

Unit - 2 : Decision Theory.

Bayesian Decision Theory:

Bayes Theorem

→ named after british mathematician Thomas Bayes,

Is a mathematical formula for determining
conditional probability.

→ where it is used? → Trying to find answer.

ie, there are only 2 possible events for given question:

A : It will rain tomorrow.
B : It will not rain tomorrow. } 50% chance to
rain tomorrow if it's raining today.

50% chance to rain tomorrow if it's raining today.

→ Probability theory is framework for making decisions under uncertainty.

→ In classification, baye's rule is used to calculate the probabilities of classes.

→ Thus, bayesian decision theory , is based on the existence of prior distributions of parameters .

Example:

- * Tossing a coin is a random process because we cannot predict at any toss whether the outcome will be head or tail - that's why we toss coin, or buy lottery tickets.
 - * We can only talk about the probability that the outcome of next toss will be heads or tails.
 - * extra piece of knowledge that we don't have access to are named unobservable variables.
 - * In coin tossing example, only observable variable is the outcome of the toss.

$$X = f(z)$$

where $f(\cdot)$ is deterministic function that defines the outcome from the unobservable piece of knowledge.

Such 'x' are bernoulli distributed where the parameter of distribution P_0 is the probability that outcome is head:

$$\Phi(x=1) = \phi_0$$

$$P(x=0) = 1 - P(x=1)$$

$$= 1 - P_0$$

Assume we want to predict the outcome of next toss if we know p_0 ,

prediction will be heads if $P_0 > 0.5$,
otherwise tails

x^t is random variable 8

$x^t = 1$ if outcome of toss 't' is head.
 $x^t = 0$ " "

Sample = {H, H, H, T, H, T, T, H, H} tail.

$$x = \{1, 1, 1, 0, 1, 0, 0, 1, 1\}$$

$$P_0 = \frac{\sum_{t=1}^N x_t}{N} = \frac{6}{9} = 0.66$$

$$P_0 = \frac{\text{No. of faces showing heads}}{\text{total No. of faces}}$$

Baye's formula:

- also called as baye's rule or baye's law.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) * P(A)}{P(B)} \rightarrow \text{prior}$$

posterior

Likelihood

Evidence

Eg: Bill is 35 years old & earns \$40000/yr. He has a very fair credit rating. Will he buy a computer?

$$\text{e.g., } P(H|x) = \frac{P(x|H) \cdot P(H)}{P(x)}$$

Where,

x = bill is 35yr old with fair credit rating & income of 40000\$ / yr.

H = Hypothesis that bill will buy the computer.

$P(H|x)$ = prob. that bill will buy computer given that his age, income & credit rating. (posterior)

$P(H)$ = prob. that bill will buy computer (regardless of knowing his age, income etc.) (prior)

$P(x|H)$ = prob. that someone is 35yr old, has fair credit rating, earns 40,000\$/yr & has bought the computer [Likelihood]

$P(x)$ = prob. that bill is 35yr old, has fair credit rating & earns 40000/yr [Evidence]

Bayesian classifiers:

- based on baye's theorem.
- are statistical classifiers, used to find the prob. that a given tuple belong to particular class.
- exhibits high accuracy & speed when applied to large DB.

$P(x|c_i)$ - is prob. of seeing x as E/P when it is known to belong to class c_i .

The posterior probability of class c_i is :

$$\begin{aligned} P(c_i|x) &= \frac{P(x|c_i)P(c_i)}{P(x)} \\ &= \frac{P(x|c_i)P(c_i)}{\sum_{k=1}^K P(x|c_k)P(c_k)} \end{aligned}$$

↓
chain rule

Baye's classifier : Prerequisites:

Random variable :

- is a function that maps a possible set of outcome to some values.

Eg: Tossing a coin so getting head H as 1 &
 " tail T as 0
 where 0 & 1 are random variable.

Prior (or) state of Nature :

- is how likely each class is going to occur.
- priors are known before training process.

Class conditional probability or likelihood :

- prob. of how likely a feature x occurs given that it belongs to particular class.
- denoted by $P(x|A)$ where x is particular feature.
- is the quantity that we have to evaluate while training the data.

Evidence :

- prob. of occurrence of particular feature. $P(x)$.
- can be calculated using chain rule,

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

Posterior probability :

- prob. of occurrence of class A when certain features are given.
- It is what we aim at computing in test phase in which we have testing T/P or features.

Naïve Bayes classification:

- is a program which predicts a class value, given a set of attributes.
- use baye's rule to derive conditional probabilities for class variable.
- Finally output the class with highest probability.

- The dataset is divided into two parts:
 - ↳ feature matrix
 - ↳ response vector.
- feature matrix, contains all the vectors (rows) of dataset in which each vector consists of value of dependent feature.
- response vector, contains the value of class variable (prediction or output) for each row of feature matrix.

Example:

- 1) Given DB: Training Data from AllElectronics Customer DB:
 [Check whether buys-computer = Yes or No based on given condition]

ID	age	Income	Student	credit-rating	class: buys.computer
1	Youth	high	no	fair	no
2	"	"	no	excellent	"
3	middle-age	"	"	fair	yes
4	Senior	medium	"	"	"
5	"	low	yes	"	"
6	"	"	"	excellent	"
7	middle-age	"	"	"	no
8	Youth	medium	no	fair	yes
9	"	low	yes	"	no
10	Senior	medium	yes	"	"
11	Youth	"	"	"	yes
12	middle-age	"	"	excellent	"
13	middle-age	high	no	"	"
14	Senior	medium	yes	fair	"
			no	excellent	no

$$x = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$$

Now,

$$P(\text{buys-comp} = \text{yes}) = \frac{9}{14} \rightarrow 9 \text{ yes}$$

$$= \frac{9}{14} \rightarrow \text{total (yes+no)}$$

$$= \underline{\underline{0.643}}$$

$$P(\text{buys-comp} = \text{no}) = \frac{5}{14} \rightarrow 5 \text{ no}$$

$$= \underline{\underline{0.357}}$$

To compute $P(x|C_i), \forall i = 1, 2, \dots$

$$P(\text{age} = \text{youth} | \text{buys_comp} = \text{yes}) = 2/9 = 0.222$$

$$P(\text{age} = \text{youth} | \text{buys_comp} = \text{no}) = 3/5 = 0.600$$

$$P(\text{income} = \text{medium} | \text{buys_comp} = \text{yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} | \text{buys_comp} = \text{no}) = 2/5 = 0.400$$

$$P(\text{student} = \text{yes} | \text{buys_comp} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{student} = \text{yes} | \text{buys_comp} = \text{no}) = 1/5 = 0.200$$

$$P(\text{credit_rating} = \text{fair} | \text{buys_comp} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{credit_rating} = \text{fair} | \text{buys_comp} = \text{no}) = 2/5 = 0.400$$

Using these,

$$\begin{aligned} P(x | \text{buys_comp} = \text{yes}) &= P(\text{age} = \text{youth} | \text{buys_comp} = \text{yes}) \times \\ &\quad P(\text{income} = \text{medium} | \text{buys_comp} = \text{yes}) \times \\ &\quad P(\text{student} = \text{yes} | \text{buys_comp} = \text{yes}) \times \\ &\quad P(\text{credit_rating} = \text{fair} | \text{buys_comp} = \text{yes}) \\ &= 0.222 \times 0.444 \times 0.667 \times 0.667 \\ &= 0.044 \end{aligned}$$

Similarly,

$$\begin{aligned} P(x | \text{buys_comp} = \text{no}) &= 0.600 \times 0.400 \times 0.200 \times 0.400 \\ &= 0.019 \end{aligned}$$

To find the class C_i , maximize $P(x | C_i) \cdot P(C_i)$,

$$P(x | \text{buys_comp} = \text{yes}) \cdot P(\text{buys_comp} = \text{yes}) = 0.044 \times 0.643 = 0.028$$

$$P(x | \text{buys_comp} = \text{no}) \cdot P(\text{buys_comp} = \text{no}) = 0.019 \times 0.357 = 0.007$$

