligned}$$

EIGEN VECTOR

$$A = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix} \quad |A - dI| = 0$$

Matrix \downarrow const Identity Matrix

$$\begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix} - d \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

$$\begin{bmatrix} -1 & 1 \\ -2 & -3-1 \end{bmatrix} = 0$$

$$1^2 + 3d + 2 = 0$$

$$d_1 = -1 \quad d_2 = -2 \Rightarrow \text{EIGEN VALUE}$$

$$(A - dI)V = 0$$

→ EIGEN VECTOR

$$d = -1$$

$$\begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ -2 & -2 \end{bmatrix}$$

COVARIANCE MATRIX IN 3D

$$C = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) & \text{cov}(x,z) \\ \text{cov}(y,x) & \text{cov}(y,y) & \text{cov}(y,z) \\ \text{cov}(z,x) & \text{cov}(z,y) & \text{cov}(z,z) \end{bmatrix}$$

eg	X	Y	$X_i - \bar{M}_x$	$Y_i - \bar{M}_y$	$(X_i - \bar{M}_x)(Y_i - \bar{M}_y)$
2.5	2.4	0.69	+0.49	0.3381	
0.5	0.7	-1.31	-1.21	1.5851	
2.2	2.9	0.39	0.99	0.3861	
1.9	2.2	0.09	0.39	0.0261	
3.1	3.0	1.29	1.09	1.4061	
2.3	2.7	0.49	0.79	0.3871	
2.0	1.6	0.19	-0.31	-0.0589	
1.0	1.1	-0.81	-0.81	0.6561	
1.5	1.6	-0.31	-0.31	0.0961	
1.1	0.9	-0.71	-1.01	0.7171	

$$\begin{bmatrix} 1 & 1 \\ -2 & -2 \end{bmatrix} \begin{bmatrix} V_{11} \\ V_{12} \end{bmatrix} = 0$$

$$V_{11} + V_{12} = 0 \quad V_{11} = -V_{12}$$

$$-2V_{11} - 2V_{12} = 0 \quad V_{11} = -V_{12}$$

$$V_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

for A_2

$$\begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \begin{bmatrix} V_{21} \\ V_{22} \end{bmatrix} = 0$$

$$2V_{21} + V_{22} = 0 \quad -2V_{21} - V_{22} = 0$$

$$V_{21} = -V_{22} \quad V_{21} = 1 \quad V_{22} = -2$$

$$V_2 = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

PROPERTIES:

- (i) Can only be found for square matrix.
- (ii) Not every square matrix has eigen vector.
- (iii) For given $N \times N$ matrix, there are N eigen vectors.
- (iv) If you scale the vectors by some amount before you multiply it, you will get same multiple of it as result.

a. Find eigen vectors and eigen values.

$$A = \begin{bmatrix} 8 & -8 & -2 \\ 4 & -3 & -2 \\ 3 & -4 & 1 \end{bmatrix}$$

$$\begin{array}{rcl} -3 - 8 & = & -11 \\ 4 + 6 & = & 10 \\ -16 + 9 & = & -7 \\ -11 - 7 & = & -18 \end{array}$$

$$\text{ANS. } |A - dI| = 0$$

$$\begin{vmatrix} 8-d & -8 & -2 \\ 4 & -3-d & -2 \\ 3 & -4 & 1-d \end{vmatrix} = 0$$

$$(8-d)(-3+d^2+2d) + 8(4-4d+6)$$

$$-2(-16+9+3d) = 0$$

$$(8-d)(d^2+2d-11) + 8(10-4d)$$

$$-2(3d-7) = 0$$

$$-d^3 - 2d^2 + 11d + 8d^2 + 16d - 88 + 80 - 32d$$

$$-6d + 14 = 0$$

$$-d^3 + 6d^2 - 11d + 6 = 0$$

$$-d^3 + d^2 + 5d^2 - 5d - 6d + 6 = 0$$

$$-d^2(d-1) + 5d(d-1) - 6(d-1) = 0$$

$$(d-1)(-d^2 + 5d - 6) = 0$$

$$(d-1)(-d^2 + 3d + 2d - 6) = 0$$

$$(d-1)(-d(d-3) + 2(d-3)) = 0$$

$$d = 1, 2, 3.$$

$$d_1 = 1$$

$$\begin{bmatrix} 7 & -8 & -2 \\ 4 & -4 & -2 \\ 3 & -4 & 0 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \end{bmatrix} = 0$$

