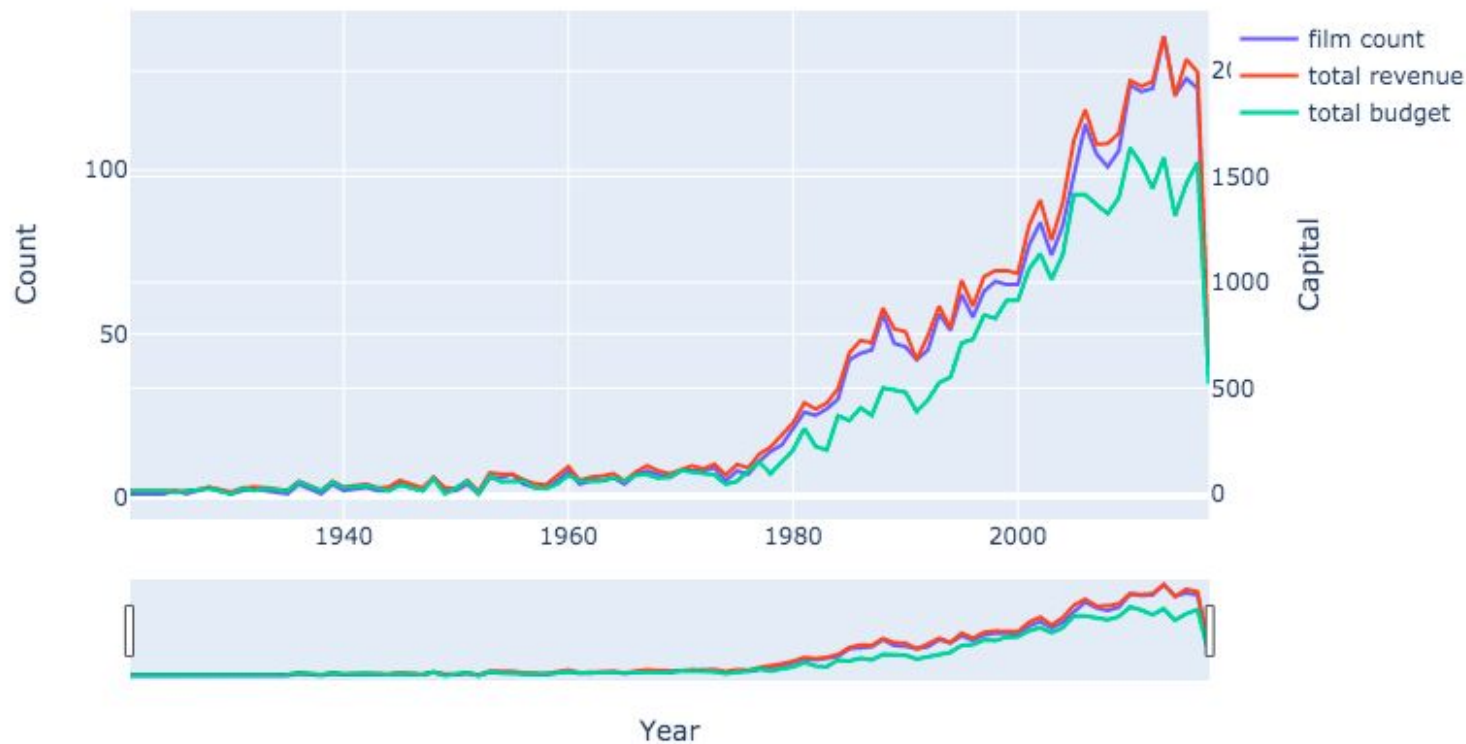# TMDB Box Office Collection: Prediction

Predicting the box office revenue of over 4000 movies based on