Among both
[yes] has a
more probab

\therefore the naive bayesian classifier predicts

buys-computer = yes for tuple x

2) Training Example : Play Tennis

Day	outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	"	"	"	Strong	"
D3	Overcast	"	"	Weak	Yes
D4	Rain	Mild	"	"	"
D5	"	Cool	Normal	"	"
D6	"	"	"	Strong	No
D7	Overcast	"	"	"	Yes
D8	Sunny	Mild	High	Weak	No
D9	"	Cool	Normal	"	Yes
D10	Rain	Mild	"	"	"
D11	Sunny	"	"	Strong	"
D12	Overcast	"	High	"	"
D13	"	Hot	Normal	Weak	"
D14	Rain	Mild	High	Strong	No

Given instance,

$$x = (\text{outlook} = \text{Sunny}, \text{Temperature} = \text{cool}, \text{Humidity} = \text{High}, \text{wind} = \text{Strong})$$

Now,

$$P(\text{Play} = \text{Yes}) = 9/14 = 0.643$$

$$P(\text{Play} = \text{No}) = 5/14 = 0.357$$

To compute $P(x|c_i)$ for $i=1, 2, \dots$

$$P(\text{outlook} = \text{Sunny} | \text{Play} = \text{Yes}) = 2/9 = 0.222$$

$$P(\text{"} | \text{Play} = \text{No}) = 3/5 = 0.600$$

$$P(\text{Temp} = \text{cool} | \text{Play} = \text{Yes}) = 3/9 = 0.333$$

$$P(\text{Temp} = \text{cool} | \text{Play} = \text{No}) = 1/5 = 0.200$$

$$P(\text{Humidity} = \text{high} | \text{Play} = \text{Yes}) = 3/9 = 0.333$$

$$P(\text{"} | \text{Play} = \text{No}) = 4/5 = 0.800$$

$$P(\text{wind} = \text{Strong} | \text{Play} = \text{Yes}) = 3/9 = 0.333$$

$$P(\text{"} | \text{Play} = \text{No}) = 3/5 = 0.600$$

Testing phase:

$$\begin{aligned}
 P(\text{Yes} | x) &= [P(\text{Sunny} | \text{Yes}) \times P(\text{cool} | \text{Yes}) \times P(\text{High} | \text{Yes}) \times \\
 &\quad P(\text{Strong} | \text{Yes})] * P(\text{Play} = \text{Yes}) \\
 &= [0.222 \times 0.333 \times 0.333 \times 0.333] * 0.643 \\
 &= 0.008 * 0.643 \\
 &= \boxed{0.005}
 \end{aligned}$$

Similarly for,

$$\begin{aligned}
 P(\text{No} | x) &= [0.600 \times 0.200 \times 0.800 \times 0.600] * 0.357 \\
 &= 0.072 * 0.357 \\
 &= \boxed{0.025}
 \end{aligned}$$

$$\therefore P(\text{Yes} | x) < P(\text{No} | x).$$

Naive Bayes classifier predicts [play-tennis = 'No'],

Advantage:

- fast for training & prediction.
- easily interpretable.
- can handle both continuous & discrete data.

Disadvantage:

- data scarcity & chances of loss of accuracy.
- Zero frequency.
i.e. if the category of any variable is not seen in training dataset then model assigns 0% zero prob. to that category & prediction cannot be made.

Applications:

- Text classification.
- Sentiment Analysis
- Recommendation System.

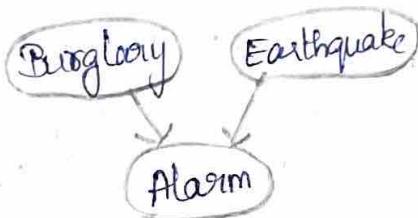
Classification:

Bayesian Network:

- capture joint prob of events represented by model.
- bayesian belief n/w describe the joint prob distribution for set of variables.
- Naive bayes is restricted form of bayesian n/w where nodes should have no parents & ~~edges~~ node corresponding to attribute variable have no edges b/w them.
- BN is made up of two parts:
 1. directed acyclic graph.
 2. set of parameters.

Directed Acyclic graph:

- nodes are random variables.
(can be discrete or continuous)
- arrows connect pair of nodes.
(x is parent of y if there is an arrow from x to y)



Set of parameters:

- parameters are prob. in these conditional prob. distribution
- each node X_i has conditional prob. distribution
 $P(X_i | \text{Parents}(X_i))$ that quantifies effect of parents on node

burglary

B	$P(B)$
F	0.999
T	0.001

earthquake

E	$P(E)$
F	0.998
T	0.002

Alarm

B	E	A	$P(A B,E)$	$P(A B,E)$
F	F	F	0.999	
F	F	T	0.001	
F	T	F	0.71	
F	F	F	0.06	
T	F	T	0.29	
T	F	F	0.94	
T	T	F	0.05	
T	T	T	0.95	

Conditional probability Tables.

Bayesian Belief network:

- "A bayesian n/w is a probabilistic graphical model which represents a set of variables & their conditional dependences using directed acyclic graph".
- Also called as bayes n/w, belief n/w, decision n/w or bayesian model.

Joint probability Distribution:

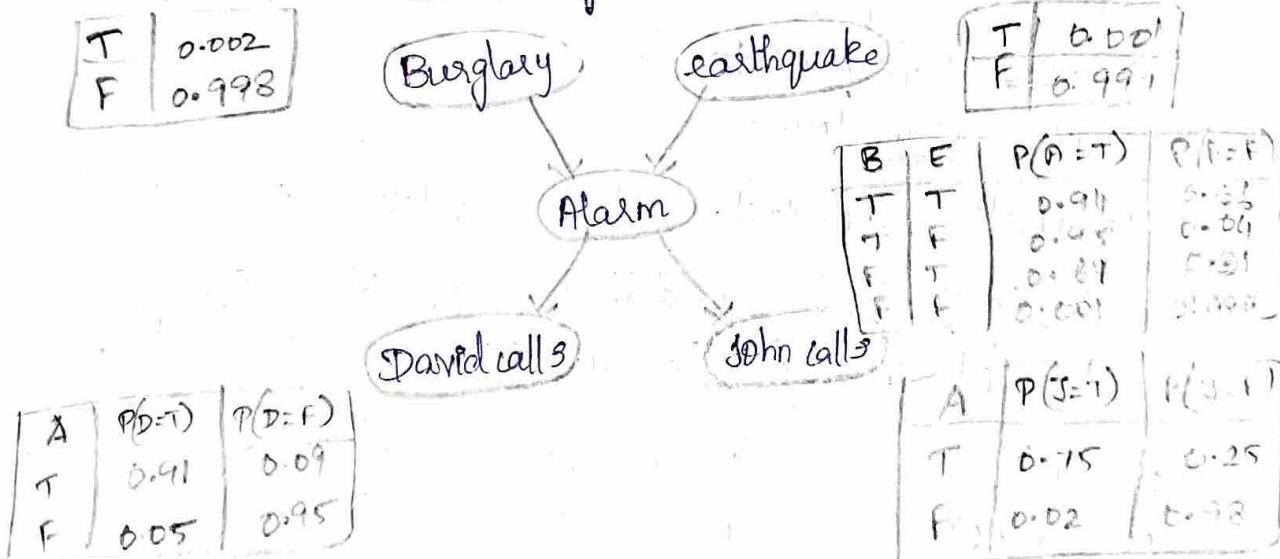
* If we have variables of x_1, x_2, \dots, x_n , then prob. of diff. combination of x_1, x_2, \dots, x_n is known as joint prob. distribution.

$$P[x_1, x_2, \dots, x_n] = P[x_1 | x_2, x_3, \dots, x_n] \cdot P[x_2, x_3, \dots, x_n]$$

$$= P[x_1 | x_2, x_3, \dots, x_n] \cdot P[x_2 | x_3, x_4, \dots, x_n] \cdots P[x_{n-1} | x_n] P[x_n]$$

Example:

① Calculate the prob. that alarm has sounded, but there is neither burglary nor earthquake occurred, & david, John both called the Harry.



