

## DSBDA

\*enlive

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

Page: \_\_\_

- \* pandas - is a python library used in data analysis. It provides various data structures for manipulating & analyzing data.
  - import pandas as pd
  - df = pd.read\_csv('path')
- \* numpy - (Numerical Python) numpy consists of multi-dimensional array objects & collections to process array. It is used for working with array. It is used to perform math operations (linear algebra, matrix, fourier transform)
  - import numpy as np
- \* Dataframe - tabular data structure. It is 2D tabular data structure with labeled axis (rows & columns).
- \* Transposing - It is used to permute the dimensions of array.
- \* shape & Reshape -
  - shape - To change the size triple of array of dimensions.
  - Reshape - To change dimensions of array into new shape.
- \* melt() - converts data from wide format to long format.

### \* Data cleaning -

It means fixing ~~that~~ bad data.

bad data - empty cells, data in wrong format  
 wrong data  
 duplicates

#### (a) Remove empty cells

- i) new\_df = df.dropna()
- ii) df.fillna(0)

#### (b) Data in wrong format

- i) pd.to\_datetime(df['Date'])

#### (c) Wrong data

```
df.loc[row_no, column_name] = value,
```

#### (d) Duplicates

df.duplicated() # True / False

df.drop\_duplicates(inplace=True)

### \* Data integration -

It is a data processing technique that involves combining data from multiple data sources into a coherent data store & provide unified view of data.

#### a) Read from csv file

```
import pandas as pd
pd.read_csv('heart.csv')
```

#### b) Read from google sheet

```
import pandas as pd
new_google_sheet_url = 'url'
df = pd.read_csv(new_google_sheet_url)
```

## c) Read Data from Database

```

import pandas as pd
import sqlite3
# Read via SQLite database
con = sqlite3.connect("your-db-link")
# Read table via select stmt
player = pd.read_sql_query("SELECT * FROM Table", con)
# close the connection
con.close()

```

## \* Data transformation

The process of converting raw data into a format or structure that would be more suitable to make model. It transforms values.

## a) # Multiply each value by 10 in array.

```

import pandas as pd
import numpy as np
df = pd.DataFrame(np.array([[1, 2, 3],
                           [4, 5, 6],
                           [7, 8, 9]]),
                  columns=['a', 'b', 'c'])
df.transform(func = lambda x: x*10)

```

## b) import pandas as pd

```

df.groupby('C')['A'].mean()
mean = df.groupby('C')[["A"]].mean().rename("N").reset_index()
df_1 = df.merge(mean)
df['N3'] = df.groupby(['C'])['A'].transform('mean')

```

- \* **matplotlib**  
It is a plotting library used to plot graphs.  
`matplotlib.pyplot` provides functions.
  - \* **seaborn**  
It is a library based on `matplotlib` which provides interface to draw attractive & informative statistical graphs.
  - \* **sklearn**  
It is a tool to build statistical models including classification, regressions, clustering etc.

## 1. Data Preparation

\*enlive

Date : / /  
Page:

- \* Dataset - heart.csv
- \* Operations -
  - 1) Find shape of data
  - 2) Find missing values
  - 3) Find data type of each col
  - 4) Find out zero's
  - 5) Find mean age of patient

```
import pandas as pd  
df = pd.read_csv('heart.csv')  
1) df.shape  
2) df.isnull()  
3) df.dtypes  
4) df.isnull().sum()  
5) df['Age'].mean()
```

age: 39.0

(min=29.0, q1=37.0, median=41.0, q3=46.0, max=88.0)

(count=303, std=11.0, var=180.0)

## 4/5. Data Preparation

Date:  
Page:

- \* dataset - facebook & temperature
- \* operation -
  - 1) create data subset
  - 2) Merge data
  - 3) Sort data
  - 4) Transposing data
  - 5) Shape & Reshape data



