

# lecture no.01

03-03-2025

## \* Machine Learning:

we make the machine to learn from the data Continuously. we make the machine to predict the thing that will happen in the future, by the given data. There are many methods that we use for predicting the data.

(i) Supervised learning :

(ii) Unsupervised learning :

(iii) Reinforcement learning :

- Before Applying any algo onto the data, the dataset must be complete -

→ Supervised learning:

where we know what the output is, or the characteristic & feature are known to predict the dataset. e.g. predicting the diseases on bases of the data given. Dataset — ML — output  
Train  
|  
Test

→ Unsupervised learning:

like we don't know what the output is, we have some data to predict and learn about the dataset, but we don't actually know the answer.

Dataset — ML — Groups .

→ Reinforcement learning:

There is no dataset in reinforcement learning. There are some reward & punishment which we make the machine to learn from it.

In RL we make the agent to learn from the environment .

→ Why we need to preprocess the data?

- Data must be complete
- no missing value.
- Data must be accurate.

→ 1st preprocess the data. Secondly check whether the output is give or not.

### • Supervised learning:

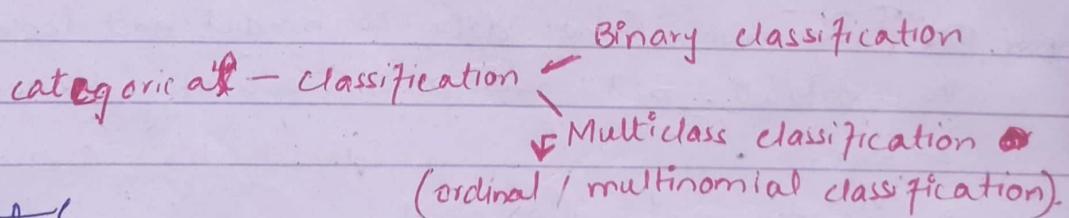
- Dataset given.
- Features                      output.  
 $x_1, x_2, x_3 \rightarrow y$ .

Feathers, no. of legs, Tail etc  $\rightarrow$  Bird, cat, cow etc.

→ If the output is categorical then it is classification.

mean the value of output is countable.

- $\rightarrow$  If the value has only two output like yes, no or 0, 1, then it is Binary Classification.
- $\rightarrow$  If the value is more than two option then it is multiclass classification, multinomial classification.



in ordinal

→ If the output is discrete then it is Regression

problem. e.g. To predict the house price, etc.

- Linear Regression  $\rightarrow$  value/output will be changed each time.

### Algorithms for classification:

- Logistic Regression
- SVM
- KNN
- Decision Tree
- Naive Bayes,

## Algorithms for Regression.

- Linear Regression .
- Ordinal Multiclass ~~Reg~~ classification :  
Specific. Order ,  
Strongly Agree , Disagree , Neutral Agree .  
Best , average , worst .
- Multinomial Multiclass classification :  
Non - Specific order .
  - cat , cow , Bird .
  - male , female .
  - Red Blue , white .
  - Hot , mild , cold .

## Assignment :

Download a dataset from Kaggle, UCI .

- Dataset type , no. of rows , columns , feature ( 15 ) ,  
rows ( 1500 )

