

## NOTES

### Ass 6-1

a) Read CSV file :-

- import pandas as pd :-
  - It import pandas library.
  - It is a powerful data manipulatio<sup>n</sup> & analysis library in python.

-  $df = pd.read\_csv(\text{"file path"})$

↓  
variable.

↓  
It is a funct<sup>n</sup> provided by pandas to read csv file into data frames.

b) df. Shape :-

- Attribute of a pandas Data Frame. it returns a tuple representing the dimensions of the Data Frame.
- 1<sup>st</sup> is row & 2<sup>nd</sup> is columns.  
          ↑                                  ↑  
          element                        element.

\* Shape :- layout of data.

c) df. dtypes :-

- It returns a series containing the data types of each column in the Data Frame.



a) - `df.isnull().sum()`   
 - Count of missing values in each column of the data frame

- `df['Thal'].mode()`   
 - Find the common values in Thal column of the Data Frame

- `df['Thal'].fillna(df['Thal'].mode(), inplace=True)`   
   
 select the 'Thal' column from DataFrame.   
 It is a method. Used to fill missing value with specified values   
 inplace = True   
 This parameter updates the DataFrame directly into string changing the 'Thal' column in the original DataFrame   
 It calculate the mode of the 'Thal' column & convert into string

- `inplace = False`   
 - It omitted, it return a new DataFrame with the filled values.

- `df['Ca'].fillna(df['Ca'].mean(), inplace=True)`   
   
 It calculates the mean value of the 'Ca' column.



\* Inconsistency :- Same data fill in the different table & diff Format

Page No.	
Date	

b) `df.drop_duplicates(inplace=True)`

- Removing Inconsistency
- We need to remove duplicates.

c) - `import matplotlib.pyplot as plt` :-  
- It import MATLAB-like plotting framework.

- `import numpy as np` :-

- It import numpy library
- It is used for mathematical operations & creating arrays.

- `plt.boxplot(df.Age)` :-

- It creates a box plot for the 'Age' column of the Data Frame

- `plt.show()` :-

- It display the box plot.

d) - `import seaborn as sns` :-

- It import seaborn library.
- It is a visualization library.

- `sns.histplot(df.Age, kde=True)`

It is used  
to plot a  
histogram

It is data that  
you want  
to plot a  
histogram

It estimate the  
data's probability  
density funct<sup>n</sup>



\* Histogram - It is a graphical representation of the distribution of numerical data.

g) - mean\_age = df['Age'].mean() - Calculate (mean)

f) zeros\_count = (df == 0).sum()

- It creates a Dataframe with 'True' if it is equal to 0 & 'False' otherwise. It contains the value 0.

→ Add the no of times each column in the Dataframe.



## Ass. 2.

Read the file

```
- import pandas as pd  
df = pd.read_csv('file path')  
print(df)
```

- df.dtypes # find the data types

- import numpy as np

- numeric\_df = df.select\_dtypes(include=[np.number])  
- print(numeric\_df)

include only  
numeric data

This line creates  
a new Dataframe

a) # Find Std. dev & Variance of every numerical attribute.

```
- Std_dev = numeric_df.std() → Calculate Std. dev  
print(Std_dev)
```

```
- Var = numeric_df.var() → Calculate Var  
print(Var)
```

b) # Find Covariance & perform Correlation analysis the using Correlation Coefficient.

```
- Covariance_matrix = numeric_df.cov()  
print(Covariance_matrix)
```

```
- Correlation_matrix = numeric_df.corr()  
print(Correlation_matrix)
```



c) # How many independent features are present in the given dataset?

- independent\_features = np.linalg.matrix\_rank  
(numeric\_df.corr())

- It is a function.  
- It calculates the rank of a matrix, which is the maximum no. of independent columns or rows in the matrix.

It calculates the correlation matrix of the numeric columns in the DataFrame.

