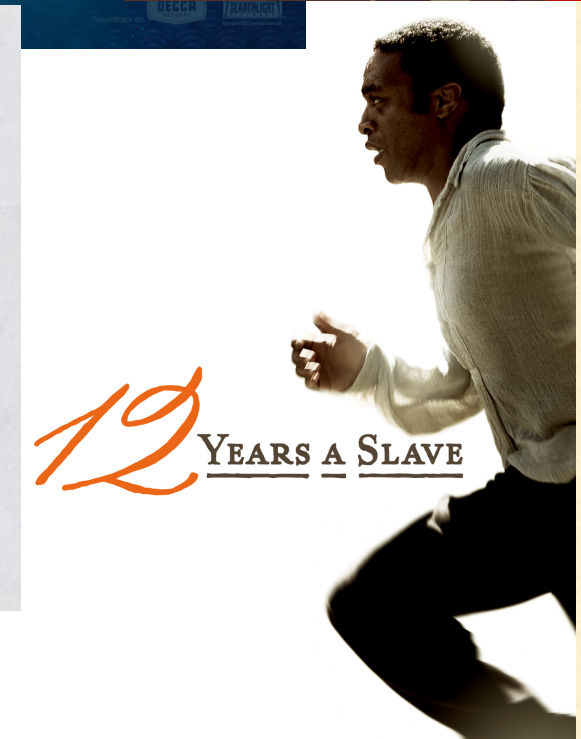
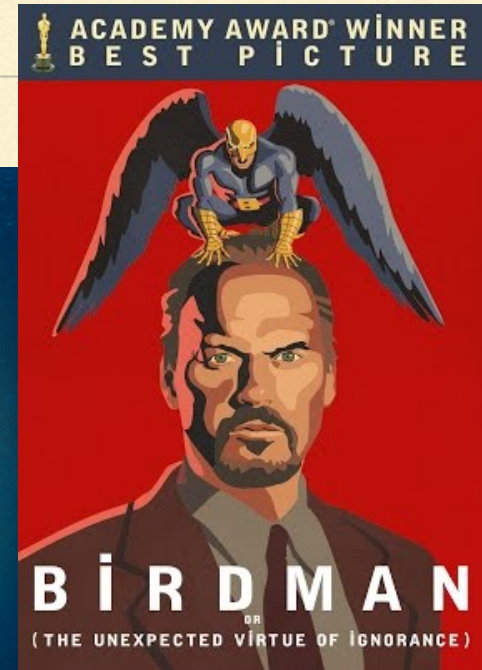

