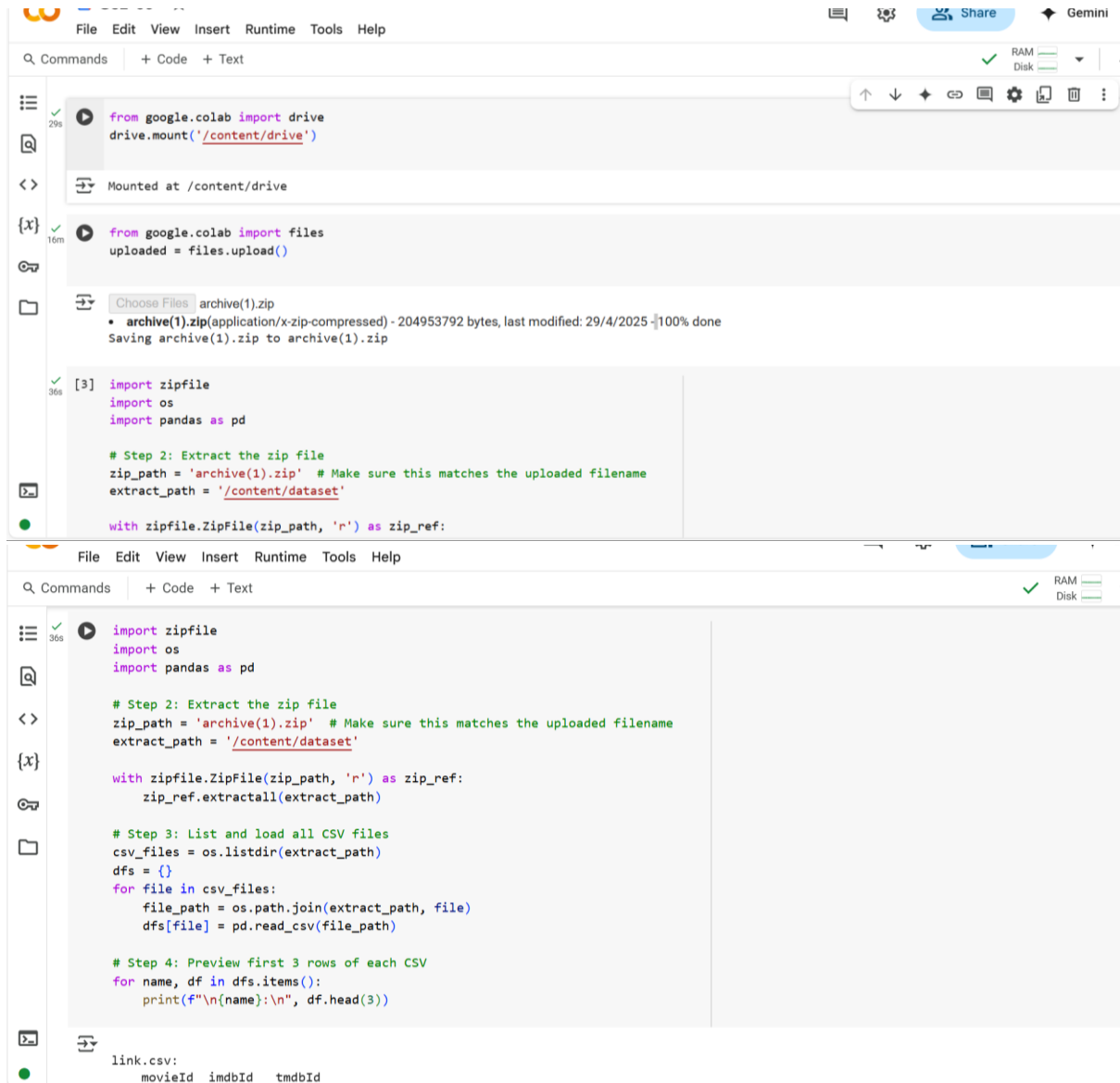


Name:-Priyal S. Mahajan

Roll No.:-CS3-76

PRN:- 202401040203

## • EDS ACTIVITY-2



The image displays two screenshots of a Google Colab notebook interface, showing the steps to mount a drive, upload a file, and extract it.