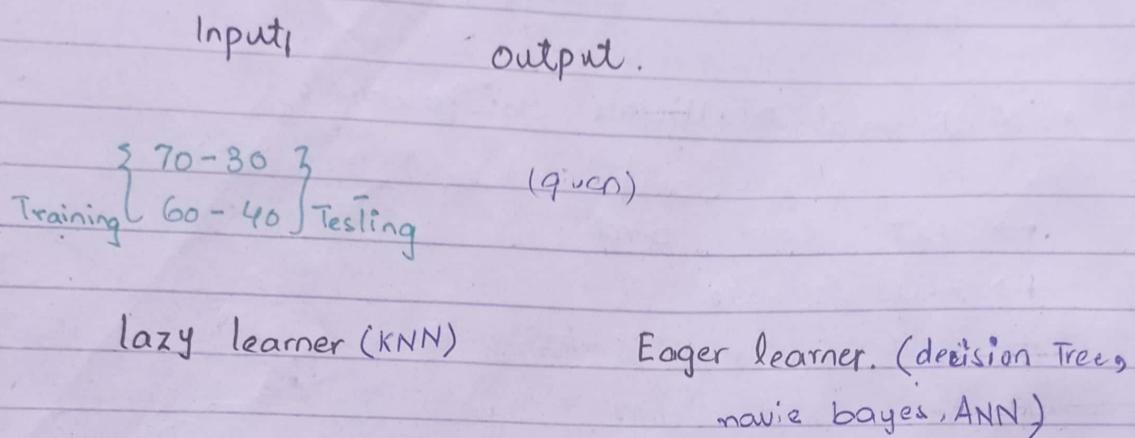
# Lecture no. 02

10-03-2024.

## \* Supervised Learning:

### → Classification:

- There must be some input, feature in classification.
- In classification we divide the data into training & testing data
- Lazy learner. That why we call it learner.
- Eager Learner.



→ KNN (k-nearest neighbour), lazy learner.

example: (binary classification.)

Height cm	Weight kg	T-shirt-Size.	Distance	S-D
158	58	M	4.24	9
158	59	M	3.6	6
158	63	M	3.6	6
160	59	M	2.22	3
160	60	M	1.4	1
163	60	M	2.2	4
163	61	M	2	2
160	64	L	3.016	5
163	64	L	3.6	6
165	61	L	4	7
165	62	L	4.1	8
165	65	L	5.0	10

168	62	L	7.06	11
168	63	L	7.28	12
x168	66	L		
x170	63	L		
x170	64	L		
x170.	68	L		

Q: identify shirtsize of customer having ht 161 and wt 61kg.

\* Step 1: calculate Similarity Based on Distance:

(i) Euclidean formula:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

(ii) Manhattan Distance.

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

$$\begin{aligned} d(x, y) &= \sqrt{(161 - 158)^2 + (61 - 58)^2} \\ (161, 61) &(158, 58) = \sqrt{9 + 9} \\ &= \sqrt{18} \\ &= 4.24. \end{aligned}$$

- if the value are discrete and numerical then avoid using decision tree.

- if there is only one attribute in the dataset

that is textual then convert it into the numerical value like

if age is given like old, young or child then convert old into numbers

like 0 for old, 1 for young and 2 for child

so the no of wash-B can't be 2.5 so the data is invalid

\* Step 2: Find k-Nearest Neighbour.

$$k=3 \Rightarrow M \quad M \quad M$$

$$k=5 \Rightarrow M \quad M \quad M \quad M \quad M$$

• Before applying any algo check if the data is valid or invalid like the attribute in no of washroom in uni is given 2, 2.5 & 3

• if the attribute is given like names of person etc so we cannot convert it into number so here we can drop the attribute

if age is given like old, young or child then convert old into numbers

like 0 for old, 1 for young and 2 for child

so the no of wash-B can't be 2.5 so the data is invalid

# Lecture no. 03

17-03-2025

## \* Logistic Regression - Classification:

↓ value will predict the class (ordinal, binary)

→ first we define a threshold.

like some value 50.

→ Threshold can be multiple.

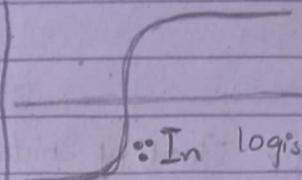
→ Identification of constants.

Formula:

Logistic Regression with one variable

$$y = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

or attribute or feature.



In logistic regression the dependent variable ( $y$ ) is categorical.

$$= \frac{1}{1 + e^{(\beta_0 + \beta_1 x)}}$$

make a sigmoid.

Linear Regression - Regression

∴  $Bx \rightarrow x$  is the feature or attribute: In linear regression it is given.

