# CP1 - What movies make the most money at the box office
# Milestone Report 1

## Problem Statement

In a world where movies made an estimated $41.7 billion in 2018, the film industry is more popular than ever. But what movies make the most money at the box office? How much does a director matter? Or the budget?
Can we build models, which will be able to accurately predict film revenue?

This project will help film production companies understand key features of having high revenue.

## Data Source

The major data source comes from the public dataset uploaded to Kaggle.com (https://www.kaggle.com/c/tmdb-box-office-prediction/overview).

This dataset with metadata on over 7,000 past films from The Movie Database. Data points provided include cast, crew, plot keywords, budget, posters, release dates, languages, production companies, and countries.

## Data Cleaning

There are 8 JSON-style columns. We will parse them and create categorical and dummy variables. For example, column "belongs_to_collection":

**belongs_to_collection**

```
for i, e in enumerate(master['belongs_to_collection'][:5]):
    print(i, e)
```
```
0 [{'id': 313576, 'name': 'Hot Tub Time Machine Collection', 'poster_path': '/iEhb00TGPucF0b4joM1ieyY026U.jpg', 'backdrop_p
ath': '/noeTVcgpBiD48fDjFVic1Vz7ope.jpg'}]
1 [{'id': 107674, 'name': 'The Princess Diaries Collection', 'poster_path': '/wt5AMbxPTS4Kfjx7Fgm149qPfZl.jpg', 'backdrop_p
ath': '/zSEtYD77pKRJlUPx34BJgUG9v1c.jpg'}]
2 nan
3 nan
4 nan
```