- print (independent\_features)

d) # Unwanted Features. - creates an empty set to hold names of features that might be redundant due to high correlation.

```
loop iter. unwanted_features = set()
ates      → for i in range(len(correlation_matrix.columns)):
overs column
indices in    for j in range(i):
correln      if abs(correlation_matrix.iloc[i,j]) > 0.8:
matrix        colname = correlation_matrix.columns[i]
              unwanted_features.add(colname)
              print(unwanted_features)
```

Loop checks previous column indices to avoid repeating comparison. → Checks if correlation coefficient bet<sup>n</sup> columns exceeds 0.8

→ Assigns the name of the column to a variable, representing a highly correlated feature.

Adds the name of the highly correlated features.



e) # Perform data discretization using equi-Frequency binning method on the age attribute

num - bins = 5  
df['age - bins'] = pd.qcut(df['Age'], q = num - bins, labels = False)

*(Perform equi-Frequency binning on attribute of the data-frame)*

- Set the Variable
- It indicate that we want to divide to data into 5 bins

Accesses or creates a column named 'age - bins' in DataFrame 'df' to store discretized value

The no. of quantiles or bins into which to divide the data

o/p. bin indices, not labels, as integers starting from 0

Func<sup>n</sup> from Pandas library Perform quantile-based discretization

```
import matplotlib.pyplot as plt
plt.figure(figsize=(5,3))
plt.hist(df['Age'], bins=num - bins, edgecolor='black', alpha=0.7)
plt.title('Equi-Freq. Binning for Age attribute')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.show()
```



- Data bet<sup>n</sup> 0 & 1 by subtracting the minimum value & dividing by the range.

F) From sklearn.preprocessing import MinMaxScaler, StandardScaler

It transforms data to have mean 0 & std. dev. 1

attributes\_to\_normalize = ['RestBP', 'Chol', 'MaxHR']

# Min-Max normalization

min\_max\_scaler = MinMaxScaler() <sup>object</sup>

create duplicate of the DataFrame  $\Leftarrow$  df\_min\_max\_normalized = df.copy()

df\_min\_max\_normalized[attributes\_to\_normalize] = min\_max\_scaler.fit\_transform(df[attributes\_to\_normalize])

df\_min\_max\_normalized[attributes\_to\_normalize]

Apply Min-Max normalization to selected columns & stores the normalized values in 'df\_min\_max\_normalized'

Selects specific columns from 'df\_min\_max\_normalized'



# # Z - Score Normalization

Z - Score - Scaler = StandardScaler ()

df - z - score - normalized = df . copy ()

df - z - score - normalized [attributes - to - normalize] =  
z - score - scaler . fit - transforms  
(df [attributes - to - normalize])

df - z - score - normalized [attributes - to - normalize]

$10^{**}(\text{len}(\dots)-1)$  Calculate decimal  
Scaling Factor.

decimal - Scaling - Factor =  $10^{**}(\text{len}(\text{str}(\text{int}(\text{df}$

Counts  $\leftarrow$  [attributes - to - normalize].abs().max().max()))  
digits in

the

max

value

by

converting

it to

integers &

a string,

& calculates

the length

of that

string.

df - decimal - scaled = df . copy ()

df - decimal - scaled [attributes - to - normalize] =

df [attributes - to - normalize] / decimal - Scaling - Factor

df - decimal - scaled [attributes - to - normalize]



## ASS 4

### \* WEKA :-

- It stand for Waikato Environment For knowledge Analysis.
- It is open source sw.
- It provides tools for data preprocessing.
- It is used to preprocess the big data, apply different machine learning algo. on big data & compare various o/p.
- This sw is easy to work with big data & train a machine using machine learning algo.

### \* Advantage :-

- Testing new ideas quickly.
- Easy to learn & solve data mining problem.
- Code is not required; everything can be done through the UI.

### \* Disadvantage :-

- It can only handle small datasets.
- Lack of algo. for deep learning.
- The complication in converting the string into numeric math.
- Alternative algo. are not available.

### \* Clustering :-

- It is groups similar data points together based on characteristics.

Counts digits in the max. value by converting it to integer & a string, & calculate the length of that string.



