

AUGUST 2020



ANALYZING MOVIE DATA WITH PYTHON

By: Luke Bravo, Matt Jiang, Ali Daanesh Sayyed

MASTERS OF SCIENCE BUSINESS ANALYTICS 2021

The University of Texas at Austin



Overview

- Description of Project Goals
- Exploratory Analysis
- Solution and Insights



Description - What We're Investigating

- Dataset on movies released between 1927 - 2016
- Includes information
 - movie's production
 - reviews, genres
 - actors involved
 - directors involved
 - post production value
 - variety of other variables



Steps We Took

- Initially we did an Exploratory Analysis
- Realized there was significant amount of data on production of movies
- Questions we asked about the data:
 - Can we predict movie rating?
 - Can we construct a reliable model to predict revenue & budget based on the variables we have in the data?
 - Which movies, actors and directors are most notable in the features we are using to analyze?
 - Can We Explore various predictive modeling methods to analyze our data?

Importance of the Problem

- Netflix, Hulu, and HBO Max want to understand what drives movies success
- Our analyses provide helpful insights for big media companies to create the most attractive content for their audience
- Potentially help increase profit margins for media companies

Importance of the Problem Continued

- Looking for:
 - trends within the industry that provide economic insights to how movies generate revenue
 - factors that contribute movie success
 - help media companies select the best cast
 - popular genres in order to maximize gross profit

Exploratory Analysis

- Adjusting prices to 2020 USD
 - Pulled data on monthly CPI and inflation
 - Used PANDAS to calculate annual averages and scale each movie based on release date
- Simple Explorations
- Regression Based Explorations

Longest & Shortest Movie in Dataset

```
df['duration'].max()
```

```
511.0
```

```
df['duration'].min()
```

```
7.0
```

How many unique movie titles have The in their name

```
df['movie_title'].dropna().drop_duplicates()\n.map(lambda s: s.split(' ')[0] == 'The').value_counts()
```

```
False    3934
```

```
True       982
```

```
Name: movie_title, dtype: int64
```

Top 10 directors in dataset

```
df['director_name'].value_counts()[:10]
```

Steven Spielberg	26
Woody Allen	22
Martin Scorsese	20
Clint Eastwood	20
Spike Lee	16
Ridley Scott	16
Steven Soderbergh	15
Renny Harlin	15
Tim Burton	14
Oliver Stone	14

```
Name: director_name, dtype: int64
```

Top 10 Genres

```
df['genres'].value_counts()[:10]
```

Drama	233
Comedy	205
Comedy Drama	189
Comedy Drama Romance	185
Comedy Romance	157
Drama Romance	150
Crime Drama Thriller	98
Horror	67
Action Crime Drama Thriller	65
Crime Drama	62

Name: genres, dtype: int64

```
def splitem(s):  
    return s.split('|')[0]  
  
df['onegenre'] = df['genres'].map(splitem)  
df['onegenre'].value_counts()
```

Comedy	1313
Action	1113
Drama	944
Adventure	439
Crime	340
Biography	250
Horror	221
Documentary	84
Animation	61
Fantasy	48
Mystery	32
Thriller	21
Sci-Fi	13
Western	12
Family	11
Romance	6
Musical	4
Film-Noir	1
Game-Show	1
Music	1
History	1