$$\begin{aligned} 7V_1 - 8V_2 - 2V_3 &= 0 \\ 4V_1 - 4V_2 - 2V_3 &= 0 \\ 3V_1 - 4V_2 &= 0 \end{aligned}$$

$$3V_1 - 4V_2 = 0$$

$$V_1 = 4 \quad V_2 = 3$$

$$16 - 12 - 2V_3 = 0$$

$$2V_3 = 4 \quad V_3 = 2$$

$$\lambda_1 = 2$$

$$\left[\begin{array}{ccc|c} 6 & -8 & -2 & V_1 \\ 4 & -5 & -2 & V_2 \\ 3 & -4 & -1 & V_3 \end{array} \right] = 0$$

$$4V_1 - 5V_2 - 2V_3 = 0$$

$$3V_1 - 4V_2 = V_3 = 0$$

$$\rightarrow 6V_1 - 8V_2 - 2V_3 = 0$$

$$- 4V_1 - 5V_2 - 2V_3 = 0$$

$$2V_1 - 3V_2 = 0$$

$$2V_1 = 3V_2$$

$$V_1 = 3 \quad V_2 = 2$$

$$3(3) - 4(2) - V_3 = 0$$

$$V_3 = 9 - 8 = 1$$

$$\lambda_1 = 3$$

$$\left[\begin{array}{ccc|c} 5 & -8 & -2 & V_1 \\ 4 & -6 & -2 & V_2 \\ 3 & -4 & -2 & V_3 \end{array} \right] = 0$$

$$V_1 = \begin{bmatrix} 4 \\ 3 \\ 2 \end{bmatrix}$$

$$\begin{aligned} 7V_1 - 8V_2 - 2V_3 &= 0 \\ 4V_2 - 6V_2 - 2V_3 &= 0 \end{aligned}$$

$$V_1 - 2V_2 = 0 \quad V_1 = 2V_2$$

$$V_1 = 2 \quad V_2 = 1$$

$$5(2) - 8(1) - 2V_3 = 0$$

$$10 - 8 = 2V_3 \quad V_3 = 1$$

$$\begin{aligned} 3V_1 - 4V_2 - 2V_3 &= 0 \\ 3(2) - 4(1) - 2(1) &= 0 \quad \checkmark \end{aligned}$$

$$V_3 = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$$

FINDING EIGEN EQU

$$\lambda^3 - [\text{Sum of diagonal elements}] \lambda^2 + [\text{Sum of diag elements}] \lambda - \text{minor of } A_{11}$$

$$- |A| = 0$$

minor of A_{11}

+ minor of A_{22}

$$= | -3 & -2 | + | 8 & -2 | + | 8 & -8 |$$

$$= -4 + 3 + 4 = 11$$

→ PCA is statistical procedure that uses an orthogonal transformation to convert a set of observation of possibly correlated variables into a set of values of linearly uncorrelated variables called PRINCIPLE COMPONENTS

- It is a way of identifying pattern in data and expressing data in such a way as to highlight their similarity and differences.
- Since pattern in data can be hard to find in data of high dimension, where luxury of geographical location is not available, PCA is powerful tool to analyse that data.
- Once you have found this pattern in data, then you can compress the data by reducing no. of dimensions without much loss of info.

PROCEDURE:

- collect the data
- subtract the mean
- Calculate co-variance matrix
- Calculate eigen vectors and corresponding eigen values of co-variance matrix
- choose first k-eigen vectors (sorted in decreasing order) that will be ~~fit~~ near k-dimensions.
- transform the original and dimensional data points into k-dimensional.

→ This final step in PCA where we have to choose component ~~for~~ eigen vector in such a way that loss of info is min we wish to keep in our data and form a feature vector.

→ Final data = raw feature vector \times raw data adjusted matrix with eigen vector in column transpose so that eigen vector are now in row with most significant eigen vector at top.

