

Box Office > 10,000,000 \$ vs. IMDb Rating (log Y-Achse)

Y-axis: Box Office (\$)

X-axis: IMDb Rating

Legend: Daten (green dots), Regressionslinie (log Y) (red line)

A hand-drawn arrow points to the right side of the plot.

Director: Number of Movies vs. Average IMDb Rating

Average IMDb Rating

Number of Movies

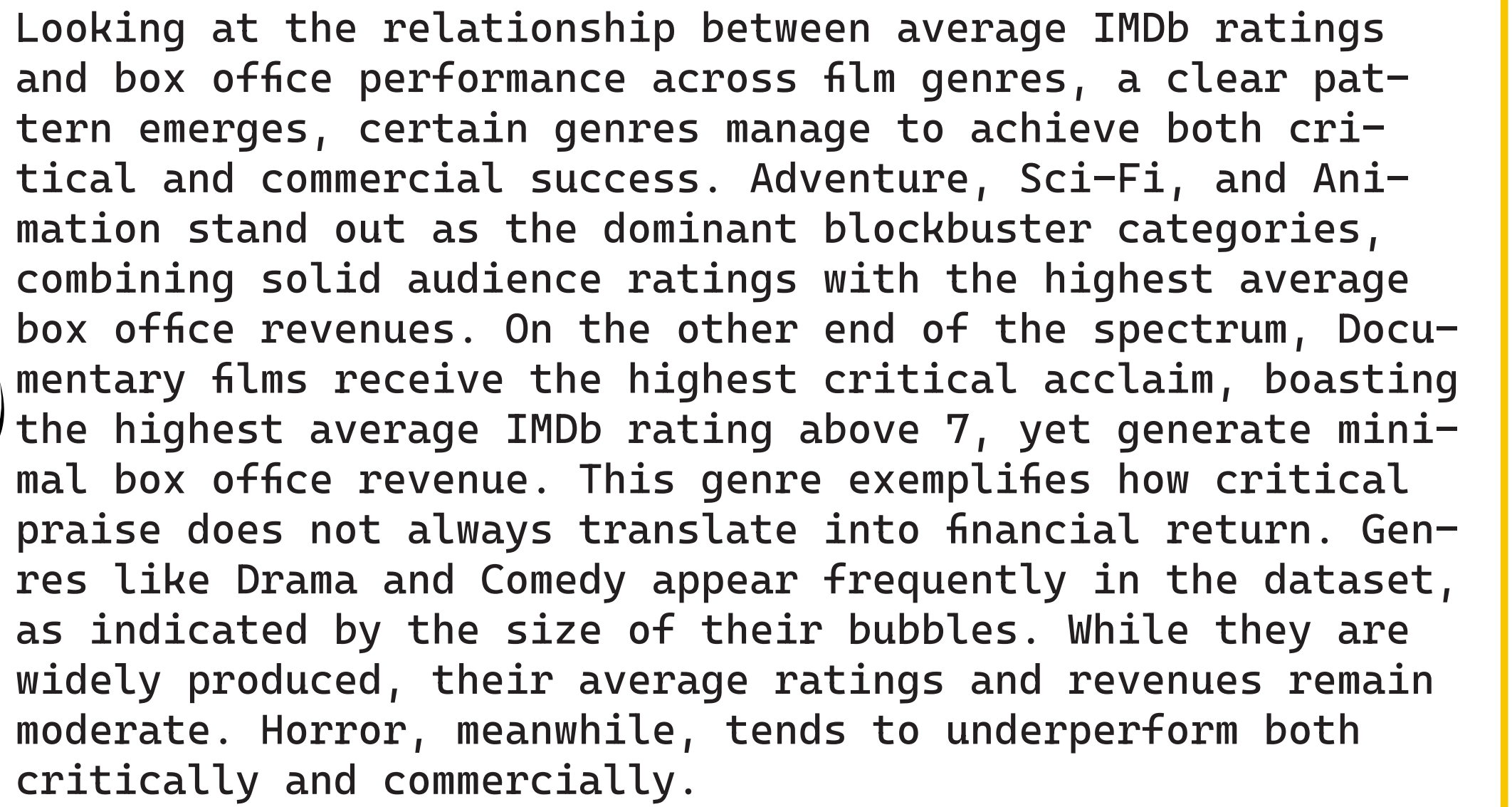
Number of Movies

Now we point our attention towards the „Plot“ field of our dataset. There, the films story is outlined in a few sentences ranging from a minimum of X to a maximum of Y words. According to the IMDb guidelines, the text is sourced from user submissions. There, writing style and rules are also specified. The plot data from OMDb seems to match the IMDb values. For these wordclouds, we have separated the movies into the upper (green