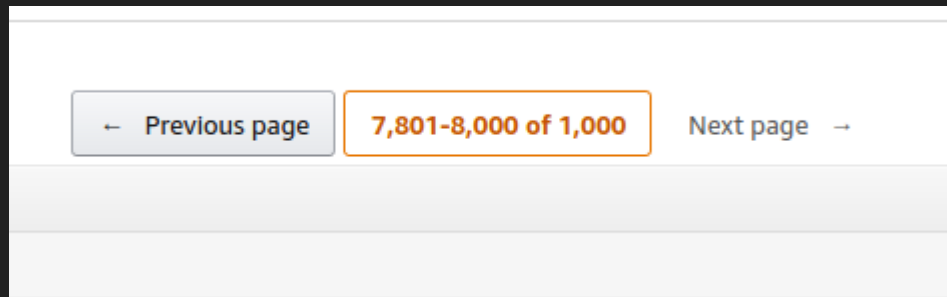# Webscraping

# Top 10,000 Movies from Box Office Mojo

**<u>Flow:</u>**

- Collect list of movies
- Request tabs of cast and crew
- Parse  information

**<u>Interesting Aspects:</u>**

- Going out of bounds
- Inconsistent/missing information



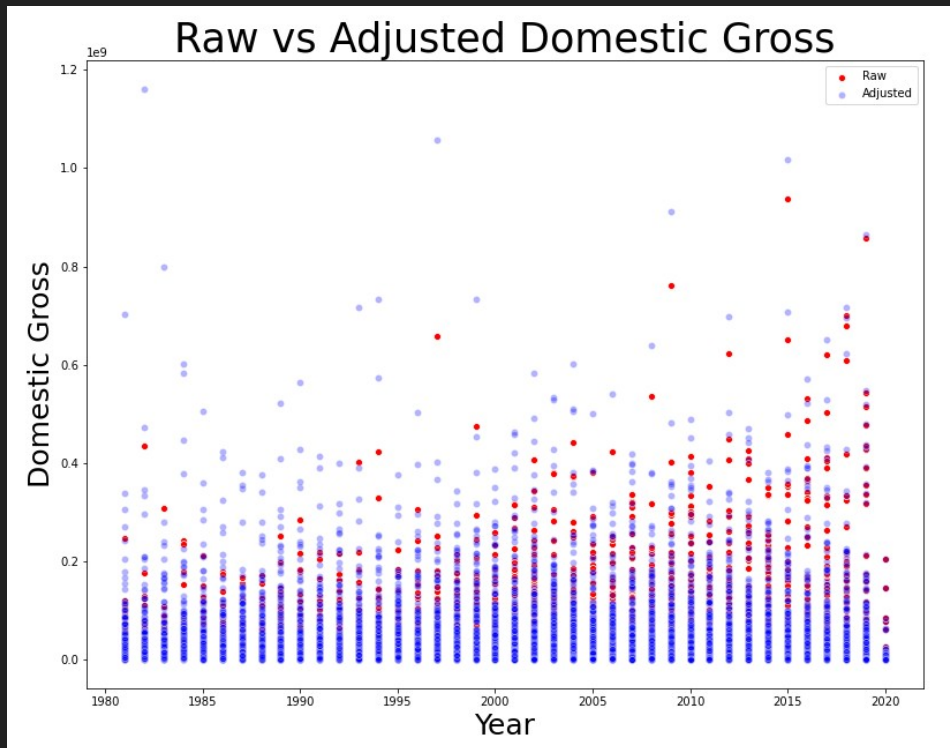← Previous page | **7,801-8,000 of 1,000** | Next page →

