

Data

Analysis

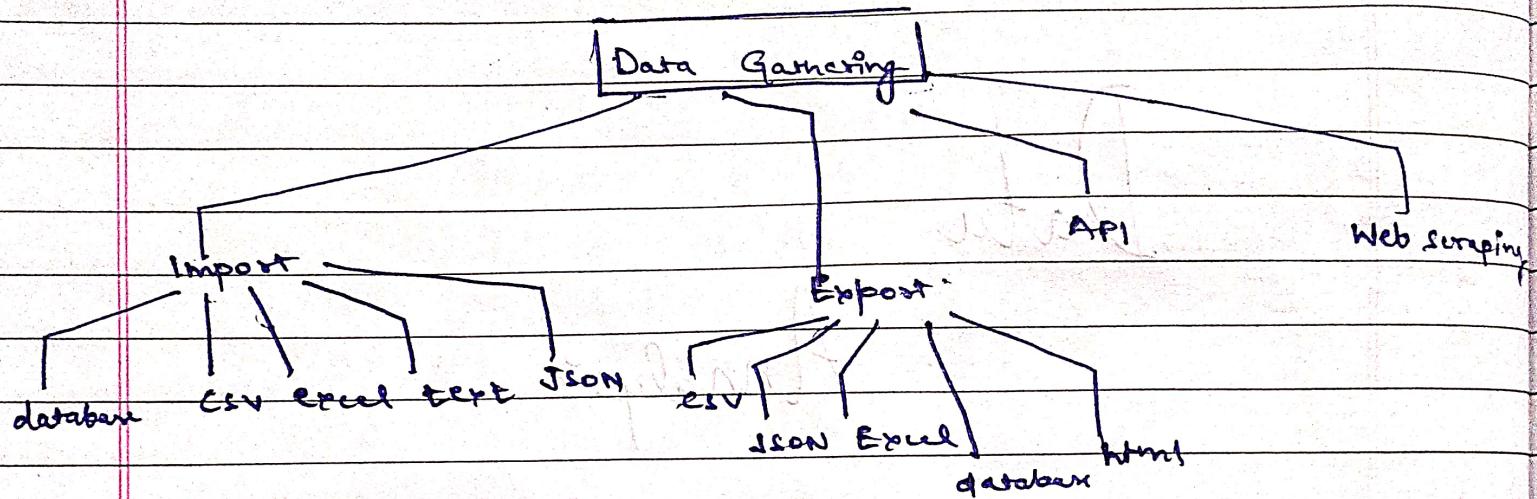
Process

Data Analysis process

1) Data Gathering.

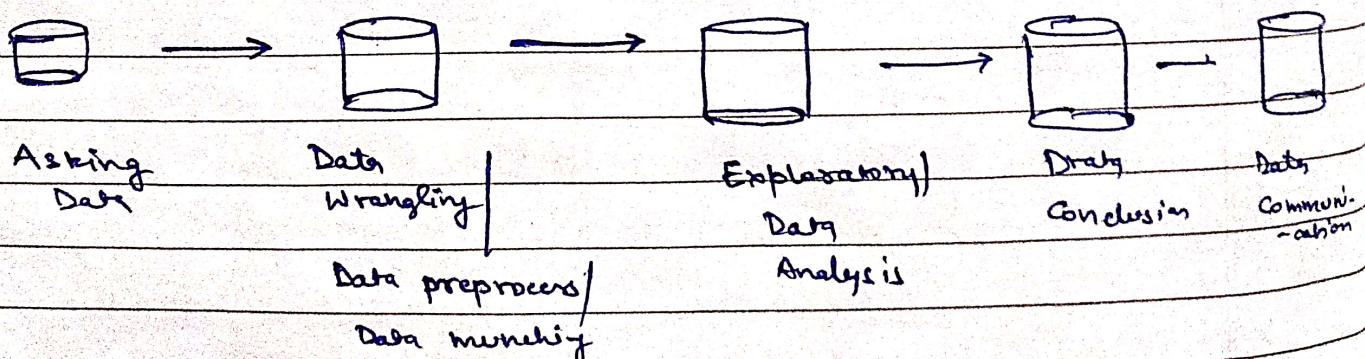
2) Data Cleaning

3) EDA -



Data Analysis is a process of inspecting, cleansing, transforming and modelling data with a goal to discover useful information.

Data Analysis process



Way 1 → get the data

Step 1: Asking Ques'

Way 2 → problem statement is given

Ques1

What feature will contribute to my analysis

Ques2

What feature are not important for my analysis

Ques3

Which of feature have strong correlation?