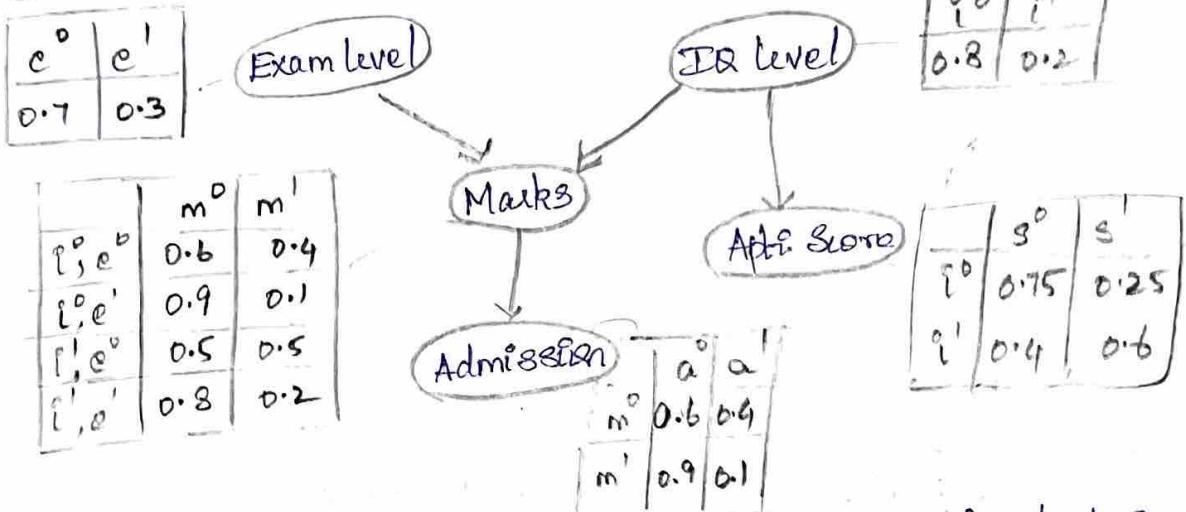
List of all events in this n/w:

- Burglary (B) • earthquake (E)
- Alarm (A) • David calls (D) • John calls (J)

$$P(J, D, A, \neg B, \neg E) = ?$$

$$\begin{aligned}
 P(J, D, A, \neg B, \neg E) &= P(J|A) * P(D|A) * P(A | \neg B \wedge \neg E) * P(\neg B) * P(\neg E) \\
 &= 0.75 * 0.91 * 0.001 * 0.998 * 0.999 \\
 &= \underline{\underline{0.00068045}}
 \end{aligned}$$

②



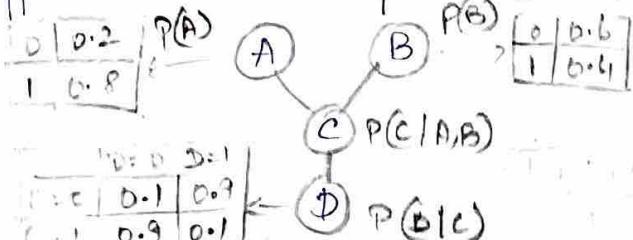
- i. Calculate the prob. that despite of exam level being difficult, the student having a low IQ & low Aptitude score, manages to pass the exam & secure admission.

$$\begin{aligned}
 P[A=1, M=1, I=0, E=1, S=0] &= ? \\
 &= P(a^1 | m^1) \times P(m^1 | i^0 e^0) \times P(i^0) \times P(e^0) \\
 &= 0.1 \times 0.1 \times 0.8 \times 0.3 \times 0.75 \\
 &= \underline{\underline{0.0018}}
 \end{aligned}$$

- ii. Calculate the prob that the student has a high IQ level & Aptitude score, the exam being easy yet fails to pass & doesn't secure admission.

$$\begin{aligned}
 P[A=0, M=0, I=1, E=0, S=1] &= ? \\
 &= P(a^0 | m^0) \times P(m^0 | i^1 e^0) \times P(i^1) \times P(e^0) \\
 &= 0.6 \times 0.5 \times 0.2 \times 0.7 \times 0.6 \\
 &= \underline{\underline{0.0252}}
 \end{aligned}$$

- 3) Suppose there are four boolean random variables:



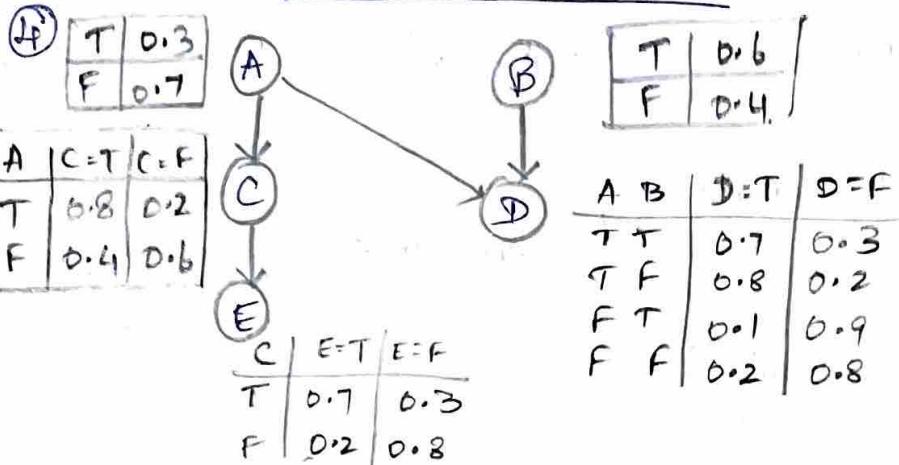
A	B	C = 0	C = 1
0	0	0.2	0.8
0	1	0.6	0.4
1	0	0.3	0.7
1	1	0.5	0.5

Joint prob. distribution of random variable is given by:

$$P(A, B, C, D) = P(A) \cdot P(B) \cdot P(C|A, B) \cdot P(D|C)$$

i) Find $P(A=0, B=1, C=0, D=1)$

$$\begin{aligned}
 &= P(A=0) * P(B=1) * P(C=0 | A=0, B=1) * P(D=1 | C=0) \\
 &= 0.2 * 0.4 * 0.6 * 0.9 \\
 &= \underline{\underline{0.0432}}
 \end{aligned}$$



Given : $P(A=T) = 0.3$, $P(C=T | A=T) = 0.8$

$$P(B=T) = 0.6, P(C=T | A=F) = 0.4$$

$$P(D=T | A=T, B=T) = 0.7$$

$$P(D=T | A=T, B=F) = 0.8$$

$$P(D=T | A=F, B=T) = 0.1$$

$$P(D=T | A=F, B=F) = 0.2$$

$$P(E=T | C=T) = 0.7$$

$$P(E=T | C=F) = 0.9$$

5, $P(D=T) = ?$

$$\begin{aligned}
 &= P(D=T | A=T, B=T) + P(D=T | A=T, B=F) + P(D=T | A=F, B=T), \\
 &\quad P(D=T | A=F, B=F)
 \end{aligned}$$

$$\begin{aligned}
 &= P(D=T | A=T, B=T) * P(A=T, B=T) + P(D=T | A=T, B=F) * P(A=T, B=F) \\
 &\quad + P(D=T | A=F, B=T) * P(A=F, B=T) + P(D=T | A=F, B=F) * P(A=F, B=F)
 \end{aligned}$$

$$\begin{aligned}
 &= P(D=T | A=T, B=T) * P(A=T) * P(B=T) + P(D=T | A=T, B=F) * P(A=T) * P(B=F) \\
 &\quad + P(D=T | A=F, B=T) * P(A=F) * P(B=T) + P(D=T | A=F, B=F) * P(A=F) * P(B=F)
 \end{aligned}$$

$$= (0.7 * 0.3 * 0.6) + (0.8 * 0.3 * 0.4) + (0.1 * 0.7 * 0.6) + (0.2 * 0.7 * 0.4)$$

$$= \underline{\underline{0.32}}$$

6, $P(A=T | C=T) = ?$

$$= \frac{P(C=T | A=T) * P(A=T)}{P(C=T)}$$

$$now, P(C=T) = P(C=T, A=T) + P(C=T, A=F)$$

$$= P(C=T | A=T) * P(A=T) + P(C=T | A=F) * P(A=F)$$

$$= (0.8 * 0.3) + (0.4 * 0.7)$$

$$= 0.52$$

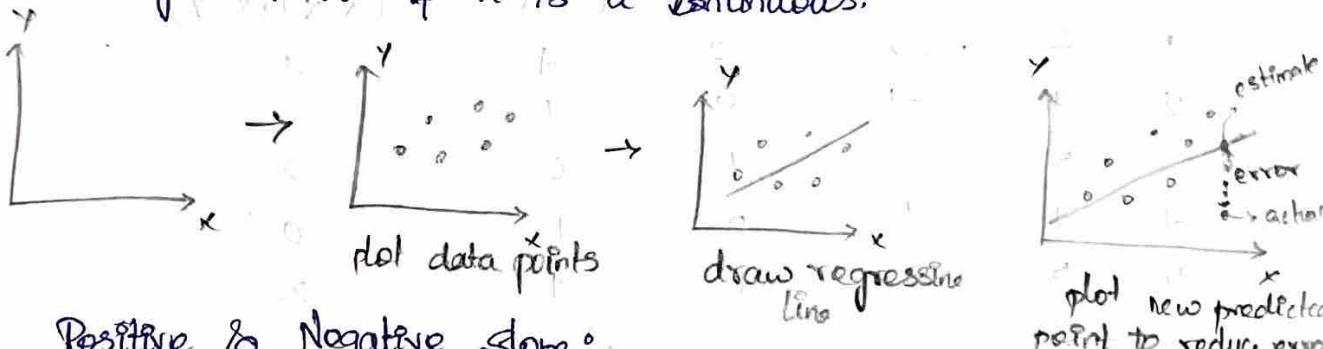
$$\frac{P(C=T | A=T) \cdot P(A=T)}{P(C=T)} = \frac{0.8 \times 0.3}{0.52} = \underline{\underline{0.46}}$$