border) and lower (red border) .25 percentile to outline potential difference in the plots. We then used a open source sentiment dataset from the Text-Machine-Lab at UMass Lowell to color the words. In addition, we also trained models using RandomForest and XGBoost with varying text vectorizers, to try to predict the rating for a given Plot.

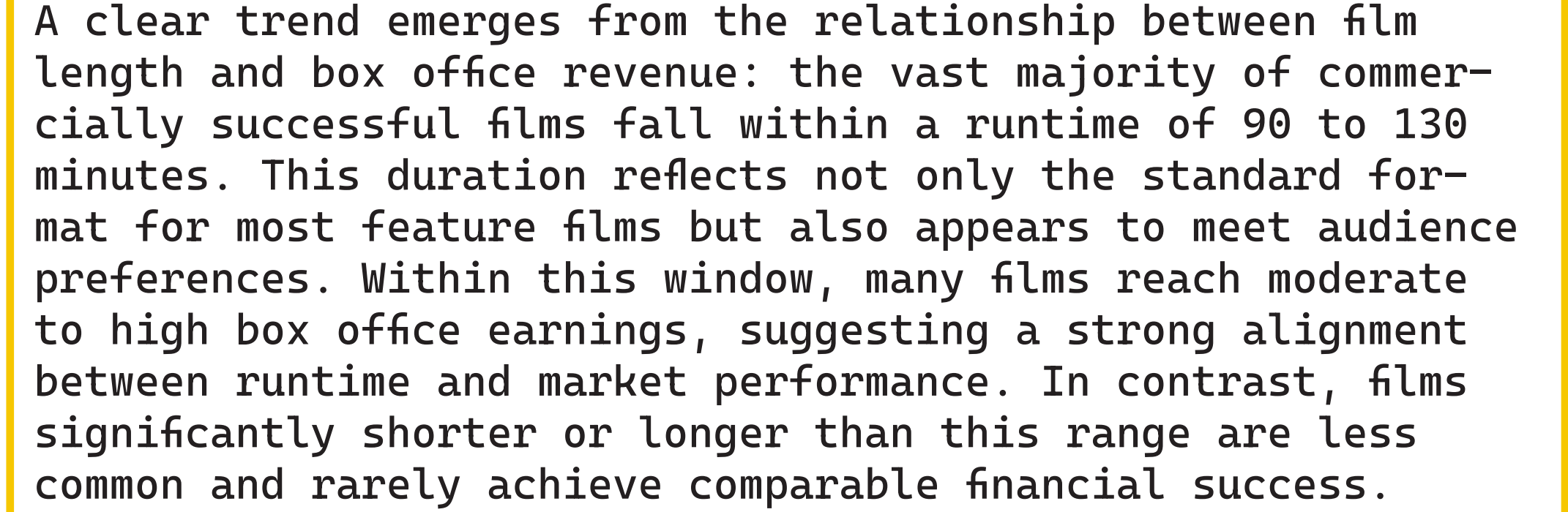
This shows that there is at least some information to be extracted from a movies plot

A tour into which **IMDb** features effect a movies community rating and box office

Once we had the data, we decided to split up and each dive into a separate aspect and subset of features.



Even the most subjective narrative elements—like plot summary and the filmmaking process, our findings suggest that strategic da

[illegible]