



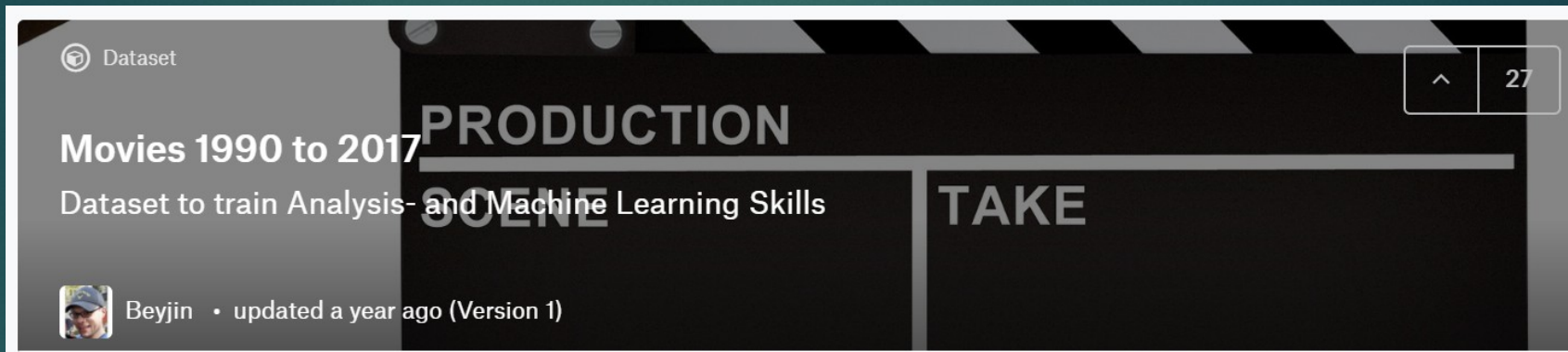
# PREDICTING KEY METRICS FOR MOVIE PRODUCTION

BY NAVEEN ARUMANDLA

# Key Questions Movie Producers Need Answers For

- ▶ Why this project?
- ▶ Why now?
- ▶ Why with this group of talent?
- ▶ With this investment of money?
- ▶ **What's the genre?**
- ▶ **Who is the audience?**
- ▶ **Will this win an award?**
- ▶ **Will the audience like this?**
- ▶ **What is the probability of success?**

# Dataset



```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 import datetime
        6 movies = pd.read_csv('movies-1990-to-2017/Movie_Movies.csv') #Movies dataset with movies from 1900 to 2018
        7 movies.head(5)
```

C:\Users\Jarvis\Anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3049: DtypeWarning: Columns (14) have mixed types. Specify dtype option on import or set low\_memory=False.  
interactivity=interactivity, compiler=compiler, result=result)

Out[1]:

	Awards	Country	DVD	Director	Language	Plot	Poster	Production	Rated	Released	Runtime	Title	Type	Website	Year	imdbID	i
0	NaN	USA	NaN	Rose Cummings	English	Rachel constantly hears her baby cry from the ...	NaN	NaN	NaN	26 Apr 2012	20 min	Baby's Breath	movie	NaN	2012	tt2268369	
1	NaN	USA	NaN	James Byrne	NaN	The struggle against unfortunate circumstances...	NaN	NaN	NaN	NaN	9 min	Winter Trees	movie	NaN	2008	tt1560760	
2	NaN	USA	NaN	Dimitri Buchowetzki	NaN	NaN	NaN	NaN	NaN	27 Mar 1926	50 min	The Crown of Lies	movie	NaN	1926	tt0016750	

# Data Cleaning

- This is mostly a string dataset with no numeric values. So these would have to be converted accordingly

```
In [2]: 1 movies.describe()
```

Out[2]:

	imdbRating
count	62073.000000
mean	6.416236
std	1.360472
min	1.000000
25%	5.600000
50%	6.500000
75%	7.300000
max	10.000000

```
In [3]: 1 movies.dtypes
```

Out[3]:

Awards	object
Country	object
DVD	object
Director	object
Language	object
Plot	object
Poster	object
Production	object
Rated	object
Released	object
Runtime	object