LINEAR REGRESSION ON INDEPENDENT FILM DATA



DATA SET

Box Office **Mojo**

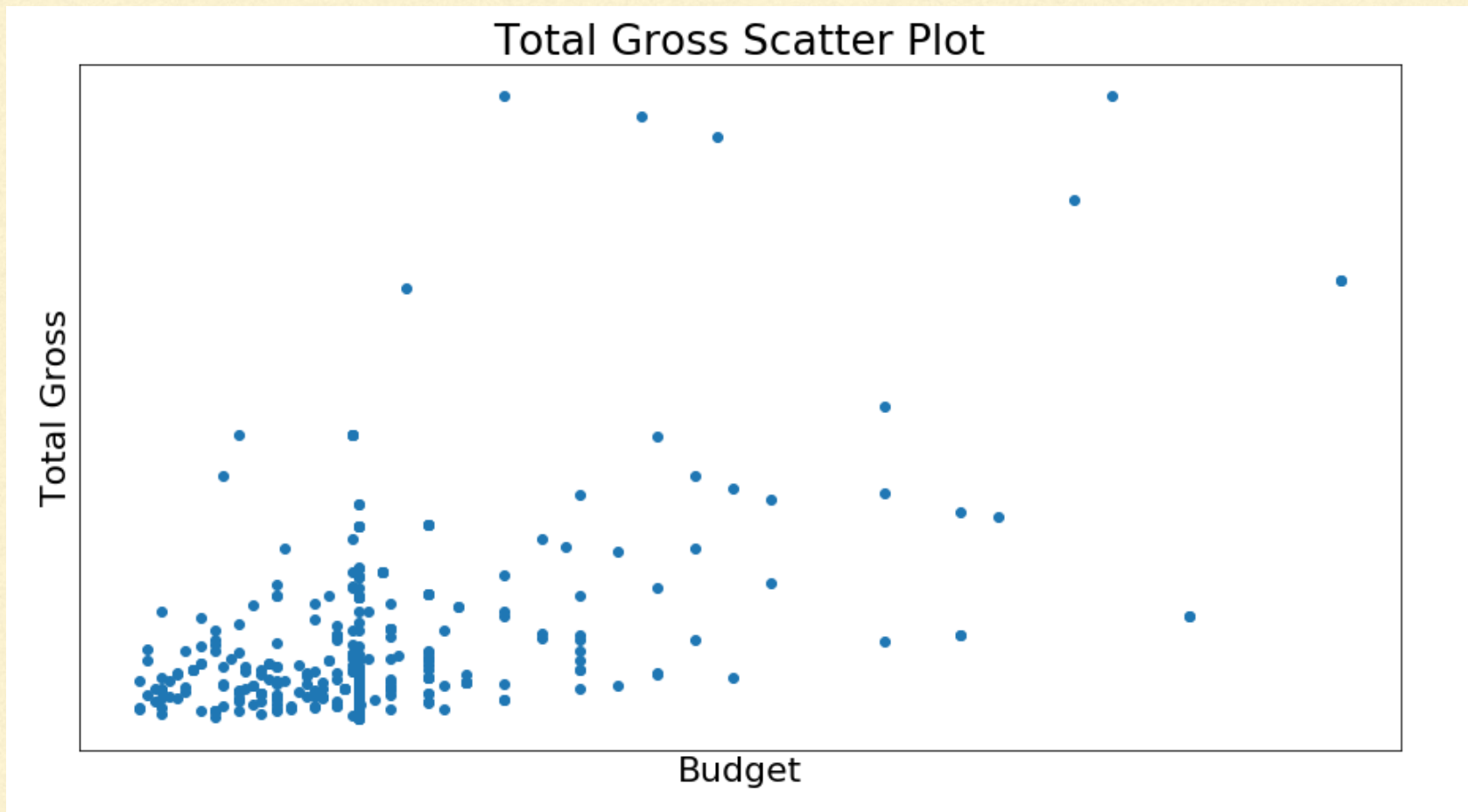


- Movies between 2008-2017 that are not “Big Six”
- Movie Revenue = theater ticket sales + disc sales
- Features:
 - Numerical - budget, runtime, days in release
 - Categorical - genre, MPAA rating, release date