Regression:

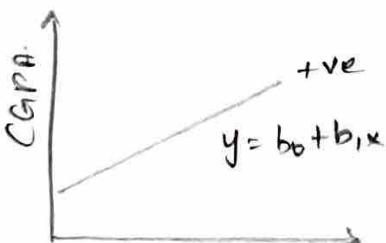
- is a statistical modeling to show the relationship b/w two variables with linear regression.
- is form of predictive modeling technique which investigate the relationship b/w a dependent & independent variable.

Simple Linear Regression:

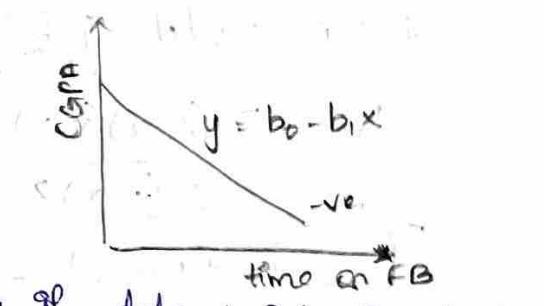
- data is modeled using straight line.
- $y = b_1 x + b_0$ is a straight line. y where
 - x = independent variable.
 - y = dependent variable.
 - b_1 = slope
 - b_0 = y -intercept.
- every value of x there is a corresponding value of y if it is a continuous.



Positive & Negative slope:



- If data point in x axis increases then increase in y axis value is called a positive slope.



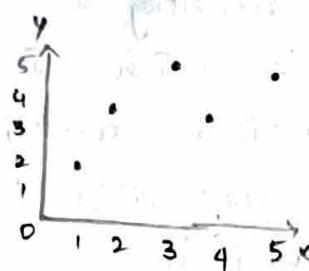
- If data point in x increasing & dependent variable y decreasing is called negative slope.

Example:
1) Given:

independent x	dependent y
1	2
2	4
3	5
4	4
5	5

Find the regression line
& R^2 co-efficient.

Step 1: plot data



Step 2: calculate Slope.

Straight line formula: $y = b_0 + b_1 x$

where slope $b_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$

and $\bar{x} = \text{mean of } x$

$\bar{y} = \text{mean of } y$

$$\bar{x} = \frac{1+2+3+4+5}{5} = [3]$$

$$\bar{y} = \frac{2+4+5+4+5}{5} = [4]$$

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x})(y - \bar{y})$
1	2	-2	-2	4	4
2	4	-1	0	1	0
3	5	0	1	0	0
4	4	1	0	1	0
5	5	2	1	4	2
$\bar{x} = 3$ $\bar{y} = 4$				$\sum = 10$	$\sum = 6$

Linear model: $\hat{y} = b_0 + b_1 x$

$$b_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{6}{10} = [0.6]$$

$$\therefore \hat{y} = b_0 + 0.6x \rightarrow ①$$

$$\text{Now, at mean } (3, 4), ① \Rightarrow b_0 = \bar{y} - 0.6\bar{x} \\ = 4 - 0.6(3)$$

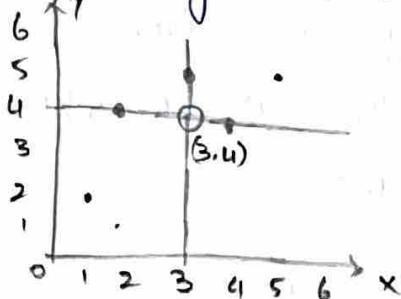
$$[b_0 = 2.2]$$

Sub b_0 in equ ①.

$\therefore \hat{y} = 0.6x + 2.2$, ps a linear model equation.

83:

Find regression line.



At mean (3, 4)

given $x = \{1, 2, 3, 4, 5\}$

$$y = 0.6x + 2.2$$

Sub x in y equ.

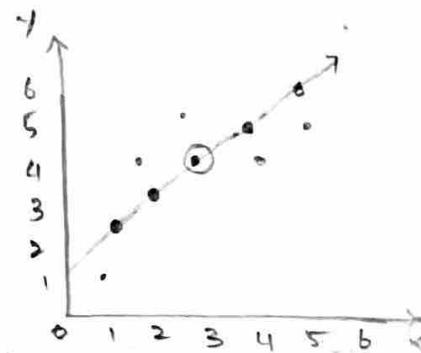
$$y = 0.6(1) + 2.2 = 2.8$$

$$y = 0.6(2) + 2.2 = 3.4$$

$$y = 0.6(3) + 2.2 = 4$$

$$y = 0.6(4) + 2.2 = 4.6$$

$$y = 0.6(5) + 2.2 = 5.2$$



∴ Regression line points are:

$$(1, 2.8) (2, 3.4) (3, 4) (4, 4.6) (5, 5.2)$$

Mean Squared error:

* task ps find the distance b/w actual & predicted value
to reduce the error or reduce the distance.

* line with least error ps line of regression.

* R^2 value is a statistical measure of how close the data to be fitted in regression line

$$R^2 = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2}$$

where,

\hat{y} = predicted value of y

\bar{y} = mean value of y .

∴ from the given points,

$$\begin{aligned} R^2 &= 1 - \frac{(2-2.8)^2 + (4-3.4)^2 + (5-4)^2 + (4-4.6)^2 + (5-5.2)^2}{(2-4)^2 + (4-4)^2 + (5-4)^2 + (4-4)^2 + (5-4)^2} \\ &= 1 - \frac{2.4}{6} \Rightarrow 1 - 0.4 \end{aligned}$$

$$R^2 = 0.6$$

Condition:

If $R^2 \approx 0.3$ \Rightarrow data points are far away from regression line.

If $R^2 \approx 0.7$ \Rightarrow data points are closer to the regression line.

Hence for the given values the $R^2 = 0.6$ which is closer to regression line.

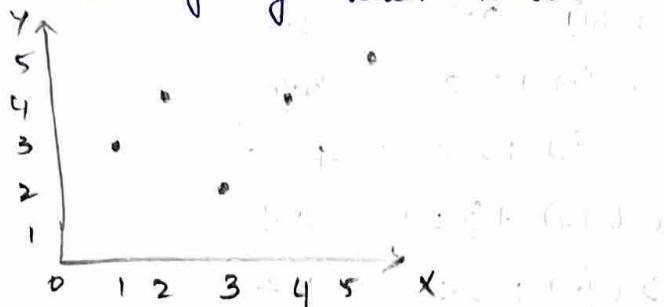
2) Given:

x	1	2	3	4	5
y	3	4	2	4	5

a) find least square regression line $y = ax + b$

b) Estimate value of y when $x = 10$

S1:



S2:

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

$$\bar{y} = \frac{3+4+2+4+5}{5} = 3.6$$

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x})(y - \bar{y})$
1	3	-2	-0.6	4	-1.2
2	4	-1	0.4	1	-0.4
3	2	0	-1.6	0	0
4	4	1	0.4	1	0.4
5	5	2	1.4	4	2.8
$\bar{x} = 3$ $\bar{y} = 3.6$				$\sum = 10$	$\sum = 4$