We create two new columns from column "belongs_to collection", first one is collection name and second one has collection or not. We assume that other information from this column we can't use for future prediction.

```
master['collection_name'] = master['belongs_to_collection'].apply(lambda x: x[0]['name'] if x != {} else 0)
master['has_collection'] = master['belongs_to_collection'].apply(lambda x: len(x) if x != {} else 0)

master = master.drop(['belongs_to_collection'], axis=1)
```
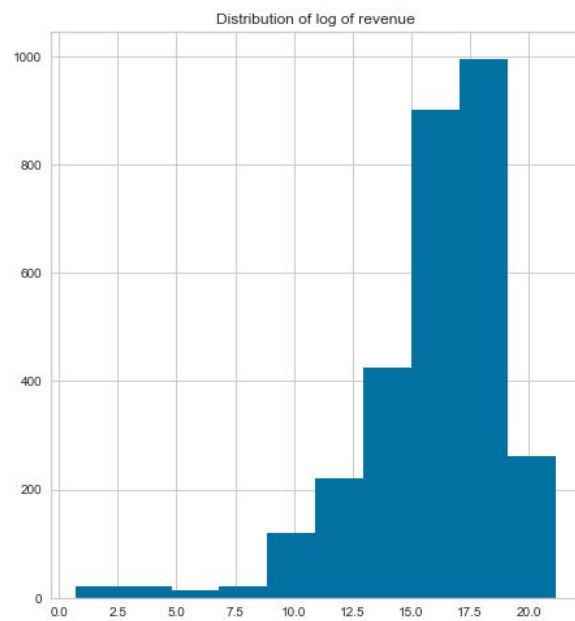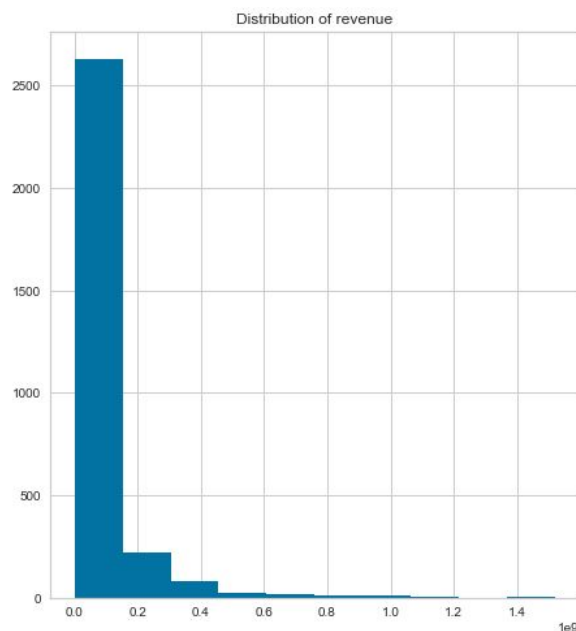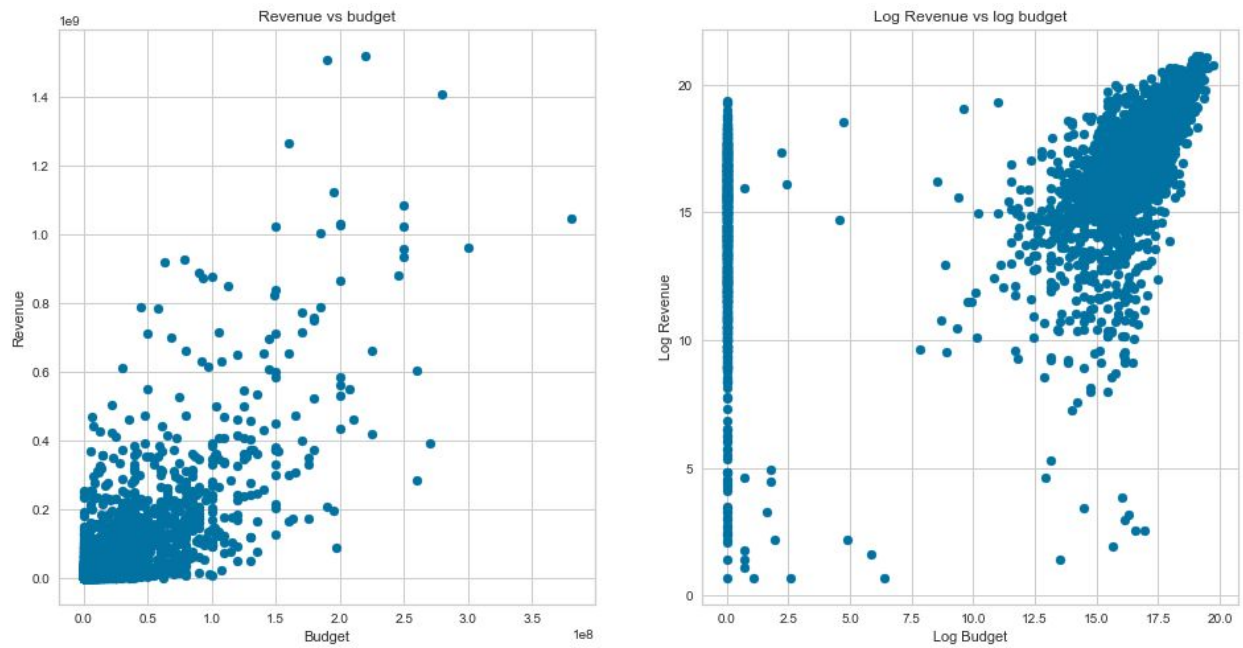
## Exploratory Analysis

**SOME INSIGHTS**

- Average length of a movie is about **107 minutes**
- The **longest** movie: *"Carlos"* is **5hrs 38min** long
- The **most popular** movies: *"Wonder Woman"* and *"Beauty and the Beast"*
- Movies with the **biggest budget**: *"Pirates of the Caribbean: On Stranger Tides"* and *"Pirates of the Caribbean: At World's End"*
- The **biggest** movie **producers**: *Warner Bros.* and *Universal Pictures*
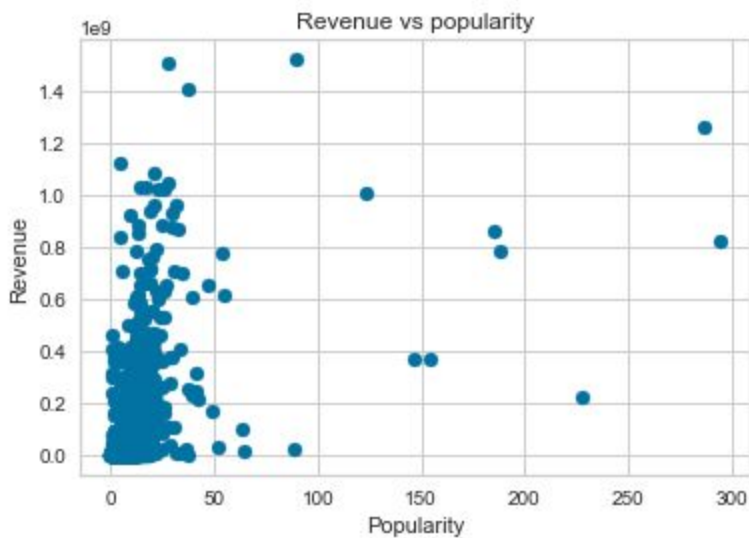- The most popular **genres**: *Drama* and *Comedy*