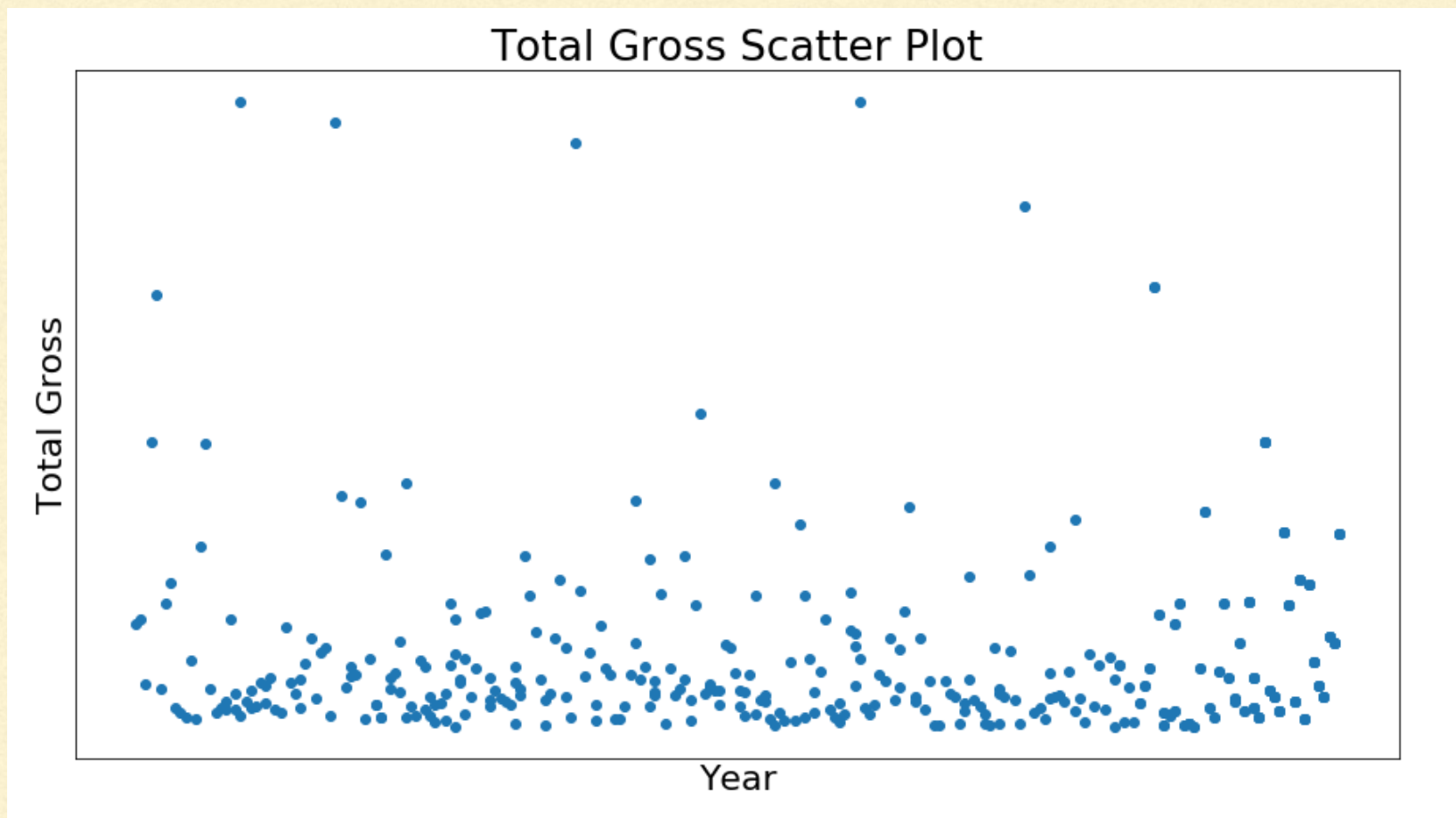
DATA CLEANING

- Adjust monetary values for inflation (2017)
 - Impute mean values for NaNs in numerical features
 - Merge DVD/Blu-ray sale data
-