# Data Cleaning

# Correcting and Dropping Inconsistencies

## Issues:

- No dates
    - Drop rows
- No budgets
    - Drop budget column
- Wrong lengths
    - Check for min/hr
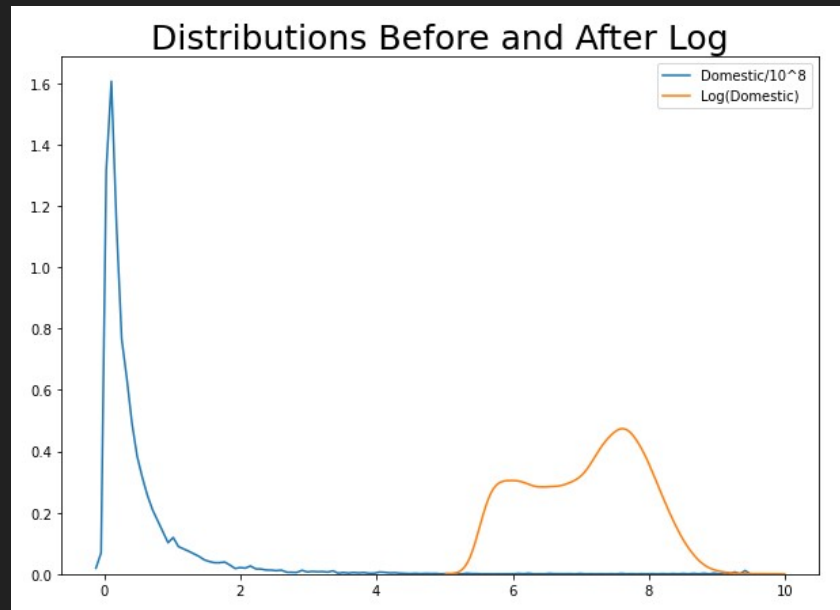- Inflation adjustment
    - CPI data from BLS

# Feature Engineering

# How to deal with so much categorical data?

- 1000s of actors, directors, etc
  - Sort by date
  - Find experience levels by date
  - Record maximum experience
- Multiple genres
  - Dummies
  - Form combinations
  - Disregard strange genres

# What about distributions?

- Target (Domestic Box Office Gross)~
  - log-normal
  - Same for experience levels (not-shown)
- Year and experience
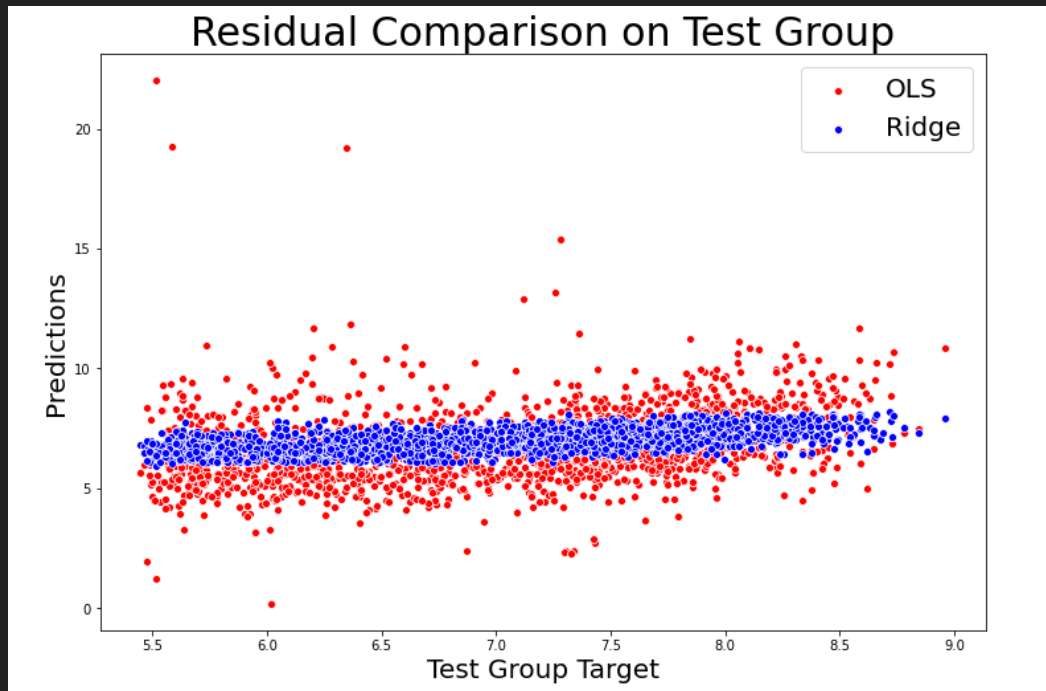  - Unintended interaction
  - Limit to 1980
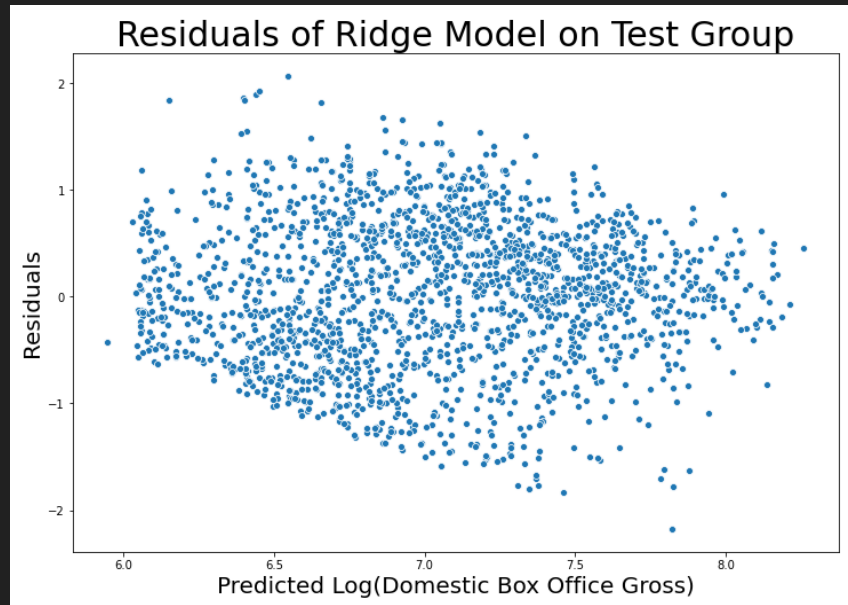


