

5/3/25

Lab 0

70 00

- ① Initialise values directly into dataframe
import pandas as pd

```
data = {
    'USN': ['IBM22CS100', 'IBM22CS200',
            'IBM22CS300', 'IBM22CS400',
            'IBM22CS500'],
    'Name': ['Alice', 'Bob', 'Charlie', 'David',
              'Eve'],
    'Marks': [85, 78, 92, 88, 76]
}
```

df = pd.DataFrame(data)

print(df)

O/P	USN	Name	Marks
0	IBM22CS100	Alice	85
1	IBM22CS200	Bob	78
2	IBM22CS300	Charlie	92
3	IBM22CS400	David	88
4	IBM22CS500	Eve	76

- ② Importing datasets from sklearn.datasets

from sklearn.datasets import load_diabetes

dia = load_diabetes()

df = pd.DataFrame(dia.data, columns=dia.feature_names)

df.head()

O/P	age	sex	bp	s1	s2	s3	s4	s5	s6
0	0.03	0.05	0.06	0.02	-0.04	-0.4	-0.0	-0.01	-0.017

③ Importing datasets from a specific .csv file
 file = 'sample-sales-data.csv'
 df = pd.read_csv(file)
 df.head()

	Product	Quantity	Price	Sales	Region
O/P					
0	laptop	3	1000	5000	North
1	mouse	15	20	300	West

④ Download dataset from repositories

file = "dataset of diabetes.csv"

df = pd.read_csv(file)

df.head()

O/P	ID	No. of children	Gender	AGE	Urea	G	HbA1c	Chol	TG
0	302	17975	F	30	4.7	46	4.9	4.2	0.9
	HbC	LDL	VLDL	BMI	CLASS				
	2.9	1.4	0.5	24.0	N				

To Do : Stock Market Data Analysis

import yfinance as yf

import pandas as pd

import matplotlib.pyplot as plt

tickers = ['HDFCBANK.NS', 'ICICIBANK.NS',
 'KOTAKBANK.NS']

data = yf.download(tickers, start = "2024-01-01",
 end = "2024-12-31", group_by = 'tickers')

data.head()

print(data.shape)

print(data.columns)

```
hdfc_data = data ['HDFC BANK. NS']  
print (hdfc_data.describe ())  
hdfc_data ['Daily Return'] = hdfc_data ['Close']. pct -  
change ()  
plt.figure (figsize = (12, 8))  
plt.subplot (2, 1, 1)  
hdfc_data ['Close']. plot (title = "HDFC - Closing Price",  
color = 'blue')  
plt.ylabel ('Closing Price')  
plt.subplot (2, 1, 2)  
hdfc_data ['Daily Returns']. plot (title = "HDFC -  
Daily Returns", color = "red")  
plt.ylabel ('Daily Returns')  
plt.tight_layout ()  
plt.show ()
```

Repeat for

ICICI

and

KOTAK

Banks

BSE

EOS

in NSE

5/3/23

Lab 1

To do

```
import pandas as pd  
data = pd.read_csv('housing.csv')  
print(data.info())  
print(data.describe())  
print(data['Area Proximity'].value_counts())  
print(data.isnull().sum() [data.isnull().sum() > 0])
```

Op area - proximity

ISLAND	9136
INLAND	6551
NEAR OCEAN	2658
NEAR BAY	2290

Op total - bedrooms - 207

Sum 0.63 m

To Do - Diabetes Dataset

① df = pd.read_csv('diabetes.csv')

② Missing values :

val = df.isnull().sum()

O/P: Series ([], dtype: int64)

③ Handling missing values

num-col = df.select_dtypes(include=['float64', 'int64'])
columns

imputer = SimpleImputer(strategy='mean')

df[num-col] = imputer.fit_transform(df[num-col])

cot-col = df.select_dtypes(include=['object']).
columns

imputer-cot = SimpleImputer(strategy='most-frequent')

df[cot-col] = imputer-cot.fit_transform(df[cot-col])