EXPLORATORY DATA ANALYSIS



EXPLORATORY DATA ANALYSIS

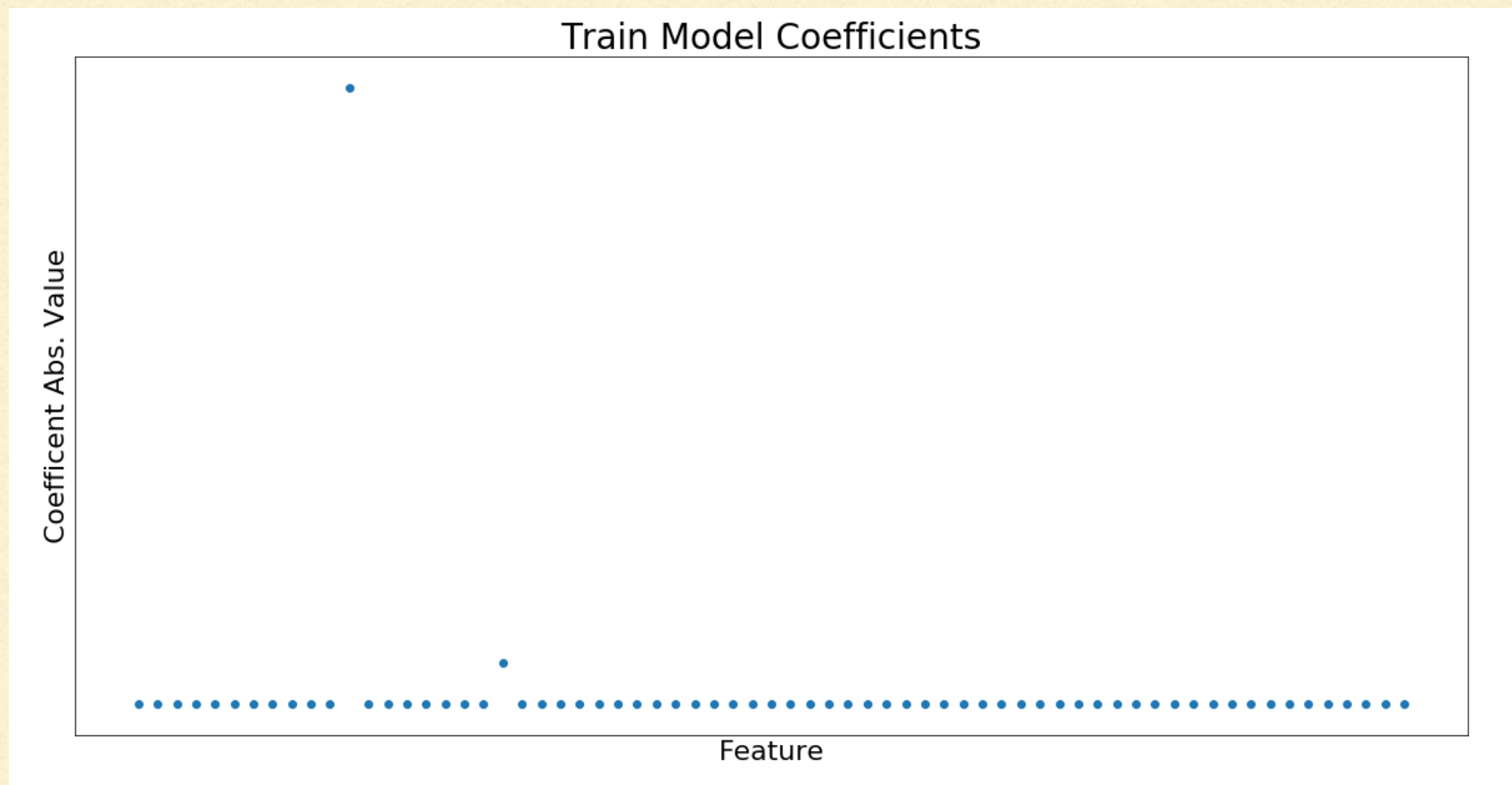


MODELING

