

3/3/23

## Lab-1



Write python code, consider filename as "housing.csv"

- (i). To load .csv file into the data frame.
- ii. To display information of all columns
- iii. To display statistical information of all numerical.
- iv. To display the count of unique labels for "Ocean Proximity" column.
- v). To display which attributes (columns) in a dataset have missing values count greater than zero.

(i). `pd. df = pd.read_csv("housing.csv")`

(ii) `print(df)`

(iii) `print(df.describe())`

(iv). `print(df[["Ocean Proximity"]].value_count())`

(v). ~~`missing_value = df.isnull().sum()`~~  
~~`missing_column = df.missing_values[missing_value]`~~  
~~`print(missing_column)`~~

1. Which column in the dataset had missing values? How did you handle them?

Dropped the missing values with the help of dropna since the number of rows are less with missing value is less compatibility.

2. What were the categorical ~~data~~ columns did you identify in the dataset? How did you encode them?

diabetes.csv : ["Gender", "Class"]

adult.csv : ["workclass", "education", "marital-status", "occupation", "relationship", "gender", "native-country", "income"]

and we used LabelEncoder to encode it.

3. What is the difference between Min-Max scaling and standardization? When would you use one over the other?

Both MinMax Scaling and standardization are feature scaling techniques.

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

$$x' = \frac{x - \bar{x}}{s}$$

We use max when we do not have significant outliers. Then we use standardization when we have outliers.

~~Q1~~

negative effect from a break off without a

no evidence to exist

(E. "outlier test") no break, no significant

break, no significant change

break, (1) significant, (2) not significant

break, p = unknown: Insignificant

chaos & endo

end of diag not important fail

most important not important

initially no more than one result -> and

multiple result & multiple test

standard deviation - robust - standard dev

multiple standard deviations - a rule

very few even standard deviation test

no significant difference, normal distribution

normal & equal var and equal shape

normal & homoscedastic error

normal & independent

normal & homoscedastic error

Demonstrate the steps to build a machine - Learning model that predicts the median housing price using the California housing price dataset.

1. Perform the describe and info steps:

```
import pandas as pd  
housing = pd.read_csv("content/drive/..../")  
print(housing.info())  
print(housing.describe())
```

Output : column = 9  
rows = 20640.

2. Plot histogram for each feature

The histogram for median-income & house-median-age can give us insight into the distribution of these features. For example, median-income might show a right-skewed distribution, indicating that most households have lower median income, while house-median-age might show how the ages of houses are distributed.

3 Demonstrate process of creating set & differences b/w random & stratified test.

Random Stratified.

Split data randomly into training & testing

Both are representative of the overall distribution & a feature.

May not preserve the distribution of key feature in training & testing set. Preserves the distribution across training & testing set.

4 List geographical feature.

This graph visualizes housing prices in relation to their geographic location, with the color representing median house value and size representing population.

5 Which feature correlates to maximum:

median income usually shows the highest correlation with the median house value.

6 List the features that could improve correlation.

By creating rooms per household & bedrooms per room.

The median house value is improved.

7 total bedrooms list the features that needs to be cleaned & demonstrate cleaning process.

(2) Categorical data to numerical data.

Yes, ocean-proximity is categorical feature and method used to convert is OneHotEncoder.

9 Important Scaling feature :-

- ensure all feature contribute equally to result.

Technique used StandardScaler & MinmaxScaler.

⑩

17/03/25

LAB-3

Q) Solving the following Linear Regression Problem using Matrix approach. Find linear Regression of the data of week and product - sales.

$x_i$  (week)       $y_i$  (Sales in thousands).

$$\begin{matrix} 1 & 2 \\ 2 & 4 \\ 3 & 5 \\ 4 & 9 \end{matrix}$$

$$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$$

$$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$$

$$\begin{matrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{matrix}$$

$$\beta = ((x^T x)^{-1} x^T) y$$

$$x^T x = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{pmatrix} \times \begin{bmatrix} 1 & 2 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

$$(x^T x)^{-1} = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}^{-1} = \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix}$$

$$(x^T x)^{-1} x^T = \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix} \times \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0.5 & 0 & -0.5 \\ -0.3 & -0.1 & 0.1 & 0.3 \end{bmatrix}$$