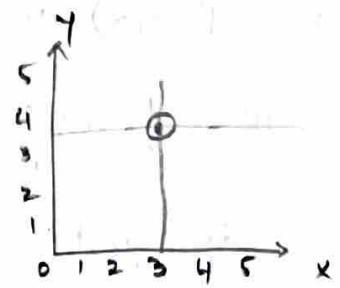
$$\therefore \bar{y} = ax + b$$

$$a = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{4}{10} = 0.4$$

$$\therefore y = 0.4x + b \rightarrow ①$$

Now, at mean $(3, 3.6)$

$$\begin{aligned} \textcircled{1} \text{ equ } \Rightarrow b &= y - 0.4 \cdot x \\ &= 3.6 - 0.4(3) \\ \boxed{b = 2.4} \end{aligned}$$

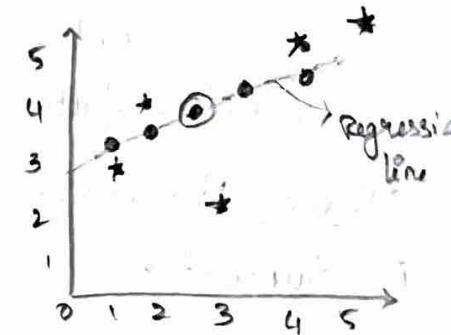


Thus, $y = 0.4x + 2.4$ is a linear model

Q3: find regression line,

$$x = \{1, 2, 3, 4, 5\}$$

$$\begin{aligned} y &= 0.4(1) + 2.4 \Rightarrow 2.8 \\ y &= 0.4(2) + 2.4 \Rightarrow 3.2 \\ y &= 0.4(3) + 2.4 \Rightarrow 3.6 \\ y &= 0.4(4) + 2.4 \Rightarrow 4.0 \\ y &= 0.4(5) + 2.4 \Rightarrow 4.4 \end{aligned}$$



∴ points for Regression line are } $(1, 2.8), (2, 3.2), (3, 3.6), (4, 4.0), (5, 4.4)$

(b) when $x = 10$,

$$y = 0.4(10) + 2.4$$

$$\boxed{y = 6.40}$$

Association Mining:

• finding regularities in data.

i.e. what products were often purchased together.

* Frequent patterns are patterns that appear frequently in a data set

↳ Set of items, such as milk & bread, that appear frequently together in a transaction dataset is a frequent itemset.

Association Rules:

• is an implication of the form $x \rightarrow y$ where 'x' is an antecedent & 'y' is consequent of rule.

• used to find dependency b/w two items x & y.

Eg: $P(y|x) \rightarrow$ is prob. that somebody who buys x also buys y .

i.e. 90% of transaction that purchase bread & butter also purchase milk. $(x \rightarrow y)$

"if" part = antecedent [before \rightarrow part i.e., x]

"then" part = consequent [after \rightarrow part i.e., y]

Here antecedent & consequent are disjoint.

(i.e., have no items in common)
Antecedent = bread & butter

Consequent = milk

Confidence factor = 90%.

Three Measures:

① Support:

$$\text{support}(x,y) = P(x,y) = \frac{\text{No. of cust. who bought } x \text{ and } y}{\text{No. of customers}}$$

② Confidence:

$$\text{Conf}(x \rightarrow y) = P(y|x) = \frac{P(x,y)}{P(x)} \xrightarrow{\text{support}}$$

$$\therefore \text{Conf}(x \rightarrow y) = \frac{\text{No. of customers who bought } x \text{ & } y}{\text{No. of customers who bought } x}$$

③ Lift (or) Interest:

$$\text{Lift}(x \rightarrow y) = \frac{P(x,y)}{P(x) \cdot P(y)} = \frac{P(y|x)}{P(x)}$$

Market basket Analysis:

- is the earliest form of frequent pattern mining & association rules.

↳ Apriori Alg.:

- proposed by Agrawal et.al in 1996.

2 Steps:

- finding frequent itemset. i.e. those which have enough support
- Converting them to rules with enough confidence.

Concepts:

- An Item : is an item/article in basket.
- I : the set of all items sold in store.
- transaction : items purchased in basket, it may have \downarrow transaction ID
- transactional dataset : set of all transactions.

Example:

- 1) find all frequent itemset using Apriori with min. support of 2 and confidence of 50% for following transactions:

TID	Set of Items ID's
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

1-item Set:

1-item Set	freq.
I1	6
I2	7
I3	6
I4	2
I5	2

→ take only items
with freq. ≥ 2
i.e., min. supp = 2

Freq. 1-item Set	freq.
I1	6
I2	7
I3	6
I4	2
I5	2

2-item Set:

2-item Set	freq.
I1, I2	4
I1, I3	4
I1, I4	1
I1, I5	2
I2, I3	4
I2, I4	2
I2, I5	2
I3, I4	0
I3, I5	1
I4, I5	0

→

Freq. 2-item Set	freq.
I1, I2	4
I1, I3	4
I1, I5	2
I2, I3	4
I2, I4	2
I2, I5	2

3-Item Set:

3-Item Sets	freq.
I ₁ , I ₂ , I ₃	2
I ₁ , I ₂ , I ₄	1
I ₁ , I ₂ , I ₅	2
I ₁ , I ₃ , I ₄	0
I ₁ , I ₃ , I ₅	1
I ₁ , I ₄ , I ₅	0
I ₂ , I ₃ , I ₄	0
I ₂ , I ₃ , I ₅	1
I ₂ , I ₄ , I ₅	0
I ₃ , I ₄ , I ₅	0



freq. 3-item sets	freq.
I ₁ , I ₂ , I ₃	2
I ₁ , I ₂ , I ₅	2

4-Item Set:

4-Item Set	freq.
I ₁ , I ₂ , I ₃ , I ₅	1



Not possible

Now,

freq. 3-Item set :

$$I = \{1, 2, 3\} \cup \{1, 2, 5\}$$

Non-empty subsets 'S' are :

- {1} {2} {3} {1, 2} {1, 3} {2, 3}
- {1} {2} {5} {1, 2} {1, 5} {2, 5}

How to form Association Rule?

for every non-empty subset 'S' of I,

$$S \rightarrow (I - S)$$

If $\text{Supp}(I) / \text{Supp}(S) \geq \text{min_conf(50%)}$

① for Item Set I : {1, 2, 3}, Subsets: {1}, {2}, {3}, {1, 2}, {1, 3}, {2, 3}

Rule 1: {1} \rightarrow {2, 3}

$$S \rightarrow (I - S)$$

$$\text{Supp} = \frac{2}{9} \quad (1, 2, 3) \text{ occur 2 times}$$

Total Transaction

$$= 22.22\%$$

$$\text{Conf} = \frac{\text{Supp}(1, 2, 3)}{\text{Supp}(1)} = \frac{2/9}{6/9} = \frac{2}{6} = \boxed{33.33\%} < 50\%$$

\therefore Invalid Rule.

Rule 2: $\{2\} \rightarrow \{1, 3\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \frac{\text{Supp}(1, 2, 3)}{\text{Supp}(2)} = \frac{2}{7} = \boxed{28.57\%} < 50\% \quad \times$$

\therefore Invalid Rule.

Rule 3: $\{3\} \rightarrow \{1, 2\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \frac{\text{Supp}(1, 2, 3)}{\text{Supp}(3)} = \frac{2}{6} = \boxed{33.33\%} < 50\% \quad \times$$

\therefore Invalid Rule.

Rule 4: $\{1, 2\} \rightarrow \{3\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \frac{\text{Supp}(1, 2, 3)}{\text{Supp}(1, 2)} = \frac{2}{4} = \boxed{50\%} \geq 50\% \quad \checkmark$$

\therefore Valid Rule.

Rule 5: $\{1, 3\} \rightarrow \{2\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \frac{\text{Supp}(1, 2, 3)}{\text{Supp}(1, 3)} = \frac{2}{4} = \boxed{50\%} \geq 50\% \quad \checkmark$$

\therefore Valid Rule.

Rule 6: $\{2, 3\} \rightarrow \{1\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \frac{\text{Supp}(1, 2, 3)}{\text{Supp}(2, 3)} = \frac{2}{4} = \boxed{50\%} \geq 50\% \quad \checkmark$$

\therefore Valid Rule.

(2) For Itemset 2: $\{1, 2, 3\}$, Subsets: $\{\{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}\}$

Rule 1: $\{1\} \rightarrow \{2, 3\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \text{Supp}(1, 2, 3) / \text{Supp}(1) \Rightarrow 2/6 = \boxed{33.33\%} < 50\% \quad \times$$

\therefore Invalid Rule.