(y) dependent variable is not categorical

logistic Regression with multi Variable.

make a straight line.

$$\text{Formula.} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}$$

$$y = \beta_0 + \beta_1 x$$

$$= \frac{1}{1 + e^{-\beta_0 - \beta_1 x_1 - \beta_2 x_2}}$$

$$\beta_0 = 3, \beta_1 = 0.8$$

input	output
no 5	0
1.2	0
2.3	0
3.1	1
4.5	1
5.0	1
5.8	1
6.3	1
7.1	1
8.0	1

$$= \frac{1}{1+e^{-(3+(0.8)(0.5))}} \\ = 0.97$$

$$2) \quad \frac{1}{1+e^{-(3+(0.8)(1.2))}} \\ = 0.98$$

$$3) \quad \frac{1}{1+e^{-(3+(0.8)(2.3))}} \\ = 0.99$$

$$4) \quad \frac{1}{1+e^{-(3+(0.8)(3.0))}} \\ = 0.99$$

$$5) \quad \frac{1}{1+e^{-(3+(0.8)(4.5))}} \\ = 0.99$$

x	y_p	y_p (-3)
0.5	0.97	(0) 0.0691
1.2	0.98	(0) 0.115
2.3	0.992	(0) 0.238
3.1	0.995	(0) 0.372
4.5	0.998	(1) 0.645
5.0	0.999	(1) 0.731
5.8	0.9995	(1) 0.837
6.3	0.996	(1) 0.84
7.1	0.9998	(1) 0.935
8.0	0.9999	(1) 0.967

accuracy  $\rightarrow 7/10 = 70\%$        $9/10 = 90\%$

$$\beta_0 = \frac{\sum x^2 \sum y - \sum x \sum xy}{n \sum x^2 - (\sum x)^2}$$

$$\beta_i = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

Assignments:

1) Write a program in python and read the dataset / display the data

2) Library for logistic regression  
which library is used to import logistic regression

when the values of  $B$  is not known.

$X$	$Y$	$X^2$	$XY$
0.5	0	0.25	0
1.2	0	1.44	0
2.3	0	5.29	0
3.1	1	9.61	3.1
4.5	1	20.25	4.5
5.0	1	25	5.0
5.8	1	33.64	5.8
6.3	1	39.69	6.3
7.1	1	50.41	7.1
8.0	1	64	8.0

$$\sum x = 43.80 \quad \sum y = 7 \quad \sum x^2 = 250 \quad \sum xy = 39.8$$

$$B_0 = -3 \quad 390\% \\ B_1 = 0.8 \quad \text{accuracy}$$

$$B_0 = 3 \quad 370\% \\ B_1 = 0.8 \quad \text{accuracy}.$$

$$B_0 = \frac{\sum x^2 \sum y - \sum x \sum xy}{n \sum x^2 - (\sum x)^2}, \quad B_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \\ = 0.157$$

$$= \frac{(250)(7) - (43.80)(39.8)}{10(250) - (43.80)^2}$$

$$= 0.157, 0.006.$$

e.g. if value of  $y$  is not given.

$$x = 5.2$$

$$y = ?$$

check the better accuracy ( $B_0$  &  $B_1$ )

then

Put values in formula.

$$* Y = \frac{1}{1 + e^{-(B_0 + B_1 x_1)}}$$

$$Y = \frac{1}{1 + e^{(-0.3 + (0.8)(5.2))}}$$

Given a dataset of students containing their study hours to identify whether they pass or fail the exams. If the student studied 3.2 hours for exams, you need to identify whether he/she will pass the exam using the appropriate algo based on the dataset.

Students	Study hours ( $x$ )	Pass(1)/Fail(0), ( $y$ )	Identify the choice of algo. Also identify the accuracy of your algorithm.
1	1.5	0	
2	2.0	0	
3	2.5	0	
4	3.0	0	
5	3.5	1	
6	4.0	1	
7	4.5	1	
8	5.0	1	
9	5.5	1	
10	6.0	1	
	$\sum x = 37.5$	$\sum y = 6$	