Let first look at the histograms of the target feature - **revenue**. As we can see revenue distribution has a high skewness. It is better to use log transformation of revenue.
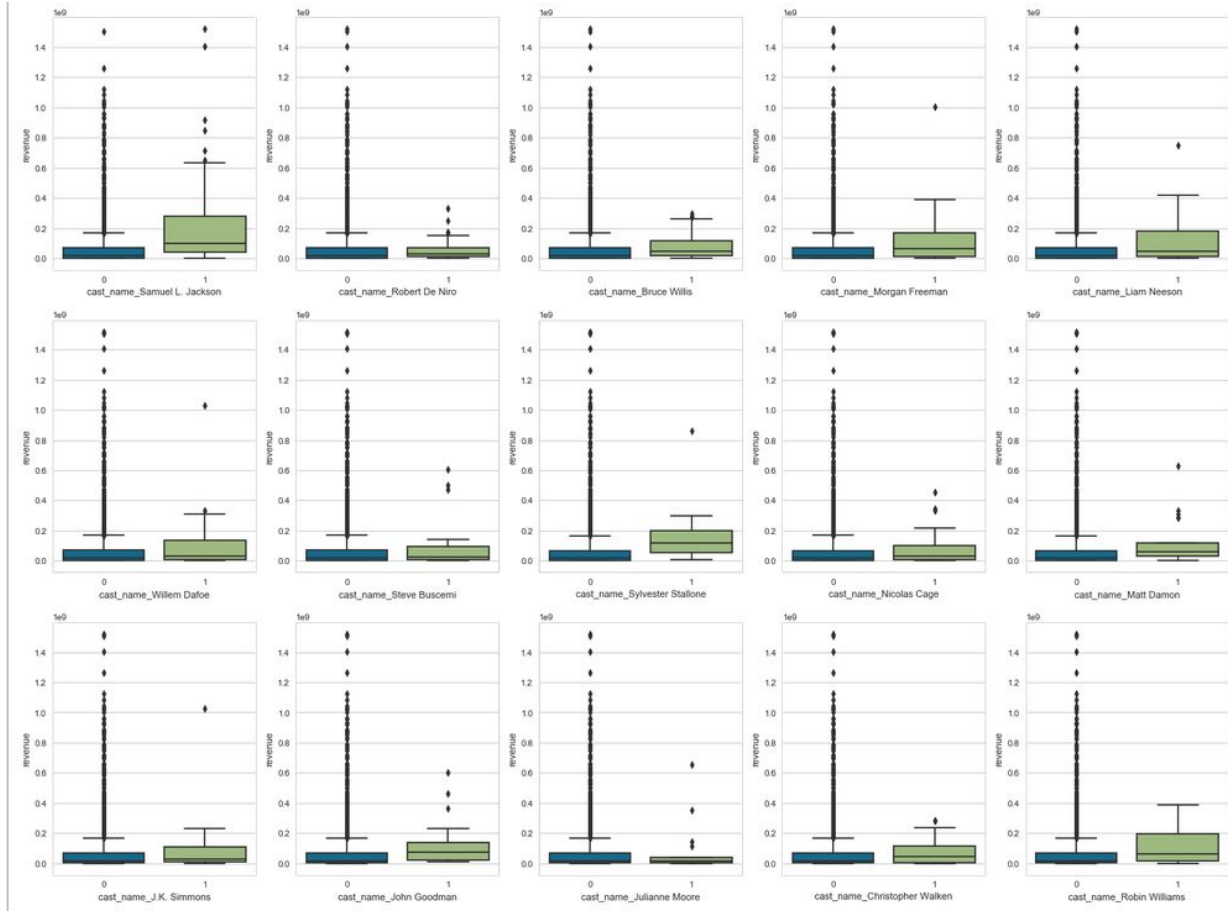
First feature we will look at is **budget**, intuitively the budget should be somehow correlated with revenue.