the data and features we have for 3000 movies.
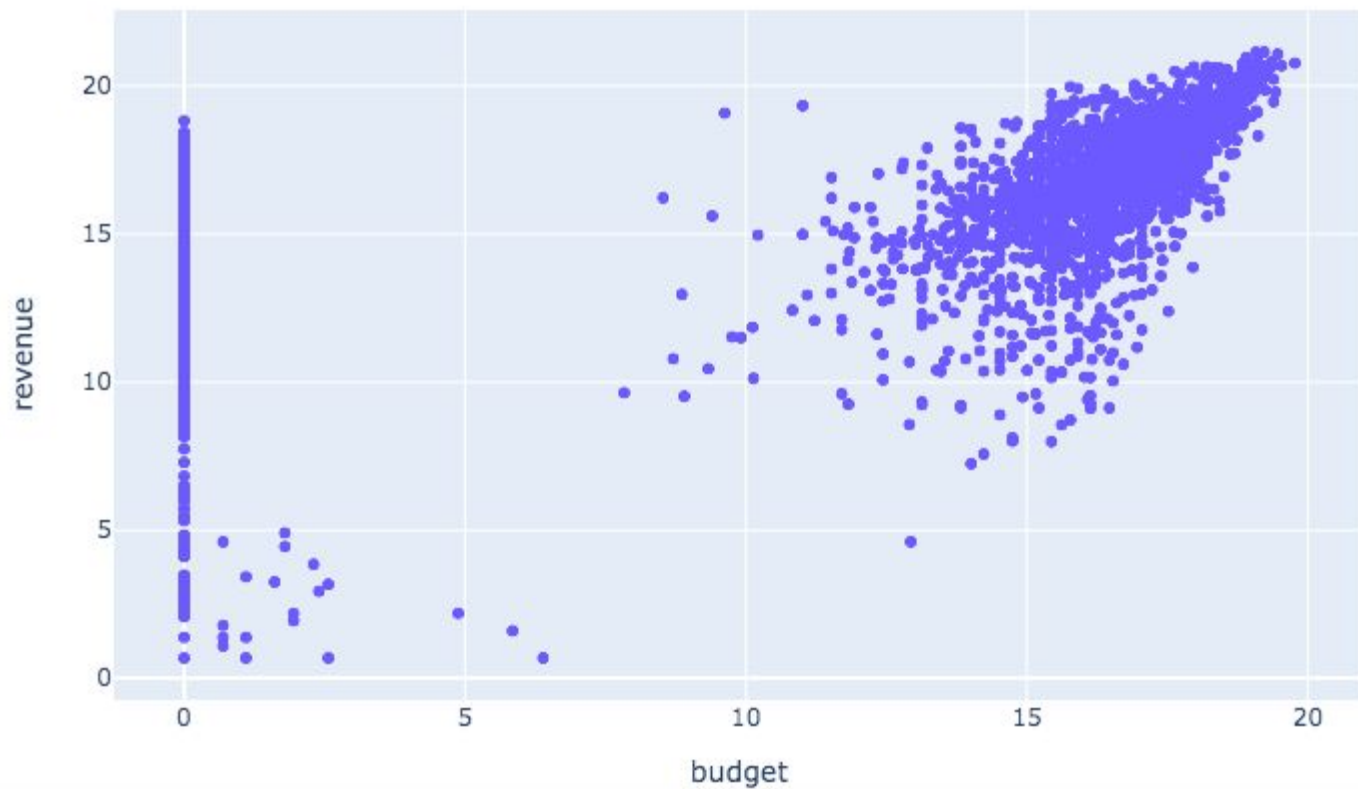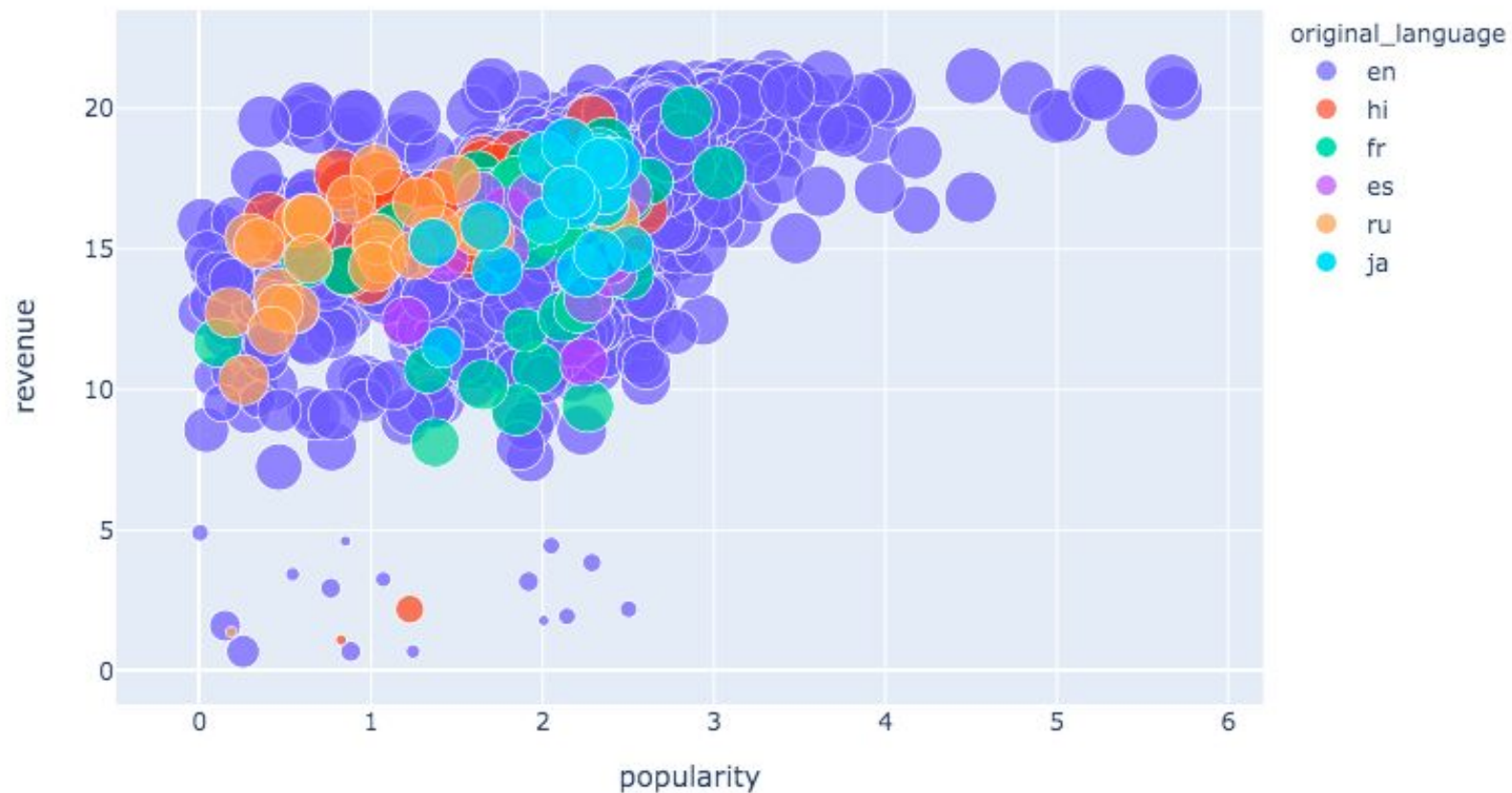
Adamya Nayyar

