

```
In [129]: # Import the numpy and pandas packages
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

Task 1: Reading and Inspection

```
## Subtask 1.1: Import and read

# Import and read the movie database. Store it in a variable called movies.

In [64]: #Write your code here
movies = read_csv('IMDb_Movies.csv')
OrgData = movies

Out[64]:
```

	color	director_name	num_critics_for_reviews	gross	genres	actor_1_facebook_likes	gross	genres	num_user_for_reviews	language	country	rating	budget	title_year	actor_2_facebook_likes
0	Color	James Cameron	723.0	178.0	0.0	855.0	Joel David Moore	10000.0	76050547.0	Action/Adventure/Fantasy/Sci-Fi	USA	PG-13	23700000.0	2009.0	936.0
1	Color	Gore Verbinski	302.0	149.0	563.0	100.0	Orlando Bloom	40000.0	390404152.0	Action/Adventure/Fantasy	USA	PG-13	30000000.0	2007.0	5000.0
2	Color	Sam Mendes	602.0	168.0	0.0	161.0	Rory Kinnear	11000.0	200470415.0	Action/Adventure/Thriller	UK	PG-13	24500000.0	2015.0	393.0
3	Color	Christopher Nolan	813.0	164.0	22000.0	23000.0	Christian Bale	27000.0	448130642.0	Action/Thriller	USA	PG-13	25000000.0	2010.0	23000.0
4	NaN	Doug Walker	NaN	NaN	111.0	NaN	Rob Doyle	NaN	NaN	Documentary	NaN	NaN	NaN	NaN	12.0
...
5038	Color	Scott Smith	1.0	87.0	2.0	318.0	Daphne Jung	637.0	NaN	Comedy/Drama	Canada	NaN	NaN	2013.0	470.0
5040	Color	Benjamin Roberts	13.0	76.0	0.0	0.0	Howell Mossy	0.0	NaN	Crimedrama/Mystery/Thriller	USA	TV-14	NaN	NaN	593.0
5041	Color	Daniel Hsia	14.0	10.0	0.0	489.0	Daniel Henney	94.0	10443.0	Comedy/Drama/Romance	USA	PG-13	NaN	2013.0	0.0
5042	Color	Jon Gunn	43.0	90.0	16.0	16.0	Brian Herzberg	86.0	85222.0	Documentary	UK	PG	1100.0	2004.0	23.0

5043 rows × 16 columns

```
## Subtask 1.2: Inspect the dataframe

Inspect the dataframe's columns, shapes, variable types etc.

In [65]: #Write your code here
movies.info()

Out[65]:
Out[66]:
movies.info()
Name: movies, dtype: object
RangeIndex: 5043 entries, 0 to 5042
Data columns (total 16 columns):
color                5043 non-null object
director_name        5043 non-null object
num_critics_for_reviews 4928 non-null float64
duration             5028 non-null float64
director_facebook_likes 4939 non-null float64
actor_2_facebook_likes 5020 non-null float64
actor_3_name         5030 non-null object
actor_1_facebook_likes 5038 non-null float64
gross               4939 non-null object
genres              5043 non-null object
actor_1_name        5035 non-null object
movie_title         5043 non-null object
num_voted_users     5043 non-null int64
cast_total_facebook_likes 4943 non-null float64
director_facebook_likes 5020 non-null object
facecnumber_in_poster 5043 non-null object
plot_keywords       4890 non-null object
movie_imdb_link     5023 non-null object
num_user_for_reviews 5033 non-null object
language            5033 non-null object
country            5038 non-null object
content_rating      4951 non-null float64
title_year         4930 non-null float64
actor_2_facebook_likes 5030 non-null float64
actor_3_facebook_likes 5043 non-null float64
aspect_ratio        5043 non-null float64
movie_facebook_likes 5043 non-null int64
dtype: object
float64(12), int64(3), object(13)
memory usage: 1.1+ MB
```

Task 2: Cleaning the Data

```
## Subtask 2.1: Inspect Null values

Find out the number of Null values in all the columns and Also, find the percentage of Null values in each column. Round off the percentages upto two decimal places.