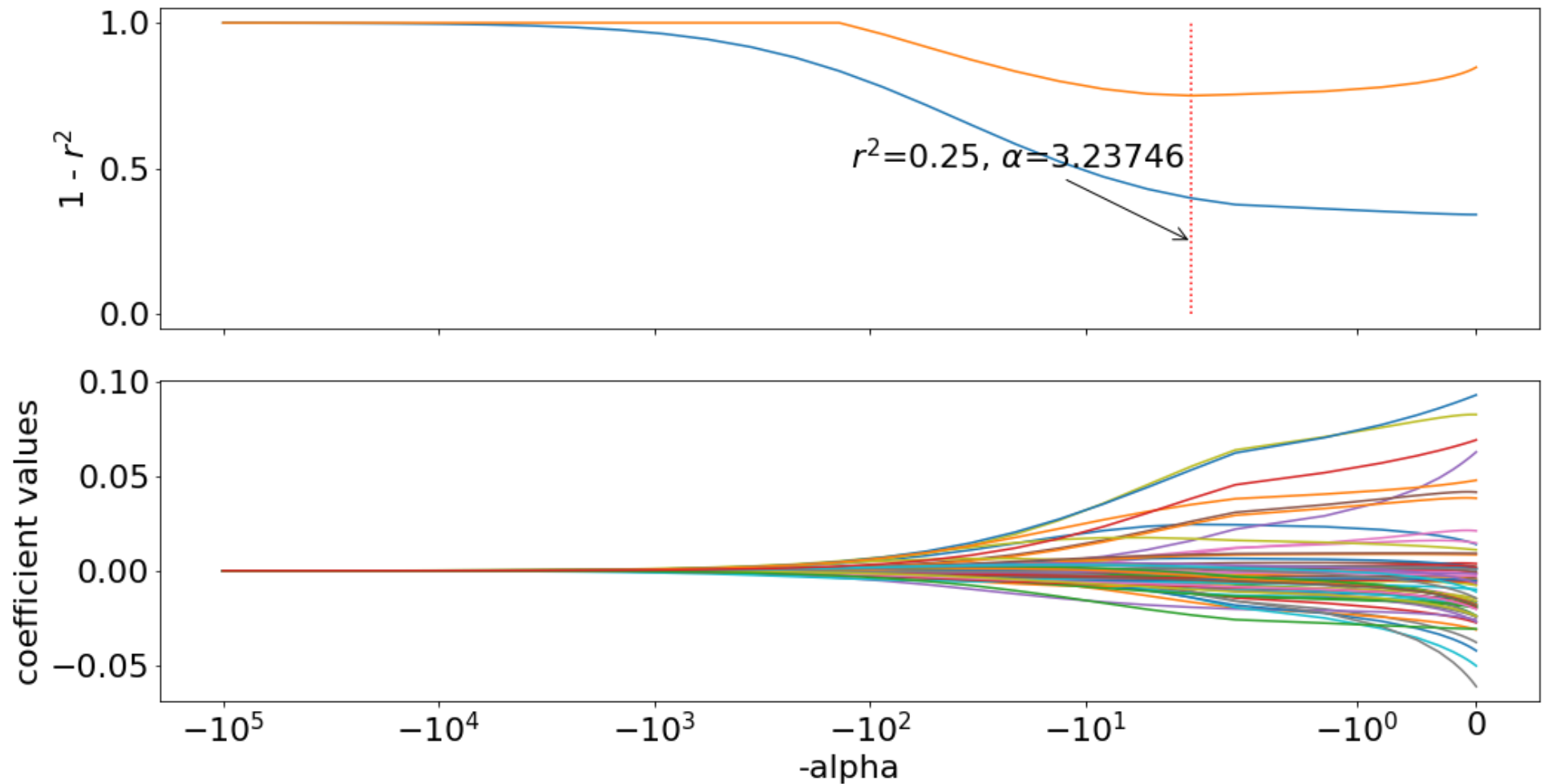
- 67 total features
 - Scikit-learn modeling - using “normalize” scaling
-