# Data Exploration
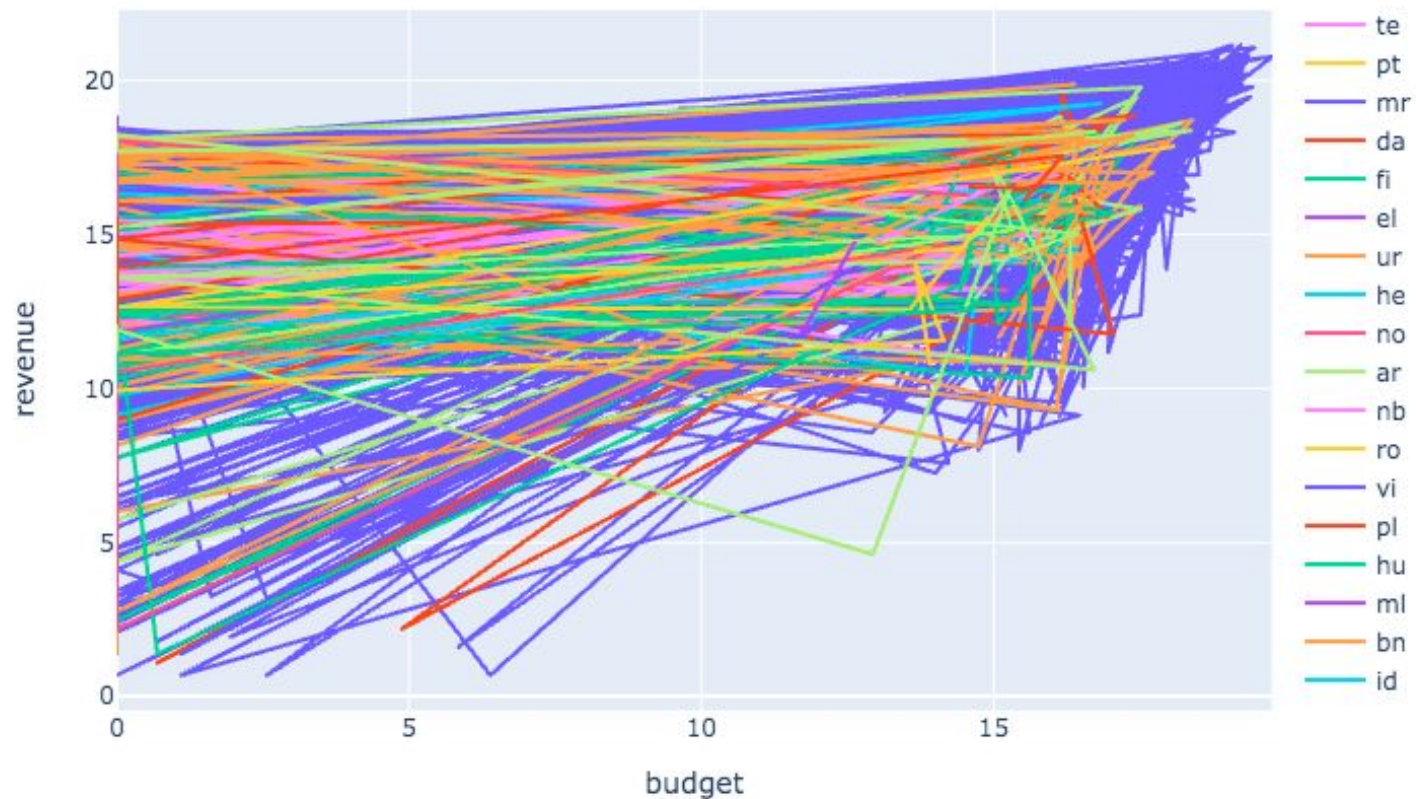


Number of films and total