In [67]: # Write your code for column-wise null count here
movies.isnull().sum(axis=0).sort_values(ascending=False)

Out[67]:
gross                894
budget              492
aspect_ratio        393
plot_keywords       329
title_year         108
director_name       104
director_facebook_likes 104
num_critics_for_reviews 23
actor_3_name        23
actor_2_facebook_likes 23
num_user_for_reviews 19
color              19
duration           13
facecnumber_in_poster 13
actor_1_name       13
actor_2_facebook_likes 13
language           12
actor_3_facebook_likes 7
movie_facebook_likes 7
genres             0
movie_title        0
num_voted_users    0
movie_imdb_link    0
jmdp_score         0
cast_total_facebook_likes 0
dtype: object

## Subtask 2.2: Drop unnecessary columns

For this assignment, you will mostly be analyzing the movies with respect to the ratings, gross collection, popularity of movies, etc. So many of the columns in this dataframe are not required. So it is advised to drop the following columns.
```

- color
- director_facebook_likes
- actor_1_facebook_likes
- actor_2_facebook_likes
- actor_3_facebook_likes
- actor_2_name
- actor_3_name
- cast_total_facebook_likes
- actor_3_name
- facecnumber_in_poster
- content_rating
- country
- movie_imdb_link
- aspect_ratio
- plot_keywords

```
In [78]: # Write your code for dropping the columns here. It is advised to keep inspecting the dataframe after each set of operations
movies = movies.drop()
```

```
In [71]:
movies
Out[71]:
```

	director_name	num_critics_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language	budget	title_year	imdb_score	movie_facebook_likes
0	James Cameron	723.0	76050547.0	Action/Adventure/Fantasy/Sci-Fi	CCI Pounder	Avatar	886204	3054	English	23700000.0	2009.0	7.9	33000
1	Gore Verbinski	302.0	390404152.0	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	1238	English	30000000.0	2007.0	7.1	0
2	Sam Mendes	602.0	200470415.0	Action/Adventure/Thriller	Christopher Waltz	Spectre	275868	994	English	24500000.0	2015.0	6.8	85000
3	Christopher Nolan	813.0	448130642.0	Action/Thriller	Tom Hardy	The Dark Knight Rises	1144337	2701	English	25000000.0	2012.0	8.5	164000
4	Doug Walker	NaN	NaN	Documentary	Doug Walker	Star Wars: Episode VII - The Force Awakens	8	NaN	NaN	NaN	NaN	7.1	0
...
5038	Scott Smith	1.0	NaN	Comedy/Drama	Eric Mabius	Signed Sealed Delivered	629	6	English	NaN	2013.0	7.7	84
5039	NaN	43.0	NaN	Crimedrama/Mystery/Thriller	Nadine Zie	The Following	73897	239	English	NaN	NaN	7.5	3000
5040	Benjamin Roberts	13.0	NaN	Drama/Comedy/Thriller	Eric Boboche	A Plague So Few	38	3	English	1400.0	2013.0	6.3	16
5041	Daniel Hsia	14.0	10443.0	Comedy/Drama/Romance	Alan Rick	Shanghai Calling	1235	9	English	NaN	2012.0	6.3	560
5042	Jon Gunn	43.0	85222.0	Documentary	John August	My Date with Drew	4285	84	English	1100.0	2004.0	6.6	456

5043 rows × 13 columns

```
## Subtask 2.3: Drop unnecessary rows using columns with large percentage of null values

Now, on inspection you might notice that some columns have high percentage (greater than 5%) of Null values. Drop all the rows which have Null values in any of those columns.

