

# Pandas

18 February 2025 07:36

## Advanced Concepts in Python

- 1) Decorators
- 2) generators
- 3) iterators

## Pandas

- 1) Data Analyst interview
- 2) PySpark →
- 3) Data scientist →



### What is Pandas?

Pandas is an **open-source Python library** used for data manipulation and analysis. It provides **data structures** like **Series** (1D) and **DataFrame** (2D) to handle structured data efficiently.

### Where is Pandas Used?

Pandas is widely used in **data science, analytics, and machine learning** for:

- ✓ Data Cleaning & Preprocessing – Handling missing values, duplicates, and formatting data.
- ✓ Data Analysis – Performing statistical calculations and data summarization.
- ✓ Data Visualization – Plotting data using built-in functions and integration with Matplotlib & Seaborn.
- ✓ Data Transformation – Merging, reshaping, and filtering large datasets.
- ✓ Time Series Analysis – Handling time-indexed data like stock prices or IoT data.
- ✓ Big Data Handling – Works well with CSV, Excel, SQL databases, and JSON files.

## Data Manipulation

Empid	Emp Name	Empsal
123	Suman Singh	10000
345	Kanna Jithu	20000
987	Ravi Kohli	30000

=>

Empid	Emp Name	Empsal
123	Purn	10000
345	Uma	20000
987	Leela	30000

### Structured :-

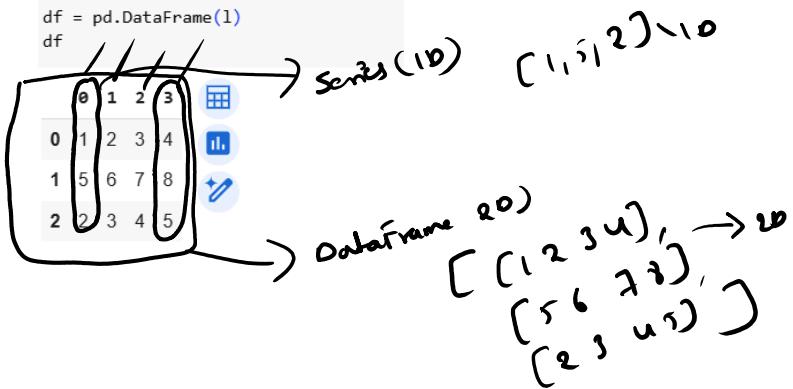
- 1) Excel ✓
- 2) CSV ✓
- 3) Text in DB (SQL)
- 4) MySQL ✓

### Unstructured :-

- 1) MongoDB
- 2) Cassandra
- 3) AWS S3
- 4) Video ✓
- 5) Images

{ - - -  
y

```
l=[[1,2,3,4],[5,6,7,8],[2,3,4,5]]  
df = pd.DataFrame(l)
```



# Few Important Methods

19 February 2025 07:28

.loc & .iloc

.set\_index()

.apply()

.reset\_index()

----> Modify the Runtime column

----> Modify the Gross column

.sort\_values()

.rename()

.value\_counts()

.De-nesting Genre for top 10 movies based on Meta\_score

.duplicated()

.

.isnull()

.notnull()

.fillna()

.dropna()

.drop\_duplicates()

.drop()

.loc & .iloc

→ Both are used to filter / Select / update row3 | columns in a dataframe.

diff .loc & .iloc

→ .loc uses labels

→ .iloc uses indexes / positions

Syntax

df.loc[  
  ↓   ↓  
  rows   columns]