revenue per year
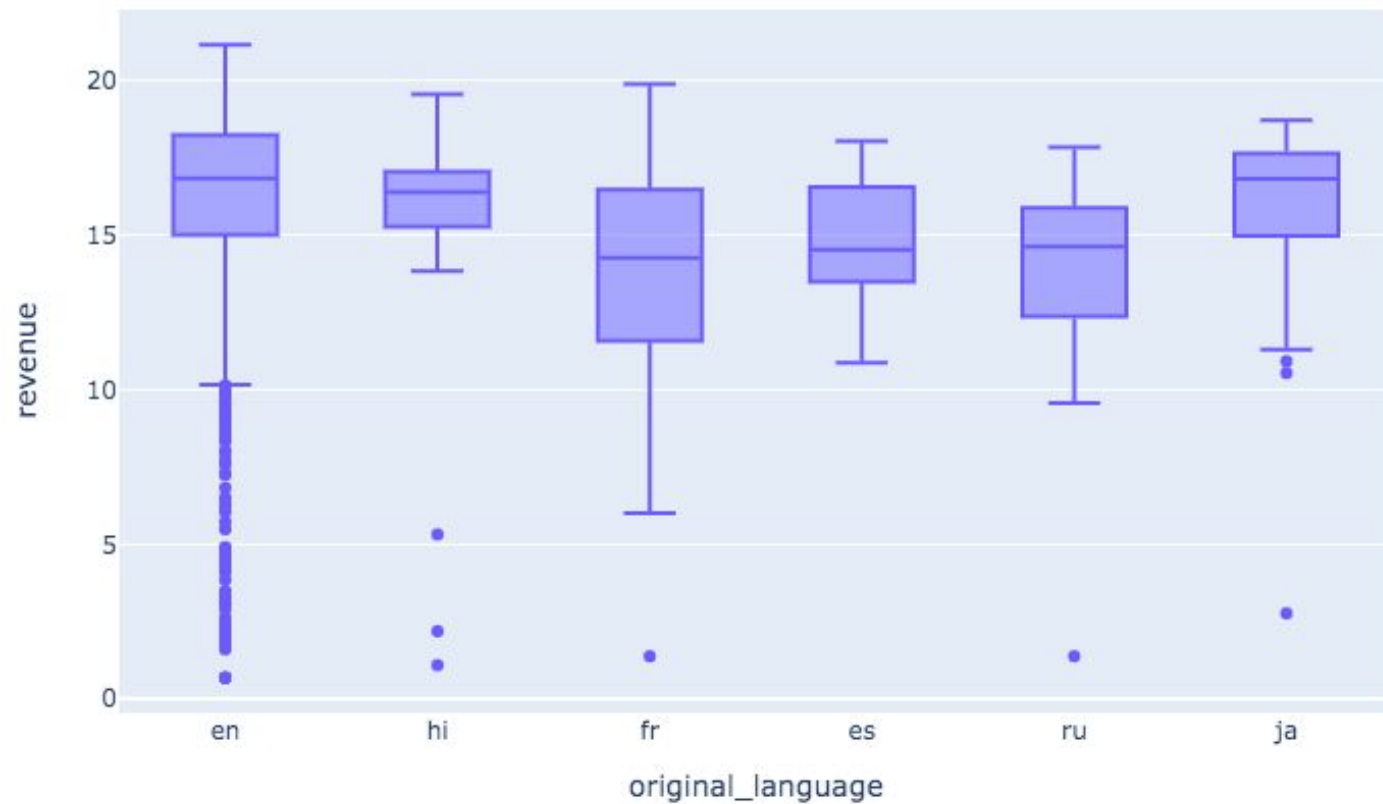
# Log Budget vs Log Revenue

