

# PROBLEM STATEMENT AND HYPOTHESIS

Can a movie's financial success be predicted?

This project creates a linear machine-learning model that does just that!

This is an example of a supervised learning problem with a continuous response



- Haxan Films
- \$600K Production Budget
- \$140M Domestic Box Office



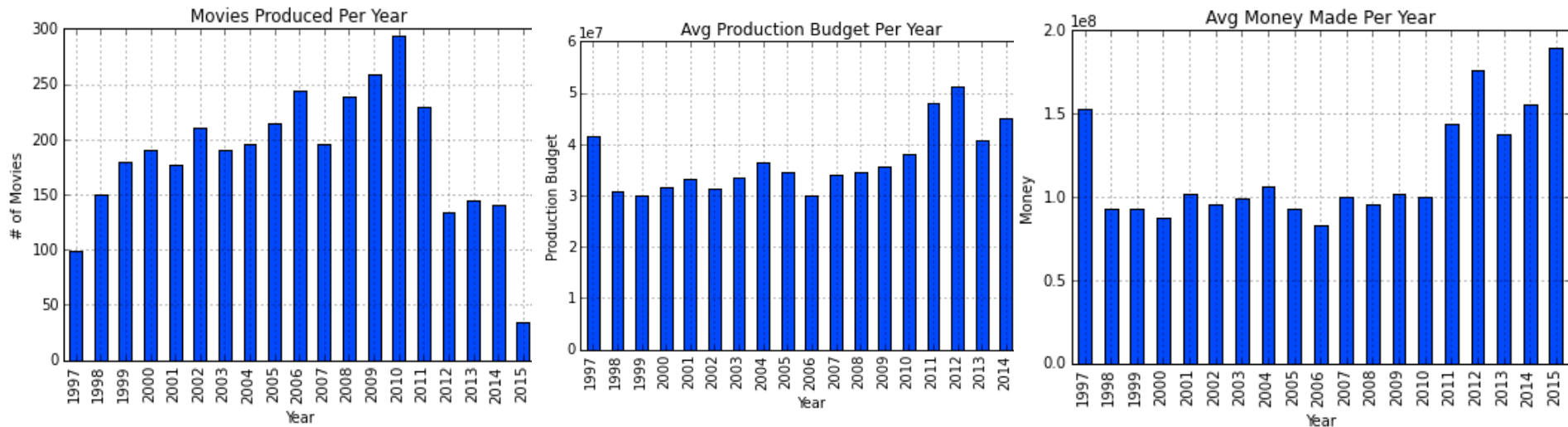
- Disney Studios
- \$270M Production Budget
- \$90M Domestic Box Office

# ***DATA SET DESCRIPTION & PREPROCESSING WORK***

- The data for this project was obtained from OpusData, a movie database company behind the popular site *The Numbers*. The dataset included ten csv files of which seven were used for this project, all linked through an Opus Data Identification Number (ODID).
  - CSV Files: movie\_summary, acting\_credits, technical\_credits, movie\_ratings, production\_companies, releases, and keywords.
- Preprocessing took the majority of time spent on this project and a huge thanks to Ramesh, Patrick, Liam, and Sinan for all your help! It involved:
  - Creating new columns to be used as features and for analysis
  - Sorting the csv files and then merging them into one large final file with **3525 ROWS AND 901 COLUMNS!**
  - Example: The actor Leonardo Dicaprio represents one actor column out of 250, and a “1” is placed in each row of that column for a movie he is in, and a “0” for each movie he is not in.