In [72]: # Write your code for dropping the rows here
# Round with 2 for 2 decimal points
round(movies.isnull().sum()) / len(movies) * 100, 2)

Out[72]:
gross                1.73
budget              0.97
aspect_ratio        0.76
plot_keywords       0.64
title_year         0.21
director_name       0.02
director_facebook_likes 0.02
num_critics_for_reviews 0.46
num_user_for_reviews 0.40
language           0.24
actor_3_name       0.45
movie_facebook_likes 0.00
jmdp_score         0.00
num_voted_users    0.00
movie_imdb_link    0.00
dtype: float64

In [73]: movies = movies[movies['gross'] > 10000000]
movies = movies[movies['budget'] > 10000000]

In [74]: round(movies.isnull().sum()) / len(movies) * 100, 2)

Out[74]:
language           0.08
actor_3_name       0.08
num_critics_for_reviews 0.43
movie_facebook_likes 0.00
jmdp_score         0.00
title_year         0.00
budget            0.00
num_user_for_reviews 0.40
num_voted_users    0.00
movie_title        0.00
genres            0.00
director_name      0.00
dtype: float64

## Subtask 2.4: Drop unnecessary rows

Some of the rows might have greater than five NaN values. Such rows aren't of much use for the analysis and hence, should be removed.

In [75]: # Write your code for dropping the rows here
movies.isnull().sum(axis=0).sort_values(ascending=False) > 5).sum()

Out[75]:
0

In [76]: movies = movies[movies.isnull().sum(axis=0).sort_values(ascending=False) <= 5]
movies

In [77]: movies[movies['language'] == 'English']

Out[77]:
```

	director_name	num_critics_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language	budget	title_year	imdb_score	movie_facebook_likes
0	James Cameron	723.0	76050547.0	Action/Adventure/Fantasy/Sci-Fi	CCI Pounder	Avatar	886204	3054	English	23700000.0	2009.0	7.9	33000
1	Gore Verbinski	302.0	390404152.0	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	1238	English	30000000.0	2007.0	7.1	0
2	Sam Mendes	602.0	200470415.0	Action/Adventure/Thriller	Christopher Waltz	Spectre	275868	994	English	24500000.0	2015.0	6.8	85000
3	Christopher Nolan	813.0	448130642.0	Action/Thriller	Tom Hardy	The Dark Knight Rises	1144337	2701	English	25000000.0	2012.0	8.5	164000
5	Andrew Stanton	462.0	73056879.0	Action/Adventure/Sci-Fi	Daryl Sabara	John Carter	212024	738	English	26370000.0	2012.0	6.6	24000
...
5033	Shane Carruth	143.0	624276.0	Drama/Sci-Fi/Thriller	Shane Carruth	Primer	72639	371	English	7000.0	2004.0	7.0	19000
5034	Neill Dea Lusa	35.0	0.070071	Thriller	Ian Gamston	Caitie	589	35	English	7000.0	2005.0	6.3	74
5035	Robert Rodriguez	56.0	2406920	Action/Crime/Drama/Romance/Thriller	Carlos Gallardo	El Mariachi	52055	130	Spanish	0.0070	1992.0	6.9	0
5037	Edward Burns	14.0	0.004584	Comedy/Drama	Kerry Bishi	Newlyweds	1338	14	English	9000.0	2010.0	6.4	413
5042	Jon Gunn	43.0	85222.0	Documentary	John August	My Date with Drew	4285	84	English	1100.0	2004.0	6.6	456

3917 rows × 13 columns

```
## Subtask 2.5: Fill NaN values

You might notice that the 'language' column has some NaN values. Here, on inspection, you will see that it is safe to replace all the missing values with 'English'.