- Duplicate rows with same movie title removed



- Column 'Runtime' cannot be converted to string because it has values followed by strings "h" and ",",

- Slice function used to remove "min" string at the end of column value



- Two columns created to identify which rows have "h" and ",". Identified rows were removed from Dataframe

```
In [10]: 1 movies = movies.drop_duplicates(subset = 'Title')
2 movies.head(5)
```

Out[10]:

Country	Director	Language	Plot	Production	Released	Runtime	Title	Year	imdbID	imdbRating	imdbVotes	Award?	Unnamed: 0	Genre
USA	Rose Cummings	English	Rachel constantly hears her baby cry from the ...	NaN	2012-04-26	20 min	Baby's Breath	2012	tt2268369	NaN	NaN	0	0	Short
USA	Julia Hechler	English	A Gift introduces Samuel Green, Washington Sta...	NaN	2013-05-27	2 min	A Gift	2013	tt3405286	NaN	NaN	0	6	Documentary
Sri Lanka	Udara Siriruwan	Sinhalese	NaN	NaN	2014-03-20	23 min	Journey	2014	tt3816698	NaN	NaN	0	8	Short

In [13]:

```
1 sub1 = 'h'
2 sub2 = ','
3 sub3 = '.'
4 movies['Runtime'] = movies['Runtime'].str.slice(0,-4) #remove last four characters from 'Runtime' column which has 'min' str
5 movies['Runtime_index1'] = movies['Runtime'].str.find(sub1)
6 movies['Runtime_index2'] = movies['Runtime'].str.find(sub2)
7 movies = movies.dropna(subset=['Runtime'])
8 movies.head(5)
```

Out[13]:

Plot	Production	Released	Runtime	Title	Year	imdbID	imdbRating	imdbVotes	Award?	Unnamed: 0	Genre	Runtime_index1	Runtime_index2
Rachel constantly hears her baby cry from the ...	NaN	2012-04-26		Baby's Breath	2012	tt2268369	NaN	NaN	0	0	Short	-1.0	-1.0
A Gift introduces Samuel Green, Washington Sta...	NaN	2013-05-27		A Gift	2013	tt3405286	NaN	NaN	0	6	Documentary	-1.0	-1.0
NaN	NaN	2014-03-20		Journey	2014	tt3816698	NaN	NaN	0	8	Short	-1.0	-1.0

- Dataset has been cleaned to contain necessary variables

- 1) Runtime (in float)
- 2) Movie Rating (in float)
- 3) Votes (in int)
- 4) Identifier for award winning movie



```
In [20]: 1 movies['Runtime'] = pd.to_numeric(movies['Runtime']).astype(float)
2 movies['imdbRating'] = pd.to_numeric(movies['imdbRating']).astype(float)
3 movies['imdbVotes'] = movies['imdbVotes'].fillna('0',inplace=True)
4 movies['Votes_index1'] = movies['imdbVotes'].str.replace(',', '')
5 movies['imdbVotes'] = pd.to_numeric(movies['imdbVotes']).astype(float)
6 movies.dtypes
```

```
Out[20]: Awards          object
Country          object
Director          object
Language          object
Plot              object
Production        object
Released          datetime64[ns]
Runtime           float64
Title             object
Year             object
imdbID            object
imdbRating        float64
imdbVotes         float64
Award?           int32
Unnamed: 0        int64
Genre             object
Runtime_index1    int32
Runtime_index2    int32
Votes_index1      object
dtype: object
```