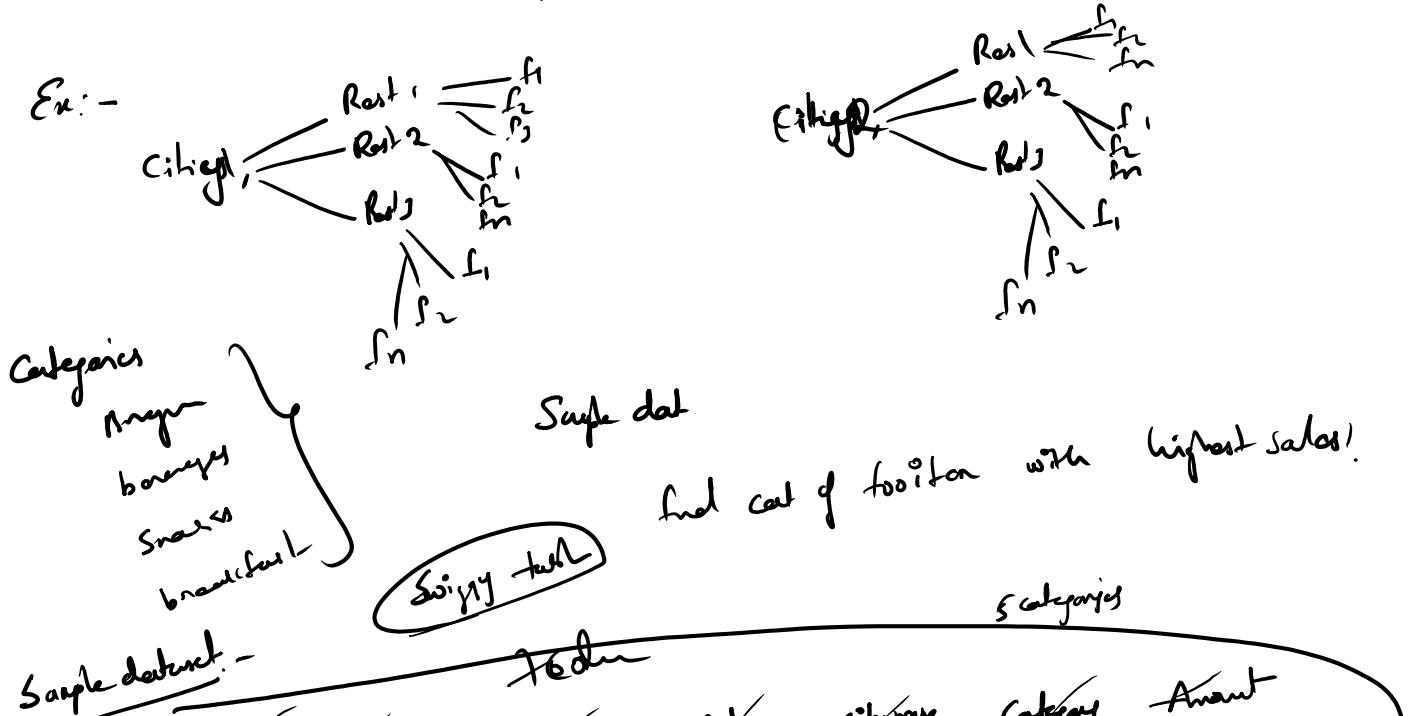
df.iloc[  
  ↓   ↓  
  row   columns]

## groupby

24 February 2025 08:02

→ helps in dividing dataset into groups, on top of it you can perform transformations/operations on each group.

Ex:-

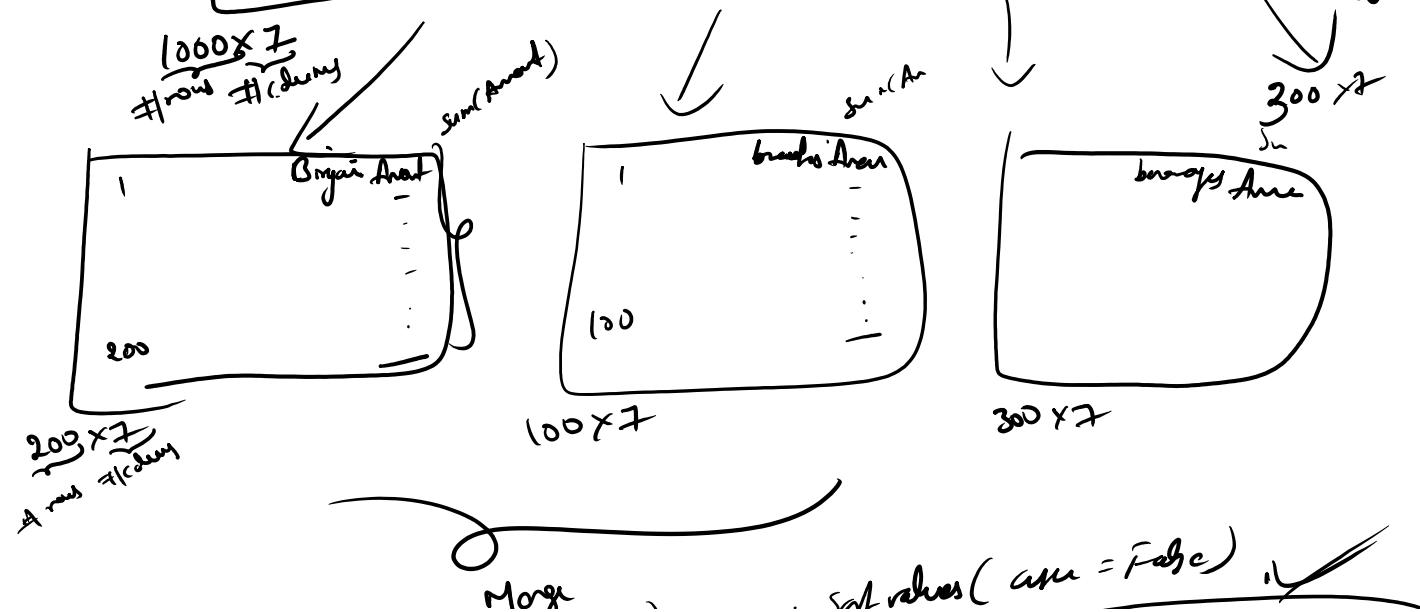


Sample dataset:-

7 columns

S.No	Food Item	Rest Name	Rating	city name	Category	Amount
1	$f_1$	$R_1$	9.5	Kolkata	Biryani	250
2	$f_2$	$R_1$	8.5	Mumbai	breakfast	90