**Top Screenshot:**

- Code Cell 1:** `from google.colab import drive; drive.mount('/content/drive')`. The output shows "Mounted at /content/drive".
- Code Cell 2:** `from google.colab import files; uploaded = files.upload()`. The output shows a file "archive(1).zip" (application/x-zip-compressed) - 204953792 bytes, last modified: 29/4/2025, 100% done, being saved to "archive(1).zip".
- Code Cell 3:** `import zipfile; import os; import pandas as pd`. The output shows the file path and the extraction path: `zip_path = 'archive(1).zip' # Make sure this matches the uploaded filename` and `extract_path = '/content/dataset'`.

**Bottom Screenshot:**

- Code Cell 4:** `import zipfile; import os; import pandas as pd`. The output shows the file path and the extraction path: `zip_path = 'archive(1).zip' # Make sure this matches the uploaded filename` and `extract_path = '/content/dataset'`.
- Code Cell 5:** `with zipfile.ZipFile(zip_path, 'r') as zip_ref: zip_ref.extractall(extract_path)`. The output shows the file path and the extraction path: `zip_path = 'archive(1).zip' # Make sure this matches the uploaded filename` and `extract_path = '/content/dataset'`.
- Code Cell 6:** `# Step 3: List and load all CSV files; csv_files = os.listdir(extract_path); dfs = {}; for file in csv_files: file_path = os.path.join(extract_path, file); dfs[file] = pd.read_csv(file_path)`. The output shows the file path and the extraction path: `zip_path = 'archive(1).zip' # Make sure this matches the uploaded filename` and `extract_path = '/content/dataset'`.
- Code Cell 7:** `# Step 4: Preview first 3 rows of each CSV; for name, df in dfs.items(): print(f"\n{name}:\n", df.head(3))`. The output shows the file path and the extraction path: `zip_path = 'archive(1).zip' # Make sure this matches the uploaded filename` and `extract_path = '/content/dataset'`.

**Final Output:**

```
link.csv:
movieId  imdbId  tmdbId
```

	movieId	title	genres
8555	26018	Chase a Crooked Shadow (1958)	Crime Film-Noir Mystery Thriller
8933	26580	Park Is Mine, The (1986)	Action Drama Thriller
9249	27249	Trumpet of the Swan, The (2001)	Animation Drama Musical
9315	27396	Gentleman's Game, A (2002)	Drama
9770	31797	White Banners (1938)	Drama

```
# Q2. Identify movies with no ratings.
rated_movie_ids = ratings["movieId"].unique()
movies_no_ratings = movies[~movies["movieId"].isin(rated_movie_ids)]
movies_no_ratings
```

	movieId	title	genres
	8555	26018 Chase a Crooked Shadow (1958)	Crime Film-Noir Mystery Thriller
	8933	26580 Park Is Mine, The (1986)	Action Drama Thriller
	9249	27249 Trumpet of the Swan, The (2001)	Animation Drama Musical
	9315	27396 Gentleman's Game, A (2002)	Drama
	9770	31797 White Banners (1938)	Drama
...	...	...	...
	26818	128886 Love at the Top (1974)	Comedy Drama
	26872	129201 The Time Being (2012)	Mystery
	26933	129443 Thank You a Lot (2014)	Drama
	27004	129820 Spare Parts (2015)	Children Drama
	27056	130040 Mo (1983)	Horror

534 rows x 3 columns

```
[6] # Q3. Find the standard deviation of all ratings
std_rating = np.std(ratings_np)
std_rating
```

np.float64(1.051988892994865)

```
[7] # Q4. Get the minimum and maximum ratings.
min_rating = np.min(ratings_np)
max_rating = np.max(ratings_np)
(min_rating, max_rating)
```