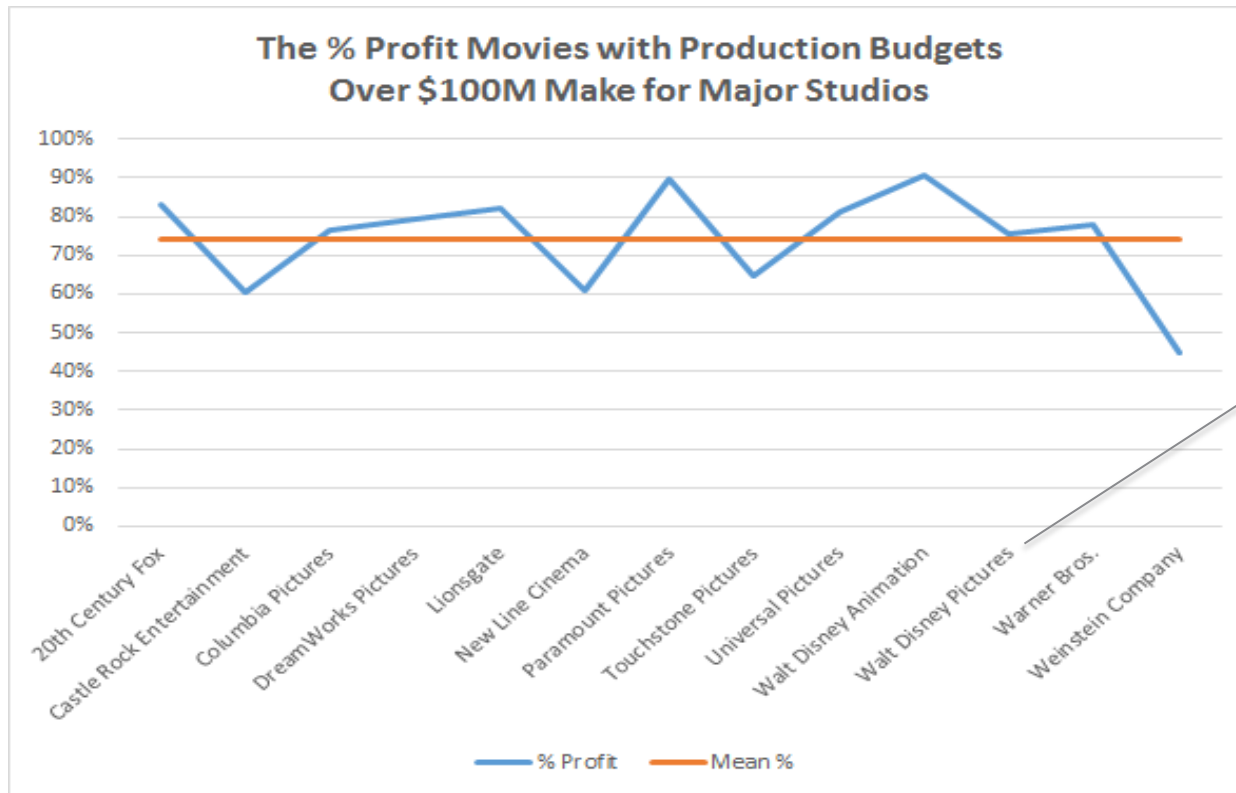
# DATA EXPLORATION & ANALYSIS

- **SEQUELS:** While sequels represented 9% of the analyzed movies, they on average made \$240 million dollars more than original films.
- Dramas and comedies are the most prevalent genres, contemporary fiction the most prevalent creative type, and “R” being the most frequent rating.
- The data corroborates the “blockbuster” theory towards movie production – In order to make money, studios have to be willing to spend it!



# DATA EXPLORATION & ANALYSIS (CONT)

Analysis of movie production budgets and production companies revealed that for the top movie studios, on average 75% of their profits come from the movies they produced with over \$100 million production budgets.



This should make Disney feel better about Lone Ranger

# FEATURE SELECTION

**FEATURE\_COLS** = 859 feature columns

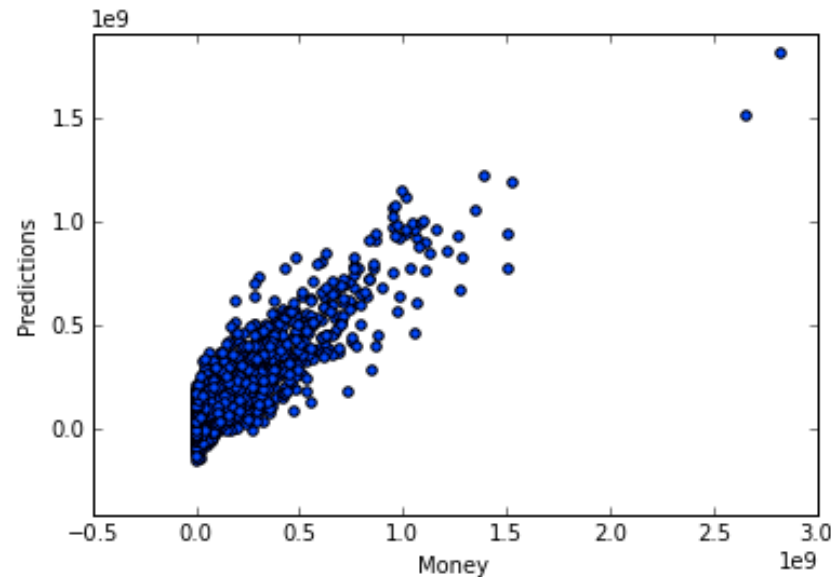
- 854 rows were from merged csv files relating to actors, directors, production companies, ratings, releases, and keywords
- 5 rows were from the original summary csv: production budget, runtimesquared, sequel, release month, and Christmas releases
  - Note: Removed features that contained metrics after a movie's release!