# Who is the audience? What Genre?

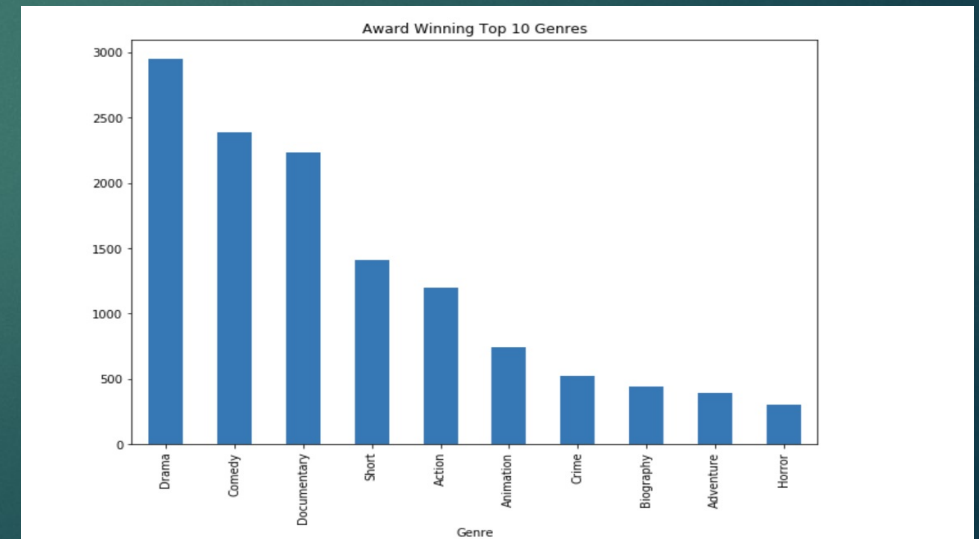
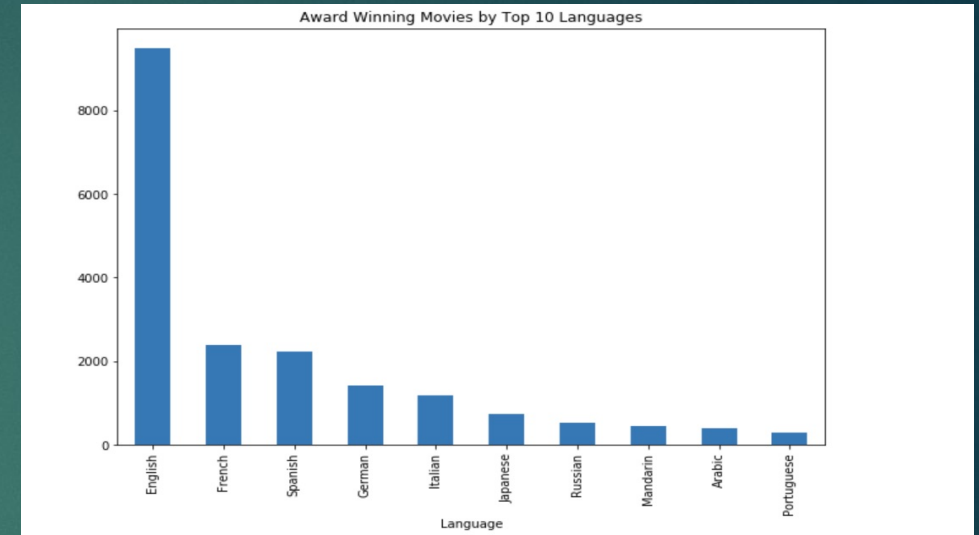
```
In [24]: 1 def new_language(df):
2         new_df = \
3         (df.set_index(df.columns.drop('Language',1).tolist())
4          .Language.str.split(' ',expand=True)
5          .stack()
6          .reset_index()
7          .rename(columns={0:'Language'})
8          .loc[:,df.columns])
9         return new_df

In [28]: 1 val1 = input("Enter preferred movie type to choose a genre and language:\n1. Award winning movie\n2. Highest rated movies (n
2
3 if val1 == '1':
4     award_movies = movies[movies['Award?']==1]
5     print(award_movies.groupby('Genre').size().sort_values(ascending=False).head(10).plot(kind='bar',figsize=(10,7),title="A
6     reshaped_award = new_language(award_movies)
7     x=reshaped_award.groupby(['Language']).size().sort_values(ascending=False).head(10)
8     print("Award Winning Movies by Top 10 Languages")
9     print(x)
10
11 elif val1 == '2':
12     rating = movies[movies['Award?']==0]
13     rating['imdbVotes'] = rating['imdbVotes']>10000
14     rating['imdbRating'] = rating['imdbRating'] > 8
15     print(rating.groupby('Genre').mean()['imdbRating'].sort_values(ascending=False).head(10).plot(kind='bar',figsize=(10,7),
16     reshaped_rating = new_language(rating)
17     x=reshaped_rating.groupby(['Language']).size().sort_values(ascending=False).head(10)
18     print("High Rated Movies (>10000 votes and >8.0 rating) by Top 10 Languages")
19     print(x)
20
```