④ Handling Categorical Attributes

df1 = df.copy()

df1['Gender'] = df1['Gender'].str.upper()

ordinal-encoder = OrdinalEncoder(categories=[['M', 'F']])

df1['Gender-encoded'] = ordinal-encoder.fit_transform
(df1[['Gender']])

onehot-encoder = OneHotEncoder()

encoded-data = onehot-encoder.fit_transform
(df1[['CLASS']])

encoded-array = encoded-data.toarray()

encoded-df = pd.DataFrame(encoded-array,
columns=onehot-encoder.
get_feature_names_out
(['CLASS']))

df-encoded = pd.concat ([df1, encoded df], axis=1, ignore_index=True)

df-encoded . drop (['order', 'class'], axis=1, inplace=True)

⑤ Data Normalization

norm = MinMaxScaler()

df-encoded [['area']] = norm.fit_transform (df-encoded [['area']])

scales = StandardScaler()

df-encoded [['age']] = scales.fit_transform (df-encoded[['age']])

⑥ Removing Outliers

df1 = df-encoded

Q1 = df1['area'].quantile(0.25)

Q3 = df1['area'].quantile(0.75)

IQR = Q3 - Q1

lower = Q1 - 1.5 * IQR

upper = Q3 + 1.5 * IQR

df1['area'] = np.where (df1['area'] > upper, upper,

np.where (df1['area'] < lower, lower, df1['area']))

df2 = df-encoded

df2['area-zscore'] = stats.zscore (df2['area'])

df2['area'] = np.where (df2['area-zscore'].abs() > 3, np.nan,

df2['area'])

$df^3 = df - \text{encoded}$

$df^3['\text{area-zscore}'] = \text{stats.zscore}(df^3['\text{area}'])$

$\text{median-area} = df^3['\text{area}'].\text{median}()$

~~$df^3['\text{area}'] = \text{np.where}(\text{abs}(df^3['\text{area-zscore}']) > 3, \text{median-area}, df^3['\text{area}'])$~~

(Median value is same as mean)

• this way (Median and Mean) = median-area
[0] median-area

• step 4) generation of area + color - feature
or color + intensity of gray

then add total back color - feature after
step area + color

• (center + std) is not
of standard deviation

area (std) point in positive - direction

(positive [center]) std = center - center
value - mean

0 = positive - below

• center - center is center as
standard deviation

[center - [center] std] std = standard
deviation

([center] - [center] std) = positive - below

(below) positive

positive - below = positive - above

12/3/25

Lab 2

⑤ import pandas as pd
from collections import Counter

def entropy(data):
 labels = data['label'].value_counts()
 counts = Counter(labels)
 probabilities = [count / len(labels) for count in counts.values()]

(6) entropy - value = - sum(p * math.log2(p))
 for p in probabilities if p > 0

return entropy - value

def gain(data, feature):

initial - entropy = entropy(data)

feature - values = data[feature].unique()

weighted - entropy = 0

for value in feature - values:

subset = data[data[feature] == value]

weighted - entropy += (len(subset) / len(data)) * Entropy(subset)

return initial - entropy - weighted - entropy

```
def id3(data, features, target_attribute):  
    if len(data['label'].unique()) == 1:  
        return data['label'].iloc[0]  
  
    if len(features) == 0:  
        return data['label'].value_counts().index[0]  
  
    best_feature = max(features, key=lambda feature:  
        gain(data, feature))  
  
    tree = {best_feature: {}}  
    features = [f for f in features if f != best_feature]  
    for value in data[best_feature].unique():  
        subset = data[data[best_feature] == value].  
            drop(columns=[best_feature])  
  
        if len(subset) == 0:  
            tree[best_feature][value] = data['label'].  
                value_counts().index[0]  
        else:  
            tree[best_feature][value] = id3(subset,  
                features, target_attribute)  
  
    return tree
```

import math

~~data = {}~~

df = pd.DataFrame(data)

features = ['outlook', 'temperature', 'humidity', 'wind']

target_attribute = 'label'

decision_tree = id3(df, features, target_attribute)