100 rows



Merge  
 about (sum)  
 5,00,000

oranges  
 2,00,000

breakfast  
 6,00,000

Snacks  
 1,00,000

CS  
 7,50,000

(Sum => copy => subg) => transformation

sort values ( axis = False )

breakfast 6,00,000

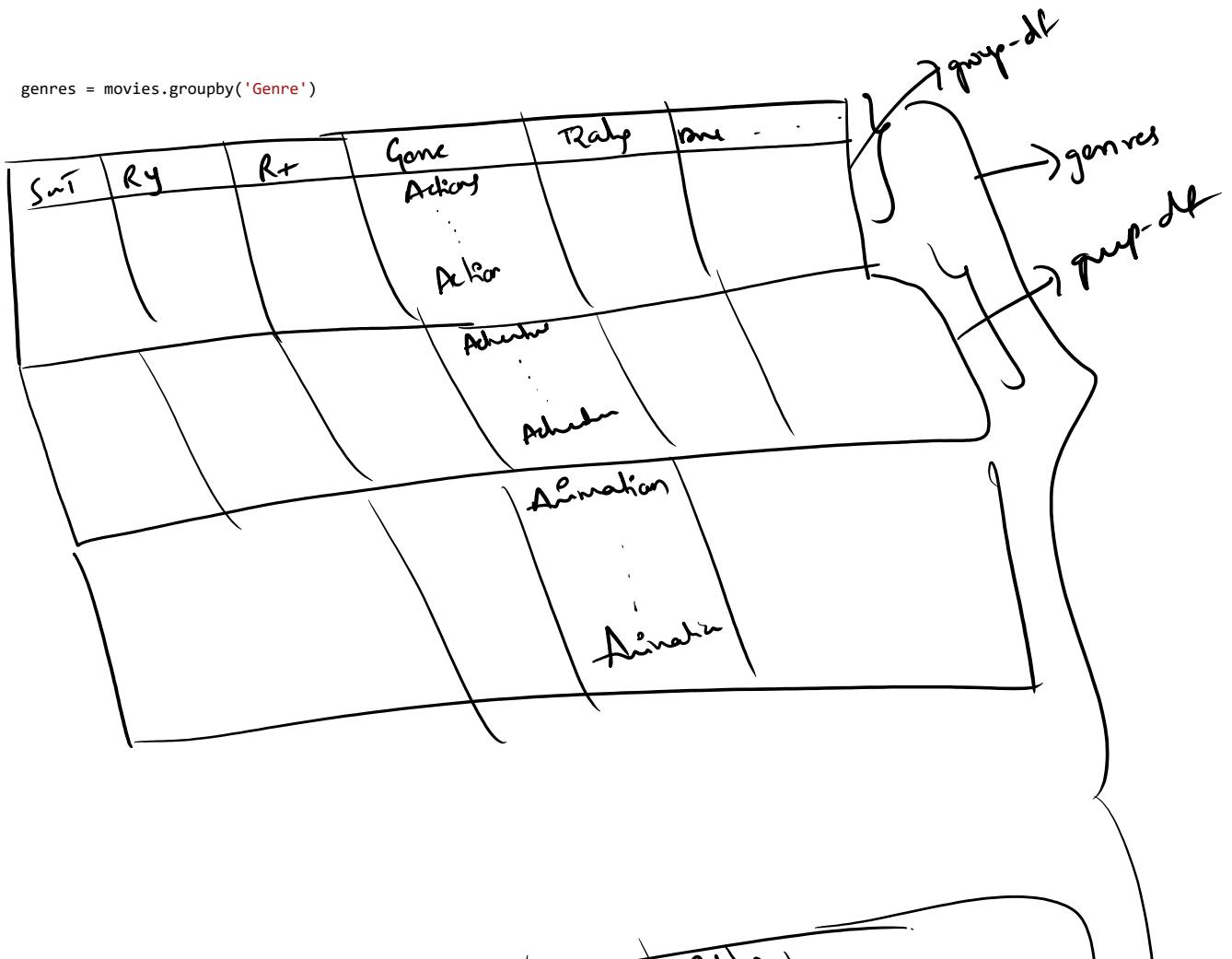
oranges 2,00,000

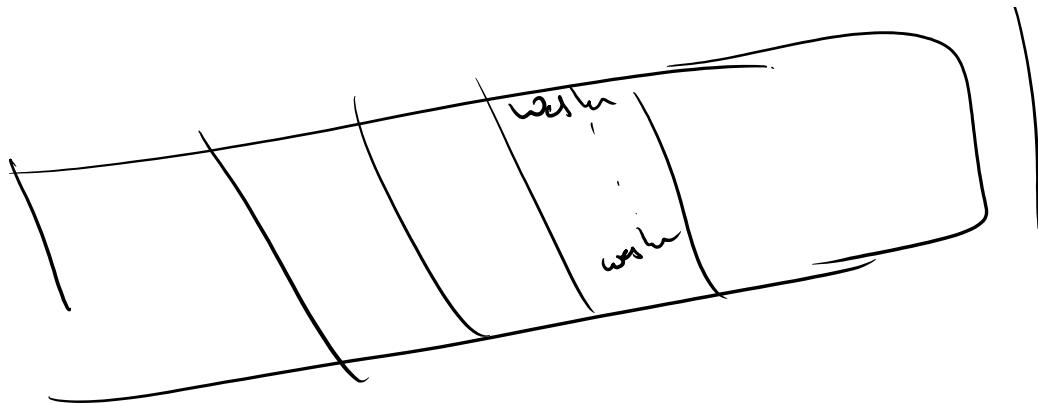
Snacks 1,00,000

CS 7,17,000

	Series_Title	Released_Year	Runtime	Genre	IMDB_Rating	Director	Star1	No_of_Votes	Gross	Metascore	
0	The Shawshank Redemption	1994	142	Drama	9.3	Frank Darabont	Tim Robbins	2343110	28341469.0	80.0	info
1	The Godfather	1972	175	Crime	9.2	Francis Ford Coppola	Marlon Brando	1620367	134966411.0	100.0	
2	The Dark Knight	2008	152	Action	9.0	Christopher Nolan	Christian Bale	2303232	534858444.0	84.0	
3	The Godfather: Part II	1974	202	Crime	9.0	Francis Ford Coppola	Al Pacino	1129952	57300000.0	90.0	
4	12 Angry Men	1957	96	Crime	9.0	Sidney Lumet	Henry Fonda	689845	4360000.0	96.0	

```
genres = movies.groupby('Genre')
```





`movies.groupby("Director")["No_of_Votes"].sum()`

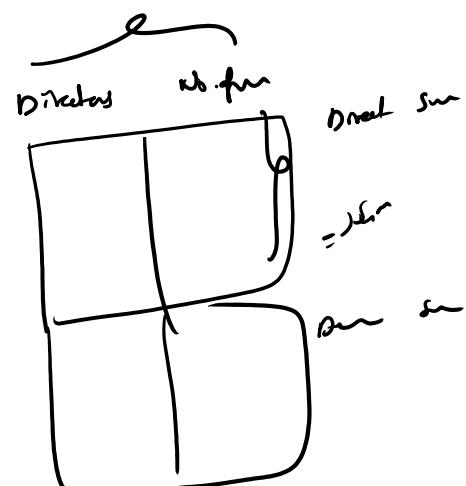
$1000 \times 2$

10 Directors

			No_of_Votes	Director
1				...
2				...
				...
				...
				...
				...
				...
				...
				...