(np.float64(0.5), np.float64(5.0))

```
[8] # Q5. Calculate the median rating.
median_rating = np.median(ratings_np)
median_rating
```

np.float64(3.5)

```
# Q6. Calculate the standard deviation of ratings for each movie and find the top 10 with highest variance.
rating_std = ratings.groupby("movieId")["rating"].std().dropna().sort_values(ascending=False)
movies.set_index("movieId").loc[rating_std.head(10).index]
```

	movieId	title	genres	genre_set
	119569	Quatsch und die Nasenbärbande (2014)	Children	{n, i, l, e, d, r, C, h}
	112577	Willie & Phil (1980)	Comedy Drama Romance	{m, n, e, l, d, o, a, r, y, c, C, R, D}
	87948	Bright Victory (1951)	Drama	{m, a, r, D}
	94027	Uwasa No Onna (The Woman in the Rumor) (Her Mo...	Drama Romance	{m, n, e, l, o, a, r, c, R, D}
	128173	Beethoven's Big Break (2008)	Children Comedy	{n, m, i, l, e, l, d, o, r, y, C, h}
	99012	Gay Bed and Breakfast of Terror, The (2007)	Comedy Horror	{m, H, e, l, d, o, r, y, C}
	34240	Karol: A Man Who Became Pope (Karol, un uomo d...	Drama	{m, a, r, D}
	126403	Apostle Peter and The Last Supper (2012)	Drama	{m, a, r, D}
	6253	Down and Out with the Dolls (2001)	Comedy	{m, e, d, o, y, C}
	126959	The Epic of Everest (1924)	Documentary	{m, n, e, o, t, a, r, c, y, u, D}



- —
- —
- —



&lt; &gt;

 $\{x\}$ 

0s



- —
- —
- —



&lt; &gt;

 $\{x\}$ 

183



```
array([False, False, False, False, False, False, False, False, False,
       False])
```



File Edit View Insert Runtime Tools Help

🔍 Commands + Code + Text

- 
- 
- 



&lt; &gt;

$$\{x\}$$


```
[14] # Q13. Find the number of unique users in the dataset
num_users = dfs['rating.csv']['userId'].nunique()
num_users
```

→ 138493

```
[37] # Q14. List the top 5 users who have given the most 5-star ratings.
      top_5star_users = ratings[ratings["rating"] == 5.0]["userId"].value_counts().head(5)
      top_5star_users
```

```

count
userId
72008    1540
131894    1300
119661    933
48498     927
82418     868
dtype: int64

```



```
[16] # Q.15 Count how many movies are rated by more than 500 users.
```

```
{x}
⇒ np.int64(4483)
```

8374

np.int64(118205)



```
[21] # Q.18 What is the average number of ratings per user?

avg_ratings_per_user = ratings.shape[0] / ratings['userId'].nunique()
avg_ratings_per_user
```

144.4135299257002

⇒ [('Drama', 13344)]

	count
rating	
0.5	239125
1.0	680732
1.5	279252
2.0	1430997

```
[27] # Q21. Find the 10 movies that received the most ratings.
      most Rated = ratings["movieId"].value_counts().head(10)
      movies.set_index("movieId").loc[most Rated.index]
```



	title	genres
movieId		
296	Pulp Fiction (1994)	Comedy Crime Drama Thriller
356	Forrest Gump (1994)	Comedy Drama Romance War
318	Shawshank Redemption, The (1994)	Crime Drama
593	Silence of the Lambs, The (1991)	Crime Horror Thriller
480	Jurassic Park (1993)	Action Adventure Sci-Fi Thriller
260	Star Wars: Episode IV - A New Hope (1977)	Action Adventure Sci-Fi
110	Braveheart (1995)	Action Drama War
589	Terminator 2: Judgment Day (1991)	Action Sci-Fi
2571	Matrix, The (1999)	Action Sci-Fi Thriller
527	Schindler's List (1993)	Drama War

```
# Q22. Identify users who have rated more than 100 movies.
active_users = ratings["userId"].value_counts()
active_users[active_users > 100]
```



	count
userId	
118205	9254
8405	7515
82418	5646
121535	5520
125794	5491
...	...
74293	101
20433	101
20427	101
3440	101
115563	101

51869 rows × 1 columns

dtype: int64







🔍 Commands + Code + Text

<> 7484

```

min    1995-01-09 11:46:44
max    2015-03-31 06:40:02

dtype: object

```



🔍 Commands + Code + Text

	year	rating
	1891	3.000000
	1893	3.375000
	1894	3.071429
	1895	2.833333
	1896	3.282609
	...	...
	2011	3.519526
	2012	3.588706
	2013	3.490962
	2014	3.524484
	2015	2.920181

118 rows × 1 columns

dtype: float64