Rule 2: $\{2\} \rightarrow \{1, 3\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \text{Supp}(1, 2, 3) / \text{Supp}(2) \Rightarrow 2/7 = \boxed{28.57\%} < 50\% \quad \times$$

\therefore Invalid Rule.

Rule 3: $\{5\} \rightarrow \{1, 2\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \text{Supp}(1, 2, 5) / \text{Supp}(5) = 2/2 \Rightarrow [100\%] > 50\% \quad \checkmark$$

\therefore Valid Rule.

Rule 4: $\{1, 2\} \rightarrow \{5\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \text{Supp}(1, 2, 5) / \text{Supp}(1, 2) = 2/4 = [50\%] \geq 50\% \quad \checkmark$$

\therefore Valid Rule.

Rule 5: $\{1, 3\} \rightarrow \{2\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \text{Supp}(1, 2, 3) / \text{Supp}(1, 3) = 2/4 = [50\%] \geq 50\% \quad \checkmark$$

\therefore Valid Rule.

Rule 6: $\{2, 3\} \rightarrow \{1\}$

$$\text{Supp} = 2/9 = 22.22\%$$

$$\text{Conf} = \text{Supp}(1, 2, 3) / \text{Supp}(2, 3) = 2/3 = [66.67\%] > 50\% \quad \checkmark$$

\therefore Valid Rule.

Hence all the valid Rules are the association rule for the given transaction.

- Q2) find the frequent itemset for following Transaction DB with min. Supp = 50%.

Tid	Items
10	A, C, D
20	B, C, E
30	A, B, C, E
40	B, E

Since total Tid = 4 & min. support

$$\therefore \text{min. Supp} = 2$$

i.e., half of total Tid
Also since min. Conf. is not given, we take all possible association rules.

1-item Set:

1-item set	freq
A	2
B	3
C	3
D	1
E	3



freq. 1-item set	freq
A	2
B	3
C	3
E	3

2-item sets

2-item set	freq.
{A, B}	1
{A, C}	2
{A, E}	1
{B, C}	2
{B, E}	3
{C, E}	2



freq. 2-item set	freq.
{A, C}	2
{B, C}	2
{B, E}	3
{C, E}	2

3-item set:

3-item set	freq.
A, C, E	1
A, B, E	1
B, C, E	2



freq. 3-item set	freq.
B, C, E	2

Now possible Association Rule for {B, C, E} are:

$$I = \{B, C, E\}, S = \{B\} \{C\} \{E\} \{B, C\} \{B, E\} \{C, E\}$$

$$\therefore [S \rightarrow I - S]$$

i.e., $B \rightarrow \{C, E\}; C \rightarrow \{B, E\}; E \rightarrow \{B, C\}$

$$\{B, C\} \rightarrow E; \{B, E\} \rightarrow C; \{C, E\} \rightarrow B$$

Losses And Risks:

- caused due to decisions that are not equally good or only

Eg: A financial institution when making a decision for loan applicant should take into account the potential gain & loss as well.

Actions: α_i - as decision to assign the Q/P to class C_i

Loss: λ_{ik} - is loss incurred for taking action α_i when the Q/P actually belongs to C_k .

\therefore Expected risk for taking action α_i e.g.:

$$R(\alpha_i | x) = \sum_{k=1}^K \lambda_{ik} P(C_k | x)$$

Decision: Choose α_i with minimal expected risk.

i.e. choose α_i , $R(\alpha_i | x) = \min_k R(\alpha_k | x)$

Loss & Risk : 0/1 Loss:

Correct decision : 0 loss

Incorrect decision : 1 loss.

Let us define 'k' actions α_i , where $i=1, 2, \dots, k$.

Special case of 0/1 loss is:

$$\lambda_{ik} = \begin{cases} 0, & \text{if } i=k \\ 1, & \text{if } i \neq k. \end{cases}$$

(Correct decision has no loss)

Now, risk for taking action α_i is:

$$R(\alpha_i | x) = \sum_{k=1}^k \lambda_{ik} P(C_k | x)$$

$$= \sum_{k \neq i} 1 \cdot P(C_k | x)$$

• If no loss then risk is 0 since $\lambda_{ik} = 0$
• If $k \neq i$ then $\lambda_{ik} = 1$

$$[R(\alpha_i | x) = 1 - P(C_i | x)]$$

Decision: To minimize the risk, take the most probable class $P(C_i | x)$

Loss & Risk : Reject:

- In some applications, wrong decisions namely misclassification may have very high cost.

Eg: manual decision is made if automatic system has low certainty of its decision.

Here, consider an extra option: $(k+1)^{\text{st}}$ class, i.e., reject

$$\lambda_{ik} = \begin{cases} 0, & \text{if } i=k \\ \lambda, & \text{if } i=k+1 \\ 1, & \text{otherwise} \end{cases} \quad 0 < \lambda < 1$$

Risk of reject is,

$$R(\alpha_{k+1} | x) = \sum_{k=1}^k \lambda P(C_k | x)$$

$$= \lambda$$

Risk of misclassification is,

$$R(\alpha_i | x) = \sum_{k \neq i} P(C_k | x)$$

$$[R(\alpha_i | x) = 1 - P(C_i | x)]$$

Decision: Choose c_i which min. $R(c_i|x)$, $1 \leq i \leq k$.

O/P: { Choose c_i , if $P(c_i|x) > P(c_k|x)$ for all $k \neq i$ &
 $P(c_i|x) > 1 - \lambda$
 reject, otherwise.

Hence, If $\boxed{\lambda=0} \rightarrow$ we always reject.

i.e., correct classification.

If $\boxed{\lambda \geq 1} \rightarrow$ never reject

i.e., costlier than error

Discriminant Functions: (function of several variables used to assign items into one or more groups)

- classification can be seen as implementing a set of discriminant functions, $g_i(x)$, $i = 1, 2, \dots, k$.

choose c_i if $\boxed{g_i(x) = \max g_k(x)}$

$$g_i(x) = \begin{cases} -R(c_i|x), & \text{max. corresponds to min. risk.} \\ P(c_i|x), & \text{using 0/1 loss function.} \\ P(x|c_i) P(c_i), & \text{ignore } p(x) \end{cases}$$

This divides the feature space into ' k ' decision regions of R_1, R_2, \dots, R_k where

$$R_i = \{x | g_i(x) = \max g_k(x)\}$$

Where there is 2 classes, define a single discriminant as:

$$\boxed{g(x) = g_1(x) - g_2(x)}$$



Choose $\begin{cases} c_1, & \text{if } g(x) > 0 \\ c_2, & \text{otherwise} \end{cases}$

When $\boxed{k=2} \Rightarrow$ it is dichotomizer classification system.

& $\boxed{k \geq 3} \Rightarrow$ it is polychotomizer.

Parametric Methods: (how we estimate prob from given dataset)

- statistic is any value that is calculated from given sample.
- In statistical inference, we make decision using the information provided by sample.
- first approach is parametric where we assume that sample is drawn from some distribution that obeys a known model, for eg. Gaussian
- advantage of parametric approach is that model is defined to small No. of parameters.
Eg: mean, variance - sufficient statistics of distribution
- so we estimate parameters of distribution from samp. use these estimates in assumed model, get an estimated distribution \rightarrow which we use to make decision.
- method used to estimate parameters of distribution is "maximum likelihood estimation".

Parametric Estimation:

* Density estimation is general case of estimating $p(x)$.

\hookrightarrow for classification:

estimated densities are class densities $P(x|c_i)$
prior $[P(c_i)]$, posteriors, $[P(c_i|x)]$

\hookrightarrow for regression:

estimated density is $P(y|x)$

where 'x' is 1-dimensional & hence the densities are univariate.

for an independent & identically distributed sample,

$$X = [x^t]_{t=1}^N$$

x^t is instance drawn from some known prob. density family.

$$x^t \sim p(x|\theta)$$

want to find θ that makes sampling x^t from $p(x|\theta)$ as likely as possible.

Maximum Likelihood Estimation:

- Since x^t are independent, likelihood of parameter θ given sample x is the product of likelihood of individual points.

$$l(\theta|x) = p(x|\theta) = \prod_{t=1}^N p(x^t|\theta)$$

- Search for θ that maximizes the likelihood which is denoted by $l(\theta|x)$.
- we use log of likelihood where it takes max. value.