Let's look at popularity. We can see some clear trends that an increase in popularity tends to lead to higher revenue.



We have columns with most common cast members, so let's plot boxplots.

As you can see mostly films with these actors tend to have higher revenue.
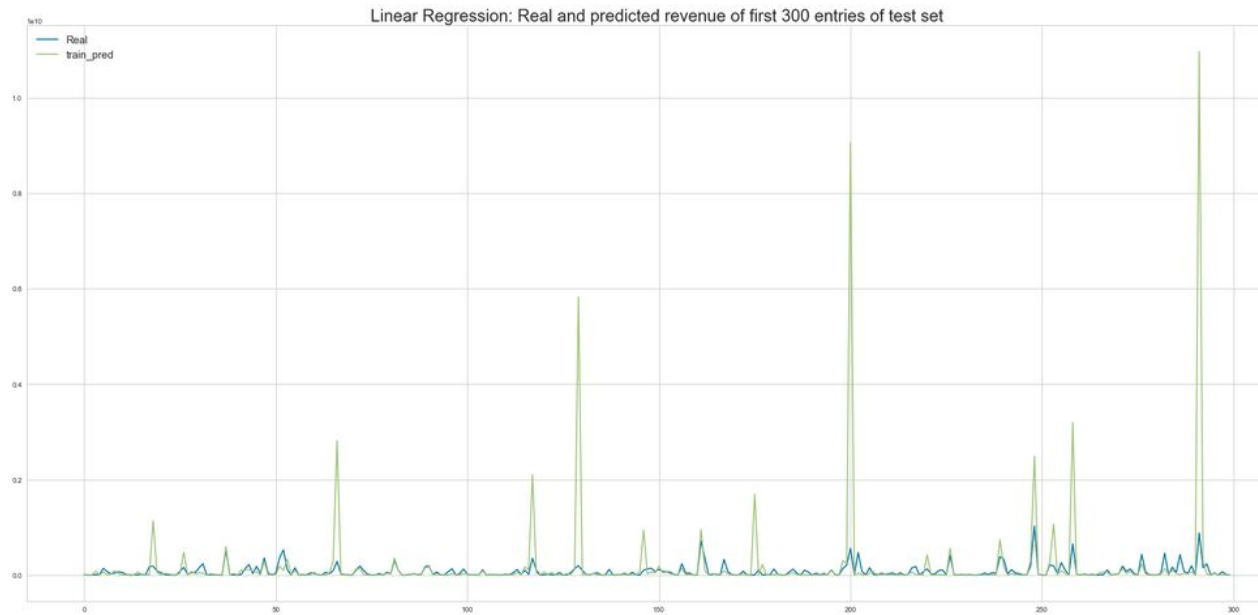
# Machine Learning