Enter preferred movie type to choose a genre and language:

1. Award winning movie
2. Highest rated movies (no awards)
3. Award winning and high rating

- Typing 1, gives me a visual for the top 10 languages in which award winning movies were made and top 10 genres they were made in
- Similar graphs will be displayed for choices 2, 3





# Probability of success?

```
In [*]: 1 val2 = input("Enter preferred genre:")
2
3 prob_award = movies[movies['Award?'] == 1]
4 prob_award['imdbRating'] = prob_award['imdbRating'] < 8
5 count_genre_award = prob_award.loc[prob_award['Genre'] == val2].count()
6 probability1 = ((count_genre_award['Genre']) / (prob_award['Genre'].count())) * 100
7 print("Probability of '%s' genre winning an award is %.0f%%" % (val2, probability1))
8
9 probRated = movies[movies['Award?'] == 0]
10 probRated['imdbVotes'] = probRated['imdbVotes'] > 10000
11 probRated['imdbRating'] = probRated['imdbRating'] > 8
12 count_genreRated = probRated.loc[probRated['Genre'] == val2].count()
13 probRated['Genre'].count()
14 probability2 = ((count_genreRated['Genre']) / (probRated['Genre'].count())) * 100
15 print("Probability of '%s' genre being rated above 8 without any award is %.0f%%" % (val2, round(probability2)))
16
17 probBoth = movies[movies['Award?'] == 0]
18 probBoth['imdbVotes'] = probBoth['imdbVotes'] > 10000
19 probBoth['imdbRating'] = probBoth['imdbRating'] > 8
20 count_genreBoth = probBoth.loc[probBoth['Genre'] == val2].count()
21 probBoth['Genre'].count()
22 probability3 = ((count_genreBoth['Genre']) / (probBoth['Genre'].count())) * 100
23 print("Probability of '%s' genre winning an award and being rated above 8 is %.0f%%" % (val2, round(probability3)))

Enter preferred genre: 
```

Enter preferred genre:Drama

```
C:\Users\Jarvis\Anaconda3\lib\site-packages\ipykernel_launcher.py:4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing
after removing the cwd from sys.path.
C:\Users\Jarvis\Anaconda3\lib\site-packages\ipykernel_launcher.py:10: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing
# Remove the CWD from sys.path while we load stuff.
C:\Users\Jarvis\Anaconda3\lib\site-packages\ipykernel_launcher.py:11: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing
# This is added back by InteractiveShellApp.init_path()
```

```
Probability of 'Drama' genre winning an award is 23%
Probability of 'Drama' genre being rated above 8 without any award is 6%
Probability of 'Drama' genre winning an award and being rated above 8 is 6%
```

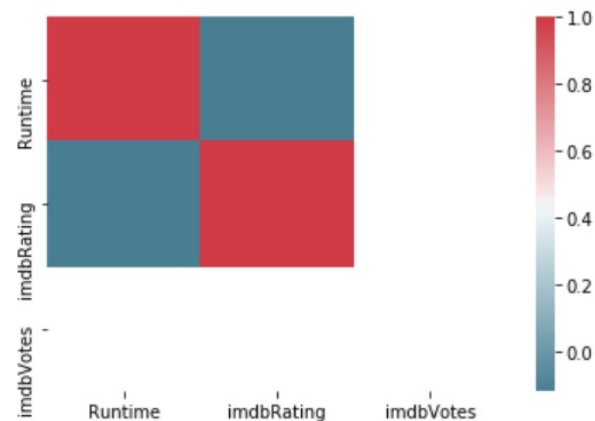
Inputting the preferred type of genre, will provide the probability of the genre:

- 1) Winning an award
- 2) Being rated 8.0 or more on IMDB
- 3) Both being awarded and being rated above 8.0

# Relationship between Variables

```
In [21]: 1 corr = movies.loc[:,movies.dtypes == 'float64'].corr()  
        2 sns.heatmap(corr, xticklabels=corr.columns, yticklabels=corr.columns, cmap=sns.diverging_palette(220, 10, as_cmap=True))
```

Out[21]: <matplotlib.axes.\_subplots.AxesSubplot at 0x160f79aa128>



“Number of votes” on IMDB and “IMDB rating” have a positive but weak correlation.

**THEREFORE, more voting population means better movie ratings**



# Evidence

## ▶ DATA MODEL VALIDATED

The movie poster for 'The Dark Knight' (2008) features Batman standing in front of a burning Gotham City skyline.

1. **The Dark Knight** (2008)

PG-13 | 152 min | Action, Crime, Drama

★ 9 ☆ Rate **84** Metascore

When the menace known as The Joker emerges from his mysterious past, he wreaks havoc and chaos on the people of Gotham. The Dark Knight must accept one of the greatest psychological and physical tests of his ability to fight injustice.

Director: Christopher Nolan | Stars: Christian Bale, Heath Ledger, Aaron Eckhart, Michael Caine

Votes: 2,070,814 | Gross: \$534.86M

[On Disc at Amazon](#)

The movie poster for 'The Dark Knight Rises' (2012) features Batman in a dark, dramatic pose with a fiery background.

2. **The Dark Knight Rises** (2012)

PG-13 | 164 min | Action, Thriller

★ 8.4 ☆ Rate **78** Metascore

Eight years after the Joker's reign of anarchy, Batman, with the help of the enigmatic Catwoman, is forced from his exile to save Gotham City, now on the edge of total annihilation, from the brutal guerrilla terrorist Bane.

Director: Christopher Nolan | Stars: Christian Bale, Tom Hardy, Anne Hathaway, Gary Oldman

Votes: 1,389,868 | Gross: \$448.14M

[On Disc at Amazon](#)

The movie poster for 'Batman Begins' (2005) features Batman in a dark, atmospheric setting with a city skyline in the background.

3. **Batman Begins** (2005)

PG-13 | 140 min | Action, Adventure

★ 8.2 ☆ Rate **70** Metascore

After training with his mentor, Batman begins his fight to free crime-ridden Gotham City from corruption.

Director: Christopher Nolan | Stars: Christian Bale, Michael Caine, Ken Watanabe, Liam Neeson

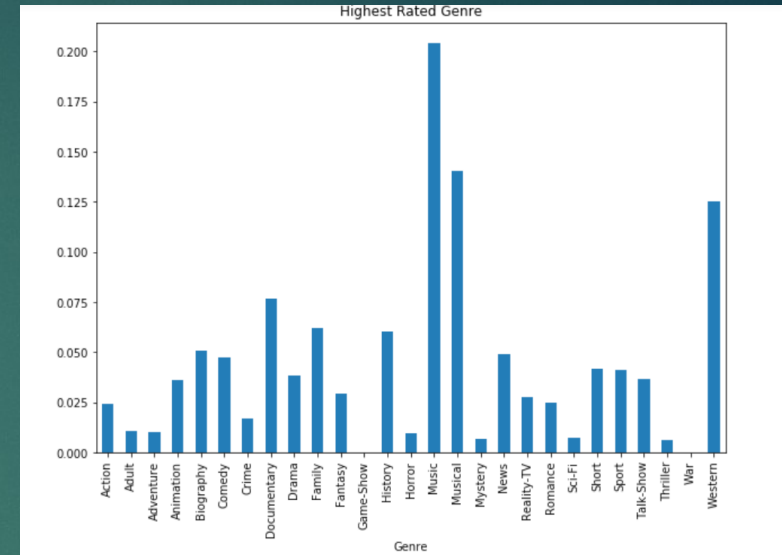
Votes: 1,196,382 | Gross: \$206.85M

[On Disc at Amazon](#)