Finally,

$$((X^T X)^{-1} X^T) Y = \begin{bmatrix} 1 & 0.5 & 0 & -0.5 \\ -0.5 & -0.1 & 0.1 & 0.3 \end{bmatrix}^{-1} \times \begin{bmatrix} 2 \\ 4 \\ 5 \\ 9 \end{bmatrix}$$

$$\geq \begin{bmatrix} -0.5 \\ -0.6 - 0.4 + 0.5 + 2.7 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} -0.5 \\ 2.2 \end{bmatrix}$$

~~$$y = -0.5 + 2.2x$$~~

→ Normal:

$$\sum x = 10 \quad \sum x^2 = 30.$$

$$\sum y = 20. \quad \sum xy = 61.$$

$$\bar{x} = 2.5 \quad \bar{y} = 5, \quad \bar{x}^2 = 7.5, \quad \bar{xy} = 15.25.$$

$$\frac{\bar{y} - \bar{xy}}{\bar{x}^2 - (\bar{x})^2}$$

$$\Rightarrow \frac{15.25 - 5 \times 2.5}{7.5 - 6.25} = 2.2$$

$$x_0 = \bar{y} - \bar{x}\bar{x} = 5 - 2.2 \times 2.5 = -0.5$$

$$y = 2.2x - 0.5.$$



1 Did you perform any data processing steps?

Yes, handled missing value by filling them with the column mean. Also applied label encoding to categorical (like "State") and scaled numerical features for 1000-companies.csv to normalize the data.

2 Did you visualize the regression line for canada-per-capita-income.csv?

Yes, the regression line was plotted. The plot shows a strong linear relationship between year & per-capita income, meaning that as the year increases, per capita income also rises.

3 Predicted salary for (12 years experience, 10 test score, 10 interview score)?

By

The predicted salary is printed in the script and depends on the trained model's coefficient.

4 Did you encode categorical variables for 1000-companies.csv?

Yes, the "State" column was encoded using LabelEncoder().

By /



Did you scale the feature? If yes, why?

Yes, because of R & D Spend, Administration, & Marketing Spend have different units, feature Scaling (using Standard Scaler()) was applied to improve model performance.

21/03/25

## Lab-4

Logistic Regression :- : (1)

- (1). Student - passed or fail based on the study  
 Intercept  $a_0 = -5$  Inefficient  $a_1 = 0.8$

(a)  $P(z) = \frac{1}{1+e^{-z}}$   
 $z = a_0 + a_1 x$   
 $z = -5 + 0.8x$

Student has studied for 5 hours

$P(z) = \frac{1}{1+e^{-(-5+0.8 \cdot 5)}}$   $\rightarrow 0.6456$

- (b) Probability that a student who studies for 5 hours will pass the exam

$P(z) = \frac{1}{1+e^{-5+0.8x}}$   $\rightarrow 0.6456$

- (c) Determining predicted class ( $P(F)$ ) for student based on threshold 0.5

As  $0.6456 > 0.5$  this statement will Pass.

- Q.2. Consider  $z \in [2, 1, 0]$  for 3 classes. Apply softmax function to find probability value of 3 classes.

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^3 e^{z_j}}$$

$$P(z) = \frac{e^2}{e^0 + e^1 + e^2} = 0.6456$$

$$P(1) = \frac{e^1}{e^0 + e^1 + e^2} = 0.0900$$

$$P(0) = \frac{e^0}{e^0 + e^1 + e^2} = 0.2448$$

For - HR comma - sep. csv.

Q.1 Which variable had clear and direct impact on employee retention? Why?

Sol.: satisfaction level, time spent company, number project, average monthly hours.

Reason: They showed strong correlation with retention based on EDA.

Q.2 Accuracy of logit regression model? Good?

Accuracy (eg n 79%)

Yes it has good accuracy, it is above baseline, through improvement may be possible.

For Zoo dataset.

i) Did you perform data processing steps?

Dropped animal name as it's not useful for classification.

ii) Were there many missing values?

No missing values detected.

iii) What does the confusion matrix tell you about the performance?

Shows class-wise performance.

It shows how many instances of each

Class were correctly or incorrectly classified.

It helps assess accuracy and identify area of miscalculation.

iv) Which class type were most frequently misclassified? Why?

Class types with similar features (e.g. mammal & humans) were most frequently misclassified. This happened due to overlapping features and insufficient distinction between classes.

