

Topic Modeling Report

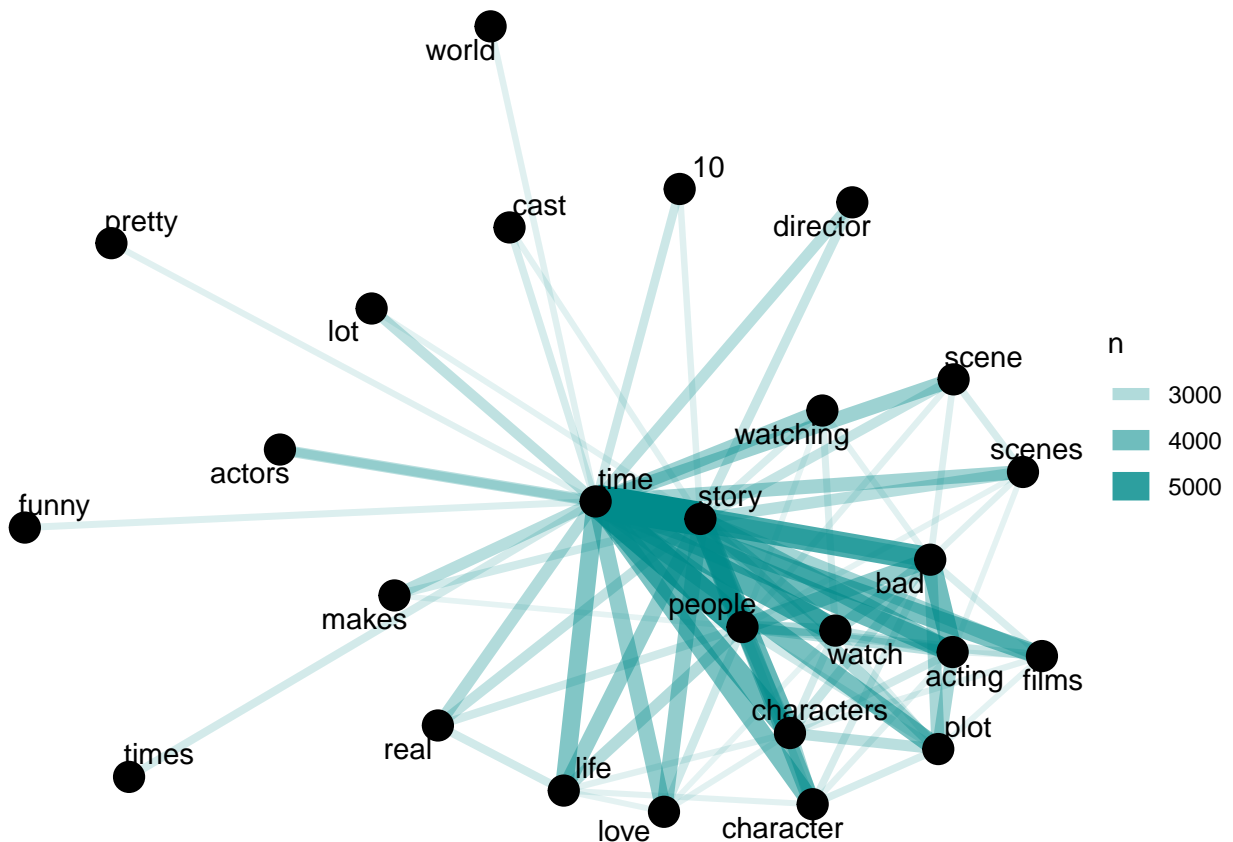
Group 7

2022-11-12

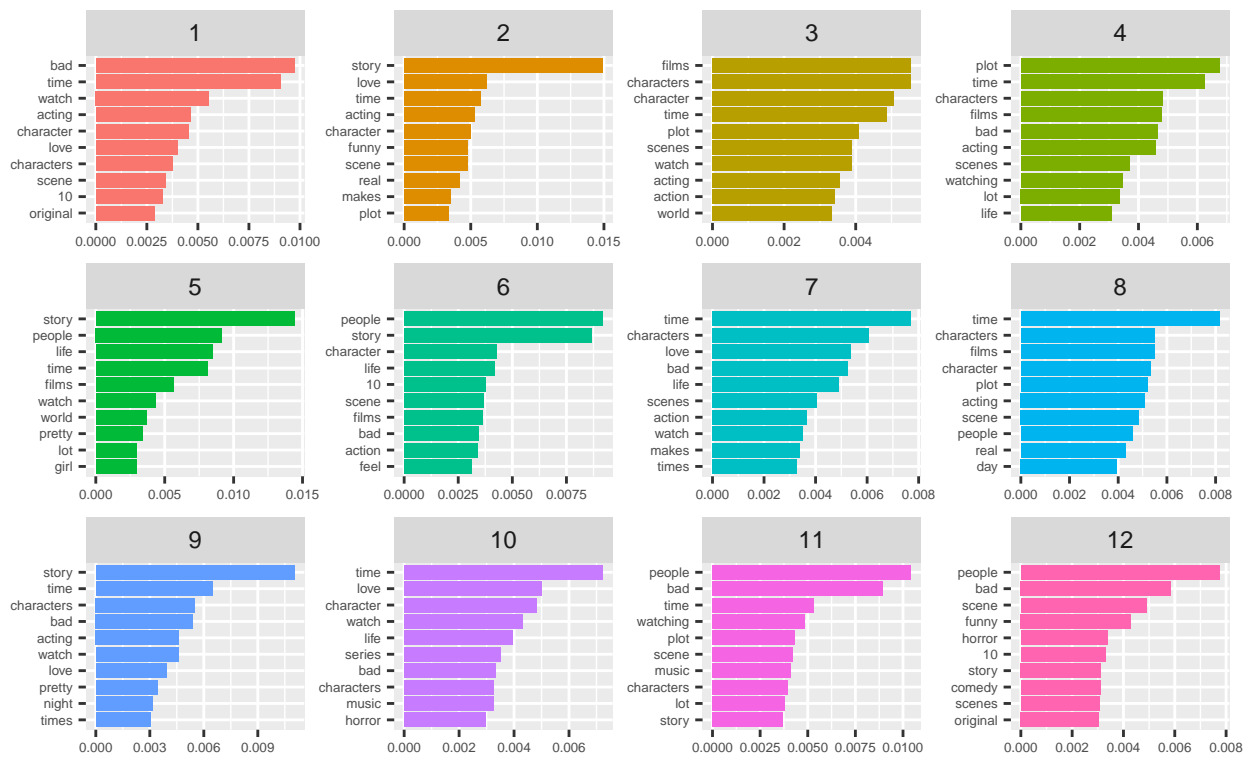
We use IMDB reviews as our data and do the following analysis. First, we set up three stop words, movie, movies and film, since these three words appear most frequently and do not make much sense in movie reviews. Then, we use `unnest_tokens()` function to both break the text into individual tokens and transform it to a tidy data structure.

Afterwards, we look at the terms' inverse document frequency to measure if a word is important in our analysis. As a result, considering both the number of occurrences and term frequency(n/total), we have top 5 terms, 'trivialboring', 'mad', 'stop.oz', 'bad', and 'bought'. Three of them are negative, so that we may conclude that there are more negative review than the positive. In order to prove it, we further draw some plots.

Here, we draw a plot showing how words correlate to each other. We filter the words by saying that the total number of occurrence is larger than 2200 and get 26 words, which is centered around 'time', 'story', 'bad' and 'people'.



Top 10 terms in each LDA topic



β