decision_tree

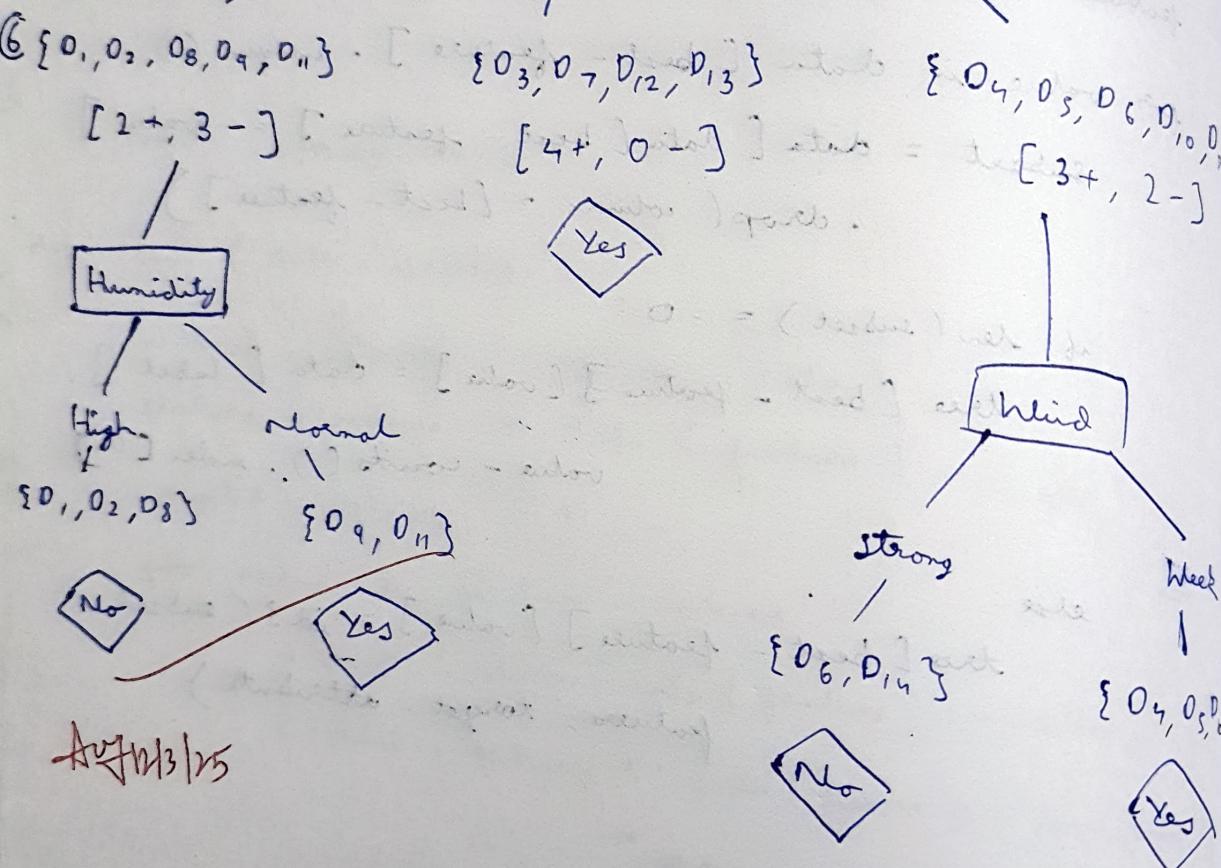
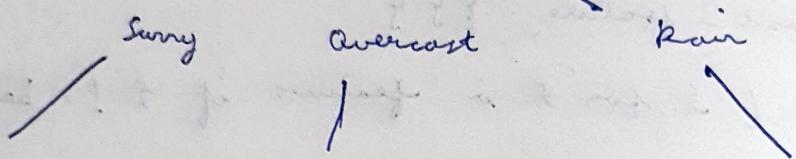
Q11

{'outlook': {'sunny': {'humidity': {'high': 'no', 'normal': 'yes'}}},
 'overcast': 'yes';
 'rain': {'wind': {'weak': 'yes', 'strong': 'no'}}}

(5)

O_1, O_2, \dots, O_{14}

Outlook



Aug 12/3/25

Lab 3

Linear Regression

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
x = np.array([1, 2, 3, 4, 5])
```