Distributions Before and After Log

# Modeling

# OLS Model

- Quick sense of fit
- Initial reduction of variables
- Ultimately over-fit
  - High R^2 on training
  - Low overall P
  - Higher P for features
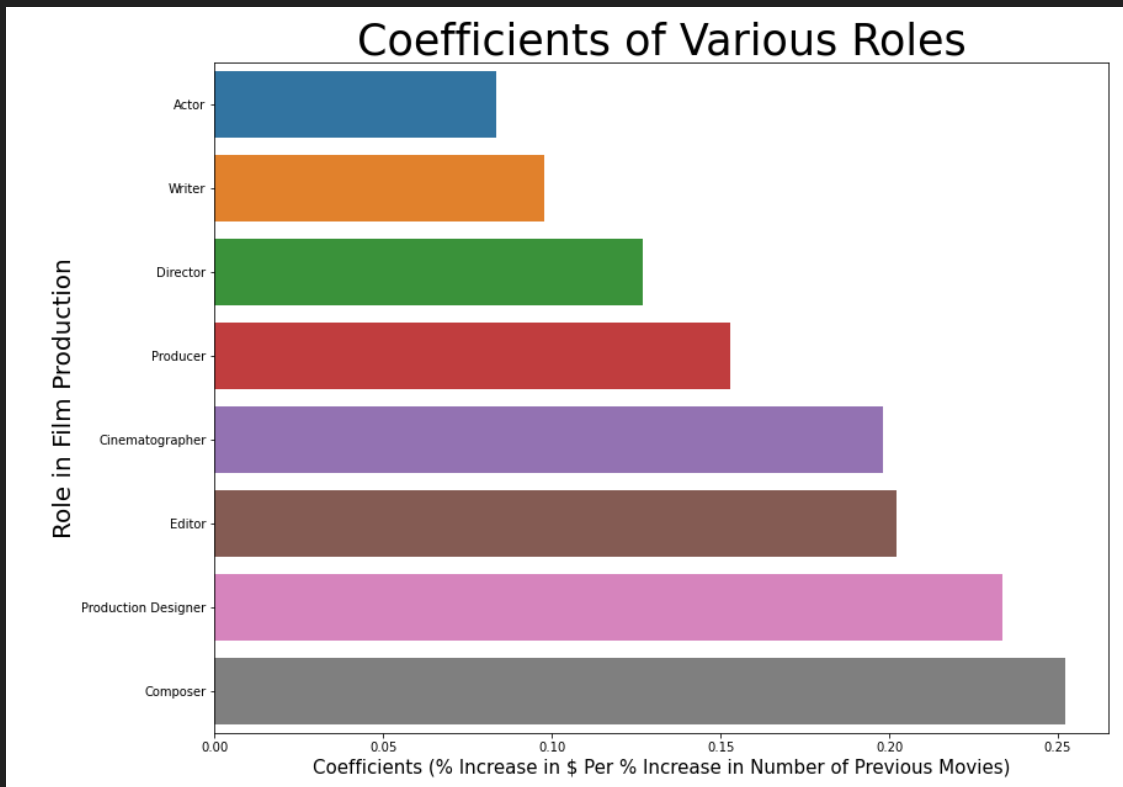  - Variance on Test Set