## 1. Linear Regression

Once the data was transformed and usable for machine learning, the data was split into training and test sets with the neighbor players appended to the test set to make sure predictions applied to them.

```python
from sklearn.model_selection import train_test_split
```
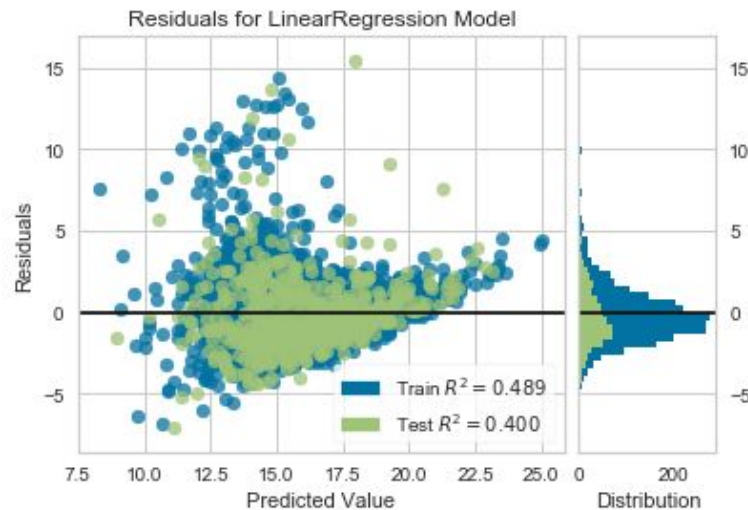
```python
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```python
#Linear Regression
regressor = LinearRegression()
regressor.fit(X_train, y_train)
y_pred=regressor.predict(X_test)
```

Linear Regression: Real and predicted revenue of first 300 entries of test set

```
#Linear Regression metrics
pd.set_option('display.float_format', lambda x: '%.3f' % x)
print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))
print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_pred))
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
print('Root Mean Squared Log Error:',np.sqrt(mean_squared_log_error( y_test, y_pred )))
```

```
Mean Absolute Error: 1.6500602047686255
Mean Squared Error: 5.356984834769914
Root Mean Squared Error: 2.314516112445518
Root Mean Squared Log Error: 0.20010788070298655
```
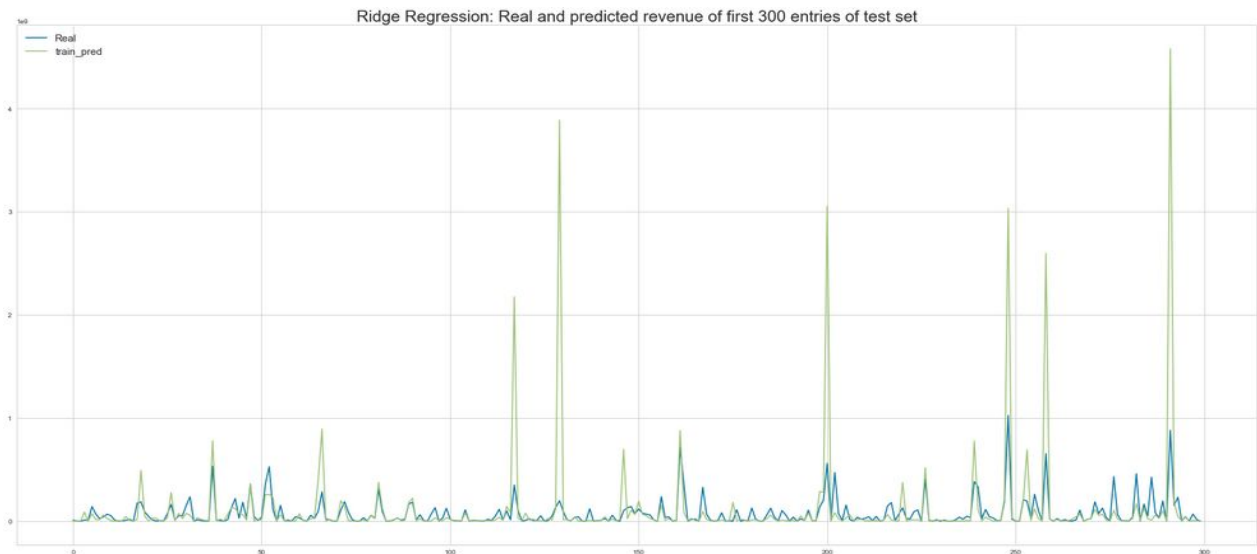