Churn

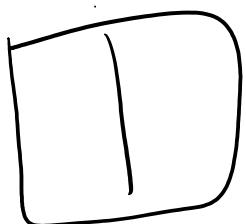
=>



10

Days

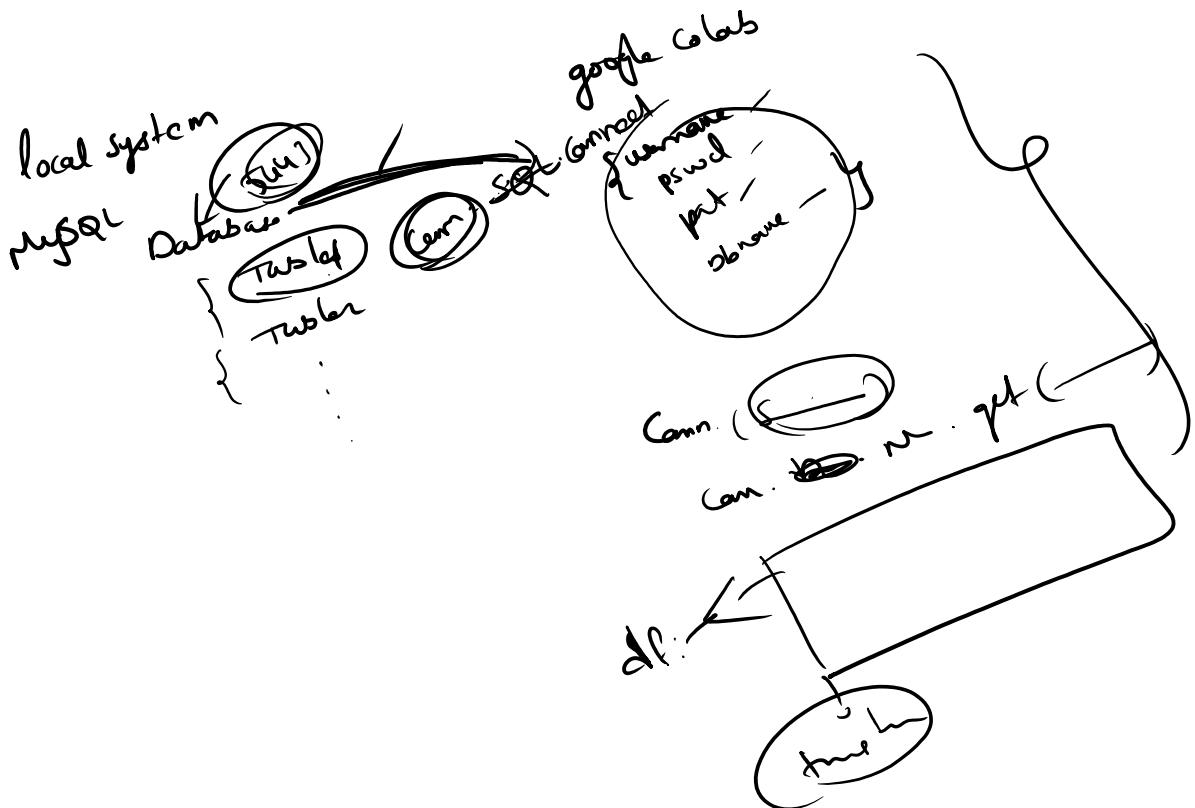
=>



$10 \times 1$

den	No_of_Votes
D <sub>1</sub>	-
D <sub>2</sub>	-
	-
D <sub>10</sub>	-

dmr	Not dm
D_1	-
D_2	-
D_3	-



```
=====agg
max
min
sum
mean
mode
std
var
count
size - group wise count
```

first - first row in every group

last - last row in every group

nth(num) - if num = 5, 5th row in every group

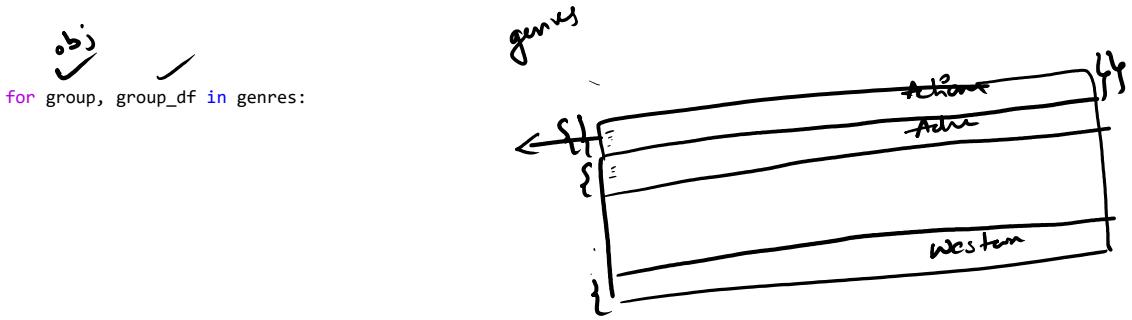
get\_group("Horror") - returns Horror group related data

groups - returns dict with keys as group names and values as index positions

describe() - similarly like pandas df on groups it will summarize every column(count, mean, std, 25% etc)

sample(2) - randomly picks 2 rows on each group

nunique() - gives column wise unique count



```
# looping on groups
# find highest rated movie from each genre
df=pd.DataFrame(columns=movies.columns)
for group, group_df in genres:
    df = pd.concat([df, group_df[group_df['IMDB_Rating'] == group_df['IMDB_Rating'].max()]])
```

Divide & conquer  
group\_df['IMDB\_Rating'].max()

