# Assignment no.04

Name :- Nandini Thorat
Div :- C
Roll no:- 384
PRN:- 202201070117

<u>Code:</u>

```python
import pandas as pd

df = pd.read_csv("/content/movie_data.csv")
#print all records of dataset
print(df)
```





```python
#1 print Names of all employees
print(df['director_name'])
```

```python
#2 print name and duration
print(df[['director_name','duration']])
```

```
          director name  duration
0         James Cameron     178.0
1         Gore Verbinski    169.0
2         Sam Mendes        148.0
3      Christopher Nolan    164.0
4         Andrew Stanton    132.0
...                  ...       ...
5037        Scott Smith      87.0
5038               NaN       43.0
5039   Benjamin Roberds     76.0
5040        Daniel Hsia    100.0
5041           Jon Gunn     90.0

[5042 rows x 2 columns]
```

```python
#1 Data cleaning
#check for missing values
print(df.isnull())

# #drop rows with missing values
df.dropna(inplace=True
```



```python
#2 convert string to upper case
df['director_name'].str.upper()
```

```
0          JAMES CAMERON
1         GORE VERBINSKI
2             SAM MENDES
3      CHRISTOPHER NOLAN
4         ANDREW STANTON
               ...
1691       JAMES BIDGOOD
1692          DARYL WEIN
1693        JAFAR PANAHI
1694     KIYOSHI KUROSAWA
1695       SHANE CARRUTH
Name: director_name, Length: 1696, dtype: object
```

```
#3. print movie title along with their year of release
df1 = df[['movie_title','title_year']]
print(df1)
```

```
                                      movie_title  title_year
0                                          Avatar      2009.0
1        Pirates of the Caribbean: At World's End      2007.0
2                                         Spectre      2015.0
3                           The Dark Knight Rises      2012.0
4                                     John Carter      2012.0
...                                           ...         ...
1691                                Pink Narcissus      1971.0
1692                              Breaking Upwards      2009.0
1693                                    The Circle      2000.0
1694                                      The Cure      1997.0
1695                                        Primer      2004.0

[1696 rows x 2 columns]
```

```
#4 calculate the total budget of all the movies
totalBudget = df['budget'].sum()
print("Total budget of all movies = ", totalBudget)
```

```
Total budget of all movies =  174826107781.0
```

```
#5 calculate mean, median, mode imdb rating
meanImdb = df['imdb_score'].mean()
medianImdb = df['imdb_score'].median()
modeImdb = df['imdb_score'].mode()
print("Mean IMDB score = ", meanImdb)
print("Median IMDB score = ", medianImdb)
print("Mode IMDB score = ", modeImdb)
```

```
Mean IMDB score =  6.467471143756558
Median IMDB score =  6.6
Mode IMDB score =  0    6.7
Name: imdb_score, dtype: float64
```

```
#6 describe gross of all movies
print(df['gross'].describe())
```

```
count    3.812000e+03
mean     5.204686e+07
std      7.016457e+07
min      1.620000e+02
25%      7.682030e+06
50%      2.922370e+07
75%      6.648842e+07
max      7.605058e+08
Name: gross, dtype: float64
```

```python
#7 minimum and maximum duration movie
minimumDuration = df['duration'].min()
maximumDuration = df['duration'].max()
print("Minimum duration movie: ", minimumDuration)
print("Maximum duration movie: ", maximumDuration)
```

```
Minimum duration movie:  37.0
Maximum duration movie:   330.0
```

```python
#8 count number of movies which are released after 2010
released_after_2010 = df[df['title_year'] > 2010]
print("Number of movies released after 2010: ",
released_after_2010['title_year'].count())
```

```
Number of movies released after 2010:   430
```

```python
#9 print count of movies released in each year
print(df.groupby('title_year').count())
```



```python
#10 correlation
print(df.corr())
```



```python
#11 covariance
print(df.cov())
```



```python
#12 print details of first 10 movies
print(df.head(10))
```

```python
#13 print details of movies with duration above 300minutes
print(df.loc[df['duration']>300])
```

```
     level_0  index   director_name  num_critic  duration    profit  \
495      495   1143  Michael Cimino       102.0     325.0  1500000.0

                     genres    lead_actor    movie_title  num_voted_users  \
495  Adventure|Drama|Western  Jeff Bridges  Heaven's Gate             9830

     num_user_for_reviews language country      budget  title_year  \