# Normalizing Data

- While calculating average rating of genres, values had massive variance

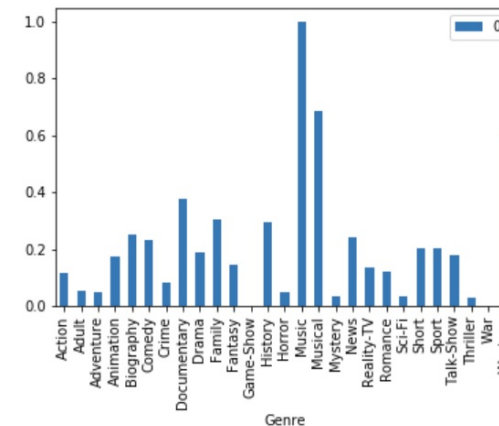


- So values were normalized to test for accuracy of analysis



```
In [63]: 1 #normalize means to check if result remains same
2
3 norm_mean = movies[movies['Award']==0]
4 norm_mean['imdbVotes'] = norm_mean['imdbVotes'] > 10000
5 norm_mean['imdbRating'] = norm_mean['imdbRating'] > 8
6 norm_mean = norm_mean.groupby('Genre').mean()['imdbRating']
7 norm_mean = pd.DataFrame(norm_mean)
8 scaler = preprocessing.MinMaxScaler()
9 norm_mean_scaled = scaler.fit_transform(norm_mean)
10 scaled_mean = pd.DataFrame(norm_mean_scaled)
11 scaled_mean.set_index(norm_mean.index).plot(kind='bar')
```

Out[63]: <matplotlib.axes.\_subplots.AxesSubplot at 0x12b4eadc0f0>



# Additional Guidelines

```
In [*]: 1 val3 = input("Enter preferred language:")
2
3 run_award = movies[movies['Award?'] == 1]
4 avg_run_award = run_award.loc[run_award['Language']== val3].mean()
5 avg_run_award = avg_run_award['Runtime'].mean()
6 print("Average runtime of award winning movie in '%s' language is %.0f minutes" % (val3,avg_run_award))
7
8 run_rated = movies[movies['Award?'] == 0]
9 run_rated['imdbVotes'] = run_rated['imdbVotes']>10000
10 run_rated['imdbRating'] = run_rated['imdbRating']>8
11 avg_run_rated = run_rated.loc[run_rated['Language']== val3].mean()
12 avg_run_rated = avg_run_rated['Runtime'].mean()
13 print("Average runtime of high rated (>8.0) movie with no awards in '%s' language is %.0f minutes" % (val3,avg_run_rated))
14
15 run_both = movies[movies['Award?'] == 0]
16 run_both['imdbVotes'] = run_both['imdbVotes']>10000
17 run_both['imdbRating'] = run_both['imdbRating']>8
18 avg_run_both = run_both.loc[run_both['Language']== val3].mean()
19 avg_run_both = avg_run_both['Runtime'].mean()
20 print("Average runtime of award winning high rated (>8.0) movie in '%s' language is %.0f minutes" % (val3,avg_run_both))
```

Enter preferred language:

```
Enter preferred language:Japanese
Average runtime of award winning movie in 'Japanese' language is 104 minutes
```

```
C:\Users\Jarvis\Anaconda3\lib\site-packages\ipykernel_launcher.py:9: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
if __name__ == '__main__':
C:\Users\Jarvis\Anaconda3\lib\site-packages\ipykernel_launcher.py:10: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
# Remove the CWD from sys.path while we load stuff.
```

```
Average runtime of high rated (>8.0) movie with no awards in 'Japanese' language is 89 minutes
Average runtime of award winning high rated (>8.0) movie in 'Japanese' language is 89 minutes
```

Inputting the preferred type of language, will provide the average runtime of the three choices by language

## Other improvements which can be made?

- Budget required
- Choosing a cast member
- Correlation between box office collection and budget