# Ridge Regression Model

- K folds for alpha
- Reduced correlated coefficients
- Better MAE on test set
  - 0.54 vs 1.04
- Allowed to keep more variables



Residuals of Ridge Model on Test Group

# Results

# Feature Coefficients of Interest

- Pair correlations < 0.5
- Three additional features
  - PG-13
  - Drama
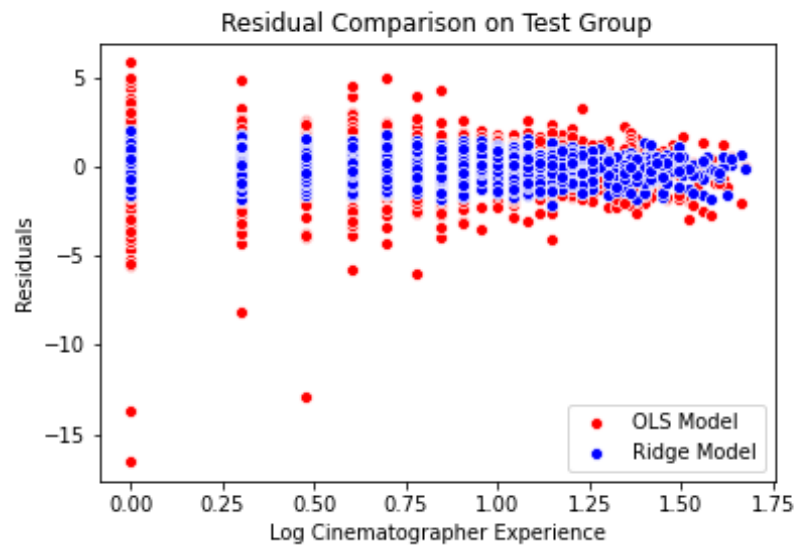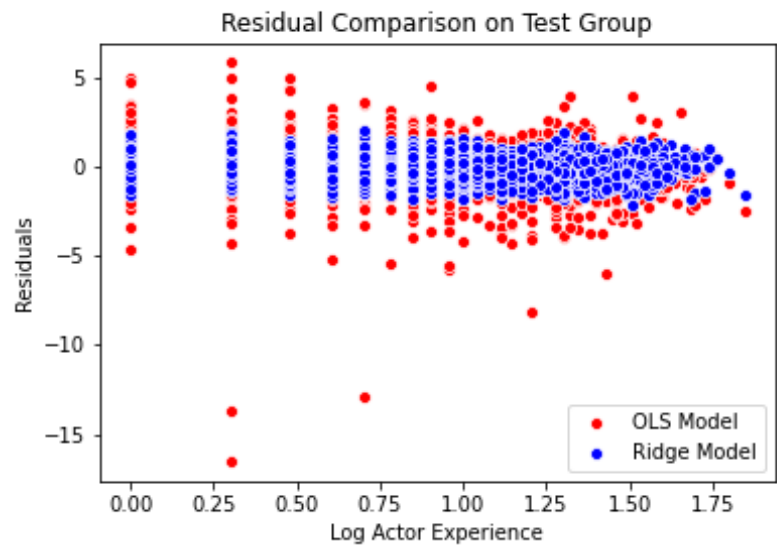- Composer experience top
- Actor experience bottom