$x^2$	$xy$	
2.25	0	
4	0	$B_0 = \frac{\sum x^2 \sum y - \sum x \sum xy}{n \sum x^2 - (\sum x)^2}$
6.25	0	
9	0	
12.25	3.5	$B_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$
16	4.0	
20.25	4.5	
25	5.0	
30.25	5.5	
36	6.0	
$\sum x^2 = 161.25$	$\sum xy = 28.5$	

$$B_0 = \frac{(161.25)(6) - (37.5)(28.5)}{10(161.25) - (37.5)^2} = -0.4909$$

$$B_1 = \frac{10(28.5) - (37.5)(6)}{10(161.25) - (37.5)^2} = 0.2909$$

$$\beta_0 = -0.4909, \quad \beta_1 = 0.2909$$

$$y = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$$y = \frac{1}{1 + e^{(-0.49 + 0.29)(1.5)}}$$

accuracy is  $\frac{7}{10} = 70\%$

x	y <sub>P</sub>	y ( )
1.5	0.486 ≈ 0	0
2.0	0.522 ≈ 1	0
2.5	0.558 ≈ 1	0
3.0	0.593 ≈ 1	0
3.5	0.628 ≈ 1	1
4.0	0.661 ≈ 1	1
4.5	0.691 ≈ 1	1
5.0	0.721 ≈ 1	1
5.5	0.751 ≈ 1	1
6.0	0.771 ≈ 1	1

# lecture no. 04

30-04-2024

## \* Regression:

- Linear Regression

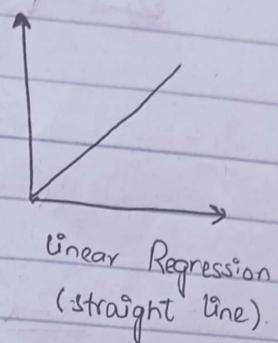
→ Formula:

- Linear Regression with one variable:

$$y_p = \beta_0 + \beta_1 x_1$$

Constant                          attribute.

∴ output can  
be continuous



- Linear Regression with multi variable.

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

- e.g:-

linear Regression for one variable.

Predict the obtained marks of a student who studies for 65 hr.

$x$ (H-Studies)	$y$ (Marks obtained)	$x^2$	$xy$	$y_p'$	error ( $E_p$ )	$(y_A - y_p')$
1	50	1	50	51.27	1.27	1.27
2	55	4	110	56.11	1.11	1.11
3	60.5	9	195	60.95	0.05	0.05
4	70	16	280	65.7	4.3	4.3
5	65	25	325	70.63	5.63	5.63
6	75	36	450	75.47	0.47	0.47
7	80	49	560	80.31	0.30	0.30
8	85	64	680	85.15	0.15	0.15
9	90	81	810	89.99	0.01	0.01
10	95	100	950	94.83	0.17	0.17
11	100	121	1100	99.67	0.33	0.33
12	105	144	1260	104.5	0.49	0.49
$\sum x = 78$	$\sum y = 935$	$\sum x^2 = 650$	$\sum xy = 6770$		$\sum E_p^2$	
					$\sum (y_A - y_p) =$	
						18.29

$$\beta_0 = \frac{\sum x^2 \sum y - \sum x \sum xy}{n(\sum x^2) - (\sum x)^2}$$

$$\beta_0 = \frac{(650)(935) - (78)(6770)}{12(650) - (78)^2}$$

$$\beta_0 = 46.43$$

$$\beta_{1,2} = \frac{n(\sum xy) - \sum x \sum y}{n(\sum x^2) - (\sum x)^2}$$

$$\frac{12(6770) - (78)(935)}{12(650) - (78)^2}$$

$$\beta_1 = 4.84$$

$$\therefore y_p = \beta_0 + \beta_1 x_1$$

$$y_p = 46.43 + 4.84(1)$$

$$\therefore E_p = (y_A - y_p)$$

$$E_{p1} = 50 - 51.27$$

# \* Performance Calculation of Linear Regression:

• Mean Square ERROR: (1st Iteration)

$$\left( \frac{\sum E_p^2}{n} \right) = \frac{(18.29)^2 / 12}{12} = 27.87$$

$(\sum E_p^2 / n) = (\sum (Y_A - Y_P)^2 / n)$

• Root mean Square:

$$\sqrt{\sum (Y_A - Y_P)^2 / n} =$$

$$\Rightarrow \sqrt{27.87} = 5.2$$

$(Y_A - Y_P^2)$

$Y_P^2$	$E_P^2$	$Y_P^3$	$E_P^2$
50.43	0.43	50.5	0.5
54.43	0.57	55	0
58.43	6.57	59.5	5.5
64.43	7.57	64	6
66.43	1.43	68.5	3.5
70.43	4.80	73	2
74.43	5.57	77.5	2.5
78.43	6.57	82	3
82.43	7.57	86.1	3.5
86.43	8.57	91	4
90.43	9.57	95.5	4.5
94.43	10.57	100	5