495                 189.0  English     USA  44000000.0      1980.0

     imdb_score  aspect_ratio  movie_likes  num_voted_reviews
495         6.8          2.35         1000            10019.0
```

```python
#14 print the quantile of movie likes
print(df['movie_likes'].quantile([0.25, 0.5, 0.75]))
```

```
0.25        0.0
0.50      225.5
0.75    11000.0
Name: movie_likes, dtype: float64
```

```python
#15 data preparation

#strip leading and trailing whitespaces if any
df['director_name'].str.strip()

#filter rows based on condition
imdb_above_8 = df[df['imdb_score'] > 8.5]
print(imdb_above_8)

#filter rows based on query
title_year_above_2008 = df.query('title_year > 2008')

#adding a new column
df['num_voted_reviews'] = df['num_voted_users'] +
df['num_user_for_reviews']

#get dummies
dummy_countries = pd.get_dummies(df['country'])
```

```
      level_0  index    director_name  num_critic  duration     profit  \
1183     1183   3174        Tony Kaye       162.0     101.0  6712241.0
1560     1560   4426  Charles Chaplin       120.0      87.0   163245.0

                   genres          lead_actor        movie_title  \
1183          Crime|Drama        Ethan Suplee  American History X
1560  Comedy|Drama|Family  Paulette Goddard        Modern Times

      num_voted_users  num_user_for_reviews language country     budget  \
1183           782437                1420.0  English     USA  7500000.0
1560           143086                 211.0  English     USA  1500000.0

      title_year  imdb_score  aspect_ratio  movie_likes  num_voted_reviews
1183      1998.0         8.6          1.85        35000           783857.0
1560      1936.0         8.6          1.37            0           143297.0
```

```python
#16  data aggregation
#renaming our gross column as profit
df.rename(columns={'gross':'profit'},inplace=True)
df
```



```python
#17 Datatype conversion
df['duration'] = df['duration'].astype('float')
print(type(df['duration'][0]))
```

```
<class 'numpy.float64'>
```

```python
#18 data wrangling

newdf1 = pd.DataFrame(df[['director_name', 'duration', 'movie_title']])
newdf2 = pd.DataFrame(df[['movie_title', 'title_year', 'imdb_score']])

# merge dataframes
merged_df = pd.merge(newdf1, newdf2)
print(merged_df.head())

#concat dataframes
concatenated_df = pd.concat([newdf1, newdf2], axis=1)
print(concatenated_df.head())
```

```
    director_name  duration                                   movie_title  \
0   James Cameron     178.0                                        Avatar
1   Gore Verbinski    169.0  Pirates of the Caribbean: At World's End
2      Sam Mendes     148.0                                       Spectre
3 Christopher Nolan   164.0                        The Dark Knight Rises
4   Andrew Stanton    132.0                                  John Carter

   title_year  imdb_score
0     2009.0         7.9
1     2007.0         7.1
2     2015.0         6.8
3     2012.0         8.5
4     2012.0         6.6
    director_name  duration                                   movie_title  \
0   James Cameron     178.0                                        Avatar
1   Gore Verbinski    169.0  Pirates of the Caribbean: At World's End
2      Sam Mendes     148.0                                       Spectre
3 Christopher Nolan   164.0                        The Dark Knight Rises
4   Andrew Stanton    132.0                                  John Carter

                                movie_title  title_year  imdb_score
0                                    Avatar       2009.0         7.9
1 Pirates of the Caribbean: At World's End    2007.0         7.1
2                                   Spectre       2015.0         6.8
3                     The Dark Knight Rises    2012.0         8.5
4                               John Carter       2012.0         6.6
```

```
#19 Data transformation

#convert duration into hours
df['duration_in_hrs'] = round(df['duration']/60, 1)
print(df['duration_in_hrs'].head(10))
```

```
0    3.0
1    2.8
2    2.5
3    2.7
4    2.2
5    2.6
6    1.7
7    2.4
8    2.6
9    3.0
Name: duration_in_hrs, dtype: float64
```

```
#20 display name of movie and director's name of first 5 movies
selected_Data = df.iloc[[1, 2, 3, 4, 5], [1, 7]]
print(selected_Data)
```

```
   index        lead_actor
1      1       Johnny Depp
2      2   Christoph Waltz
3      3         Tom Hardy
4      4      Daryl Sabara
5      5       J.K. Simmons
```