Movies Go Abroad

# Objective

How well can we predict international box office based off of movie pre-release features and domestic box office information?

What features have a stronger correlation to box office success?

# Data Collected



- 1,750 movies
- US wide released from 2000-2020
- Top grossing

# Model Selection and Workflow

| Cross Validation Tests | R² |
|---|---|
| Linear Regression | 0.535 |

↓

| Poly LASSO | 0.626 |
|---|---|

Runtime

Budget

Release Month

Metacritic Score

"Adventure" Genre
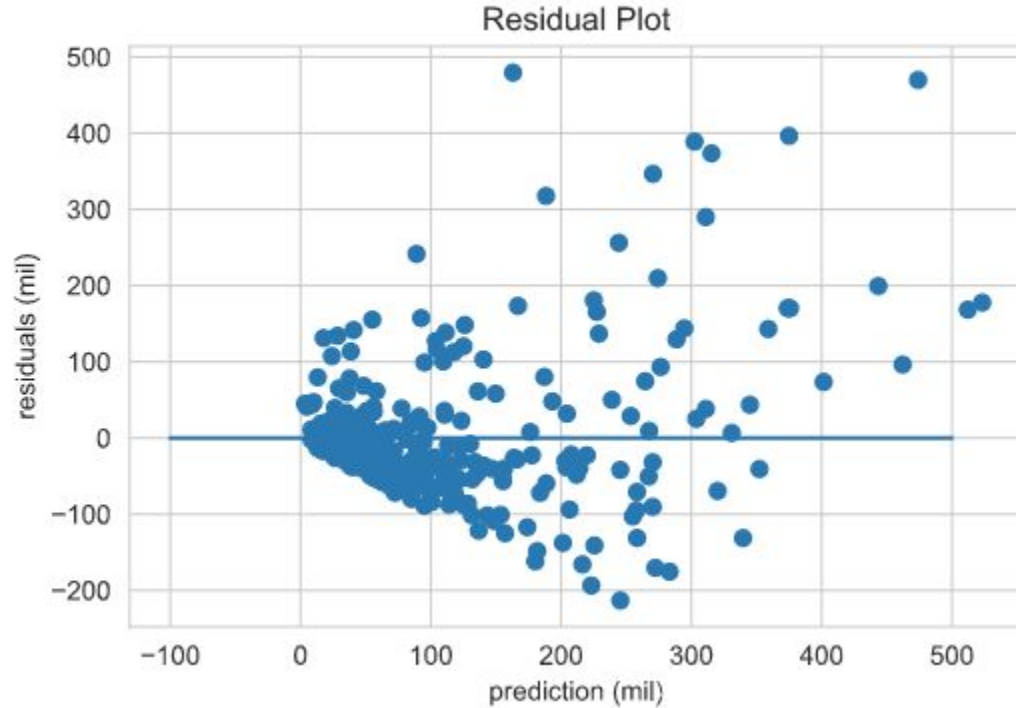
_____

Director

Dropping Outliers

# Movie Predictions

| Film Title | Actual (mil) | Predicted (mil) | Error (mil) |
|---|---|---|---|
| The 40 Year Old Virgin | $67.93 | $66.16 | $1.77 |
| Rat Race | $28.88 | $25.70 | $3.18 |
| Transformers | $691.28 | $302.40 | $388.88 |
| Avengers | $943.80 | $474.03 | $469.78 |

R² = 0.626
RMSE: 95.75

# Final Model Residual Plot



R² = 0.626
RMSE: 95.75

# Future Work

- Features:
  - Franchise
  - Popular int'l actors or directors
  - Popular int'l genres
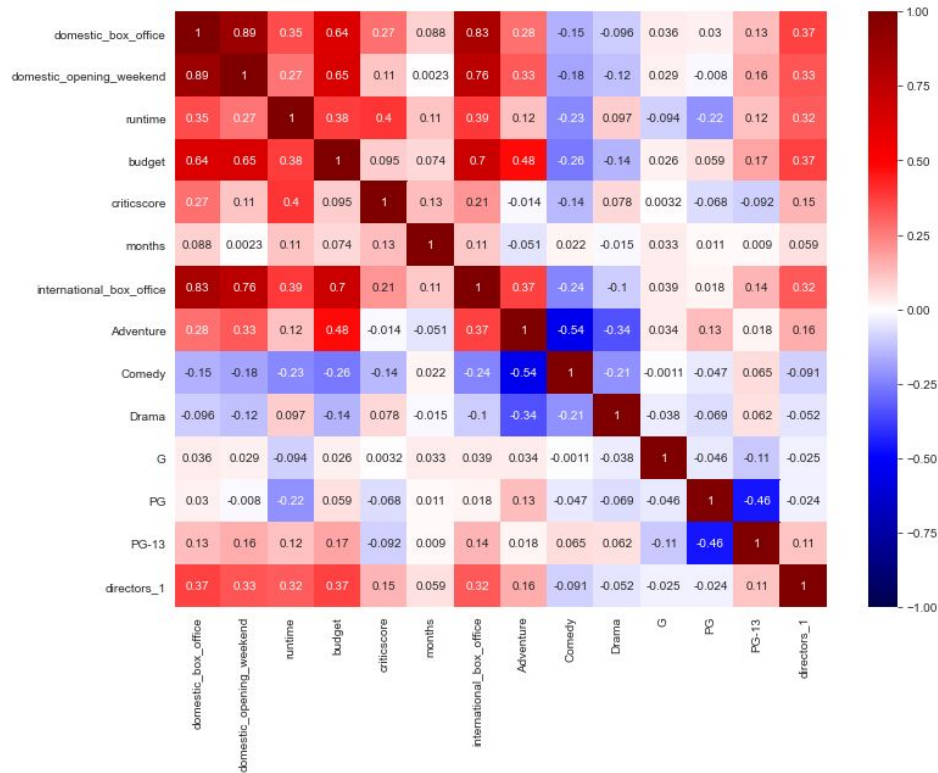- Log Regression
  - Too skewed?