Decision Tree:

branch out of P

$$\text{Entropy}(T) = \sum_{P_i} P_i \log_2 P_i$$

$$\text{Entropy}(T, A) = \sum_{i=1}^n \text{Information entropy}(A_i)$$

$$\Delta G(A) = F(T) - F(T, A)$$

$$\Delta S_k = k [F(T) - F(T, A)]$$

$$E(S) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5}$$

information - store - save - recall

loss of consistency occurs over course of

process of minimization of entropy over time

initially loss occurs every step of

~~initially~~  $\frac{1}{5}$  ~~loss~~ of entropy

information loss right after initial step  
is called initial information

however if we wait until after next step  
information loss is zero

information loss is zero after next step

- Q. Consider following dataset,  $k=3$  and test data  $(x, 35, 100)$  is ari (Person, Age, Salary) also solve using kNN classifier and predict target value given

Person	Age	Salary	k	Target	Dist
A	35	100	3	N	52.81
B	23	150	3	N	46.87
C	24	120	3	N	31.95
D	41	60	3	Y	40.44
E	43	70	3	Y	31.04
F	38	100	3	Y	60.07

Nearest neighbours  $\rightarrow$  G, C, D

Majority class  $y_i$  given

if scored close  $\in f(Y, f(x))$

then weight with  $\delta(f(Y, f(Y(i))) + \delta(f(Y, f(I)))$

else given toward  $f(Y, f(D))$

combine  $\delta(f(Y, Y))$   $\delta(f(Y, N))$   $\delta(f(Y, Y))$

repeat times until 10 iterations then

working at 10 have  $\Sigma D$  in to repeat

iterations number into both, predict

initially zero weight for all

Scoring Note:  $\delta(N, Y) + \delta(N, N)$

then  $\delta(N, Y) = f(N, Y)$

total minimum  $0 + 1 + 0.2 = 1$

### Target - 7

Ques. Explain how to choose the k value for Iris dataset. Which test can be used to choose the k value? Demonstrate using accuracy rate and error rate.

Ans. Cross validation is used to find the best k value depending on accuracy or number of misclassification.

5 folds chosen.

Ques. For diabetes dataset.  
What is the purpose of feature scaling?  
How to perform it?

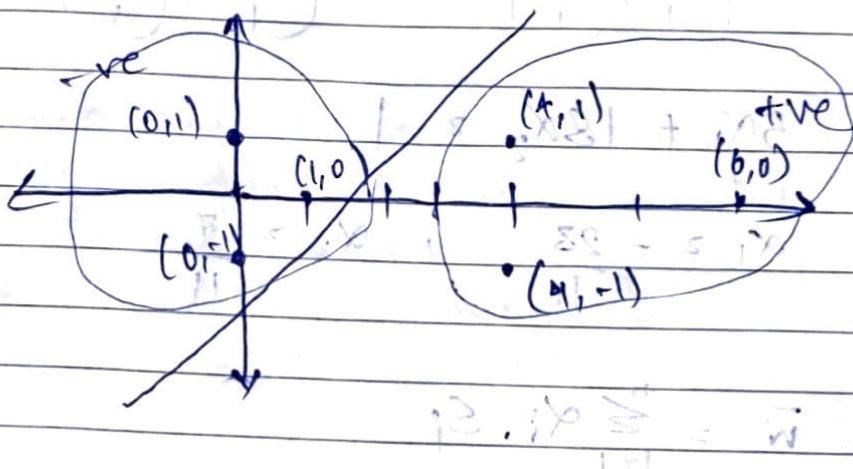
Ans. Feature scaling is used in KNN to have the same scale because if we use different scale, the features with larger scale get more weightage so the following lower scale features will get suppressed and won't have an impact on the model. To perform scaling, first the frequency distribution of all features were checked. The one having normal distributions were scaled using standard scaler and the rest using minmax scaler.

8/1/25  
8/1/25

# ab-6. - Support Vector Machines

Date: \_\_\_\_\_  
Page No.: \_\_\_\_\_

Points  $(4, 1)$ ,  $(4, -1)$  and  $(6, 0)$  belong to positive class 1 and points  $(1, 0)$ ,  $(0, 1)$  and  $(0, -1)$  belong to negative class -1. Draw an optimal hyperplane.



$$\begin{pmatrix} S_1 \\ S_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 4 \\ 1 \end{pmatrix}$$

$$\alpha_1 \begin{pmatrix} S_1 \\ S_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \alpha_2 \begin{pmatrix} S_2 \\ S_2 \end{pmatrix} = \begin{pmatrix} 4 \\ 1 \end{pmatrix} \Leftarrow \alpha_2$$