$\beta_1 = 4$

$$\sum E_p^2 = 69.49$$

↓

$$\beta_1 = 4.5$$

↓

$$\sum E_p^2 = 40$$

∴ To find  $y_p^2$ , suppose, any random value near to  $\hat{y}_1$ , like here we are supposing 4 which is nearest to the value of  $\beta_1$ . After finding the value of  $E_p^2$  we check for the distance b/w actual  $y_p$  and new  $y_p^2$ , we can change both  $\beta_0$  &  $\beta_1$ .

$$\beta_1 = 4$$

∴ Iteration:

we do the iterations for the algo by changing the value of (mean square error),

The value of mean square error must be near to 0 it can't be actual value if it is then the calculation are wrong.

→ • mean square error (for 2nd iteration)

$$\text{Error} = \frac{\sum (Y_A - Y_p^2)^2 / n}{(69 - 49)^2 / 12}$$

$$= 403.2$$

→ Root • mean Square:

$$\sqrt{\sum (Y_A - Y_p^2)^2 / n}$$

$$= 20.07$$

→ Mean Square error : (3rd iteration)

$$\text{Error} = \frac{\sum (Y_A - Y_p^3)^2 / n}{(40)^2 / 12}$$

$$= 133$$

• root mean Square:

$$= \sqrt{\sum (Y_A - Y_p^3)^2 / n}$$

$$= \sqrt{133} = 11.5$$

Q: If the student studied for 6.5 h, then how many marks will be obtain.

So, lower mean square value is 5.2

$$= 46.43 + 4.84(6.5)$$

$$= 77.89$$

• Linear Regression for multi variable.

$X_1$	$X_2$	$Y$	$Y_p^1$	$E_p^1$	$Y_p^2$	$E_p^2$
1	0.10	1.5	1.6	0.1	1.4	0.1
2	0.12	2	2.68	0.68	2.4	0.4
3	0.14	3.5	3.71	0.21	3.4	0.1
4	1.16	4.0	4.74	0.74	4.4	0.4
5	30.18	5.5	5.77	0.27	5.4	0
6	0.20	6.0	6.8	0.8	6.4	0.5

Suppose the value of  $\beta_0$  &  $\beta_1$  &  $\beta_2$

$$\beta_0 = 0.5, \beta_1 = 1, \beta_2 = 1.5$$

After find the values of  $\beta$ 's calculate mean square error and root mean error. Then change the value of  $\beta_2$ .

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 \\ &= 0.5 + (1)(1) + 1.5(0.10) \\ &= 1.6 \quad \text{so on} \end{aligned}$$

$$\begin{aligned} E_p^1 &= (Y_A - Y_p^1) \\ &= 1.5 - 1.6 \\ &= 0.1 \quad \text{so on.} \end{aligned}$$

Change the value of  $\beta_0 = 0.25$

$$\begin{aligned} \text{mean square root} &= \sqrt{\frac{\sum (Y_A - Y_p^1)^2}{n}} \\ &= \sqrt{(1.5)^2 / 6} \end{aligned}$$

$$\begin{aligned} \text{Root mean error} &= \sqrt{\frac{\sum (Y_A - Y_p^1)^2}{n}} \\ &= \sqrt{0.375} \\ &= 0.6 \end{aligned}$$

Formula of  $\beta_0$  for multi-variables.

$$\beta_0 = \frac{\sum (x_1 x_2)^2 \bar{y} - \sum (x_1 x_2) \sum (x_1 x_2) \bar{y}}{n \sum (x_1 x_2)^2 - (\sum x_1 x_2)^2},$$

# Lecture no. 05

## Performance Metrics of Classification

Accuracy:

- Prediction must be correct
- $\frac{8/12}{TP+TN / TP+TN+FN+FP}$

Precision:

- $\frac{TP}{TP + FP}$

Recall:

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

F<sub>1</sub> Square:

- $2 \times \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}}$

Dataset:

S.no	Actual( $x_p$ )	Predicted ( $y_p$ )	True positive	
			TP	Actual P, predicted R
1	A <sub>1</sub> Positive	P	TP	True negative
2	A <sub>2</sub> Positive	N	FN	Actual N, predicted N
3	A <sub>3</sub> negative	N	TN	False positive
4	A <sub>4</sub> Positive	P	FP	Actual N, Predicted P
5	A <sub>5</sub> negative.	P	FP	False negative.
6	A <sub>6</sub> negative	N	TN	
7	A <sub>7</sub> Positive	P	TP	Actual P, Predicted N.
8	A <sub>8</sub> (P)	P	TP	
9	A <sub>9</sub> N	N	TN	
10	A <sub>10</sub> N	P	FP	Confusion Matrix:
11	A <sub>11</sub> P	N	FN	2x2 matrix.
12	A <sub>12</sub> N	N	TN	

Actual

Predicted	Actual	
	P	N
P	TP = 4	FP = 2
N	FN = 2	TN = 4

Precision:

Proportion of all correctly predicted +ve cases out of all predicted +ve cases.

Recall:

Proportion of Correctly predicted +ve cases out of all actual +ve cases.

F<sub>1</sub>-Score:

Harmonic mean of precision and recall. It balances both metrics especially useful when classes are imbalanced.

3x3 Matrix. If there are 3x3 matrix.

	cat	Dog	Rabbit			
cat	TP (1)	2	5	T-N	not cat -ve	no cat -ve
Dog	3	TP (4)	2	T-P	cat +ve	cat +ve
Rabbit	1	6	—	TP (10)	F-P (other) no cat -ve	cat +ve
					F-N (other) cat +ve	no cat -ve

Confusion Matrix w.r.t cat Predicted.

	cat	Others
cat	1	7
Others	4	22

4+2+6+10,

C.M w.r.t Dog predicted:

		Predicted	
		Dog	Others
Actual	Dog	9	$5^{3+2}$
	Others	$8^{6+2}$	17 $1+5+1+10$

C.M w.r.t Rabbit predicted

		Predicted	
		Rabbit	Others
Actual	Rabbit	10	$7^{6+1}$
	Others	$5^{2+1}$	10 $1+2+3+4$

:-

- To check the performance of Regression there is a suitable method named as root mean square error.
- To check the performance of classification there is a suitable method named as Confusion matrix.

## • Linear Regression for multi-variable:

example:

$x_1$	$x_2$	$y$	$x_1^2$	$x_2^2$	$x_1y$	$x_2y$	$y_{p1}$
1	0.10	1.6	1	0.01	1.6	0.16	30
2	0.12	2.68	4	0.014	5.36	0.32	44.56
3	0.14	3.71	9	0.09	9.51	0.51	59.12
4	0.16	4.74	16	0.025	18.96	0.75	73.68
5	0.18	5.77	25	0.032	28.85	1.03	88.24
6	0.20	6.8	36	0.04	40.8	1.36	
$\sum x_1 =$	$\sum x_2 =$	$\sum y =$	$\sum x_1^2 =$	$\sum x_2^2 =$	$\sum x_1y =$	$\sum x_2y =$	
21	0.9	25.3	91	0.14	105.19	4.13	

$$\beta_0 = \frac{\sum x^2 \sum y - \sum x \sum xy}{n \sum x^2 - (\sum x)^2}, \quad \beta_1 = \frac{n \sum xy - \sum x \sum y}{n (\sum x^2) - (\sum x)^2}$$

$$\beta_0 = \frac{(0.9)(25.3) - (21)(105.19)}{6(91) - (21)^2}, \quad = \frac{6(105.19) - (91)(25.3)}{6(91) - (21)^2}$$

$$\beta_0 = -20.8, \quad \beta_1 = -15.9.$$

$$\beta_2 = \frac{n \sum x_2 y - \sum x_2 \sum y}{n (\sum x_2^2) - (\sum x_2)^2}$$