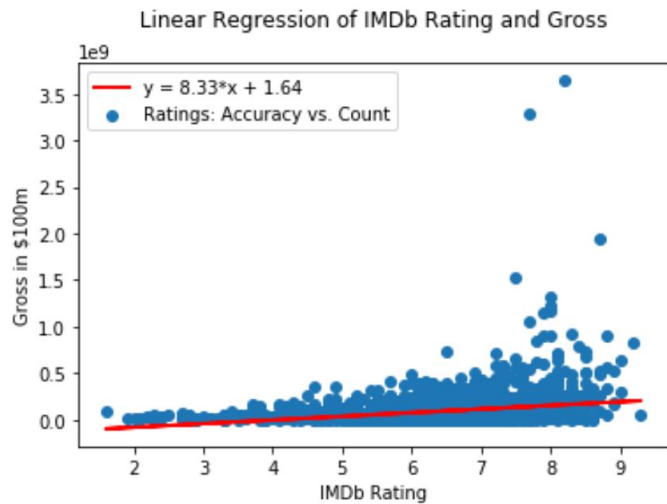
Name: onegenre, dtype: int64

Top 10 Scoring Genres

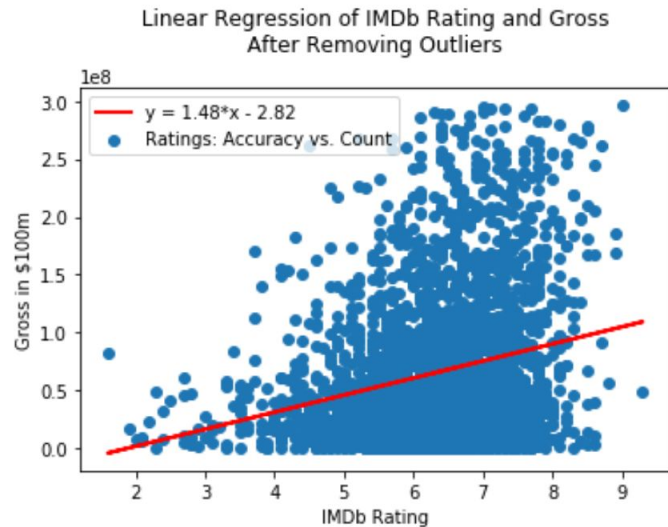
```
df.groupby('onegenre')['imdb_score'].mean()\  
    .sort_values(ascending = False)[:15]
```

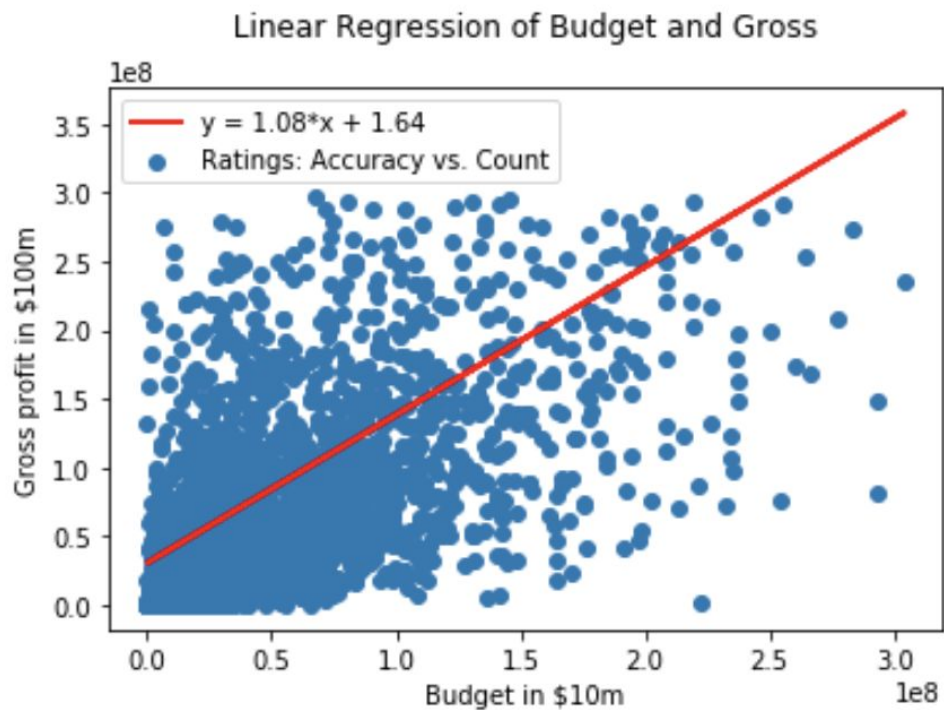
onegenre	
Film-Noir	7.600000
History	7.500000
Music	7.200000
Documentary	7.167857
Biography	7.157600
Crime	6.902941
Drama	6.767161
Animation	6.631148
Western	6.583333
Mystery	6.534375
Adventure	6.530068
Fantasy	6.381250
Action	6.231626
Comedy	6.194136
Musical	6.000000