~~$$\alpha_1 S_1 + \alpha_2 S_2 = -1$$~~

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 1 \end{pmatrix} = -1$$

$$2x_1 + 5x_2 = -1$$

→ constraint (a) has  $(1, 0)$ ,  $(0, 1)$  as its

$(1, 0)$   $\alpha_1 S_1, S_2, (0, 1) \alpha_2 S_2, S_2, \text{ and } 1$  as its

constraint is  $x_1 \leq 0$  and  $x_2 \leq 0$  of constraint

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} + x_2 \begin{pmatrix} 4 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \leq 1$$

$$5x_1 + 18x_2 = 1$$

$$x_1 = -\frac{23}{11}, \quad x_2 = \frac{7}{11}$$

$$\bar{w} = \sum_{i=1}^n \alpha_i \cdot S_i$$

$$= \alpha_1 S_1 + \alpha_2 S_2$$

$$= -\frac{23}{11} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \frac{7}{11} \begin{pmatrix} 4 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} 5/11 \\ 7/11 \\ -16/11 \end{pmatrix} \quad \text{Orientation line}$$

$$b = -\frac{16}{11}$$

$$b + \frac{16}{11} \geq 0$$

$$m = \tan \theta$$

$$\theta = \tan^{-1}(m)$$

$$\text{Orientation line} = \tan^{-1}\left(\frac{5}{7}\right)$$

line cuts  $x$ -axis at  $16/11$ .

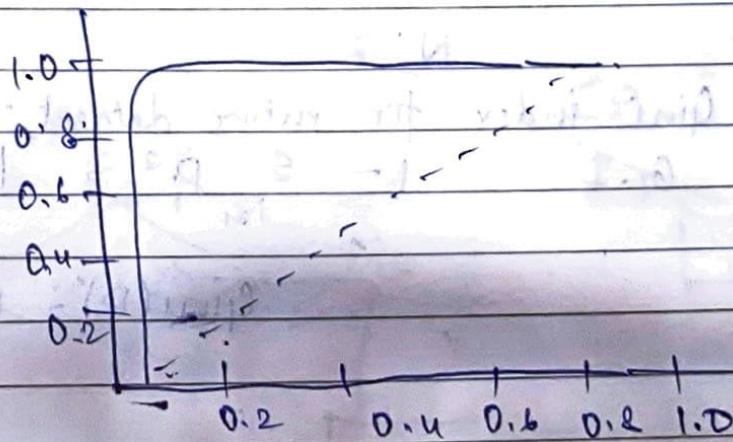
compare accuracy of kernel. both linear RBF kernel gave 100% accuracy on the test set.

Which one preferred & why?

Both preferred equally. The IRIS dataset is relatively simple & linearly separated that is why the linear kernel performs just as well as RBF.

Letter Recognition Dataset

Accuracy = 0.9055 AUC = 1.00



Oct/25

Lab - 7 ..

## Random Forest ensemble

- Q1) For sample S1 draw decision tree  
considering CGPA as root node.

S.No	CGPA	Interactions	Comm	Known	Tg
1	$\geq 9$	Y	Good	Good	Y
2	$< 9$	N	Moderate	Good	Y
3	$\geq 9$	N	Arg	Arg	N
4	$\geq 9$	N	M	M	Arg
5	$\geq 9$	Y	M	Good	Y

As root node is CGPA we can separate into unique groups i.e.,  $\geq 9$  &  $< 9$ .

For CGPA  $\geq 9$  : Y = 2, N = 2

Y : 2

N : 2

Gini index for entire dataset: Y = 3, N = 2

$$G.I = 1 - \sum_{i=1}^2 P_i^2 = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2$$

$$\begin{aligned} Gini(D) &= 1 - 0.36 - 0.16 \\ &= 0.48. \end{aligned}$$

For CGPA  $\geq 9$  : Y : 2, N : 2

$$G(\geq 9) = 1 - (0.5)^2 = (0.5)^2 = 0.5$$

for CGPA  $\leq 9 \rightarrow$  1 sample

$$y = 1, N = 5, \\ \text{Gini}(\leq 9) = 0.$$

$$\text{Weighted Gini}_{\text{CGPA}} \rightarrow \frac{4}{5} \times 0.5 + \frac{1}{5} \times 0 = 0.4$$