**RELEVANT\_FEATURES** = 119 feature columns

- Use STATS MODEL in python to trim the FEATURE\_COLS set by separating the features with p-values  $< .05$
- The p-value is the probability that the relationship we are observing is occurring purely by chance. If the p-value is less than .05, that means the confidence intervals for those coefficients do not include zero, and the null hypothesis can be rejected in favor of the alternative hypothesis - that there is a relationship between the feature and response.

# REGRESSION MODELS & EVALUATION

- Model Evaluation:
  - Mainly R-squared values but also calculated intercept, Mean Average Error (MAE), and Mean Square Error (MSE).
  - Feature importance rankings
- The R-squared value represents the variance in the observed data that is explained by the model.
- **Scikit Learn:** Best model producing the highest R-squared value of 78%.





# REGRESSION MODELS & EVALUATION (CONT)

Model	X (Features)	Y (Response)	R-squared	# of Features	Intercept	MAE	MSE
Sklearn Linear Regression Model	Feature_cols	Money	78%	859	-70,644,818	57,322,148	88,228,435
Stats Model Linear Regression	Feature_cols	Money	38%	859	NA	NA	NA
Sklearn Linear Regression Model	Relevant_Features	Money	72%	119	-18,925,571	61,393,258	99,564,775
Decision Tree Linear Regression	Feature_cols	Money	46%	859	NA	71,077,307	134,492,178
Random Forest	Feature_cols	Money	69%	859	NA	53,666,260	102,629,766

- **Stats Model:** Used mainly to create “Relevant\_Features” based on p-values
- **Decision Tree Model:** Split the data, trained on 70% of it and evaluated feature importance with a series of decisions, and then tested it on the remaining 30% it hadn’t seen.
- **Random forest model:** Ran the decision tree training test on the 70% of the data 100 different times, and averaged out the different evaluations to make predictions on the remaining 30%.



# FEATURE IMPORTANCE

## Decision Tree Features

**production\_budget 0.612**  
**runtimesquared 0.076**  
keywords3\_Cross-Class Romance 0.037  
**sequel 0.034**  
production\_method\_Live Action 0.016  
release\_month 0.015  
person\_Robert Downey, Jr. 0.013  
keywords1\_Boarding School 0.010  
keywords2\_Romance 0.008  
release\_pattern\_Special Engagement 0.006  
person\_Antonio Banderas 0.006  
keywords2\_Family Movie 0.005  
**keywords2\_Talking Animals 0.004**  
name\_Sam Raimi 0.004  
name\_Tim Burton 0.004  
name\_M. Night Shyamalan 0.004  
distributor\_Universal 0.004  
production\_method\_Digital Animation 0.003  
source\_Based on Fiction Book/Short Story 0.003  
distributor\_Warner Bros. 0.003  
creative\_type\_Science Fiction 0.003  
keywords3\_IMAX: DMR 0.003  
person\_Tom Hanks 0.003  
name\_Andrew Adamson 0.003  
production\_company1\_Columbia Pictures 0.002  
genre\_Adventure 0.002  
production\_company1\_Touchstone Pictures 0.002  
person\_Dustin Hoffman 0.002  
production\_company1\_Universal Pictures 0.002  
production\_company1\_Paramount Pictures 0.002

## Random Forest Features

**production\_budget 0.574**  
**runtimesquared 0.077**  
**sequel 0.026**  
keywords3\_Cross-Class Romance 0.022  
release\_pattern\_IMAX 0.016  
production\_method\_Live Action 0.013  
8release\_month 0.012  
distributor\_Paramount Pictures 0.008  
name\_Andrew Adamson 0.007  
production\_method\_Digital Animation 0.006  
release\_pattern\_Wide 0.006  
genre\_Thriller/Suspense 0.004  
distributor\_Walt Disney 0.004  
name\_M. Night Shyamalan 0.004  
person\_Sigourney Weaver 0.004  
release\_pattern\_Special Engagement 0.003  
keywords2\_Disaster 0.003  
genre\_Adventure 0.003  
person\_Robert Downey, Jr. 0.003  
person\_Rupert Grint 0.003  
person\_Kristen Bell 0.003  
keywords1\_Boarding School 0.003  
distributor\_Warner Bros. 0.003  
person\_Daniel Radcliffe 0.003  
keywords2\_Romance 0.003  
distributor\_Sony Pictures 0.003  
person\_Tom Hanks 0.003  
person\_John Leguizamo 0.002  
person\_Michelle Rodriguez 0.002  
source\_Based on Fiction Book/Short Story 0.002