Residuals for LinearRegression Model

It was observed that the residuals between the predicted values and actual values roughly fit the same trends according to a Linear Regression model.
It was also observed that extremely higher revenues were underpredicted while lower values were predicted closer to their actual valuation.
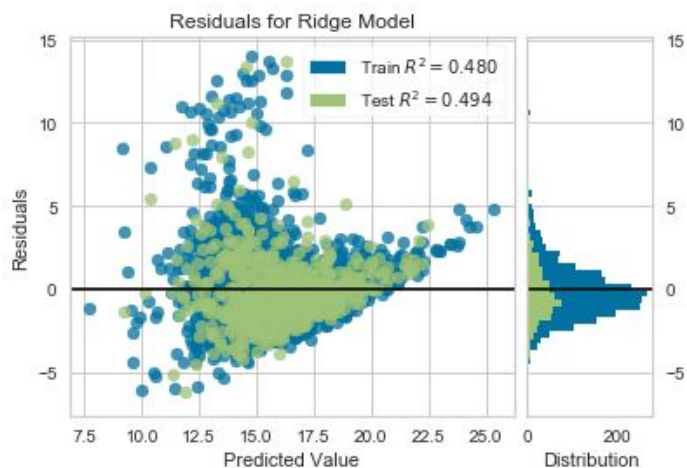
## 2. Ridge Regression

Lets try another model - Ridge Regression.



Ridge Regression: Real and predicted revenue of first 300 entries of test set

```
#Ridge Regression metrics
pd.set_option('display.float_format', lambda x: '%.3f' % x)
print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred_clf))
print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_pred_clf))
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_pred_clf)))
print('Root Mean Squared Log Error:', np.sqrt(mean_squared_log_error( y_test, y_pred_clf )))
```

```
Mean Absolute Error: 1.4990212383799935
Mean Squared Error: 4.506682175538423
Root Mean Squared Error: 2.122894763180319
Root Mean Squared Log Error: 0.1905475484817239
```



Residuals for Ridge Model

As we can see results are pretty similar to linear regression with small improvement.

## 3. XGBoost

XGBoost is a powerful machine learning algorithm especially where speed and accuracy are concerned.

If things don't go your way in predictive modeling, use XGboost. The XGBoost algorithm has become the ultimate weapon of many data scientists. It's a highly sophisticated algorithm, powerful enough to deal with all sorts of irregularities of data.

Building a model using XGBoost is easy. But, improving the model using XGBoost is difficult. This algorithm uses multiple parameters. To improve the model, parameter tuning is must.

So using the code below I am trying to fine tune XGBoost parameters.

```python
X_train, X_test, y_train, y_test = train_test_split(df_X_train, df_y_train, test_size=0.2, random_state=42)
# read in data
dtrain = xgb.DMatrix(X_train, label=y_train,feature_names=feature_names)
dtest = xgb.DMatrix(X_test, label=y_test,feature_names=feature_names)
watchlist = [(dtrain, 'train'), (dtest,'test')]
```
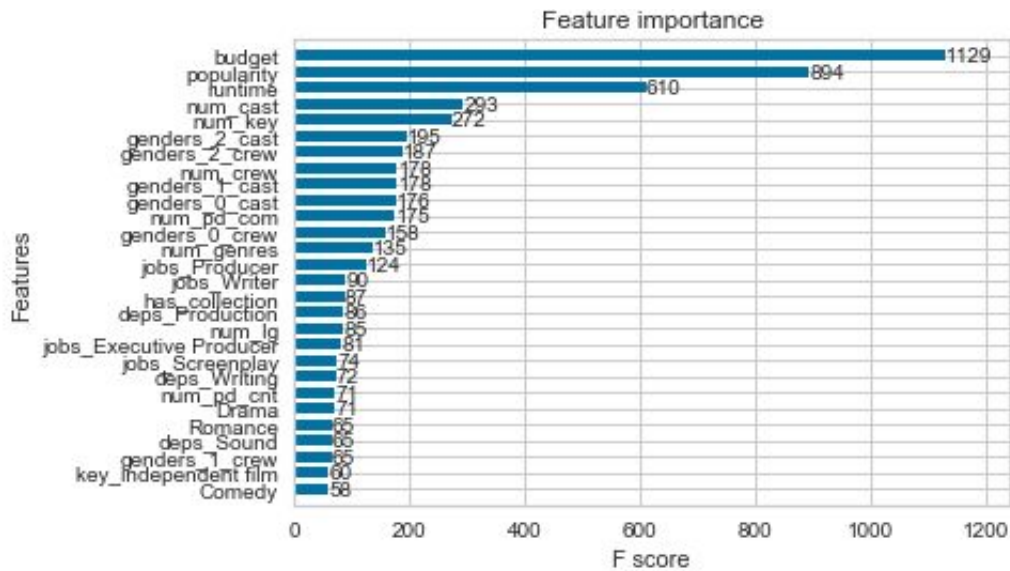
```python
md = [4,5,6]
lr = [0.01,0.05,0.1,0.3]
mcw = [1,5,20,25,30]
for m in md:
    for l in lr:
        for n in mcw:
            pd.set_option('display.float_format', lambda x: '%.3f' % x)
            print(m,l,n)
            t0 = datetime.now()
            xgb_pars = {'min_child_weight': n, 'eta': l, 'colsample_bytree': 0.9,
                        'max_depth': m,'subsample': 0.9, 'lambda': 1., 'nthread': -1, 'booster' : 'gbtree', 'silent': 1,
                'eval_metric': 'rmse', 'objective': 'reg:linear'}
            model = xgb.train(xgb_pars, dtrain, 500, watchlist, early_stopping_rounds=50,
                maximize=False, verbose_eval=0)
```