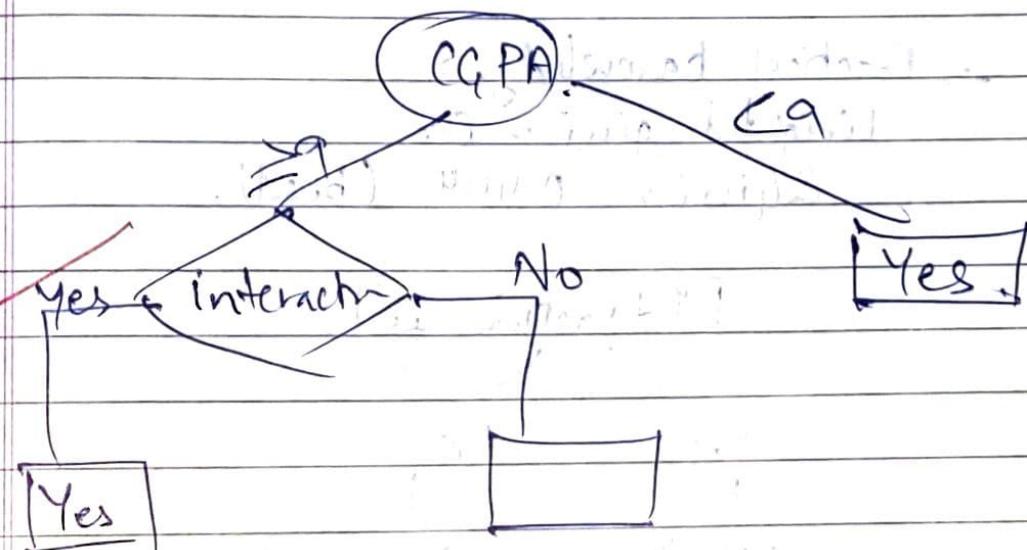
$$\text{Information Gain (Gini reduction)} \\ = 0.48 - 0.4 = \underline{0.08}$$

Split CGPA  $\geq 9$  further to get 8

Split on interaction  $y = 1, \text{Yes} = 2$   
 $N = N = 2$

$$\text{Gini} = \frac{2}{4} \times 0 + \frac{2}{4} \times 0 = 0$$

$$(\text{Gini reduction}) \approx 0.5 - 0.4 = 0.5$$



$$1. \text{ Gini} = 1 - \left(\frac{4}{8}\right)^2 - \left(\frac{4}{8}\right)^2 \Rightarrow 1 - 0.64 - 0.04 \\ \Rightarrow 0.32.$$

$$2. \text{ Gini} = \frac{2}{5} \times 0 + \frac{3}{5} \times 0.444 \Rightarrow 0.2664.$$

$$\Delta \text{Gini} = 0.32 - 0.2664 = 0.0536.$$

3. growing the "no" branch with other features.

→ CGPA :

$$\text{Weighted gini} \geq \frac{1}{3} \times 0 + \frac{2}{3} \times 0.50 \Rightarrow 0.33$$

$$\Delta \text{Gini} \geq 0.11.$$

→ Practical knowledge:

$$\text{Weighted gini} \geq 0.$$

$$\Delta \text{Gini} \geq 0.444 \text{ (best).}$$

[Interactioniveness]

Yes

No

"Yes"

[Practical knowledge]

good

Average

Yes

No



Date : \_\_\_\_\_

Page No. : \_\_\_\_\_

RF	n-estimate	mean accuracy.
1	10	0.9667
2	50	0.9667
3	100	0.9667
4	150	0.9667
5	200	0.9667
N	A	A
P	A	A

Accuracy  $\geq 1.0$

Best number of tree = 20

confusion matrix  $\rightarrow \begin{bmatrix} 10, 1, 0, 1, 0 \\ 0, 9, 0 \\ 0, 0, 11 \end{bmatrix}$

05/05/25.