Coefficients of Various Roles

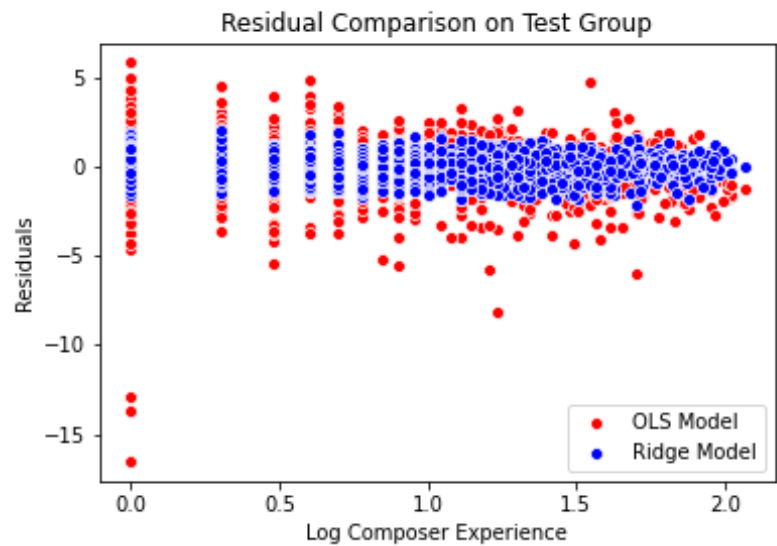# Potential Growth and Improvement

# Future Work

- Attempt to find strange interaction effects
  - Experience and genres
  - Combined experience
- Outside data