# Appendix 1



OLS Regression Results

| Dep. Variable: | international_box_office | R-squared: | 0.750 |
|---|---|---|---|
| Model: | OLS | Adj. R-squared: | 0.748 |
| Method: | Least Squares | F-statistic: | 376.7 |
| Date: | Wed, 14 Apr 2021 | Prob (F-statistic): | 0.00 |
| Time: | 21:33:56 | Log-Likelihood: | -8824.8 |
| No. Observations: | 1519 | AIC: | 1.768e+04 |
| Df Residuals: | 1506 | BIC: | 1.774e+04 |
| Df Model: | 12 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -106.8134 | 16.594 | -6.437 | 0.000 | -139.362 | -74.264 |
| domestic_box_office | 1.1888 | 0.059 | 20.126 | 0.000 | 1.073 | 1.305 |
| domestic_opening_weekend | 0.1487 | 0.178 | 0.833 | 0.405 | -0.201 | 0.499 |
| runtime | 0.6602 | 0.146 | 4.527 | 0.000 | 0.374 | 0.946 |
| budget | 0.6313 | 0.054 | 11.614 | 0.000 | 0.525 | 0.738 |
| criticscore | -0.2225 | 0.149 | -1.496 | 0.135 | -0.514 | 0.069 |
| months | 1.6078 | 0.607 | 2.647 | 0.008 | 0.416 | 2.799 |
| Adventure | 17.1199 | 6.523 | 2.625 | 0.009 | 4.325 | 29.914 |
| Comedy | -14.2784 | 6.699 | -2.131 | 0.033 | -27.419 | -1.138 |
| Drama | -0.6663 | 7.942 | -0.084 | 0.933 | -16.244 | 14.912 |
| G | 18.3626 | 20.650 | 0.889 | 0.374 | -22.143 | 58.868 |
| PG | -4.3948 | 5.951 | -0.739 | 0.460 | -16.068 | 7.278 |
| directors_1 | -14.2485 | 5.070 | -2.810 | 0.005 | -24.193 | -4.304 |

| Omnibus: | 809.317 | Durbin-Watson: | 2.218 |
|---|---|---|---|
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 34756.383 |
| Skew: | 1.789 | Prob(JB): | 0.00 |
| Kurtosis: | 26.159 | Cond. No. | 1.89e+03 |

# Appendix 2

final model regression formula

| | feature | coef |
|---|---|---|
| 0 | intercept | 1.174032e+02 |
| 1 | 1 | 2.487151e-11 |
| 2 | runtime | -2.947956e+00 |
| 3 | budget | -1.415960e-01 |
| 4 | criticscore | 7.754058e-01 |
| 5 | months | -7.347438e+00 |
| 6 | Adventure | -1.698209e+01 |
| 7 | directors_1 | -3.572943e+01 |
| 8 | runtime^2 | 1.550250e-02 |
| 9 | runtime budget | 7.583133e-03 |
| 10 | runtime criticscore | -2.119288e-02 |
| 11 | runtime months | 3.271938e-02 |
| 12 | runtime Adventure | 6.925195e-02 |
| 13 | runtime directors_1 | 5.808105e-01 |
| 14 | budget^2 | -2.438303e-03 |
| 15 | budget criticscore | 1.482026e-02 |
| 16 | budget months | 6.016295e-03 |
| 17 | budget Adventure | 5.934067e-01 |
| 18 | budget directors_1 | -2.604903e-02 |
| 19 | criticscore^2 | 1.638269e-02 |
| 20 | criticscore months | 7.981884e-03 |
| 21 | criticscore Adventure | -7.546504e-02 |
| 22 | criticscore directors_1 | -3.495226e-01 |
| 23 | months^2 | 1.601039e-01 |
| 24 | months Adventure | 1.573396e+00 |
| 25 | months directors_1 | 4.269101e+00 |
| 26 | Adventure^2 | -1.698209e+01 |
| 27 | Adventure directors_1 | 1.148067e+01 |
| 28 | directors_1^2 | -3.572943e+01 |

# Appendix 3



## International Box Office (USD)

final model regplot