Name: imdb_score, dtype: float64

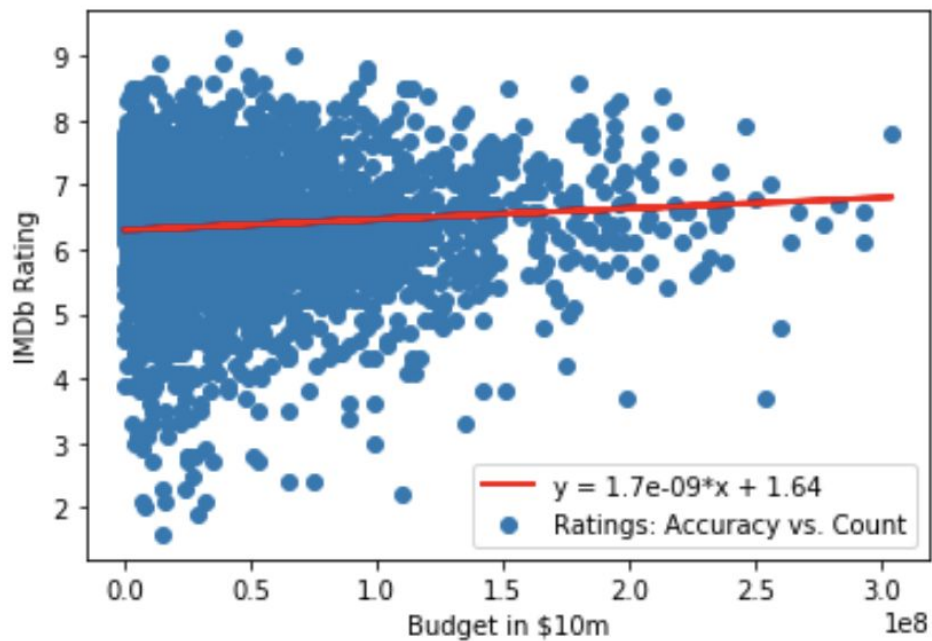


Remove outliers



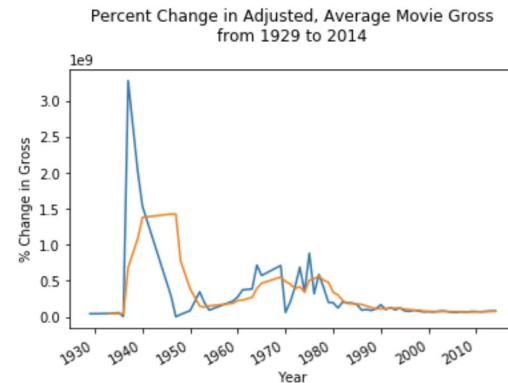
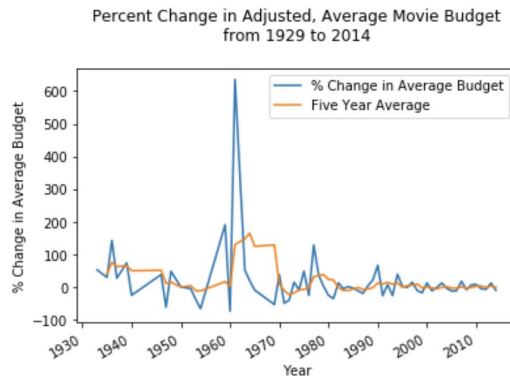
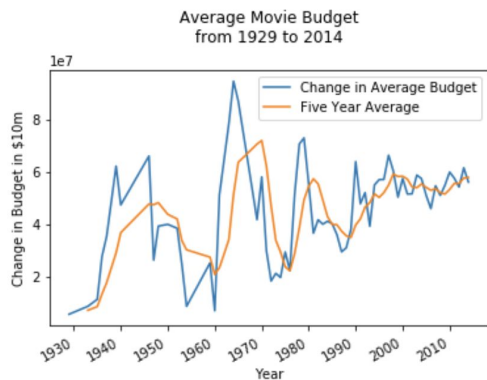