## Stats Model

'source\_Based on Ballet',  
'source\_Compilation',  
'genre\_Action',  
'genre\_Adventure',  
'genre\_Black Comedy',  
'genre\_Comedy',  
'genre\_Drama',  
'genre\_Horror',  
'genre\_Musical',  
'genre\_Romantic Comedy',  
'genre\_Thriller/Suspense',  
'person\_Amanda Seyfried',  
'person\_Anne Hathaway',  
'person\_Anthony Hopkins',  
'person\_Arnold Schwarzenegger',  
'person\_Ben Stiller',  
'person\_Bradley Cooper',  
'person\_Cameron Diaz',  
**'person\_Christina Ricci'**,  
'person\_Colin Farrell',  
'person\_Colin Firth',  
'person\_Daniel Radcliffe',  
'person\_Elizabeth Banks',  
'person\_Emma Watson',  
'person\_Freddie Highmore',  
'person\_Freida Pinto',  
'person\_Hugh Jackman',  
'person\_Jennifer Connelly',  
'person\_Jennifer Lawrence',  
'person\_Jeremy Renner',

# FEATURE IMPORTANCE (STATS MODEL CONT)

## Stats Model

'person\_John Leguizamo',  
'person\_Julia Roberts',  
'person\_Kate Winslet',  
'person\_Keanu Reeves',  
'person\_Keira Knightley',  
'person\_Kevin Spacey',  
'person\_Kirsten Dunst',  
**'person\_Leonardo DiCaprio'**,  
'person\_Michelle Rodriguez',  
'person\_Morgan Freeman',  
'person\_Natalie Portman',  
'person\_Orlando Bloom',  
'person\_Rachel McAdams',  
'person\_Rachel Weisz',  
'person\_Ralph Fiennes',  
'person\_Robert Pattinson',  
'person\_Rupert Grint',  
'person\_Sam Rockwell',  
'person\_Sam Worthington',  
'person\_Scarlett Johansson',  
'person\_Steve Carell',  
'person\_Tom Hanks',  
'person\_Vin Diesel',  
'person\_Will Smith',  
'person\_Woody Harrelson',  
'name\_Andrew Adamson',  
'name\_Ang Lee',  
'name\_Bill Condon',  
'name\_Bryan Singer',  
'name\_Chris Columbus',

## Stats Model

'name\_Chris Weitz',  
**'name\_Christopher Nolan'**,  
'name\_Jay Roach',  
'name\_Joseph McGinty Nichol',  
'name\_Louis Leterrier',  
'name\_Luc Besson',  
'name\_M. Night Shyamalan',  
'name\_Michael Bay',  
'name\_Raja Gosnell',  
'name\_Rob Cohen',  
'name\_Robert Zemeckis',  
'name\_Roland Emmerich',  
'name\_Sam Raimi',  
'name\_Steven Spielberg',  
'creative\_type\_Kids Fiction',  
'production\_method\_Hand Animation',  
'production\_method\_Live Action',  
'production\_method\_Rotoscoping',  
'production\_method\_Stop-Motion Animation',  
'production\_company1\_Gold Circle Films',  
'production\_company1\_New Line Cinema',  
'production\_company1\_Relativity Media',  
'production\_company1\_Universal Pictures',  
'production\_company1\_Walt Disney Animation',  
'production\_company1\_Warner Bros.',  
'production\_company2\_Participant Media',  
'production\_company2\_Regency Enterprises',  
'distributor\_Dreamworks SKG',  
'distributor\_Warner Bros.',  
'release\_pattern\_Expands Wide',

## Stats Model

'release\_pattern\_IMAX',  
'release\_pattern\_Limited',  
'release\_pattern\_Special Engagement',  
'release\_pattern\_Wide',  
'keywords1\_Animals Gone Bad',  
'keywords1\_Family Movie',  
'keywords1\_Marvel Comics',  
'keywords1\_Relationship Advice',  
'keywords1\_Secret Agent',  
'keywords1\_Terminal Illness',  
'keywords1\_Visual Effects',  
'keywords2\_Animal Lead',  
'keywords2\_Artists',  
'keywords2\_Disaster',  
'keywords2\_Immigration',  
'keywords2\_Kidnap',  
'keywords2\_Road Trip',  
'keywords2\_Secret Agent',  
'keywords3\_Coming of Age',  
'keywords3\_Cross-Class Romance',  
'keywords3\_End of the World',  
'keywords3\_Faulty Memory',  
'keywords3\_Heist',  
'keywords3\_IMAX: DMR',  
'keywords3\_Mistaken Identity',  
'keywords3\_Time Travel',  
'keywords3\_Visual Effects',  
'sequel',  
'production\_budget']

# BUSINESS APPLICATIONS & CONCLUSION

- It is possible to create a linear model to help predict a movie's financial performance!
- The models in this project are a great start to creating one that could be used by movie executives to help with “greenlight” decision making.
  - “production budget” of \$80 million + “runtimesquared” of 10404 minutes + “release\_pattern\_IMAX” = \$260 million movie!
- The challenge to this and other linear models is providing the very best features that will boost the R-SQUARED values (% of error covered by the model).

