# COURSERA
# Applied Data Science
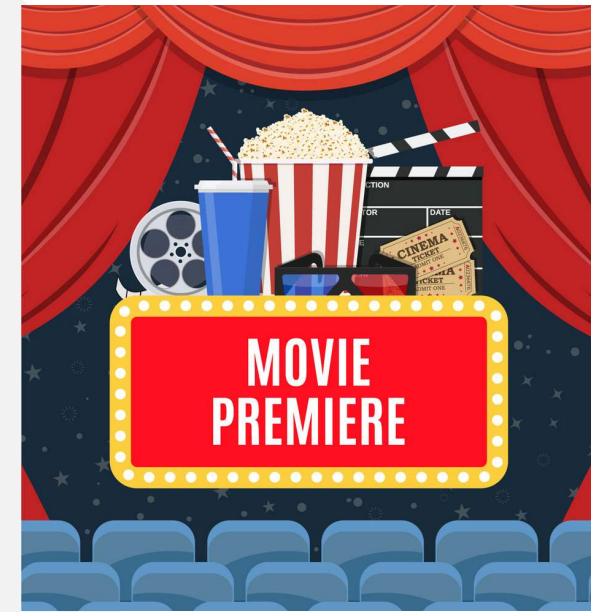
Capstone Project

Movie earnings prediction

# 1. Background

Cinematography is a money generating business.
- The producer and film studio do not only make art, they also want to make money.
- The actors who salary is percentage of the profit.
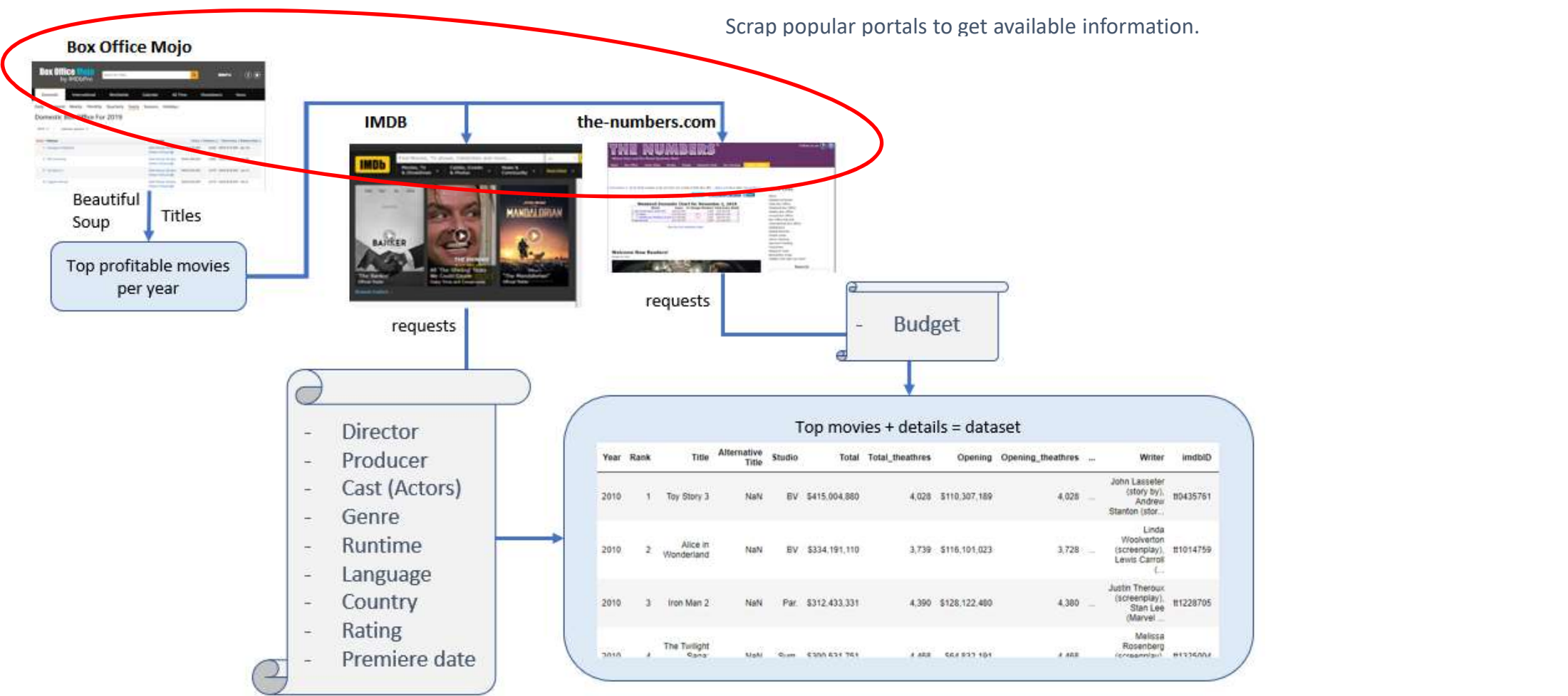- The investors who put money in the production.
They all want to ensure their investment will be profitable.