Ques4

Do I need data preprocessing

Ques5

What kind of feature manipulation/ engineering is required?

How can I ask better question

1) Subject matter expertise

2) Experience

Step 2: Data Wrangling / Munging

process of transforming and mapping data from one "raw" data form into another with an intent of making it appropriate.

1) Gathering data [CSV, API, Web scrapping, Databases]

2) Assessing data [shape, info, is_unique, describe]

3) Cleaning data [

Ass data:

1) Find no. of rows / column (shape)

2) Detg type of various columns (info())

3) Checking for missing values (isna())

4) check for duplicate data (is_unique)

5) Memory occupied by the dataset (info)

6) High level mathematical overview of the data (describ)

Cleaning Data

- 1) Missing Data (e.g. mean)
- 2) Remove duplicate data (drop duplicate)
- 3) Incorrect data type (datatype)

Step 3 Exploratory Data AnalysisExploratory

- 1) Find Correlation & Covariance
- 2) Doing Univariate and multivariate analysis
- 3) plotting graphs (data visualisation)

(feature engineering)

Segment

- 1) Removing outliers
- 2) Merging Dataframes
- 3) Adding new Column

Step 4 Drawing Conclusion

[Take a small sample
try to conclude to entire
(Inferential stats)]

Machine
learning

(prediction)

Data

Science

Conclusion
concludeDescriptive
statistics.Descriptive Stats

- 1) Is Rohit Sharma a better batsman in 2nd innings (IPL)
- 2) Does being a female increase your chance of survival (Titanic)
- 3) Is delhi the most costly place to eat out (Zomato)

Step 5 Communicating Result | Data story telling.

Report

Result

Blog
pos

Presentation

CSV → comma separated value

TSV → tab separated value

document read - CSV pandas

1) `df = pd.read_csv('name')` → Open a local CSV file.

2) Opening a CSV file from a URL

```
import requests
```

```
from io import StringIO
```

```
url = ''
```

```
headers = {"User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10.14; rv: 66.0) Gecko/20100101 Firefox/66.0"}  
req = requests.get(url, headers=headers)
```

```
data = StringIO(req.text)
```

```
pd.read_csv(data)
```

3) Sep parameter

```
pd.read_csv('file.tsv', sep='\t', names=[n1, n2, n3, ...])
```

3-12 यह column के names की हैं

4) Index_col parameter

```
pd.read_csv('file', index_col='column-name')
```

Use data: Index

5) Header parameter

```
pd.read_csv('test.csv', header=1)
```

3-12 column name की 4-5 तक row हैं

6) Use_col parameter

, 3-12 के column की 2-3 हैं।

```
pd.read_csv('file', usecols=['enrol', 'gender', 'edu'])
```

7) Squeeze parameter

```
pd.read_csv('file', usecols=['gender'], squeeze=True)
```

2-3 के column की 1 के series object हैं।

8) skiprows / nrows parameters

`pd.read_csv('file', skiprows=[0, 1])` → $nrows=100$

3112 नींवे row skip करती है।

$nrows \rightarrow$ 3112 के rows load करती है।

9) Encoding parameter

3112 error & UnicodeDecodeError
→

`pd.read_csv('romantic.csv', encoding='latin-1')`

10) skip bad lines

`pd.read_csv('file', sep=',', error_bad_lines=False)`

11) dtypes → column का

3112 इन datatype change करती है।

`pd.read_csv('file', dtype={'target': int})`

12) Handling Dates

3112 default string dates को datetime dtype
करती है।

`pd.read_csv('file', parse_dates=['date']).info()`

3112 date multiple columns को month / year

`pd.read_csv('file', parse_dates=[[1, 3]])`

Imp

13) Converts

`def rename(name):`

`if name == 'Royal Challengers Bangalore':`

`return 'RCB'`

`else:`

Global
name

`pd.read_csv('file', converters={'team': rename})`

14) na_values parameter

```
pd.read_csv('file', na_values=[None])
```

State that string value None Undefined, - will not
Value to convert test to.

15) Load a huge dataset in chunks.

```
dfs = pd.read_csv('aug', chunksize=5000)
```

for chunks in dfs:

[operations]

print(chunk.shape)

Excel

```
pd.read_excel('output.xlsx')
```

```
pd.read_excel('output-file', sheet_name='Sheet-name-2')
```

Data from text

```
pd.read_csv('file', sep='|')
```

JSON

SQL

JSON → Javascript Object Notation
(Universally accepted format)

SQL → Structured query language

documentation!

pd.read_json

pandas.read_sql_query

JSON

1) from local machine

```
data = pd.read_json('train.json')
```

2) from url

```
pd.read_json('url')
```

SQL

World-Cities, Pop, Lang, Rank, (sql create)

Step 1: download dataset

Step 2: download Xamp. Start → Apache

→ MySQL

localhost/phpmyadmin

1) Create a new database

2) upload world.sql file.