$$\beta_2 = \frac{6(4.13) - (0.9)(25.3)}{6(0.14) - (0.9)^2} \Rightarrow \beta_2 = 67$$

$$y_{p1} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$= -20.8 + (-15.9)(1) + (67)(0.10)$$

$$= -30$$

# Lecture no. 06

## Unsupervised Learning ::

### \* Clustering :

(i) Partitioning:  
K-means

(ii) Hierarchical: (Bottom-up approach)  
Agglomerative

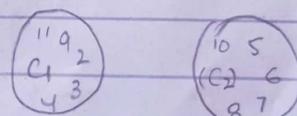
(iii) Density Based:  
DB - Scan.

### K-mean: (Partition)

example:

1st iteration.

ID	X <sub>1</sub>	X <sub>2</sub>	Centroid C <sub>1</sub>	I <sub>1</sub>	Centroid C <sub>2</sub> (1.06, 1.3) (3.68, 3.95)	I <sub>2</sub>
1	1.0	1.5	C <sub>1</sub>	I <sub>1</sub>		
2	1.2	1.3	C <sub>1</sub>	C <sub>1</sub> min(d(C <sub>1,1</sub> , d(C <sub>2,1</sub> ))		
3	0.8	1.6	C <sub>1</sub>	C <sub>1</sub>		
4	1.1	1.0	C <sub>1</sub>	C <sub>1</sub>		
5	3.5	4.0	C <sub>2</sub>	C <sub>2</sub>		
6	3.8	3.7	C <sub>2</sub>	C <sub>2</sub>		
7	3.6	4.1	C <sub>2</sub>	C <sub>2</sub>		
8	3.9	3.8	C <sub>2</sub>	C <sub>2</sub>		
9	1.3	1.2	C <sub>1</sub>	C <sub>1</sub>		
10	3.7	4.2	C <sub>2</sub>	C <sub>2</sub>		
11	1.0	1.4	C <sub>1</sub>	C <sub>1</sub>	C <sub>1</sub> = 2, 3, 4, 9, 11	
12	3.6	3.9	Centroid C <sub>2</sub>			C <sub>2</sub> = 5, 6, 7, 8, 10



$$d(C_1, 2) = \sqrt{(1.0 - 1.2)^2 + (1.5 - 1.3)^2} =$$

$$d(C_2, 2) = \sqrt{(3.6 - 1.2)^2 + (3.9 - 1.3)^2} =$$

Calculating centroid for Iteration 2.

C<sub>1</sub>

$$C_1x_1 = \frac{1.0 + 1.2 + 0.8 + 1.1 + 1.3 + 1.0}{6} = 1.066$$

$$C_1x_2 = \frac{1.5 + 1.3 + 1.6 + 1.0 + 1.2 + 1.4}{6} = 1.3$$

C<sub>2</sub>

$$C_2x_1 = \frac{3.5 + 3.8 + 3.6 + 3.9 + 3.7 + 3.6}{6} = 3.68$$

$$C_2x_2 = \frac{4.0 + 3.7 + 4.1 + 3.8 + 4.2 + 3.9}{6} = 3.95$$

# Hierarchical Clustering:

\* Agglomerative:

Example:

ID	$x_1$	$x_2$
P <sub>1</sub>	0.40	0.53
P <sub>2</sub>	0.22	0.38
P <sub>3</sub>	0.35	0.32
P <sub>4</sub>	0.26	0.19
P <sub>5</sub>	0.08	0.41
P <sub>6</sub>	0.45	0.30

Iteration 1.

	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>	P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>
P <sub>1</sub>	.					
P <sub>2</sub>	0.23	.				
P <sub>3</sub>	0.22	0.16	.			
P <sub>4</sub>	0.37	0.19	0.13	.		
P <sub>5</sub>	0.34	0.13	0.28	0.23	.	
P <sub>6</sub>	0.28	0.26	0.10	0.22	0.39	.