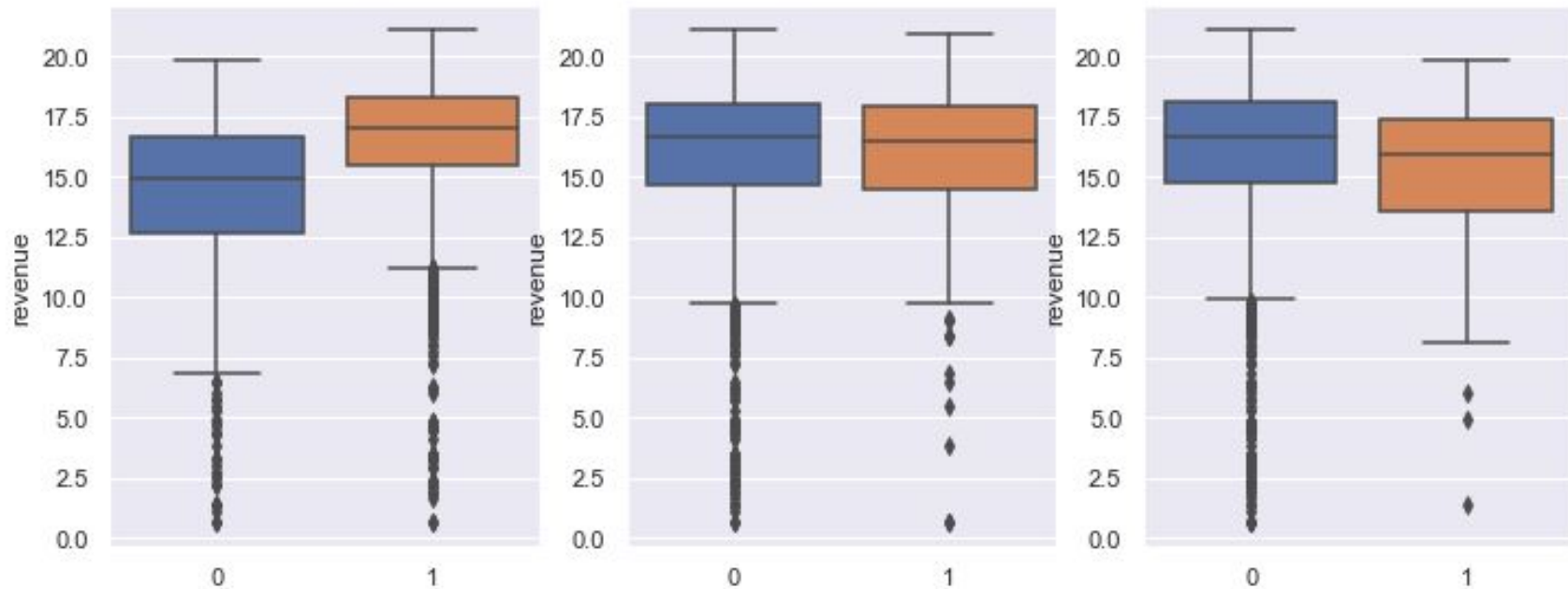Log Revenue vs Log Popularity (Buble size=Budget)