Step 3 ! pip install mysql-connector.

Step 4

import mysql.connector

```
conn = mysql.connector.connect(host='localhost', user='root',
                                password='', database='world')
```

```
data = pd.read_sql_query("SELECT * FROM city", conn)
```

Pandas Export

- o to-csv
- o to excel
- o to html
- o to json
- o to sql

→ Google
colab

CSV

```
temp-df = df.groupby('batsman')[['batsman-runs']].sum()  
reset_index()
```

```
temp-df.to_csv('batsman-runs.csv', index=False)
```

to_excel

1) `temp-df.to_excel ("batsman-runs.xlsx")`

2) Save to some another sheet.

```
temp-df.to_excel ("batsman-runs.xlsx", sheet_name=  
= 'batsman-runs')
```

3) Make two different excel sheet in same file

with pd.ExcelWriter ('output.xlsx') as writer:

```
temp-df.to_excel (writer, sheet_name='One')
```

```
temp-df2.to_excel (writer, sheet_name='Two')
```

to_html

```
df.query('batsman-runs == 6').pivot_table(index='over', columns=  
= 'balls', values='batsman-runs', aggfunc='count')  
.html ('size-heatmap.html')
```

to-JSON

df.to_json('file')

```
df.groupby('Team').sum()
df.to_json('file')
```

to-SQL

Step 1: Create a new database con using Xamp

Step 2: Create table

import pymysql

from sqlalchemy import create_engine

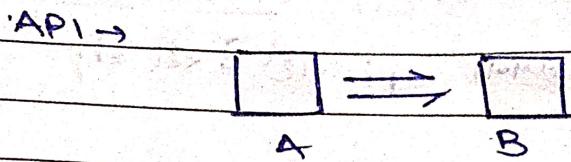
```
engine = create_engine("mysql+pymysql://root@localhost/ipl")
```

{root:{password}@mysql}/{database}

```
df.to_sql('ipl-delivery', conn=engine, if_exists='append')
```

Table name

Indirect data → इसका एक योगी विवरण वेबसाईट से मिलता है।



- 4) API pipelines जो data को A से B तक ले जाती हैं OR
 - 5) API को software के किनीय में लाते रखती हैं

website rapid api

Step-1 tndb api → google

Step-2 find wt. \rightarrow $\langle \text{aff-key} \rangle$

Create account

Profile > setting > api key

Step-3 paste full website on json viewer.

Step - 4

~~Wadsworth~~ - Hart, Chelmsford 1933 pages → 1 - 458

~~payy - 428
total movies - 8551~~

Step-5 ~~format~~ code

id (134892-358)

۴۸

Release date
10-10-2018

popularity

(cont'd.)

000175

• import pandas

import request

卷之三

```
df = pd.DataFrame()
```

for i in range (1, 429):

• `format('')`

temp_df = pd.DataFrame(response.json()['result'])
temp_df = temp_df[['id', 'tH2L-']]

~~df = df.append (temp_df, ignore_index=True)~~

elsewhere

`df_tutorialmovies.csv')`

Web scraping.

3 site की website से API होती है।

```
import pandas as pd
```

```
import requests
```

beautiful

soup → from bs4 import BeautifulSoup

```
headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/88.0.4324.150 Safari/537.36'}
```

```
webpage = requests.get('https://www.ambitionbox.com/list-of-companies?page=1', headers=headers).text
```

```
soup = BeautifulSoup(webpage, 'lxml')
```

```
company = soup.find_all('div', class_=['company-content-wrapper'])
```

```
name, rating, reviews = [], [], []
```

```
→ ctype, hq, old, employees = [], [], []
```

```
for i in company:
```

```
    name.append(i.find('h2').text.strip())
```

```
    rating.append(i.find('p', class_='rating').text.strip())
```

```
    reviews.append(i.find('a', class_='review-count').text.strip())
```

```
ctype.append(i.find_all('p', class_='info Entity')[0].text.strip())
```

```
d = {'name': name, 'rating': rating, 'reviews': reviews, 'ctype': ctype, 'hq': hq, 'old': old, 'employees': employees}
```

```
df = pd.DataFrame(d)
```

```
df
```