9.0

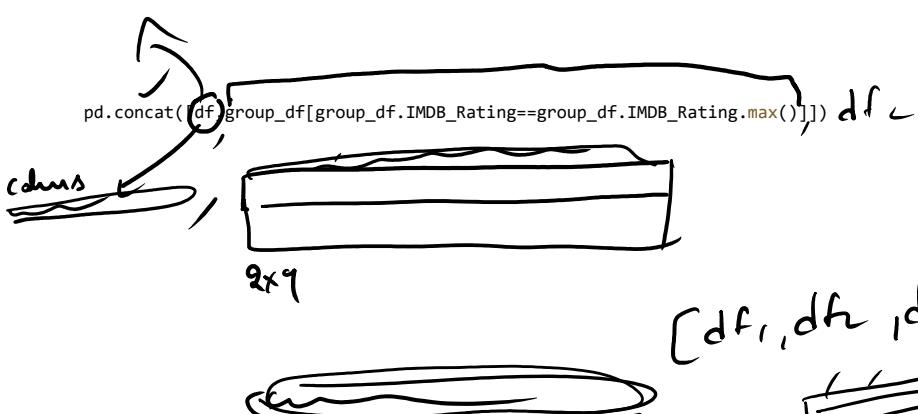
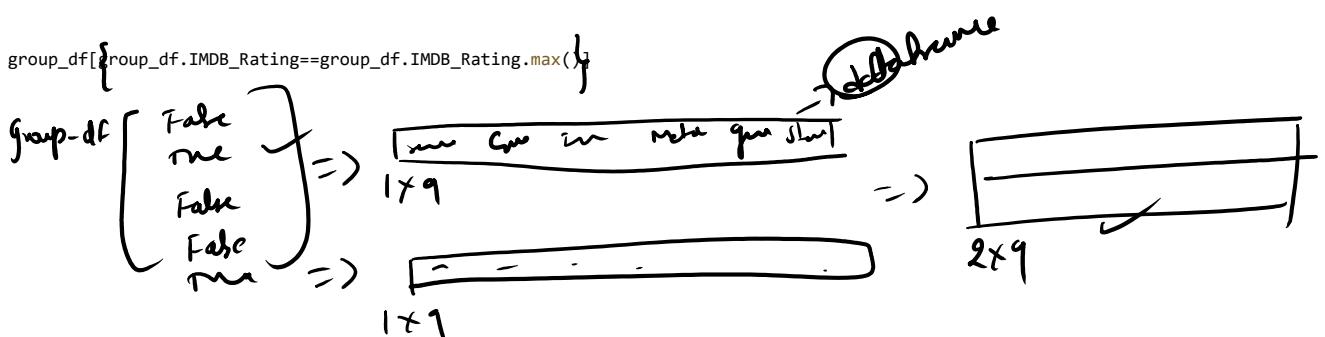
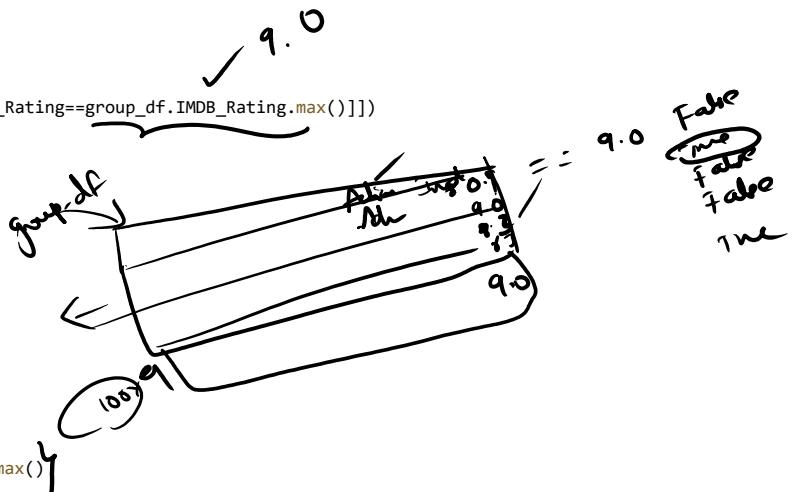
group\_df['IMDB\_Rating'] == group\_df['IMDB\_Rating'].max()