~~Lab - 8~~

## Boosting Ensemble method

CRPA	Infrac	Practi	com	Jol
$\geq 9$	Y	G	G	Y
$< 9$	N	G	M	Y
$\geq 9$	N	A	M	N
$< 9$	N	A	G	N
$\geq 9$	Y	G	M	Y
$\geq 9$	Y	G	M	Y

Initial weight  $w_0$ .

$$\sum_{j=1}^6 f_j(\text{adj}).w_j(\text{adj})$$

Pre	Ach
Y	Y
N	Y
Y	N
N	N
Y	Y

$$\rightarrow 2 \times \frac{1}{6} = 0.33$$

$$\alpha_0 = \frac{1}{2} \ln \left( \frac{1 - E_{\text{CRPA}}}{E_{\text{CRPA}}} \right)$$

$$\alpha_{\text{CRPA}} = \frac{1}{2} \ln \left( \frac{1 - 0.33}{0.33} \right)$$

$$\alpha_{\text{CRPA}} = 0.347$$

$$\rightarrow \frac{1}{6} \times 4 \times e^{-0.347} \rightarrow \frac{1}{6} \times 2 \times e^{0.347}$$

$$\rightarrow 0.9428$$

$$\text{new} = 0.1249$$

$$0.2501$$



Interaction with others - II

Activity	Y	N	with others	0.250
Interact N (F, 2)	(Y)	(Y)	with others	0.280
W	N			0.125
Y	Y			0.125
Y	Y			0.125

interaction of both methods - do

Group A (F, 2) (1, 1)

Add back with - estat 2.11

1.712

1	15.8	10	1	119182037
2	6.9	50	11.1	(4.0) 8300
3	12.8	100	12.8	(11.0) 83050
4	0.9	150	19.5	(4.0) 8322
5	2.0	200	11	(2.20) 8326.
6	20.2	12	2	(2.2.11)
7	5.9	92	11	(2.11.2.1)

confusion

[ 7028

1386 ]

1243

· (112)

turn right

8/

2.2.11, 2.2.11, 2.2.11, 2.2.11

2.11.2 + 2.11, 2.2 + 2.2 + 2.2 + 2.2 = 8

11

+

F.S. X. 89.16 =

12/05/25

Lab - 9

- Q. Compute two clusters using k-means algorithm for clustering whose cluster centers are  $(1, 1)$  &  $(5, 7)$  execute for 2 iteration.

→ Iteration 1:

Calc Euclidean dist to centroids:

Record	Close to C <sub>1</sub>	Close to C <sub>2</sub>	Assignment
$(1, 1)$	0	7.21	C <sub>1</sub>
$(1.5, 2)$	1.12	6.12	C <sub>1</sub>
$(3, 4)$	3.61	3.61	C <sub>1</sub>
$(5, 7)$	7.21	0.0	C <sub>2</sub>
$(3.5, 5)$	4.12	2.5	C <sub>2</sub>
$(4.5, 5)$	5.31	2.06	C <sub>2</sub>
$(2.5, 4.5)$	4.30	2.92	C <sub>2</sub>

New centroids:

$$C_1 = \frac{1+1.5+3}{5}, \frac{4+2+1}{3} \rightarrow 1.83, 2.33$$

$$C_2 = \frac{5+3.5+4.5+3.5}{4}, \frac{7+5+5+4.5}{4} \\ = 4.12, 5.37$$

Iteration 2 Record	(1.83, 2.33)	(4.12, 5.37)	Assignment
	Close to C <sub>1</sub>	Close to C <sub>2</sub>	
(1, 1)	1.57	5.37	C <sub>1</sub>
(1.5, 2)	0.47	4.27	C <sub>1</sub>
(2, 4)	2.04	1.77	C <sub>2</sub>
(5, 7)	5.67	1.85	C <sub>2</sub>
(2.5, 5)	3.15	0.72	C <sub>2</sub>
(4.5, 5)	3.78	0.53	C <sub>2</sub>
(3.5, 4.5)	2.74	1.07	C <sub>2</sub>

to intab with assignment

∴ New cluster are

$$C_1 = \{R_1, R_2\}, C_2 = \{R_3, R_4, R_5, R_6\}$$

New Centroids

$$C_1 = \frac{2.5}{3} , \frac{3}{3} \quad | \quad C_2 = \frac{19.5}{5} , \frac{25.5}{5}$$

For Iris dataset

The elbow plot (in the next) shows a sharp elbow at k=3, indicating that the cluster is the normal choice for petal length and width of the Iris.