```
y = np.array([1.2, 1.8, 2.6, 3.2, 3.8])
```

```
x_b = np.c_[np.ones((x.shape[0], 1)), x]
```

```
theta = np.linalg.inv(x_b.T.dot(x_b)) . dot(x_b.T) . dot(y)
```

intercept, slope = theta

```
x_input = float(input("Enter a value for x to predict y:"))
```

```
x_input_b = np.array([1, x_input])
```

```
y_pred = x_input_b . dot(theta)
```

```
plt.scatter(x, y, color='blue', label='Data points')
```

```
plt.plot(x, x_b . dot(theta), color='red',  
label='Regression Line')
```

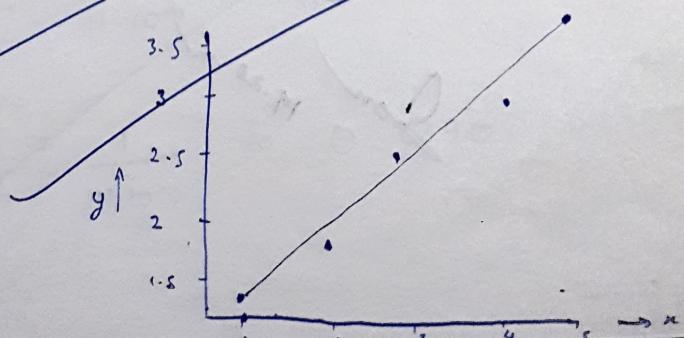
```
plt.xlabel('x')
```