Ques.

$$x = [2.5, 0.5, 9.2] \\ y = [7.4, 0.7, \dots] \text{ before}$$

$$\text{cov}(x, y) = 0.61544 \\ \text{cov}(x, x) = 0.61655$$

STEP(i)
(ii)
(iii)

$$\text{cov} = \begin{bmatrix} 0.616 & 0.615 \\ 0.615 & 0.716 \end{bmatrix} = C$$

STEP(iv) $|C - dI| = 0$

$$\begin{vmatrix} 0.616-1 & 0.615 \\ 0.615 & 0.716-1 \end{vmatrix} = \begin{vmatrix} -0.385 & 0.615 \\ 0.615 & -0.284 \end{vmatrix} \\ = -0.109 - 0.378 = -0.269$$

$$\begin{vmatrix} 0.616-1 & 0.615 \\ 0.615 & 0.716-1 \end{vmatrix} \\ = (0.616-1)(0.716-1) - 0.378 \\ = 0.441 + 1^2 - 0.441 \cdot 1.332 - 0.378 = 0$$

$$d^2 - 1.332d + 0.063 = 0$$

$$d = 1.774274 - 0.252 \\ = 1.522$$

$$\sqrt{d} = 1.234$$

$$\lambda = 1.332 \pm 1.234$$

$$= 1.283, 0, 0.49$$

$$\begin{bmatrix} 1 & 1.283 \\ 0.615 & 0.615 \\ 0.615 & 0.716-1.283 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = 0$$

$$\begin{bmatrix} -0.667 & 0.615 \\ 0.615 & -0.567 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = 0$$

$$-0.667V_1 + 0.615V_2 = 0$$

$$0.615V_1 - 0.567V_2 = 0$$

$$-0.615 \times 0.667 V_1 + (0.615)^2 V_2 = 0$$

$$0.615 \times 0.667 V_1 - (0.567)(0.667) V_2 = 0$$

$$V_1 = \frac{0.615}{0.667} V_2 = 0.922 V_2$$

\Rightarrow Raw adjusted data = Raw feature vector inverse

$\begin{matrix} X \\ \text{final data} \end{matrix}$

= raw feature vector \times final data
transpose
(raw as orthogonal,
inverse = transpose)

\Rightarrow original data = (raw feature \times final data) + Mean
vector transpose

Page No.
Date

Page No.
Date

$$\text{eg} \rightarrow a^2 + b^2 = 1$$

$$V_1^2 + V_2^2 = 1$$

V_1, V_2 are vectors

$$\begin{bmatrix} 9.9453 & 7.9876 \\ 7.9876 & -6.4153 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$-9.9453a + 7.9876b = 0$$

$$-9.9453a + 7.9876\sqrt{1-a^2} = 0$$

$$\frac{\sqrt{1-a^2}}{a^2} = \frac{9.9453}{7.9876} = 1 \quad \text{or } 1.2451$$

$$\frac{1-a^2}{a^2} = 2.4025 \quad 1.550$$

$$1-a^2 = 2.4025a^2 \quad 1.550a^2$$

$$1 = 3.4025a^2 \quad 2.550a^2$$

$$\frac{1}{a^2} = 2.550$$

$$a^2 = 0.3921$$

$$a = 0.6267$$

$$b = 0.7796$$

OR divide final
ans by $\sqrt{a^2+b^2}$

$$V_1 = 0.922 V_2 \quad V_1 = 0.922$$

$$V_2 = 1$$

An normalisation:

$$V_1 = \frac{0.922}{\sqrt{1+(0.922)^2}} = 0.677$$

$$V_2 = 0.7796$$

DECISION TREE

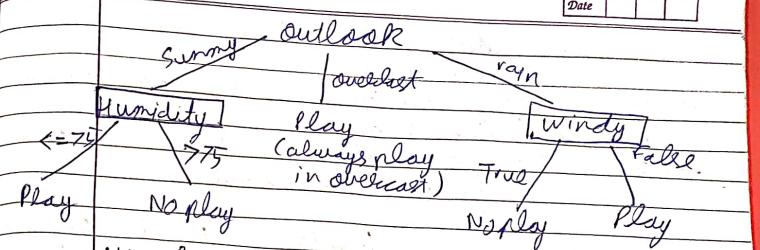
- It is a decision support tool that uses a tree like model of decision and their possible consequence including chance even outcomes, resource cost and utility.
- we can also say that decision tree is flowchart like structure in which each internal node represent test on an attribute. each branch represent outcome of test and each leaf node represent a class label.
- A path from root to leaf represent classification rules.

outlook	Temp	Humidity	Windy	Class
Sunny	79	90	True	No play
Sunny	56	70	False	Play
Sunny	79	75	True	Play
Sunny	60	90	True	No play
overcast	88	88	False	No play
overcast	63	75	True	play
overcast	84	95	False	play
rain	78	60	False	play
rain	66	70	False	No play
rain	68	60	True	No play

- which feature to be at which node is decided using formula
- last data part is anomaly.