INITIAL MODEL

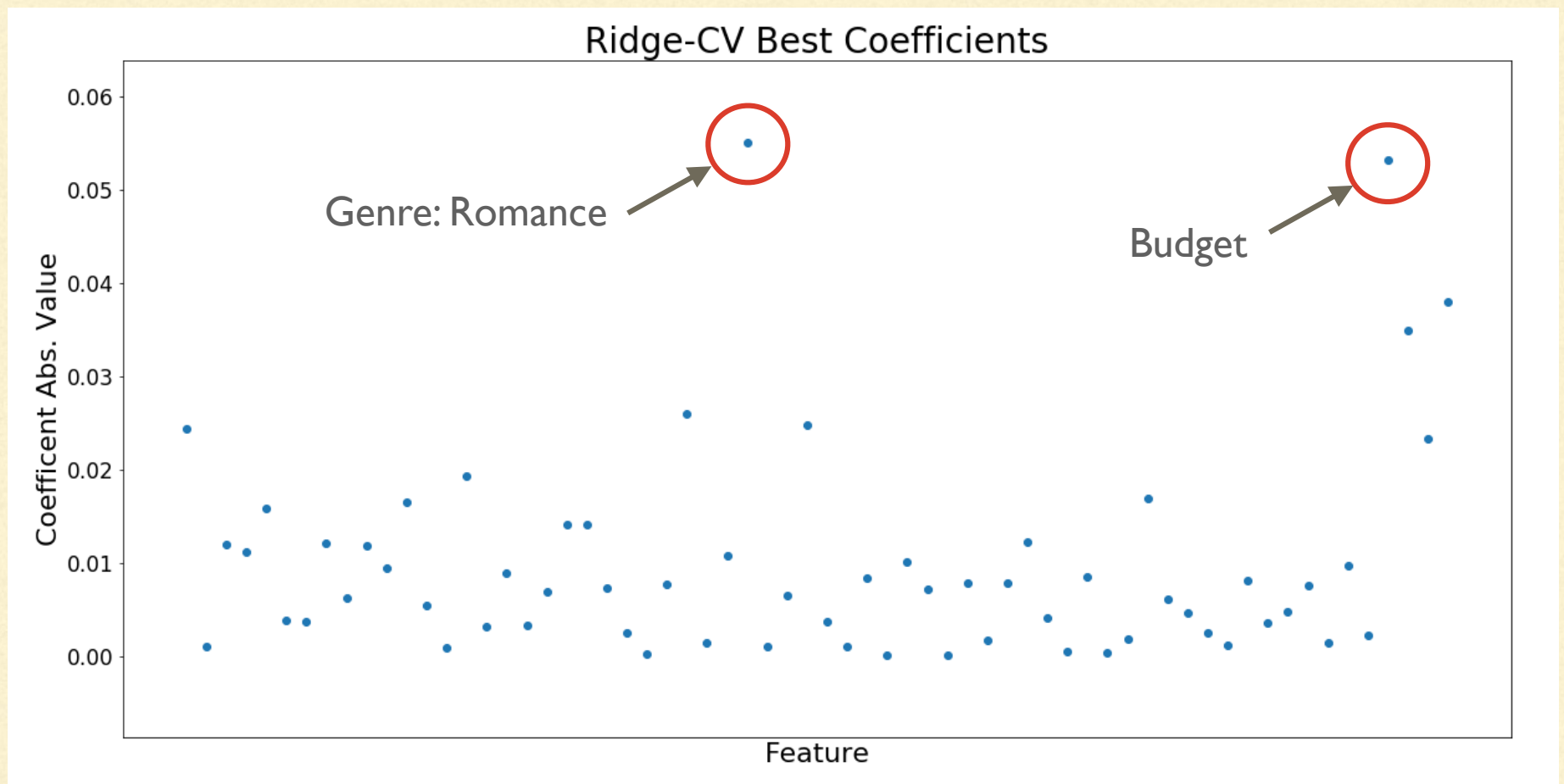
- Train R squared of 0.66 and Test R squared of $-4.06e+26$!