P<sub>1</sub>.

$$d(P_1, P_2) = \sqrt{(0.40 - 0.22)^2 + (0.53 - 0.38)^2} = 0.2$$

$$d(P_1, P_3) = \sqrt{(0.40 - 0.35)^2 + (0.53 - 0.32)^2} = 0.21$$

$$d(P_1, P_4) = \sqrt{(0.40 - 0.26)^2 + (0.53 - 0.19)^2} = 0.36$$

$$d(P_1, P_5) = \sqrt{(0.40 - 0.08)^2 + (0.53 - 0.41)^2} = 0.34$$

$$d(P_1, P_6) = \sqrt{(0.40 - 0.45)^2 + (0.53 - 0.30)^2} = 0.23$$

P<sub>2</sub>

$$d(P_2, P_3) = \sqrt{(0.22 - 0.35)^2 + (0.38 - 0.32)^2} = 0.16$$

$$d(P_2, P_4)$$

$$d(P_2, P_5)$$

$$d(P_2, P_6) :$$

P<sub>3</sub>.

$$d(P_3, P_4)$$

$$d(P_3, P_5)$$

$$d(P_3, P_6).$$

P<sub>4</sub>:

$$d(P_4, P_5)$$

$$d(P_4, P_6)$$

P<sub>5</sub>:

$$d(P_5, P_6) :$$

Close clusters:

P<sub>1</sub> P<sub>2</sub> P<sub>3</sub> P<sub>4</sub> P<sub>5</sub> P<sub>6</sub>

Iteration 2.

P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub> P <sub>6</sub>	P <sub>4</sub>	P <sub>5</sub>
P <sub>1</sub>	0			
P <sub>2</sub>	0.23	0		
P <sub>3</sub> P <sub>6</sub>	0.22	0.14	0	
P <sub>4</sub>	0.37	0.19	0.13	0
P <sub>5</sub>	0.34	0.14	0.28	0.23

$$\min((P_3 P_6), P_1) = \min(d(p_3, p_1), d(p_6, p_1))$$

$$= \min(0.22, 0.24)$$
$$= 0.22$$

$$\min((P_3 P_6), P_2) = \min(d(p_3, p_2), d(p_6, p_2))$$

$$= \min(0.14, 0.24)$$
$$0.14$$

$$\min((P_3 P_6), P_4) = \min(d(p_3, p_4), d(p_6, p_4))$$

$$= \min(0.13, 0.22)$$
$$= 0.13$$

$$\min((P_3 P_6), P_5) = \min(d(p_3 P_6), d(p_6, p_5))$$

$$= \min(0.28, 0.39)$$
$$0.28$$

Close clusters:

P<sub>3</sub> P<sub>6</sub> P<sub>4</sub>

### Iteration 3

	$P_1$	$P_2$	$P_3 P_6 P_4$	$P_5$
$P_1$	0			
$P_2$	0.23	0		
$P_3 P_6 P_4$	0.22	0.14	0	
$P_5$	0.34	0.14	0.23	0

$$\min(d(P_3 P_6 P_4), P_1) = \min(d(P_3, P_1), d(P_4, P_1))$$

$$= \min(0.22, 0.37)$$

$$= 0.22$$

$$\min(d(P_3 P_6 P_4), P_2) = \min(d((P_3 P_6), P_2), d(P_4, P_2))$$

$$= \min(0.14, 0.19)$$

$$= 0.14$$

$$\min(d(P_3 P_6 P_4), P_5) = \min(d((P_3 P_6), P_5), d(P_4, P_5))$$

$$= \min(0.28, 0.23)$$

$$= (0.23)$$

	$P_1$	$P_2 P_5$	$P_3 P_6 P_4$
$P_1$	0		
$P_2 P_5$	0.23	0	
$P_3 P_6 P_4$	0.22	0.14	0

$$\min(d(P_2, P_5), P_1) = \min(P_2, P_1), (P_5, P_1)$$

$$\min(0.23, 0.34)$$

$$= 0.23$$

$$\begin{aligned}
 d(p_3, p_6, p_4), p_2, p_5) &= \min ((p_3, p_6, p_4), (p_3, p_6, p_4), p_5) \\
 &= (0.14, 0.23) \\
 &= 0.14
 \end{aligned}$$

Close Clusters.