1. 1.1 3.0 1.0 / 1.5 1.6  
 2. 2.8 2.2 1.3 1.4 / 2.0 3.0 1.5 1.6

12/05/25

DATA

## Lab 11

(See PCA 3.1)

2. Write down

1) Mean vector  $\bar{x}$  of the data

Bar chart

2) Covariance matrix  $S$

(1,1)

Q. Reduce elimination from 2 (to 1).

1)  $\bar{x}_1 = 4 + 8 + 13 + 7 + 9 \over 5 = 9$

(4,5)

2)  $\bar{x}_2 = 11 + 4 + 5 + 14 + 8 \over 5 = 10$

(2,2)

3)  $\bar{x}_3 = 1.5 + 2.5 + 3.5 + 4.5 + 5.5 \over 5 = 3.5$

(1,2,1)

4)  $\bar{x}_4 = 1.5 + 2.5 + 3.5 + 4.5 + 5.5 \over 5 = 3.5$

(3,1,2,2)

Standardize the dataset

Mean  $\bar{x} = 9$ , Standard deviation  $s = 3$ .

$$\bar{x}_1 = 4 + 8 + 13 + 7 + 9 \over 5 = 9$$

$$\bar{x}_2 = (11 + 4 + 5 + 14) \over 5 = 10$$

$$\bar{x}_3 = 1.5 + 2.5 + 3.5 + 4.5 + 5.5 \over 5 = 3.5$$

$$\bar{x}_4 = 1.5 + 2.5 + 3.5 + 4.5 + 5.5 \over 5 = 3.5$$

$$X_C^T = \begin{bmatrix} -4 & 0.5 & 1.5 & 6.5 \\ 0.5 & -4.5 & -3.5 & 6.5 \end{bmatrix}$$

5) Compute (Standardized) mean, standard deviation

6) Compute Covariance matrix for standardized variables

7) Find the linear combination of variables such that

$$C = \left( \frac{1}{n-1} \right) X^T X$$

$$= \frac{1}{3} \begin{bmatrix} 1 & 0 & 5 & -1 \\ 2.5 & 4.5 & -3.5 & 8.5 \end{bmatrix}$$

$$B = \frac{1}{3} \begin{bmatrix} -4 & 0 & 6 & 1 & -7 \\ 2.5 & -4.5 & 1.5 & 3.5 & 5.5 \end{bmatrix} \begin{bmatrix} -4 & 1 & 2.5 \\ 0 & -4.5 & 5 & -3.5 \\ -1 & 1 & 5.5 \end{bmatrix}$$

$$C = \frac{1}{3} \begin{bmatrix} 42 & 18 & 23 \\ -33 & 1 & 69 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$C = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

$$\text{Let } (C - \lambda I) = 0, \text{ then}$$

$$\begin{vmatrix} 14 - \lambda & -11 \\ -11 & 23 - \lambda \end{vmatrix} = 0.$$

$$(14 - \lambda)^2 + 121 = 0, \text{ so } \lambda = 13$$

$$\begin{bmatrix} 14 - 2\lambda & 30.3849 & 11 \\ 30.3849 & 23 - \lambda & 0 \\ 11 & 0 & 1 - \lambda \end{bmatrix} = 0, \text{ so } \lambda = 13$$

Consider larger  $\lambda$

$$\begin{bmatrix} 14 - 30.3849 & -11 & x \\ -11 & 202.4 + 23 - 30.3849 & y \\ 11 & 0 & 1 - \lambda \end{bmatrix} \begin{bmatrix} x \\ y \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} -16.3849x & 202.4 + 11y & 0 \\ -11x & 11 + 23 - 30.3849y & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$-16.3849x - 11y = 0$$

$$x = -0.6713y$$

taking  $y = 1$

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -0.6713 \\ 1 \end{bmatrix}$$

Normalize

$$x = \frac{x}{\sqrt{x^2 + y^2}} = \frac{-0.6713}{\sqrt{(-0.6713)^2 + 1^2}}$$

$$c_1 = \begin{bmatrix} -0.5574 \\ 0.8303 \end{bmatrix}$$

Calculate principle component.

$$Z = X_0 c_1 = \begin{bmatrix} -4 & 0 & 2.5 \\ 0 & -4.5 \\ 5 & 3.5 \\ -1 & 5.5 \end{bmatrix} \begin{bmatrix} -0.5574 \\ 0.8303 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -4.305 \\ -3.736 \\ -5.693 \\ -5.124 \end{bmatrix}$$

Flora

## Accuracy before PCA

logistic regression : 0.9016

SVM : 0.8525

Random forest : 0.8361

## Accuracy After PCA

logistic : 0.8689

SVM : 0.8689.

Random forest : 0.88.