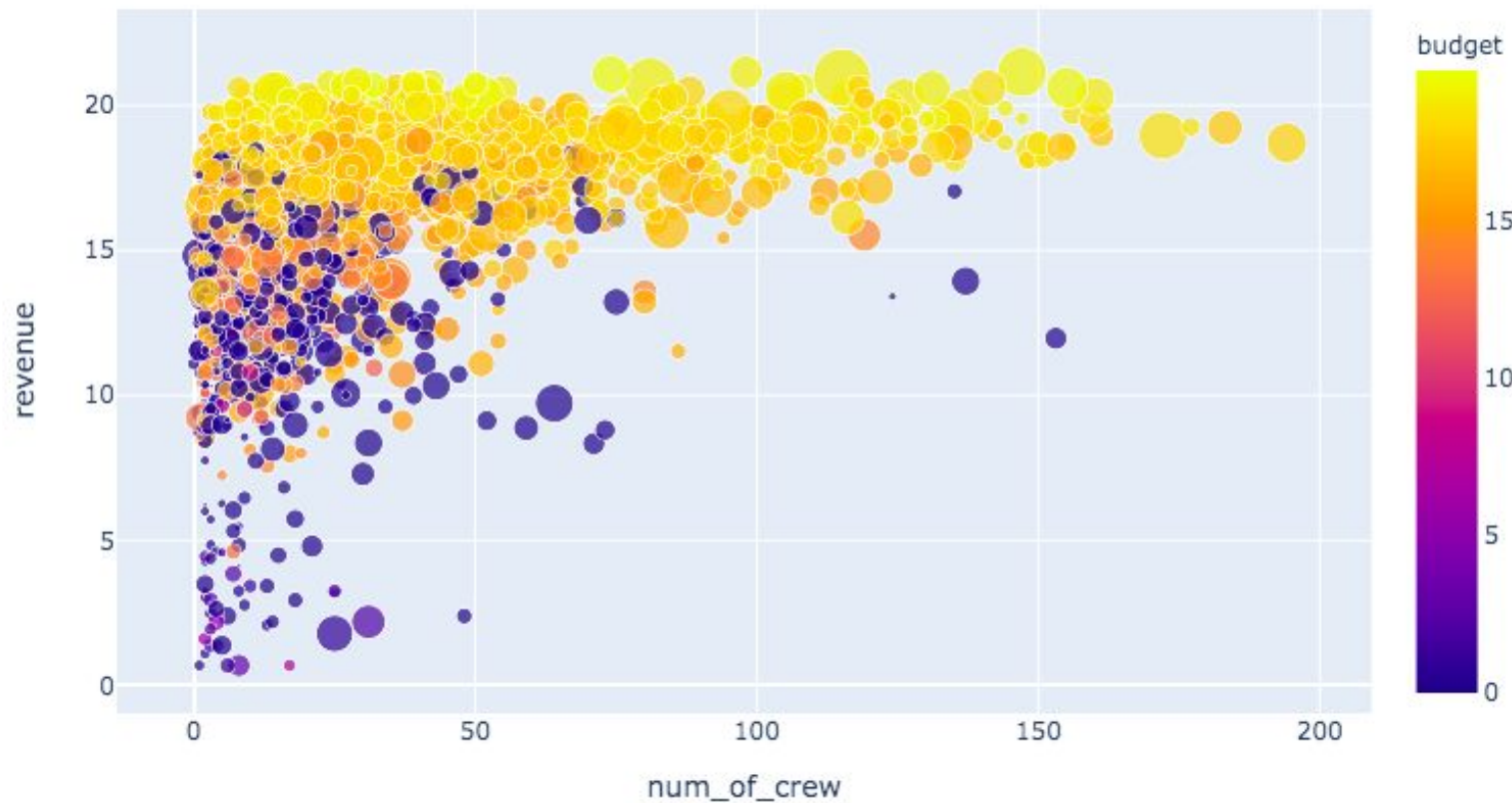Log Budget vs Log Revenue in different languages