Linear Regression of Budget and IMDb Rating



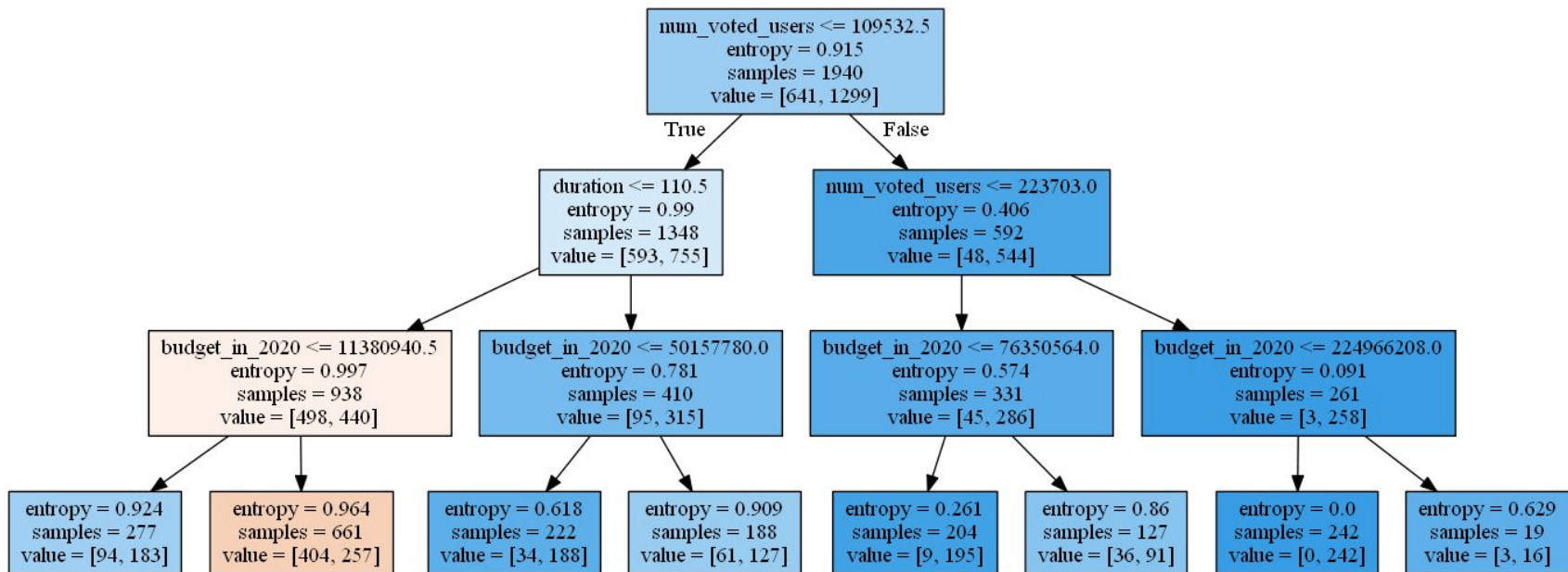
Solution and Insights

- OLS regressions showed us our data on its own cannot solve the problems on its own;
- Low R^2 suggests there are external variables that attribute to the patterns revealed by OLS.
- In light of this: used time series to get big picture of external variables' effect



Solution and Insights

- Decision trees that has no limit on max depth and number of variables overfits the data
- Stratified K-fold cross validation determines an optimal depth and limit variables complexity
- Decision trees model is not without flaws and has limitation in accuracy prediction
- Multivariate regressions to predict adjusted budget do not provide reliable data



=====

=====

Omnibus:	793.880	Durbin-Watson:	0.555
Prob(Omnibus):	0.000	Jarque-Bera (JB):	2204.359
Skew:	1.503	Prob(JB):	0.00
Kurtosis:	6.171	Cond. No.	1.89e+05

Recap

- Description of Project Goals
 - Movie Dataset; How can we help grow companies like Netflix
- Exploratory Analysis
 - Simple Analyses
 - Regression Analyses
- Solution and Insights
 - OLS Regressions
 - Decision Tree Models