So according to the code above we have following parameters with min test RMSE:
max-depth - 6
learning rate - 0.05
min child weight - 5

Feature importance

According to the feature importance plot we have next top 3 features : budget, popularity and runtime. Below I create a dataframe with these features and revenue to show how it looks like.

Also below how these features look for 5 films from the dataset

| | budget | popularity | runtime | num_cast | num_key | genders_2_cast | genders_2_crew | genders_1_cast | num_crew | genders_0_cast | num_pd_com |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3500000 | 0.556 | 90.000 | 11 | 2 | 3 | 2 | 5 | 2 | 3 | 1 |
| 1 | 0 | 2.087 | 100.000 | 7 | 7 | 2 | 5 | 1 | 13 | 4 | 2 |
| 2 | 2000000 | 1.189 | 89.000 | 3 | 1 | 2 | 0 | 1 | 1 | 0 | 1 |
| 3 | 98000000 | 7.284 | 119.000 | 31 | 6 | 20 | 12 | 2 | 16 | 9 | 4 |
| 4 | 0 | 1.219 | 101.000 | 7 | 0 | 3 | 0 | 0 | 2 | 4 | 1 |

| genders_0_crew | num_genres | jobs_Producer | jobs_Writer | has_collection | deps_Production | num_lg | jobs_Executive Producer | jobs_Screenplay | deps_Writing |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 8 | 2 | 2 | 2 | 0 | 2 | 3 | 0 | 0 | 2 |
| 1 | 3 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 2 | 2 | 4 | 0 | 0 | 6 | 2 | 1 | 2 | 4 |
| 2 | 4 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 |

| num_pd_cnt | Drama | genders_1_crew | deps_Sound | Romance | key_independent film | Comedy | revenue | predicted_revenue | title |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 0 | 1 | 16.040 | 14.407 | Ringmaster |
| 1 | 0 | 0 | 3 | 0 | 0 | 0 | 2.079 | 13.655 | He-Man and She-Ra: The Secret of the Sword |
| 1 | 1 | 0 | 0 | 0 | 1 | 1 | 10.425 | 11.509 | Cowboys & Angels |
| 4 | 0 | 2 | 0 | 0 | 0 | 0 | 16.120 | 18.579 | Cutthroat Island |
| 1 | 1 | 0 | 0 | 0 | 0 | 0 | 16.003 | 13.237 | We Are from the Future 2 |

# Further Research Work

Though I have tried pretty much works of data exploration and model selection, it is only the introduction of the whole work and the accuracy of my model is not high enough. There are still so many gaps and now I'd like to study all capabilities of improvement one by one.

Possible parts to be improved are:

- **feature engineering**
  - there are lots of features to be created.
  - missing values
  - outliers
- **feature selection**
  - feature selection methods
  - consideration of validation of selected features
- **model selection**
  - other models
  - how to improve the performance of a model(specialize for the selected model)
  - model parameter tuning