False

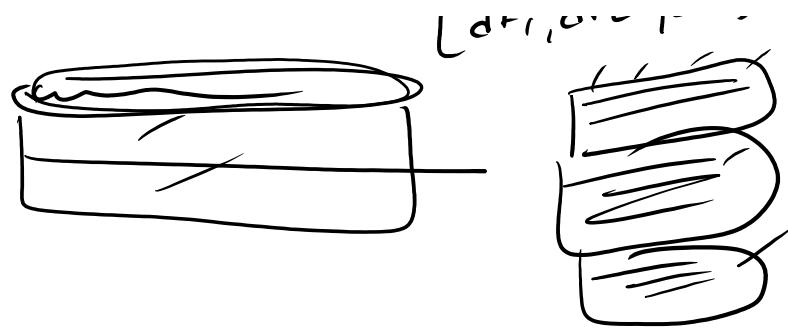
True

False

False



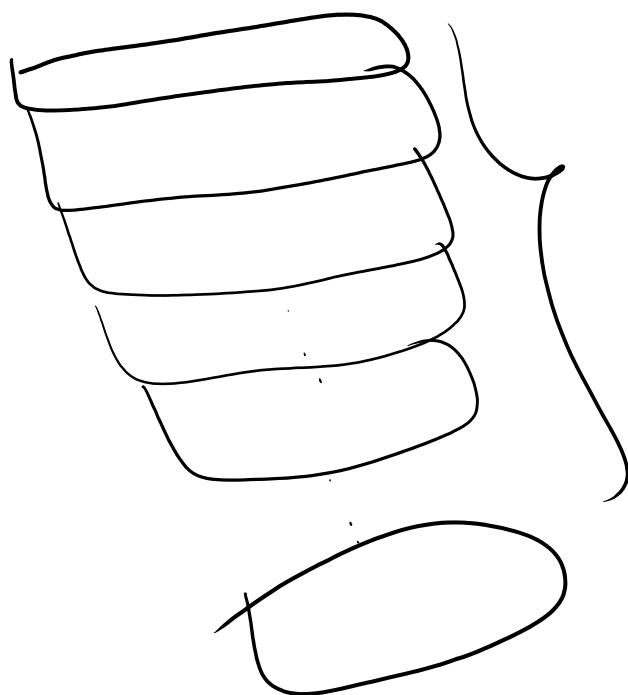
*df =*



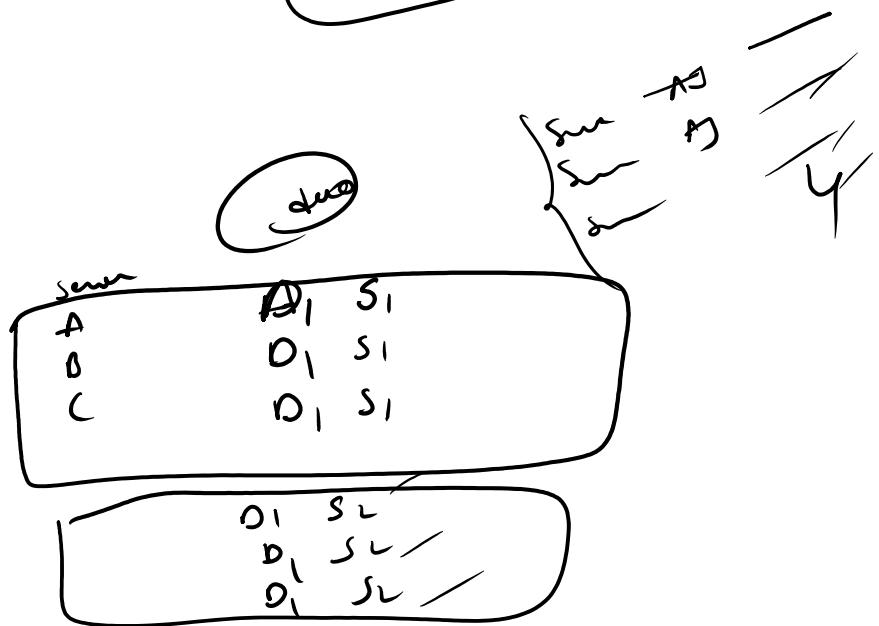
```
#using apply method
def genre_highestRated(group_df):
    return group_df[group_df.IMDB_Rating==group_df.IMDB_Rating.max()]
genres.apply(genre_highestRated)
```

*group\_by*

*step-var*



```
# groupby on multiple cols
movies.groupby(["Director", "Star1"])
```

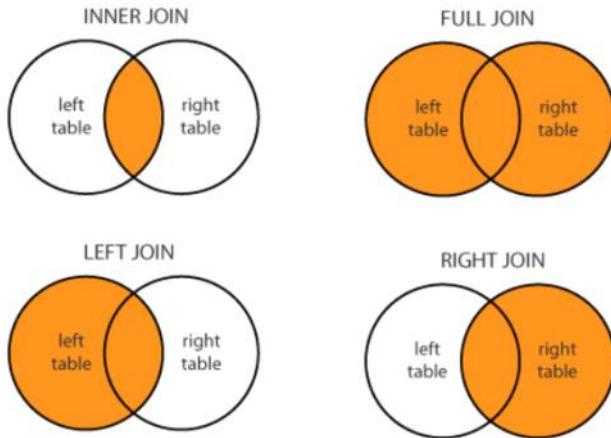


## concat, merge, visualization

28 February 2025 07:24

**Concat** - It allows us to combine 2 or more data frames along horizontal and vertical

**Merge** - It is also another way to combine data frames



## Difference Between concat() and merge() in Pandas

Feature	<code>concat()</code>	<code>merge()</code>