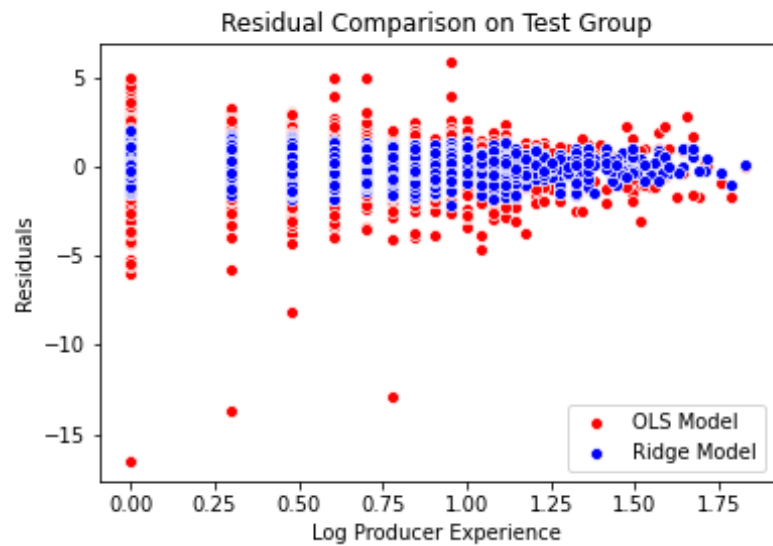  - Demographic
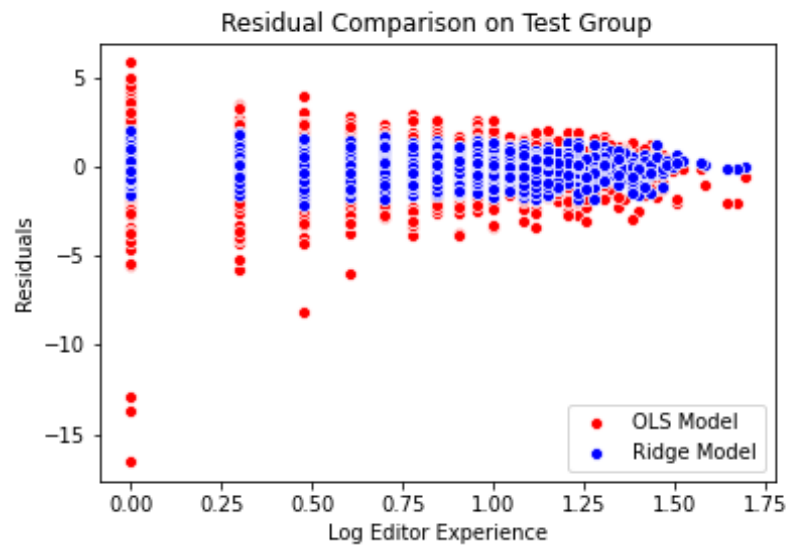  - Economic
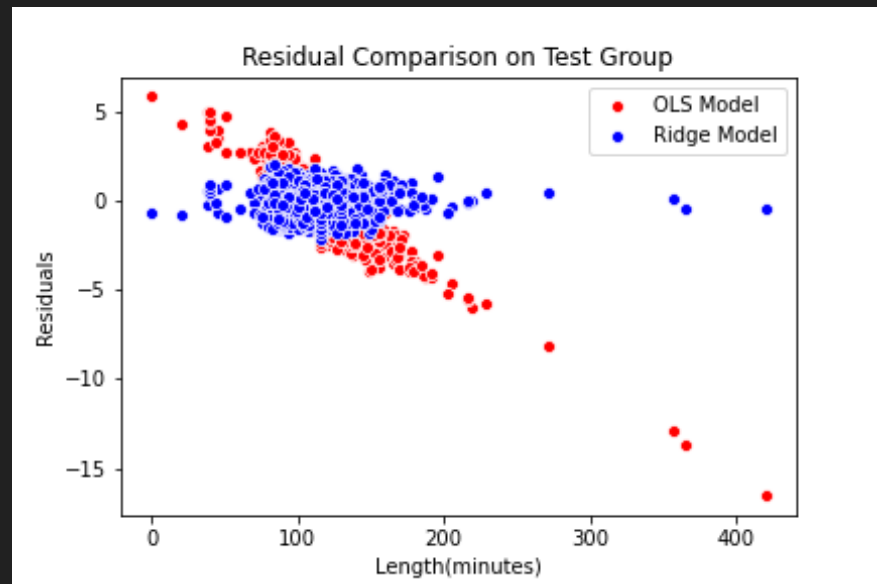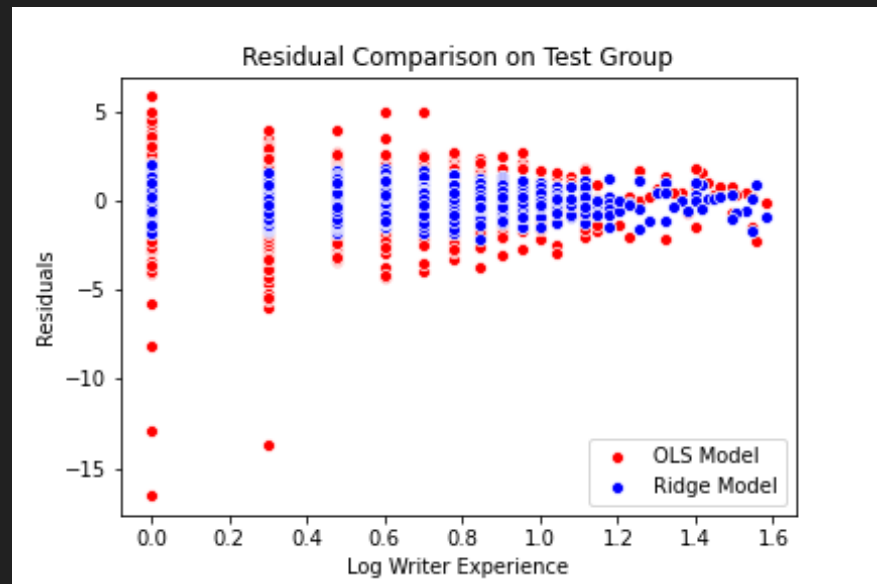- Domain research

# Appendix

KDE Plot of Degrees of Experience

# Acknowledgements