## Supervised & independent data.

Page No.	
Date	

### \* Simple k-means Clustering :-

data is dependent on oper<sup>n</sup>

- k-means clustering is an unsupervised learning algo.

- It is used to solve the clustering problems in machine learning or data science.

- k defines the no. of pre-defined clusters that need to be created in the process.

### \* Advantage :-

- Simplicity & ease to use
- Scalability
- Efficiency & speed

### \* Disadvantage :-

- Not suitable for non-linear data
- Dependence on no. of clusters
- Sensitive to outliers

### \* Hierarchical clustering :-

- It creates a hierarchical representation of the clusters in a dataset.

- It is an unsupervised learning

- It starts by treating each data points as an individual cluster

### \* Advantage :-

- Simple to implement
- Gives the best op in some cases.



\* Disadvantage :-

- Breaks the large cluster.
- Sensitive to noise & outliers
- Difficult to handle different sized clusters.



## Ass. 5

### \* Naive Bayes :-

- It is a collection of classification algo based on Bayes Theorem.
- It is a family of algo, where all of them share common principles.

### \* Random Forest :-

- It is a powerful tree learning technique in machine learning.
- It is a creating a no. of decision trees during the training phase.
- The greater no. of trees in the forest leads to higher accuracy.

### \* Decision tree :-

- It is a supervised learning technique.
- It is used both classification & Regression problem.
- But mostly it is preferred for solving classification problem.



## Ass. 6

\* Theory Same as Ass. 5

It is used to  
encode categorical  
labels with  
numerical values

# Read the dataset

module

```
- From sklearn.preprocessing import LabelEncoder  
encoder = LabelEncoder()  
data['chestPain'] = encoder.fit_transform(data['chestPain'])
```

print(data) *pick column from the Data Frame*

- It create instance.  
- Used to encode the categorical data.

- Categorical labels into  
numerical values using  
Label Encoder

- Update the Data Frame accordingly

```
- data['Thal'] = encoder.fit_transform(data['Thal'])  
print(data).
```

# Same for AHD

```
- From sklearn.model_selection import train_test_split  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.svm import SVC  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.metrics import accuracy_score
```



- data.isnull().sum()

- data['Ca'].fillna(data['Ca'].mean(), inplace=True)

- x = data.drop(columns=['AHD'])

↳ Creates new DataFrame by dropping column, named 'AHD', from DataFrame.

- y = data['AHD']

- x\_train, x\_test, y\_train, y\_test = train\_test\_split(x, y, test\_size=0.2, random\_state=42)

↳ Variable storing training & testing sets for features (x) & target variable (y).

↳ 20% of data used for testing.

Set random seed for reproducibility.

- dt\_classifier = DecisionTreeClassifier()

- svm\_classifier = Svc()

- knn\_classifier = KNeighborsClassifier()

Svc =

Support Vector Classifier

- It is a specific implement of the support vector machine algo.

- It is a supervised learning algo.

- It makes classification based on data neighbors.



- dt\_classifier.fit(x\_train, y\_train)
- svm\_classifier.fit(x\_train, y\_train)
- knn\_classifier.fit(x\_train, y\_train)
- train the decision tree classifier.
- train the svm (Support vector machine) classifier.
- train the k-nearest neighbours classifier.

- dt\_predictions = dt\_classifier.predict(x\_test)
- svm\_predictions = svm\_classifier.predict(x\_test)
- knn\_predictions = knn\_classifier.predict(x\_test)
- Generate predictions using trained decision tree.
- Same for all.

- dt\_accuracy = accuracy\_score(y\_test, dt\_predictions)
- svm\_accuracy = accuracy\_score(y\_test, svm\_predictions)
- knn\_accuracy = accuracy\_score(y\_test, knn\_predictions)
- calculate accuracy of predictions made by decision tree classifier.
- Same for all.

print(dt\_accuracy)  
print(svm\_accuracy)  
print(knn\_accuracy)

\*\*\*