Page No. _____
Date _____

Page No. _____
Date _____



We found 5 leaf node i.e. it generates 5 rules

RULES

- If it is sunny and humidity ≤ 75 then play
- If it is sunny and humidity > 75 then do not play
- If it is overcast, then play
- If it is rainy and windy then no play
- If it is rainy and not windy then play

CLASSIFICATION STEPS

- The classification of unknown input vector is done by travelling decision tree from root node to leaf node
- record entries tree at root node
- at no root, a test is applied to determine which child made the record will encounter next.
- This process is repeated until record arrives at a leaf node
- All the labeled that end up at

- a leaf node of the tree all classifier in the same way.
- There is a unique path from root to each leaf node
 - The path is a feature which is used to classify the record.

	outlook	Temp	Humid	Wind	Decision
Sunny	H	High	Weak	No	
Sunny	H	High	Strong	No	
overcast	H	High	Weak	Yes	
Rainy	Mild	High	Weak	Yes	
Rainy	C	Normal	Weak	Yes	
Rainy	C	Normal	Strong	No	
overcast	C	Normal	Strong	Yes	
Sunny	Mild	High	Weak	No	
Sunny	C	Normal	Weak	Yes	
Rainy	Mild	Normal	Weak	Yes	
Sunny	Mild	Normal	Strong	Yes	
overcast	Mild	High	Strong	Yes	
overcast	M	Normal	weak	Yes	
Rainy	Mild	High	Strong	No	

ANS Entropy = $-\sum P(I) \log_2 P(I)$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \frac{\text{MAXIMISE THIS}}{\text{TO MINIMISE SCATTERING}} \geq P(S|A) \text{Entropy}(S|A)$$

$$\text{Entropy}(\text{Decision}) = \frac{5}{14} \log_2 \frac{5}{14} + \frac{9}{14} \log_2 \frac{9}{14}$$

\downarrow

$$P(\text{No}) = 0.940$$

Gain (Decision, Wind)

$$= \text{Entropy}(\text{Decision}) - \sum P(\text{Wind}) \text{Entropy}(D|\text{Wind})$$

$$= \text{Entropy}(D) - P(\text{Weak}) \text{Entropy}(D/\text{weak}) - P(\text{Strong}) \text{Entropy}(D/\text{strong})$$

$$= \text{Entropy}(D) - \frac{8}{14} \left[\frac{2}{8} \log \frac{2}{8} + \frac{6}{8} \log \frac{6}{8} \right]$$

$$- \frac{6}{14} \left[\frac{3}{6} \log \frac{3}{6} + \frac{3}{6} \log \frac{3}{6} \right]$$

$$= 0.940 - \frac{8}{14} \left[\frac{2}{8} \log \frac{2}{8} + \frac{6}{8} \log \frac{6}{8} \right]$$

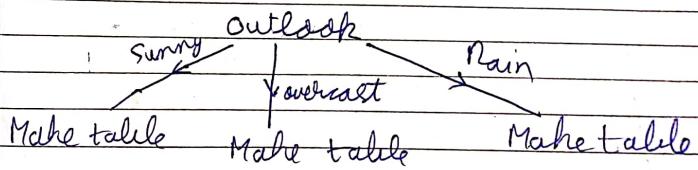
$$- \frac{6}{14} \log \frac{3}{6}$$

→ Gain (Decision, outlook),

Gain (Decision, Humid)

Gain (Decision, Temp) is calculated similarly

→ Gain (Decision, outlook) is maximum



Page No.	
Date	

Table for Sunny

Temp	Humid	Wind	Precision
H	High	Weak	No
H	High	Strong	No
Mild	High	Weak	No
C	Normal	Weak	Yes
Mild	Normal	Strong	Yes