```
plt.ylabel('y')
```

```
plt.title('Linear Regression - Model Fit')
```

```
plt.legend()
```

```
plt.show()
```



Output:

Intercept (theta 0) : 0.540000 .. 025

Slope (theta 1) : 0.6600 .. 01

Enter a value for x to predict y: 7

Predicted value for x = 7.0 : 5.1600 .. 03

Linear regression - Gradient Descent

alpha = 0.01

iterations = 1000

m = len(x)

theta 0 = 0

theta 1 = 0

for - in range (iterations):

$$y_{\text{pred}} = \theta_0 + \theta_1 * x$$

$$\text{error} = y_{\text{pred}} - y$$

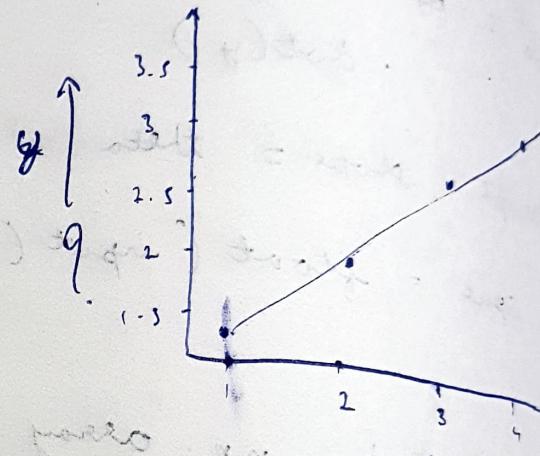
$$\theta_0 = \alpha * (1/m) * \text{np.sum(error)}$$

$$\theta_1 = \alpha * (1/m) * \text{np.sum(error * x)}$$

$$x_input = \text{float}(\text{input}())$$

$$y_{\text{pred}} = \theta_0 + \theta_1 * x_input$$

Sum
19.83 in



Lab 4 - Logistic Regression

① Solution:

a) $\alpha_0 = -5$
 $\alpha_1 = 0.8$

$$P(\text{pass}) = \frac{1}{1 + e^{-(\alpha_0 + \alpha_1 \cdot 7)}}$$

b) $\alpha_0 + \alpha_1 \cdot 7 = -5 + 0.8 \cdot 7 = 0.6$

$$P(\text{pass}) = 0.645$$

c) $P(\text{pass}) > 0.5 \rightarrow$ Here "pass"

② $P(\text{class } i) = \frac{e^{z_i}}{\sum_j e^{z_j}}, z = [2, 1, 0]$

$$\sum_j e^{z_j} = e^2 + e^1 + e^0 \approx 11.107$$

$$e^{z_1} = e^2 = 7.389$$

$$e^{z_2} = e^1 = 2.718$$

$$e^{z_3} = e^0 = 1$$

$$P(\text{class 1}) = \frac{e^2}{11.107} \approx 0.666$$

$$P(\text{class 2}) = \frac{e^1}{11.107} \approx 0.245$$

$$P(\text{class 3}) = \frac{1}{11.107} \approx 0.090$$

1. (i) Satisfaction level has highest impact due to highest correlation factor.
- (ii) It is 77.07% which is not good because it is not a proper correlation
2. (i) Yes, to get all numeric data
- (ii) No, but I dropped it from task if it was present
- (iii) It is very accurate, with almost no false predictions
- (iv) Class label 4 due to inconsistent data & patterns

Lab - 6 - KNN

02-04-25

Person	Age	Salary	target	Distance	Rank
A	18	50	N	52.8	5
B	23	55	N	46.6	4
C	24	70	N	31.9	2
D	41	60	Y	40.5	3
E	43	70	Y	31.1	1
F	38	40	Y	60.1	6

$k = 3$, Closest 3: C, E, D

2 Y's and 1 N $\rightarrow \therefore$ predicted value

$$\text{for } X = (35, 100) = Y$$

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \text{ where}$$

$$(x_1, y_1) = (35, 100)$$

① Iris :

For $k = 5$, accuracy = 0.9667

$k = 3$ error = 0.0333

For $k = 4$, accuracy = 1.000

error = 0.000

For $k = 6$, accuracy = 0.9333
error = 0.0667

1 ② Diabetes

It is important as it is distance based.
Otherwise cannot be scaled properly.

Standard scalar : $Z = \frac{x - \mu}{\sigma}$

μ = mean

σ = std.

x = feature value

Gen

value mapping \rightarrow unit normal

Conversion

value, float type (x-min) / (x-max) = 0
(0.000) = 0.000

float a = 0.000
0.000 = min
 $a = \frac{a - \text{min}}{\text{max} - \text{min}}$

0.000 = max
 $a = \frac{\text{max} - a}{\text{max} - \text{min}}$

0.000 = min
 $a = \frac{a - \text{min}}{\text{max} - \text{min}}$

16/04/25

Lab-7 : Random Forest

S1: Build a Random Forest