What will be the movie profit?

MOVIE PREMIERE

Producer

Film Studio

Actors

Producer

Investor

## 2. Data Source / Approach

Which information are available before the movie premiere to predict the earnings?
And how to get them?

Scrap popular portals to get available information.

**Box Office Mojo**

**IMDB**

**the-numbers.com**

Beautiful Soup

Titles

Top profitable movies per year

requests

requests

- Budget

requests

- Director
- Producer
- Cast (Actors)
- Genre
- Runtime
- Language
- Country
- Rating
- Premiere date

Top movies + details = dataset

| Year | Rank | Title | Alternative Title | Studio | Total | Total_theathres | Opening | Opening_theathres | ... | Writer | imdbID |
|------|------|-------|-------------------|--------|-------|-----------------|---------|-------------------|-----|--------|--------|
| 2010 | 1 | Toy Story 3 | NaN | BV | $415,004,880 | 4,028 | $110,307,189 | 4,028 | ... | John Lasseter (story by), Andrew Stanton (stor... | tt0435761 |
| 2010 | 2 | Alice in Wonderland | NaN | BV | $334,191,110 | 3,739 | $116,101,023 | 3,728 | ... | Linda Woolverton (screenplay), Lewis Carroll (... | tt1014759 |
| 2010 | 3 | Iron Man 2 | NaN | Par. | $312,433,331 | 4,390 | $128,122,480 | 4,380 | ... | Justin Theroux (screenplay), Stan Lee (Marvel ... | tt1228705 |
| 2010 | 4 | The Twilight Saga | NaN | Sum. | $300,531,751 | 4,468 | $64,832,191 | 4,468 | ... | Melissa Rosenberg (screenplay), ... | tt1325004 |

Create a dataset of most successful titles and their details in order as input for the prediction model

# 4. Exploratory data analysis

**Available features**

**Numerical**

Example: Runtime, Budget



No strong relations between movie budget and its earning

**Categorical**

Example: Actors, Director, Plot, Genre



<u>Problem:</u> how to evaluate importance of:
- Actors – 1778 enries
- Director – 539 entries
- Writer – 1478 enries
- Plot – 9814 words

**These features need to be one hot encoded which creates high dimensionality sparse matrix**

**Target**

Opening earnings



"Opening earnings" destribution with outliers

Eliminate outliers

"Opening earnings" destribution without outliers

Use opening weekend earnings instead of total earning, as the result can be evaluated quicker.

# 5. Methodology

**Preparation**

Feature engineering

**Numerical**
- Clean
- Convert to numeric type

**Categorical**
- Clean
- Hashing trick

Combine in one big data input matrix

Min Max Scaler



Dimension: 802 rows x 6443 columns

Sparse matrix

Split & apply MinMaxScaler

# 6. Modelling

Which machine learning algorithm performs best, if the input is high dimensional sparse matrix?

Use different models and evaluate the best one based on R squared coefficient.

| Model | R2 train data | R2 test data | Comment |
|-------|---------------|--------------|---------|
| Linear regression | 0.820 | -1.840 | Linear model is not appropriate for the dataset |
| Decision Tree | 0.238 | 0.185 | The result is very poor |
| Neural Network | 0.997 | 0.217 | Overfitting |
| Gradient Boosting | 0.671 | 0.435 | Overfitting, however train/test result is closer |
| Support Vector Machine | 1.0 | 0.0 | SVM does not provide any significant result |

Best performing models are Gradient Boosting and Neural Networks, however the result is still not satisfactory for model deployment. Possibly the result could be better with parameter tuning.

The reason may be the sparse data. Example:
Only 567 actors (out of 1778 total) played in more than 1 movie. This does not allow the model to use the training data, because test data contains new, different information.

# 7. Conclusion

The model does not perform well enough to allow deployment in business.
Reasons are:
- No linearly correlated features
- Sparse data: train data is not sufficient for modeling as test data contain new information.

| Movie | Actual earnings $ | Predicted earnings $ | Delta |
|---|---|---|---|
| 1 | 24'830'443 | 15'734'224 | -36.6% |
| 2 | 21'052'227 | 27'278'085 | 29.6% |
| 3 | 10'609'795 | 24'447'134 | 130.4% |

Not consistent enough.

Recommendation:

- create bigger training data set

- find a way to evaluate importance of each feature

- create hyper feature to better describe the data

- test other, more sophisticated algorithms

- fine tune model parameters