Log likelihood is defined as: $\log() \rightarrow$ converts product to sum

$$L(\theta|x) = \log l(\theta|x) = \sum_{t=1}^N \log p(x^t|\theta)$$

∴ maximum likelihood estimator (MLE):

$$\theta^* = \arg \max_{\theta} L(\theta|x)$$

* if we have a-class problem, } distribution used is } \rightarrow Bernoulli

* where there are K>2 classes, } generalization used is } \rightarrow Multinomial

* most frequently used density for } modeling class-conditional p/p densities } \rightarrow Gaussian (normal) density.

① Bernoulli Density:

- event occurs & bernoulli random variable x takes value 1 with prob. p & non-occurrence of event has prob. $1-p$ & denoted by x taking value 0.

e.g. $P(x) = p^x (1-p)^{1-x}$, $x \in (0, 1)$

Expected value & variance are:

$$E(x) = \sum_x x \cdot p(x) = 1 \cdot p + 0 \cdot (1-p)$$

$$\boxed{E(x) = p}$$

$$\text{Var}(x) = \sum_x [x - E(x)]^2 p(x)$$

$$\boxed{\text{Var}(x) = p(1-p)}$$

p is only parameter.

We need to calculate the estimator, \hat{p} :

log likelihood is:

$$\begin{aligned} L(p|x) &= \log \prod_{t=1}^N p(x_t) (1-p)^{(1-x_t)} \\ &= \sum_t x_t \log p + (N - \sum_t x_t) \log (1-p) \end{aligned}$$

∴ after solving for $\frac{dL}{dp} = 0$,

estimator: $\boxed{\hat{p} = \frac{\sum_t x_t}{N}}$

② Multinomial Density:

- Instead of two states, the outcome of random event is one of 'k' mutually exclusive & exhaustive states
- each of which has prob. of occurring p_i

with $\boxed{\sum_{i=1}^k p_i = 1}$

- x_i is 1 if outcome is state 'i', or 0 otherwise.

$$\boxed{P(x_1, x_2, \dots, x_k) = \prod_{i=1}^k p_i^{x_i}}$$

∴ MLE of p_i is

$$\boxed{\hat{p}_i = \frac{\sum_t x_t^i}{N}}$$

③ Gaussian (Normal) Distribution:

x is gaussian distributed with mean $E[x] = \mu$ and $\text{var}(x) = \sigma^2$, denoted as $N(\mu, \sigma^2)$.

density function,
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], -\infty < x < \infty$$

log likelihood is,

$$L(\mu, \sigma | x) = -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{\sum (x^t - \mu)^2}{2\sigma^2}$$

After taking partial derivatives,

$$\boxed{m = \frac{\sum x^t}{N}}, \text{ mean}$$

$$\text{and } \boxed{s^2 = \frac{\sum (x^t - m)^2}{N}}, \text{ s.d.}$$

Dimensionality Reduction:

- is the transformation of data from high-dimensional space into low-dimensional space, so that low dimensional representation retains some meaningful properties of original data.

Techniques:

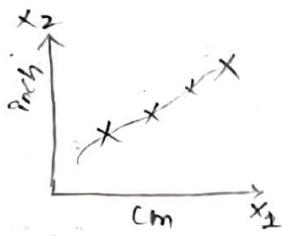
1. Principal Component Analysis (PCA)
2. Linear Discriminant Analysis (LDA).

Dimension Reduction:

- is a process of converting a dataset having vast dimensions into dataset with lesser dimensions.
- ensure that converted dataset conveys similar info.

Eg: Consider graph with 2 dimensions x_1 & x_2

- $x_1 \rightarrow$ represents measurement of object in cm
- $x_2 \rightarrow$ " " " inches



Using dimension reduction techniques:

- Convert dimension of data from 2 (x_1 & x_2) to 1 (x_1)



Benefits:

- Compress the data - which reduce storage space.
- eliminate redundant feature.
- Improve model performance.
- reduce time required for computation.

⇒ PCA :

- transforms the variables into new set of variables called as principal components.
- there can be only two principal components for two dimensional data set.

PCA Alg.:

- S1 - Get data
- S2 - Compute mean vector (μ)
- S3 - Subtract mean from given data.
- S4 - Calculate covariance matrix.
- S5 - Calculate eigen vectors & eigen values of covariance matrix.
- S6 - Choose components & form feature vector.
- S7 - Define new data set

problem:

1) Given:

$$\text{Data} = \{(2, 3, 4, 5, 6, 7); 1, 5, 3, 6, 7, 8\}$$

(or)

$$\text{Data} = (2, 1) (3, 5) (4, 3) (5, 6) (6, 7) (7, 8)$$

(or)

$$\text{Data} = \begin{cases} \text{Class 1} : x = 2, 3, 4 ; y = 1, 5, 3 \\ \text{Class 2} : x = 5, 6, 7 ; y = 6, 7, 8 \end{cases}$$

Compute the principle components using PCA alg.

Q1: Given feature vectors are:

$$x_1 = (2, 1) \quad x_2 = (3, 5) \quad x_3 = (4, 3) \quad x_4 = (5, 6) \quad x_5 = (6, 7) \quad x_6 = (7,$$

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 3 \\ 5 \end{bmatrix} \begin{bmatrix} 4 \\ 3 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix} \begin{bmatrix} 6 \\ 7 \end{bmatrix} \begin{bmatrix} 7 \\ 8 \end{bmatrix}$$

S2: calculate mean vector (μ)

$$\mu = \frac{2+3+4+5+6+7}{6}, \quad \frac{1+5+3+6+7+8}{6}$$
$$= (4.5, 5)$$

$$\therefore \text{Mean Vector } (\mu) = \begin{bmatrix} 4.5 \\ 5 \end{bmatrix}$$

S3: Subtract mean vector (μ) from given feature vectors.

$$x_1 - \mu = (2 - 4.5, 1 - 5) \Rightarrow (-2.5, -4)$$

$$x_2 - \mu = (3 - 4.5, 5 - 5) \Rightarrow (-1.5, 0)$$

$$x_3 - \mu = (4 - 4.5, 3 - 5) \Rightarrow (-0.5, -2)$$

$$x_4 - \mu = (5 - 4.5, 6 - 5) \Rightarrow (0.5, 1)$$

$$x_5 - \mu = (6 - 4.5, 7 - 5) \Rightarrow (1.5, 2)$$

$$x_6 - \mu = (7 - 4.5, 8 - 5) \Rightarrow (2.5, 3)$$

$$\begin{bmatrix} -2.5 \\ -4 \end{bmatrix} \begin{bmatrix} -1.5 \\ 0 \end{bmatrix} \begin{bmatrix} -0.5 \\ -2 \end{bmatrix} \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \begin{bmatrix} 1.5 \\ 2 \end{bmatrix} \begin{bmatrix} 2.5 \\ 3 \end{bmatrix}$$

S4: calculate covariance matrix.

$$\boxed{\text{Covariance matrix} = \frac{\sum (x_i - \mu)(x_i - \mu)^t}{n}}$$

$$m_1 = (x_1 - \mu)(x_1 - \mu)^t = \begin{bmatrix} -2.5 \\ -4 \end{bmatrix} \begin{bmatrix} 2.5 & -4 \end{bmatrix} = \begin{bmatrix} 6.25 & 10 \\ 10 & 16 \end{bmatrix}$$

$$m_2 = (x_2 - \mu)(x_2 - \mu)^t = \begin{bmatrix} -1.5 \\ 0 \end{bmatrix} \begin{bmatrix} -1.5 & 0 \end{bmatrix} = \begin{bmatrix} 2.25 & 0 \\ 0 & 0 \end{bmatrix}$$

$$m_3 = (x_3 - \mu)(x_3 - \mu)^t = \begin{bmatrix} -0.5 \\ -2 \end{bmatrix} \begin{bmatrix} 0.5 & -2 \end{bmatrix} = \begin{bmatrix} 0.25 & 1 \\ 1 & 4 \end{bmatrix}$$

$$m_4 = (x_4 - \mu)(x_4 - \mu)^t = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \begin{bmatrix} 0.5 & 1 \end{bmatrix} = \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$m_5 = (x_5 - \mu)(x_5 - \mu)^t = \begin{bmatrix} 1.5 \\ 2 \end{bmatrix} \begin{bmatrix} 1.5 & 2 \end{bmatrix} = \begin{bmatrix} 2.25 & 3 \\ 3 & 4 \end{bmatrix}$$

$$m_6 = (x_6 - \mu)(x_6 - \mu)^t = \begin{bmatrix} 2.5 \\ 3 \end{bmatrix} \begin{bmatrix} 2.5 & 3 \end{bmatrix} = \begin{bmatrix} 6.25 & 7.5 \\ 7.5 & 9 \end{bmatrix}$$

now,

$$\text{Covariance matrix} = \frac{m_1 + m_2 + m_3 + m_4 + m_5 + m_6}{6}$$
$$= \frac{1}{6} \begin{bmatrix} 17.5 & 22 \\ 22 & 34 \end{bmatrix} \rightarrow \text{after adding all values of } m_1, m_2, \dots, m_6$$

$$M = \begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix}$$

Q5: Calculate eigen value & eigen vector

λ is eigen value for matrix 'M'.