In [78]: # Write your code here
round(movies.isnull().sum()) / len(movies) * 100, 2)

Out[78]:
language           0.08
actor_3_name       0.08
num_critics_for_reviews 0.43
movie_facebook_likes 0.00
jmdp_score         0.00
title_year         0.00
budget            0.00
num_user_for_reviews 0.40
num_voted_users    0.00
movie_title        0.00
genres            0.00
director_name      0.00
dtype: float64

In [78]: movies[movies['language'] == 'English']

Out[78]:
```

	director_name	num_critics_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language	budget	title_year	imdb_score	movie_facebook_likes
0	James Cameron	723.0	76050547.0	Action/Adventure/Fantasy/Sci-Fi	CCI Pounder	Avatar	886204	3054	English	23700000.0	2009.0	7.9	33000
1	Gore Verbinski	302.0	390404152.0	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	1238	English	30000000.0	2007.0	7.1	0
2	Sam Mendes	602.0	200470415.0	Action/Adventure/Thriller	Christopher Waltz	Spectre	275868	994	English	24500000.0	2015.0	6.8	85000
3	Christopher Nolan	813.0	448130642.0	Action/Thriller	Tom Hardy	The Dark Knight Rises	1144337	2701	English	25000000.0	2012.0	8.5	164000
5	Andrew Stanton	462.0	73056879.0	Action/Adventure/Sci-Fi	Daryl Sabara	John Carter	212024	738	English	26370000.0	2012.0	6.6	24000
...
5033	Shane Carruth	143.0	624276.0	Drama/Sci-Fi/Thriller	Shane Carruth	Primer	72639	371	English	7000.0	2004.0	7.0	19000
5034	Neill Dea Lusa	35.0	0.070071	Thriller	Ian Gamston	Caitie	589	35	English	7000.0	2005.0	6.3	74
5035	Robert Rodriguez	56.0	2406920	Action/Crime/Drama/Romance/Thriller	Carlos Gallardo	El Mariachi	52055	130	Spanish	0.0070	1992.0	6.9	0
5037	Edward Burns	14.0	0.004584	Comedy/Drama	Kerry Bishi	Newlyweds	1338	14	English	9000.0	2010.0	6.4	413
5042	Jon Gunn	43.0	85222.0	Documentary	John August	My Date with Drew	4285	84	English	0.0011	2004.0	6.6	456

3917 rows × 13 columns

```
## Subtask 2.6: Check the number of retained rows

You might notice that two of the columns viz. num_critics_for_reviews and actor_1_name, have small percentages of NaN values left. You can let these columns as it is for now. Check the number and percentage of the rows retained after completing all the tasks above.

In [84]: # Write your code for checking number of retained rows here
len(movies) / len(OrgData) * 100

Out[84]:
77.15645449137418

Checkpoint 1: You might have noticed that we still have around 77% of the rows!
```

Task 3: Data Analysis

```
## Subtask 3.1: Change the unit of columns. Convert the unit of the 'budget' and 'gross' columns from $ to million $.

In [87]: # Write your code for unit conversion here
movies['budget'] = movies['budget'] / 1000000
movies['gross'] = movies['gross'] / 1000000

In [88]:
movies
Out[88]:
```

	director_name	num_critics_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language	budget	title_year	imdb_score	movie_facebook_likes
0	James Cameron	723.0	76.050547	Action/Adventure/Fantasy/Sci-Fi	CCI Pounder	Avatar	886204	3054	English	237.00000	2009.0	7.9	33000
1	Gore Verbinski	302.0	39.0404152	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	1238	English	30.000000	2007.0	7.1	0
2	Sam Mendes	602.0	20.0470415	Action/Adventure/Thriller	Christopher Waltz	Spectre	275868	994	English	24.500000	2015.0	6.8	85000
3	Christopher Nolan	813.0	44.8130642	Action/Thriller	Tom Hardy	The Dark Knight Rises	1144337	2701	English	25.000000	2012.0	8.5	164000
5	Andrew Stanton	462.0	7.3056879	Action/Adventure/Sci-Fi	Daryl Sabara	John Carter	212024	738	English	26.370000	2012.0	6.6	24000
...
5033	Shane Carruth	143.0	0.624760	Drama/Sci-Fi/Thriller	Shane Carruth	Primer	72639	371	English	7.00000	2004.0	7.0	19000
5034	Neill Dea Lusa	35.0	0.070071	Thriller	Ian Gamston	Caitie	589	35	English	0.0070	2005.0	6.3	74
5035	Robert Rodriguez	56.0	2.406920	Action/Crime/Drama/Romance/Thriller	Carlos Gallardo	El Mariachi	52055	130	Spanish	0.0070	1992.0	6.9	0
5037	Edward Burns	14.0	0.004584	Comedy/Drama	Kerry Bishi	Newlyweds	1338	14	English	0.0090	2010.0	6.4	413
5042	Jon Gunn	43.0	0.085222	Documentary	John August	My Date with Drew	4285	84	English	0.0011	2004.0	6.6	456

3917 rows × 13 columns