Purpose	Combines DataFrames along rows (axis=0) or columns (axis=1)	Joins DataFrames based on key columns (like SQL JOIN)
How it Works	Simply stacks DataFrames (row-wise or column-wise)	Matches data based on common columns or indexes
Default Behavior	Does not match values, just appends	Matches rows based on key columns
Best Use Case	When combining similar DataFrames (stacking datasets)	When you need relational joins (similar to SQL joins)

### Data Visualization(Using Matplotlib, Seaborn):

## What is Matplotlib?

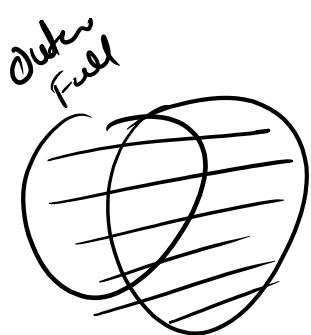
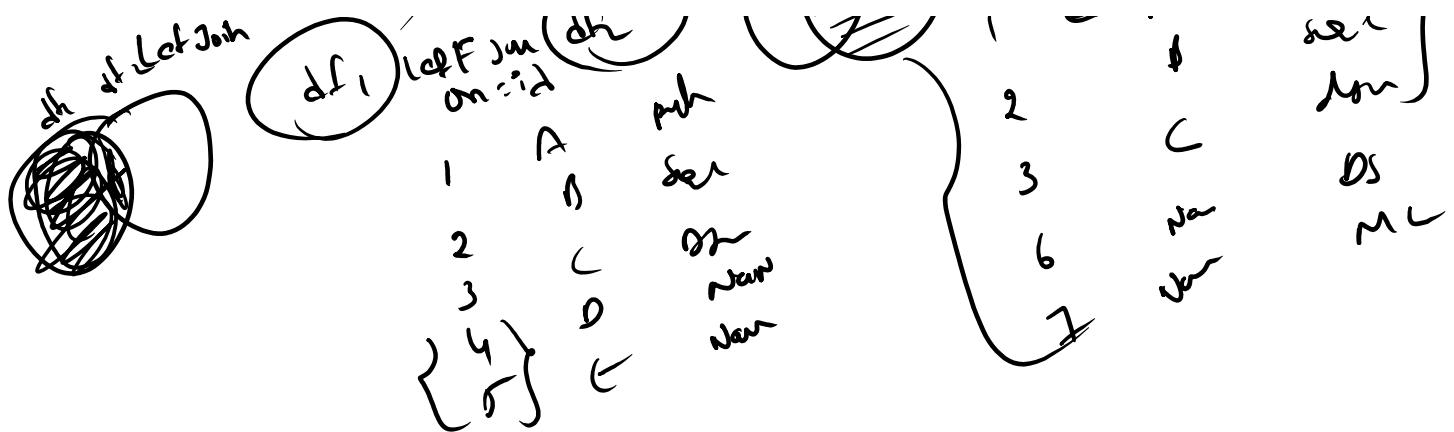
Matplotlib is a powerful Python library for creating static, animated, and interactive visualizations. It is often used for plotting line charts, bar charts, scatter plots, and histograms.

