# EDS THEORY ASSIGNMENT 1

**Name: Komal Dhanajay Gangawane**
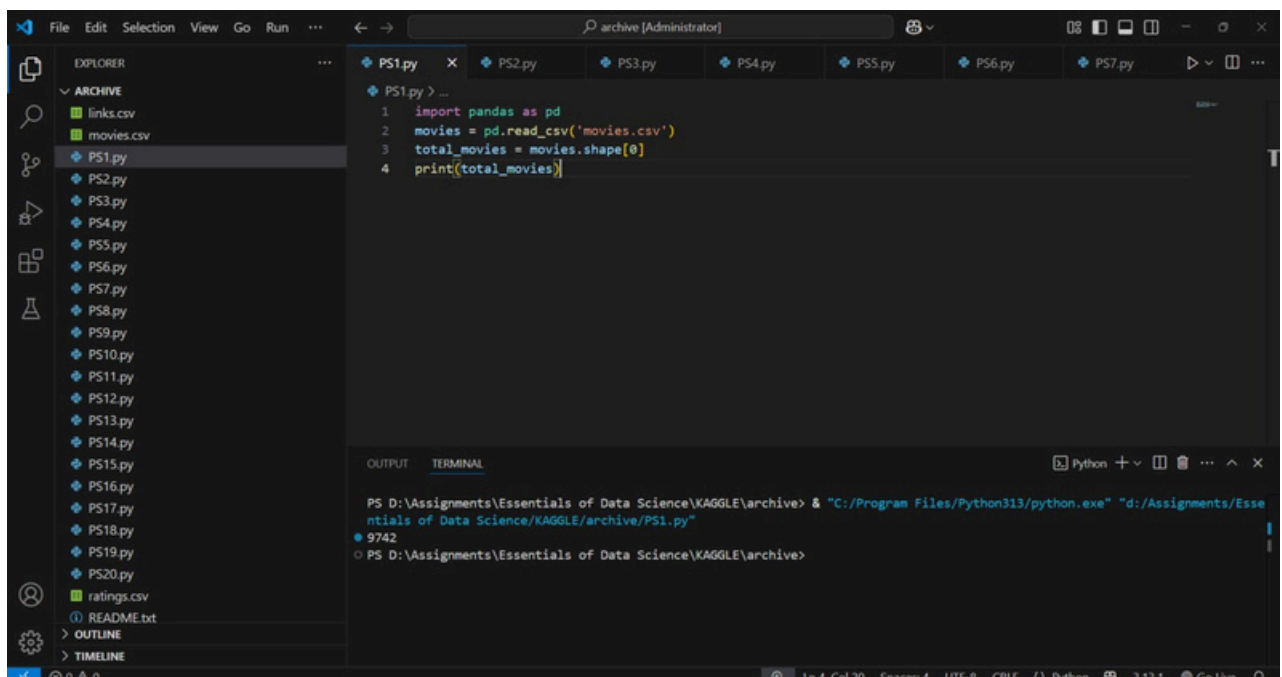
**DIV: ET2**

**Roll No: ET2-75**

**PRN: 202401070218**

URL: https://www.kaggle.com/datasets/shubhammehta21/movie-lens-small-latest-dataset

1. Count the total number of movies in the dataset.



2. Determine the number of unique genres spanning all movies.

3. Compute how many times each genre appears by splitting and exploding the genre strings.



4. Identify the five most frequent genres in the dataset.

5. Extract the release year from the movie titles.



6. Count how many movies were released in each year.

7. Derive each movie's decade (based on the release year) and show the distribution.



8. Filter out the movies that belong to the "Comedy" genre.

9. Count the number of movies that have more than one genre listed.

```python
import pandas as pd
movies = pd.read_csv('movies.csv')
comedy_movies = movies[movies['genres'].str.contains('Comedy', regex=False)]
print(comedy_movies)
```

```python
import pandas as pd
movies = pd.read_csv('movies.csv')
multiple_genre_movies = movies[movies['genres'].str.contains(r'\|')]
print(multiple_genre_movies.shape[0])
```



10. Identify the movie with the longest title.

11. Filter movies that have titles starting with the letter "A".



12. Identify movie titles that appear more than once.

13. Count how many movies are label with "(no genres listed)" in the genres field.



14. Determine the single most common word found in movie titles.

15. Calculate the percentage of movies that belong to multiple genres.



16. Count the total number of ratings

**17. Calculate the average rating overall**



**18. Identify the top 10 movies with the highest number of ratings**

```python
import pandas as pd
ratings = pd.read_csv('ratings.csv')
top10_movies = ratings.groupby('movieId')['rating'].count().nlargest(10)
print(top10_movies)
```

```
PS D:\Assignments\Essentials of Data Science\KAGGLE\archive> & "C:/Program Files/Python313/python.exe" "d:/Assignments/Esse
ntials of Data Science/KAGGLE/archive/PS18.py"
movieId
356     329
318     317
296     307
593     279
2571    278
260     251
480     238
110     237
589     224
527     220
Name: rating, dtype: int64
PS D:\Assignments\Essentials of Data Science\KAGGLE\archive>
```

## 19. Find the earliest and latest rating timestamps



```python
import pandas as pd
ratings = pd.read_csv('ratings.csv')
ratings['timestamp'] = pd.to_datetime(ratings['timestamp'], unit='s')
earliest = ratings['timestamp'].min()
latest = ratings['timestamp'].max()
print(earliest)
print(latest)
```

```
PS D:\Assignments\Essentials of Data Science\KAGGLE\archive> & "C:/Program Files/Python313/python.exe" "d:/Assignments/Esse
ntials of Data Science/KAGGLE/archive/PS19.py"
1996-03-29 18:36:55
2018-09-24 14:27:30
PS D:\Assignments\Essentials of Data Science\KAGGLE\archive>
```

## 20. Calculate the percentage of ratings that are above 4

```python
import pandas as pd
ratings = pd.read_csv('ratings.csv')
above_four = (ratings['rating'] > 4).sum()
total_ratings = ratings.shape[0]
percentage = (above_four / total_ratings) * 100
print(percentage)
```

PS D:\Assignments\Essentials of Data Science\KAGGLE\archive> & "C:/Program Files/Python313/python.exe" "d:/Assignments/Esse
ntials of Data Science/KAGGLE/archive/PS20.py"
21.5815780078S4336
PS D:\Assignments\Essentials of Data Science\KAGGLE\archive>