a) If no. of examples in the train set is N , take a sample of n -examples of random but with replacement, from original data. This sample will be training set for generating the tree.

Again out of 100, choose next 30 examples randomly for 2nd p & put them back.

S2: If there are M i/p vars., m vars. are selected at random, out of M , then best split on those m , used to split node. The value of ' m ' held constant during the generation of various trees in forest.

Ex: 10 i/p vars \rightarrow 3

subset of these vars. randomly to build tree

S3: Each tree is grown to the largest extent for new data, point final prediction of each branch & assign new data point to category part with majority votes, i.e. wherever a new data pt. is given, give that DT to the particular BT

Class
16.04.25

7/5/2025

Lab 8 - Adaboost

① Assign each item

= 1/6 weight

CGPA	Interactions	Behavioral knowledge	Communication	Job Profits
$>= 9$	Yes	Good	Spill	Yes
< 9	No	Good	Wood	Yes
< 9	No	Avg	Moderate	No
$>= 9$	No	Avg	Good	No
< 9	Yes	Good	Moderate	Yes
$>= 9$	Yes	Good	Moderate	Yes
$>= 9$	Yes	Good	Moderate	Yes

CGPA	Predicted Job Offer		Actual Job Offers	Weight
	$>= 9 \rightarrow$ Yes	$< 9 \rightarrow$ No		
$>= 9$	Yes		Yes	1/6
< 9	No		Yes	1/6
$>= 9$	Yes		No	1/6
< 9	No		No	1/6
$>= 9$	Yes		Yes	1/6
$>= 9$	Yes		Yes	1/6

$$\epsilon_{\text{CGPA}} = 2 \times \frac{1}{6} = 0.333$$

$$\alpha_{\text{CGPA}} = \frac{1}{2} \frac{\ln(1 - 0.333)}{0.333} = 0.347$$

$$\begin{aligned} Z_{\text{CGPA}} &= \frac{1}{6} \times 4 \times e^{-0.347} + \frac{1}{6} \times 2 \times e^{0.347} \\ &= 0.9928 \end{aligned}$$

$$\begin{aligned} \text{wt}(d_j)_{i+1} &= \frac{\frac{1}{6} + e^{-0.347}}{0.9928} = 0.1239 \end{aligned}$$

$$wt(\alpha_j)_{j+1} = \frac{\frac{1}{6} \times e^{0.349}}{0.9428} = 0.2501$$

<u>Interactiveness</u>	Interactiveness	Predictor Job offer	Actual Jobs Offer	Weight
Yes → Yes				0.1249
No → No	Yes	Yes	Yes	0.2501
	No	No	No	0.2501
	No	No	No	0.1249
	No	No	No	0.1249
	Yes	Yes	Yes	0.1249
	Yes	Yes	Yes	0.1249

$$\bar{x}_{\text{Interact}} = 1 \times 0.2501 = 0.2501$$

$$\alpha_{\text{Interact}} = \frac{1}{2} \ln \left(\frac{1 - 0.2501}{0.2501} \right) = 0.5490$$

$$\begin{aligned} z_{\text{Interact}} &= 0.1249 + 4 \times e^{-0.549} + 0.2501 + 1 + e^{-0.549} \\ &= 0.1249 + \frac{0.2501 + 1 + e^{0.549}}{e^{0.549}} \\ &= 0.866 \end{aligned}$$

$$wt(\alpha_j)_{j+1} = \frac{0.1249 \times e^{-0.549}}{0.866} = 0.0832$$

$$wt(\alpha_j)_{j+1} = \frac{0.2501 + e^{-0.549}}{0.866} = 0.1667$$

~~$$wt(\alpha_j)_{j+1} = \frac{0.2501 \times e^{0.549}}{0.866} = 0.5001$$~~

Practical knowledgeGood \rightarrow YesModerate \rightarrow No

Average

No instances misclassified,
no need decision stumpCommunication

	Comm Skill	Bredicted Job Offer	Actual Job Offer	Weight
Good \rightarrow Yes	Good	Yes	Yes	0.0832
Moderate \rightarrow No	Moderate	No	No	0.5001
	Moderate	No	No	0.1667
	Good	Yes	No	0.0832
	Moderate	No	Yes	0.0832
	Moderate	No	Yes	0.0832

$$\epsilon_{cs} = 1 \times 0.5001 + 3 \times 0.0832 = 0.7497$$