## What is Seaborn?

Seaborn is built on top of Matplotlib and provides a more visually appealing and statistical way to create plots. It simplifies complex visualizations like heatmaps, box plots, and violin plots.

Concat:





1	A	M
2	B	dh
3	C	nr
4	D	<del>Nar</del>
5	E	Nar
6	F	OS
7	G	ML

# Matplotlib & Seaborn

02 March 2025 19:49

## Matplotlib: A Low-Level Plotting Library

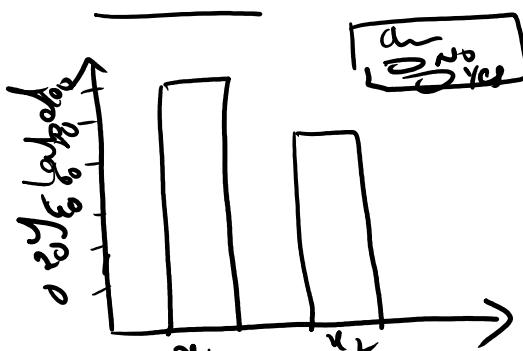
- Matplotlib is a **low-level** library that gives **full control** over every aspect of a plot.
- It works like a **canvas**, where you can modify everything, including:
  - Colors
  - Line styles
  - Axes properties
  - Legends
  - Grid styles
  - Annotations
- Since it provides raw access to plot elements, you can create **highly customized visualizations**.

## Seaborn: A High-Level Statistical Library

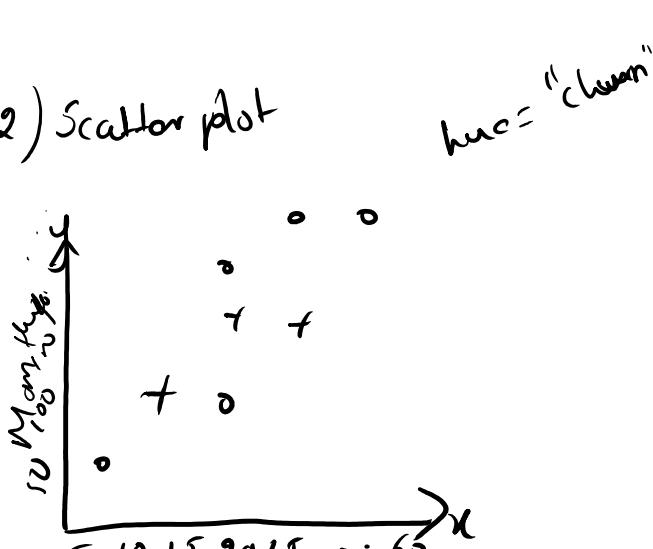
- Seaborn is built **on top of Matplotlib** but is designed to work with **structured datasets (DataFrames)**.
- It **automates** many visualizations, making plots **easier to create and visually appealing by default**.
- However, this automation **limits deep customization** because Seaborn **abstracts away** some of Matplotlib's low-level details.

- 1) Bar plot
- 2) Scatter plot
- 3) Box plot
- 4) histogram
- 5) heat map
- 6) violin plot
- 7) pie chart

### 1) Bar plot:-



### 2) Scatter plot





### 3) Box plot

Monthly-charges =  $[50, 75, 100, 120, 25, 35, 100, \dots]$

$$= [25, 35, 50, 75, 100, 120, \dots]$$

100  $\leftarrow$   $25\%$   $50\%$   $75\% \rightarrow$  200

$$IQR (\text{Inter Quartile range}) = 75\% - 25\%$$

$$\begin{cases} Q_1 = 25\% \\ Q_2 = 50\% \\ Q_3 = 75\% \end{cases}$$

$$\text{upper-whisker} = Q_3 + 1.5 \cdot IQR$$

$$= 75\% + 1.5 \cdot (75\% - 25\%)$$

$$\text{lower-whisker} = Q_1 - 1.5 \cdot (IQR)$$

$$= 25\% - 1.5 \cdot (75\% - 25\%)$$