$$|M - \lambda I| = 0$$

$$\begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

$$\begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = 0$$

$$\begin{vmatrix} 2.92 - \lambda & 3.67 \\ 3.67 & 5.67 - \lambda \end{vmatrix} = 0$$

$$\therefore (2.92 - \lambda)(5.67 - \lambda) - (3.67)(3.67) = 0$$

$$16.56 - 2.92\lambda - 5.67\lambda + \lambda^2 - 13.47 = 0$$

$$\boxed{\lambda^2 - 8.59\lambda + 3.09 = 0} \quad -b \pm \sqrt{b^2 - 4ac} \\ 2a$$

Solving this quadratic equation, $\lambda = 8.22, 0.38$.

$\therefore \boxed{\lambda_1 = 8.22}$ & $\boxed{\lambda_2 = 0.38}$ are two eigen values.

Q6: Take the highest eigen value.

The eigen vector with highest eigen value is the principle component of given data.

Equation to find eigen vector: $\boxed{MX = \lambda X}$

when M = covariance matrix.

X = eigen vector

λ = eigen value

$$\therefore \begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 8.22 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Solving these,

$$2.92x_1 + 3.67x_2 = 8.22x_1$$

$$3.67x_1 + 5.67x_2 = 8.22x_2$$

on simplification.

$$5.3x_1 = 3.67x_2 \rightarrow ①$$

$$3.67x_1 = 2.55x_2 \rightarrow ②$$

from ① & ② :

$$x_1 = 0.69x_2 \rightarrow \frac{x_1}{x_2} = \frac{0.69}{1}$$

$$\therefore \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0.69 \\ 1 \end{bmatrix}$$

Now normalize this value,

$$e_1 = \begin{bmatrix} \frac{0.69}{\sqrt{0.69^2 + 1^2}} \\ \frac{1}{\sqrt{0.69^2 + 1^2}} \end{bmatrix} = \begin{bmatrix} 0.568 \\ 0.823 \end{bmatrix}$$

Q7: Define new dataset

$$\text{Principle component } \Phi_{C_1} \quad \boxed{\begin{array}{|c|c|c|c|c|c|} \hline P_{11} & P_{12} & P_{13} & P_{14} & P_{15} & P_{16} \\ \hline \end{array}}$$

$$P_{11} = e_1^T \begin{bmatrix} x_1 - \bar{x} \\ y_1 - \bar{y} \end{bmatrix} = [0.568 \ 0.823] \begin{bmatrix} 2-4.5 \\ 1-5 \end{bmatrix}$$

$$\Rightarrow [0.568 \ 0.823] \begin{bmatrix} -2.5 \\ -4 \end{bmatrix}$$

$$= -1.42 - 3.292$$

$$\boxed{P_{11} = -4.712}$$

$$P_{12} = e_1^T \begin{bmatrix} x_2 - \bar{x} \\ y_2 - \bar{y} \end{bmatrix} \Rightarrow [0.568 \ 0.823] \begin{bmatrix} 3-4.5 \\ 5-5 \end{bmatrix}$$

$$= [0.568 \ 0.823] \begin{bmatrix} -1.5 \\ 0 \end{bmatrix}$$

$$\boxed{P_{12} = -0.852}$$

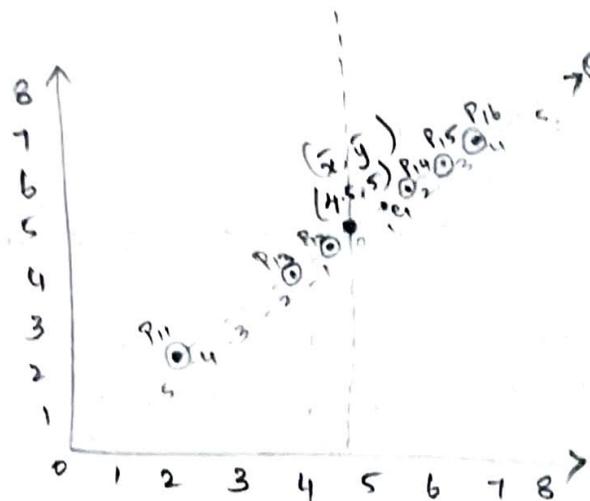
$$\boxed{P_{13} = -1.93}$$

$$\boxed{P_{14} = 1.136}$$

$$\boxed{P_{15} = 2.498}$$

$$\boxed{P_{16} = 3.889}$$

PC	P_{11}	P_{12}	P_{13}	P_{14}	P_{15}	P_{16}
	-4.712	-0.852	-1.93	1.136	2.498	3.889



project the
PC into
new subspace

Linear Discriminant Analysis: (LDA)

- was developed in the year 1936 by Ronald A. Fisher & was named Linear Discriminant or Fisher's Discriminant Analysis.
- originally was described as two-class technique later generalized by C. R. Rao as multiple discriminant analysis.
- LDA is supervised classification technique.
- used as pre-processing step in ML.

Representation of LDA models:

- ① consists of statistical properties of data, calculated for each class
- ② for single IP variable (x) there is mean & variance of variable for each class.
- ③ for multiple variables, this is same properties calculated over multivariate gaussian namely mean & co-variance
- ④ these statistical properties are estimated from data & plug into LDA equation to make predictions.

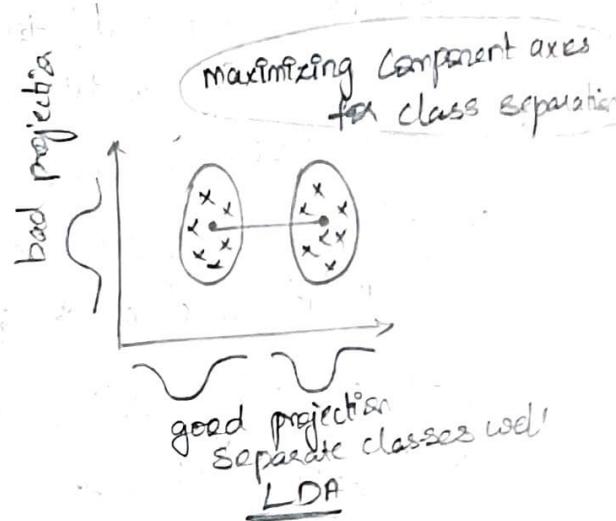
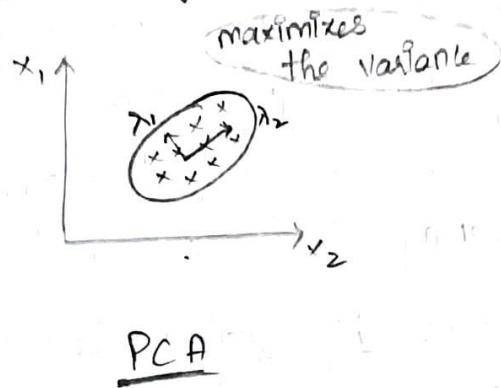
Applications:

- Medical
- Customer identification.
- facial Recognition.

Extensions to LDA:

- QDA - Quadratic Discriminant Analysis.
- Flexible Discriminant analysis. (FDA)
- Regularized Discriminant Analysis (RDA).

Comparison of PCA & LDA:



- * PCA ignores class labels & focuses on finding the principle components that maximizes variance in given data. Hence PCA is unsupervised algorithm.
- * LDA is supervised alg. that intend to find the linear discriminants that represents those axes which maximize separation b/w different classes.
- * LDA perform better multi-class classification task than PCA.
- * PCA performs better when sample size is small.
- PCA is followed by LDA.