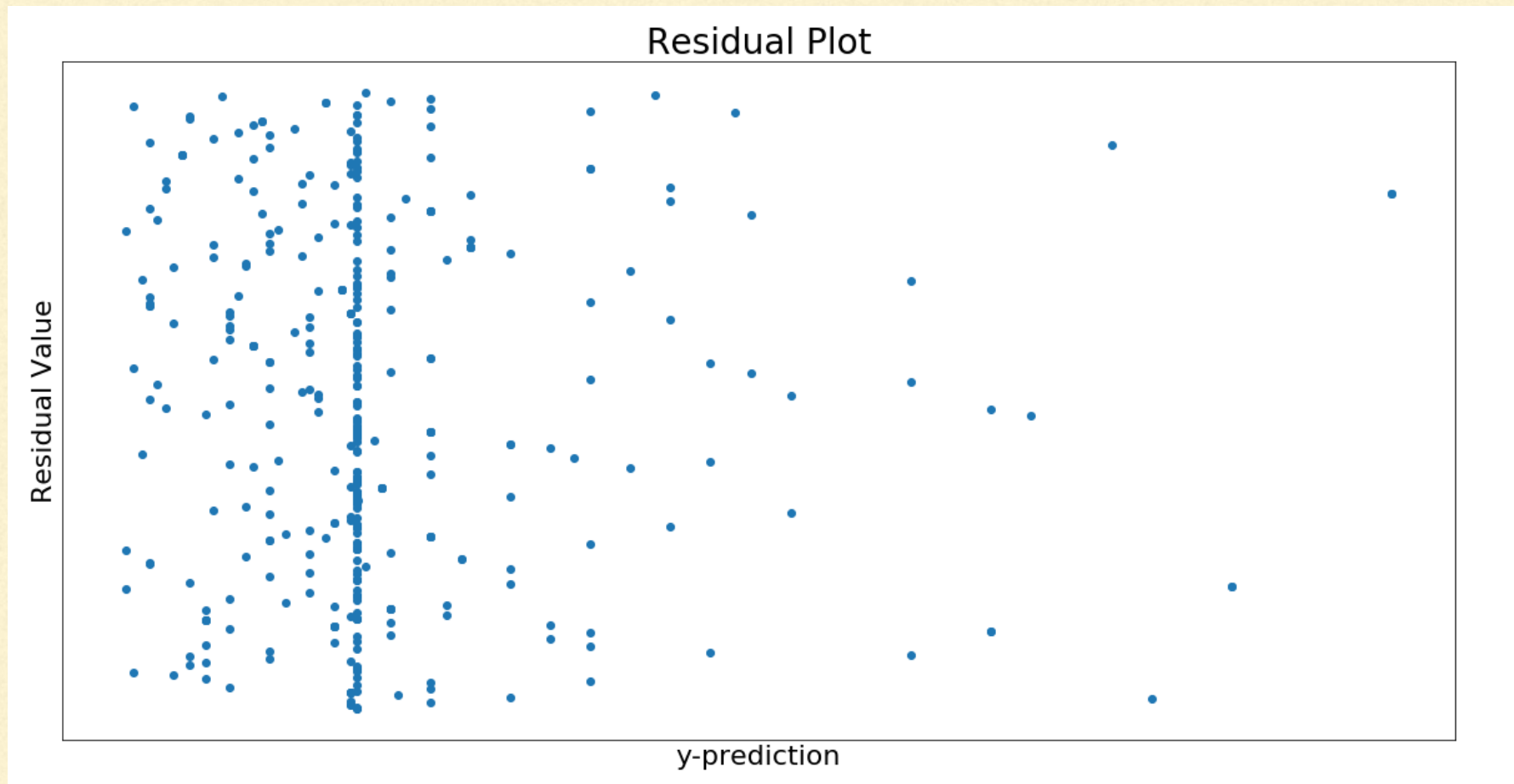
RIDGE-CV REGULARIZATION



RIDGE MODEL



RESIDUALS



NEXT STEPS

- Transform to log-scale
 - Increase data set
 - Retrieve data for viewer/critic ratings
 - Investigate retrieval of streaming data
-

LESSONS LEARNED

- Be thoughtful about feature selection
 - Don't over-filter during data-cleaning
 - Models only as good as data put into them
-

THANK YOU