```
## Subtask 3.2: Find the movies with highest profit

1. Create a new column called 'profit', which contains the difference of the two columns: 'gross' and 'budget'.
2. Sort the dataframe using the 'profit' column as reference.
3. Extract the top ten profiling movies in descending order and store them in a new dataframe - 'top10'

In [89]: # Write your code for creating the profit column here
movies['profit'] = movies['gross'] - movies['budget']
movies

Out[89]:
```

	director_name	num_critics_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language	budget	title_year	imdb_score	movie_facebook_likes	profit
0	James Cameron	723.0	76.050547	Action/Adventure/Fantasy/Sci-Fi	CCI Pounder	Avatar	886204	3054	English	237.00000	2009.0	7.9	33000	52.350547
1	Gore Verbinski	302.0	39.0404152	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	1238	English	30.000000	2007.0	7.1	0	9.0404152
2	Sam Mendes	602.0	20.0470415	Action/Adventure/Thriller	Christopher Waltz	Spectre	275868	994	English	24.500000	2015.0	6.8	85000	44.925825
3	Christopher Nolan	813.0	44.8130642	Action/Thriller	Tom Hardy	The Dark Knight Rises	1144337	2701	English	25.000000	2012.0	8.5	164000	19.8130642
5	Andrew Stanton	462.0	7.3056879	Action/Adventure/Sci-Fi	Daryl Sabara	John Carter	212024	738	English	26.370000	2012.0	6.6	24000	-19.0641321
...
5033	Shane Carruth	143.0	0.624760	Drama/Sci-Fi/Thriller	Shane Carruth	Primer	72639	371	English	7.00000	2004.0	7.0	19000	-6.375240
5034	Neill Dea Lusa	35.0	0.070071	Thriller	Ian Gamston	Caitie	589	35	English	0.0070	2005.0	6.3	74	0.063071
5035	Robert Rodriguez	56.0	2.406920	Action/Crime/Drama/Romance/Thriller	Carlos Gallardo	El Mariachi	52055	130	Spanish	0.0070	1992.0	6.9	0	2.399700
5037	Edward Burns	14.0	0.004584	Comedy/Drama	Kerry Bishi	Newlyweds	1338	14	English	0.0090	2010.0	6.4	413	-0.004416
5042	Jon Gunn	43.0	0.085222	Documentary	John August	My Date with Drew	4285	84	English	0.0011	2004.0	6.6	456	0.084122

3917 rows × 14 columns

```
## Subtask 3.3: Drop duplicate values

After you found out the top 10 profiling movies, you might have notice a duplicate value. So, it seems like the dataframe has duplicate values as well. Drop the duplicate values from the dataframe and repeat 'Subtask 3.2'.

In [92]: # Write your code for dropping duplicate values here
movies.drop_duplicates(inplace=True)

In [93]: # Write code for repeating subtask 2 here
movies

Out[93]:
```

	director_name	num_critics_for_reviews	gross	genres	actor_1_name	movie_title	num_voted_users	num_user_for_reviews	language	budget	title_year	imdb_score	movie_facebook_likes	profit
0	James Cameron	723.0	76.050547	Action/Adventure/Fantasy/Sci-Fi	CCI Pounder	Avatar	886204	3054	English	237.00000	2009.0	7.9	33000	52.350547
1	Gore Verbinski	302.0	39.0404152	Action/Adventure/Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	47							