$$\alpha_{cs} = \frac{1}{2} \frac{\ln(1 - 0.7497)}{0.7497} = -0.5485$$

$$Z_{cs} = 0.0832 \times 1 \times e^{-(-0.5485)} + 0.1667 \times 1 \times e^{-(-0.5485)} + 0.5001 \times 1 \times e^{-0.5485} + 0.0832 \times 3 \times e^{-(-0.5485)}$$

$$Z_{cs} = 0.866$$

$$wt(\alpha_j)_{j+1} = \frac{0.0832 \times e^{-(-0.5485)}}{0.866} = 0.1663$$

$$wt(\alpha_j)_{j+1} = \frac{0.1667 \times e^{-(-0.5485)}}{0.866} = 0.3331$$

$$wt(d_i)_{i+1} = \frac{0.5001 + e^{-0.5485}}{0.866} = 0.3337$$

$$wt(d_i)_{i+1} = \frac{0.0822 \times e^{-0.5485}}{0.8666} = 0.0555$$

Final

$$H_F(d_i) = \alpha_{CPA} \times \text{Yes} + \alpha_{\text{Interact}} \times \text{Yes} + \alpha_{CSL} \times \text{Yes}$$

$$= 0.347 \times 1 + 0.5490 \times 1 - 0.5485 \times 1 = 0.34$$

SI No.	α_{CPA}	α_{Interact}	α_{CSL}	Weight Avg.	Final Prediction
1	Yes	Yes	Yes	0.3475	Yes
2	No	No	No	0.0	No
3	Yes	No	No	0.347	Yes
4	No	No	Yes	-0.5485	No
5	Yes	Yes	No	0.896	Yes
6	Yes	Yes	No	0.896	Yes

(2) Best accuracy = 0.83 or 83%.

No. of trees = 160.

Confusion matrix

$$= \begin{bmatrix} 10658 & 551 \\ 2023 & 1521 \end{bmatrix}$$

7/3/23

Lab 9 - K mean clustering

Q Answer

	A	B
R1	1.0	1.0
R2	1.5	4.0
R3	3.0	7.0
R4	5.0	5.0
R5	3.5	5.0
R6	4.5	4.5
R7	3.5	

$k=2$, centroid for $C_1 = (1, 1)$,
 " " $C_2 = (5, 7)$

Iteration 1

Record No.	Close to $C_1(1, 1)$	Close to $C_2(5, 7)$	Assignment
R1(1, 1)	dist(R1, C1) = 0.0	dist(R1, C2) = 7.21	(1)
R2(3.5, 2)	dist(R2, C1) = 1.12	dist(R2, C2) = 6.12	(1)
R3(3, 4)	dist(R3, C1) = 3.61	dist(R3, C2) = 3.61	(1)
R4(5, 7)	dist(R4, C1) = 7.21	dist(R4, C2) = 0.0	(2)
R5(3.5, 5)	dist(R5, C1) = 4.12	dist(R5, C2) = 2.5	(2)
R6(4.5, 5)	dist(R6, C1) = 5.31	dist(R6, C2) = 2.06	(2)
R7(3.5, 4.5)	dist(R7, C1) = 4.30	dist(R7, C2) = 2.92	(2)

$$C_1 \{R1, R2, R3\} = \frac{1 + 1.5 + 3}{3}, \quad \frac{1 + 2 + 4}{3}$$

$$\begin{aligned} &= \frac{5.5}{3}, \quad \frac{7.0}{3} \\ &= (1.83, 2.33) \end{aligned}$$

$$C_2 = \frac{(5 + 3.5 + 4.5 + 3.5)}{4}, \quad \left(\frac{7+5+5+4.5}{4} \right)$$

$$= \frac{16.5}{4}, \quad \frac{21.5}{4}$$

$$= (4.12, 5.37)$$

Iteration 2

record No. close to $C_1 (1.83, 2.33)$ close of $C_2 (4.12, 5.37)$ Assign to cluster

$R_1 (1, 1)$ dist $(R_1, C_1) = 1.57$ dist $(R_1, C_2) = 5.37$ C_1

$R_2 (1.5, 2)$ dist $(R_2, C_1) = 0.47$ dist $(R_2, C_2) = 4.27$ C_1

$R_3 (3, 4)$ dist $(R_3, C_1) = 2.04$ dist $(R_3, C_2) = 1.77$ C_2

$R_4 (5, 7)$ dist $(R_4, C_1) = 5.64$ dist $(R_4, C_2) = 1.85$ C_2

$R_5 (3.5, 3)$ dist $(R_5, C_1) = 3.15$ dist $(R_5, C_2) = 0.72$ C_2

$R_6 (4.5, 5)$ dist $(R_6, C_1) = 3.78$ dist $(R_6, C_2) = 0.53$ C_2

$R_7 (3.5, 4.5)$ dist $(R_7, C_1) = 2.74$ dist $(R_7, C_2) = 1.07$ C_2

$$C_1 = \frac{1+1.5}{2}, \quad \frac{1+2}{2} \quad C_2 = \frac{3+5+3.5+4.5+3.5}{5},$$

$$= 1.25, 1.5$$

$$\frac{7+5+5+4.5}{5}$$

$$= (3.4, 5.1)$$

$C_1 \{ R_1, R_2 \}$

$C_2 \{ R_3, R_4, R_5, R_6, R_7 \}$

Lab 10 - PCA

heart.csv

Report accuracy:

	Before PCA	After PCA
SVM acc.	0.1648	0.1848
Logistic regression acc	0.1685	0.1739
Kardon Everest	0.1848	0.1738

Guaranteed