- ① import pandas pd  
df = pd.read\_csv('fb.csv')
- 1) # first 5 rows  
df.head()
- # last 5 rows  
df.tail()
- # with one column  
new\_df = df['like']
- # with multiple columns  
new\_df = df[['like', 'comment']]
- # condition  
new\_df = df[df['like'] > 35]
- 2) a1 = df[['comment', 'like']]  
a2 = df[['like', 'share']]  
merge = pd.merge(a1, a2)
- 3) df.sort\_values('like', ascending=True)
- 4) df.transpose
- 5) df.shape  
r1 = df.pivot\_table(index=['Type', 'category'],  
values=['comment'])

```
r2 = df.pivot_table(index=['Type', 'category'],
                     columns=[],
                     values=['comment',
                             'like'])

r3 = df.pivot_table(index=['Type', 'category'],
                     columns=[1, 1],
                     values=['comment'])

pd.melt(df)
```

(C) English = 276

(D) Spanish = 226

(E) French = 276

(F) Arabic = 166

(G) Chinese = 176

(H) German = 176

Classification of people by gender = 876

Combination of gender and language = 876 + 1326 = 2202

(I) Female = 1326

Labour share frequent contact = 2202  
 Labour share frequent contact analysis = 2202  
 Labour share frequent contact analysis = 2202

Classification of labour result = @labour <

Labour <

Female <

Male <

(Female, female) + (Male, female) <

Female, female <

Male, female <

## 6/7. Data Preparation

\* Dataset = heart.csv & airquality.csv

\* Operations → 1) Data cleaning

2) Data Integration

3) Data Transformation

• df.loc[row\_no, col\_name] = value

• df['sex'] = df['sex'].replace([0, 1], ['Female', 'Male'])

→ 4) Error correcting  
5) Data model building

1) df.fillna(0)

df1 = df.dropna()

df2 = df1.drop\_duplicates()

2) x = df2['Age'].mean()

y = df2['Age'].median()

z = df2['Age'].mode()[0]

3) df3 = df2.groupby('ca')['Age'].mean()

df4 = df3.rename("NewAge"), reset\_index()

df5 = pd.merge(df2, df4)

5) from sklearn import linear\_model

from sklearn.metrics import mean\_squared\_error

```
> model@ = linear_model.LinearRegression()
```

```
> model
```

```
> xx = df['Age']
```

```
yy = df['chol']
```

```
> model.fit(df[['chol']], df.Age)
```

```
> model.coef
```

```
> model.intercept
```

```
> y-pred = model.predict(df[['chol']])
```

```
> y-pred > print('MSE', mean_squared_error(yy, y-pred))
```

## 8 to 11. Graph

"enlivio"

Date: / /

Page:

\* Dataset: forest fire, temperature, air quality, heart dataset.

\* Operation "graph" and "plotting".

Use matplotlib & seaborn

- 1) Pie chart

2) Bar graph

3) Box plot

4) Line graph

5) Scatter plot



```
import pandas as pd
```

```
import numpy as np
```

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

```
df = pd.read_csv("heart.csv")
```

1) plt.title("Pie chart")

```
new_df = df.head()
```

```
x = new_df.Age
```

```
y = new_df.Thal
```

```
plt.pie(x, labels=y)
```

```
plt.show()
```

2) plt.title("Bar graph")

```
plt.xlabel("AGE")
```

```
plt.ylabel("CHOL")
```

```
plt.bar(x, y)
```

3) plt data = df.Chol

```
plt.title("Box Plot")
```

```
plt.xlabel("Chol")
```

```
plt.boxplot(data)
```

- 4) `x = df['Sex']`  
`y = df['Slope']`  
`plt.title("Line Graph")`  
`plt.xlabel("SEX")`  
`plt.ylabel("SLOPE")`  
`plt.plot(x, y)`
- 5) `plt.title("Scatter Graph")`  
`plt.xlabel("AGE")`  
`plt.ylabel("CHOL")`  
`plt.scatter(x, y)`
- 6) `sns.boxplot(x='Slope', y='chol', data=df)`  
`sns.scatterplot(x='Age', y='chol', data=df)`  
`sns.lineplot(data=df['Age'])`