Log Revenue Distribution for top languages

Log revenue vs Top Production Countries

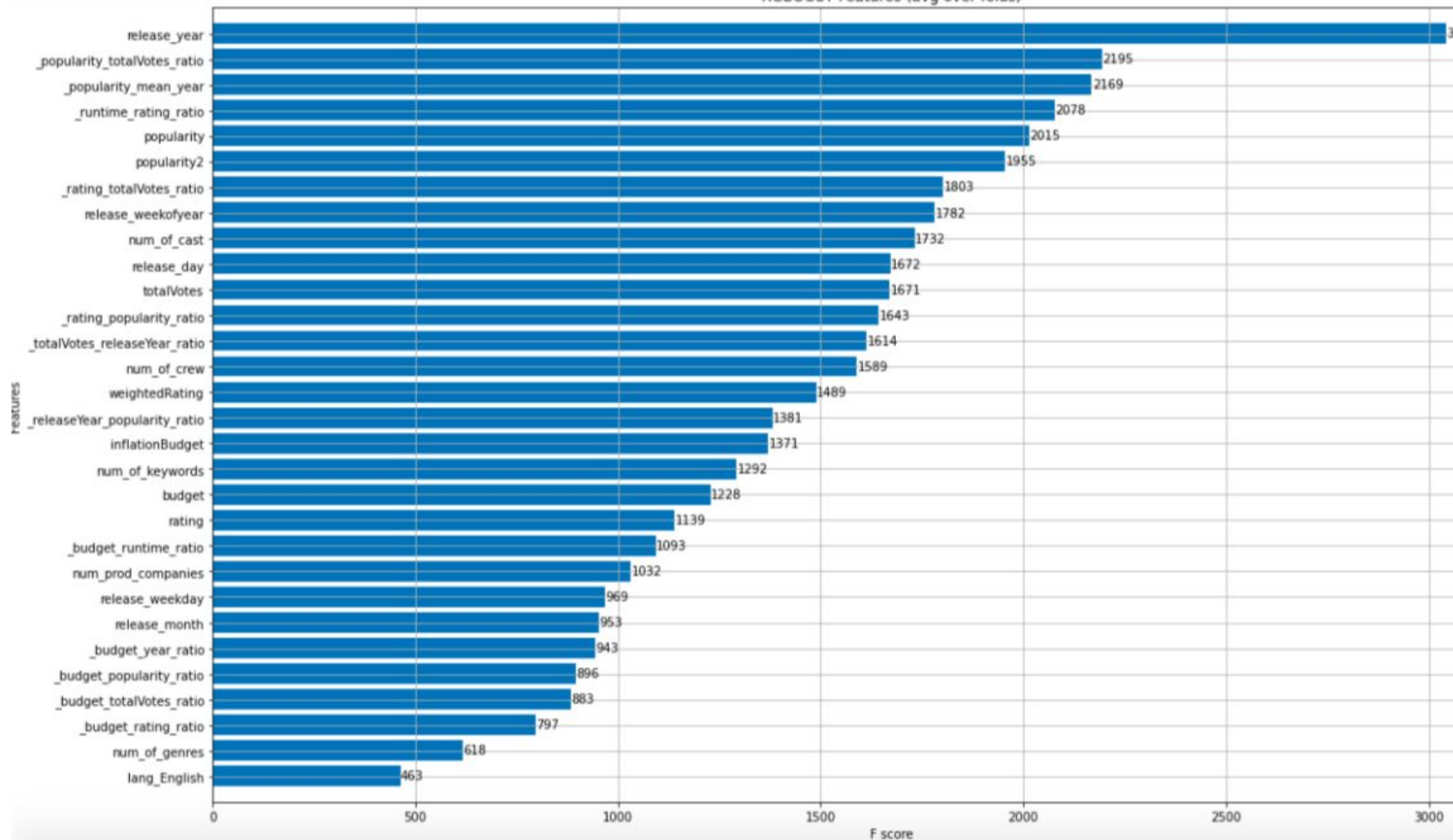Crew vs Log Revenue(Bubble size= Number of cast, color= Budget)

# Model Selection

- After fitting the data on various regression models, random forest gave us good result whereas Xgboost and Catboost model gave us the least errors.
- Random Forest model had rmse of 1.91 on validation set. Xgboost had even less error at rmse 1.83 whereas Catboost with the least rmse 1.81 .
- Since the root mean squared error of Xgboost and Catboost are very close, we can select either one as our final model. But for this project I calculated the final prediction as:

**Final prediction = 0.3*Xgboost + 0.7*Catboost**

- The above ensemble model gave the best prediction score.
- Next page have important features by Xgboost model.

XGBOOST Features (avg over folds)

# THANK YOU! ANY QUERIES?