

Text Mining South Park

Dominik Nerger (i6146759)

June 10, 2017

Contents

1	Introduction	1
2	Libraries	1
3	Data set	2
4	Pre-processing	2
5	Results and visualization	4
5.1	Wordclouds	4
5.2	Co-occurences and dendrogram	5
5.3	Sentiment and Emotion Mining	8
5.4	Non-negative Matrix Factorization	9
5.5	Structural Topic Model	11
6	Conclusion	13
7	Future work	13

1 Introduction

And tell everyone in the past for
us, that no one single answer... is
ever the answer.

UAA Leader, Ep. Go God Go XII

South Park is an animated series revolving around the four boys Stan Marsh, Kyle Broflovski, Eric Cartman and Kenny McCormick who live in South Park, Colorado. South Park has the third most episodes in regard to American animated TV series, totaling 277 episodes over 20 seasons.

Because of the huge amount of episodes, it is an interesting target for the practical assignment, which is part of the course *Information Retrieval & Text Mining*, with the possibility of acquiring new insights into a TV series that is used to satirize American culture.

At first, all libraries that were used will be presented. Afterwards, the data set and the used corpora are introduced. Then the pre-processing tasks that are performed on the data set will be explained. Next, the results and visualization will be discussed while explaining why and how everything has been implemented.

As part of the assignment is to answer combinations of W-questions, this report will also focus on the two combinations that were implemented, namely the *why-when* in the form of a sentiment timeline and the *what-when* in form of a topic timeline. In the end, a conclusion will be drawn and possible improvements for the future will be discussed.

The repository is available on GitHub¹.

2 Libraries

All scripts have been programmed in *R*. To execute the scripts, *R* needs to be installed. To view temporary files that are executed during runtime, e.g. the corpus or a TermDocumentMatrix, it is advised to install RStudio. The libraries necessary for each script are imported at the top of each script, if they are not installed they can be installed by executing:

```
install.packages("library-name")
```

In the following, all libraries that are related to Text Mining techniques will be introduced shortly.

The library **tm** is the Text Mining package of *R*, which enables pre-processing of data sets and allows to build the corpus. **RWeka** is a collection of machine learning algorithms for data mining tasks. **NMF** introduces the Non-negative Matrix Factorization to *R*. **NLP** and **OpenNLP** are libraries that provide Natural Language Processing techniques and are used for NER-Tagging. **syuzhet**

¹<https://github.com/dnerger/South-Park-Text-Mining>

uses tokenized input to extract sentiments and emotions and contains the three sentiment lexicons *bing*, *afinn* and *nrc*. The package **stm** is used for Structural Topic Modeling which is LDA with additional meta-data and can be visualized using the package **LDavis**. Libraries used for visualization include **igraph**, **ggplot2**, **ggraph**, **viridis** and **wordcloud**.

3 Data set

The data set spans from seasons **1** to **18**, adding up to 257 episodes overall with a file size of 5.41MB. It contains 70896 rows, with each row possessing information about the **season**, **episode**, **character** and **line** and is available as a *.csv* file.

The data set has been crawled by Bob Adams and is available to download on GitHub². It has been assembled by crawling the South Park Archives³. The code of this GitHub repository is not available to the public.

I made an attempt at a crawl with the R package **rvest**, which is able to harvest data from websites. However, the *html* data of the episode scripts provided on the South Park Archives differ between episodes. For 16 of 20 seasons, the whole episode needs to be pre-processed manually because the **line** and **character** are extracted individually and can not be merged because of random whitespace lines that are introduced by **rvest**.

In the scripts, the data set is used in three different variations:

1. *dialogue* - each entry holds information about season, episode, character and line
2. *by.season* - each entry holds information about season and the complete script for that season
3. *by.episode* - each entry holds information about that episode, season and script for that episode

Both *by.episode* and *by.season* are created by using the data from *dialogue*.

4 Pre-processing

In general, the standard pre-processing tasks are executed on the corpus. Depending on the later use, different packages perform these. The package **tm** is used to convert the corpus to lower case, remove punctuation and numbers as well as strip whitespaces. Furthermore, stopwords are removed using the stopword list from the Snowball stemmer project. The corpus is stemmed using Porter's stemming algorithm.⁴

²<https://github.com/BobAdamsEE/SouthParkData>

³http://southpark.wikia.com/wiki/South_Park_Archives

⁴<http://snowball.tartarus.org/algorithms/english/stop.txt>

It is used together with several DocumentTerm- and TermDocument-Matrices. From these DTMs/TDMs, sparse terms are removed. The maximum allowed sparsity depends on the data set and the intended use, with the allowed percentage being higher for the *by.episode* data set in comparison to *by.season* because there are more documents in the *by.episode* data set.

The pre-processing for the Structural Topic Model is done with the same tasks, however those are performed by the package **stm** itself.

For the Named Entity Recognition and Sentiment Annotation with package **openNLP**, no pre-processing tasks are performed. They can be executed and might improve the overall runtime, but it does not change the overall result. Especially stemming and the lower case conversion could decrease the results because stemmed or lower case entities, as well as the sentiments that are mined, will not be recognized.

5 Results and visualization

To sufficiently present the results and the corresponding visualization, each relevant script will present them in their own subsection.

5.1 Wordclouds

To start to gain an insight into South Park, I started with building wordclouds that display the most frequent words to see what South Park is about. However, even with stopwords removed, the words that are most frequent mostly include words such as *get*, *just* and *now*. From these words, nothing about the TV show can be derived.

To get meaningful answers to what is important in and characteristic for South Park, the corpus is built with a term frequency-inverse document frequency (tf-idf) weight, which reflects how important each term is to the corresponding document in the corpus.



Figure 1: Terms with highest TF-IDF score - *by.episode*



Figure 2: Terms with highest TF-IDF score - *by.season*

With this weighting, the terms that are most important over all 18 seasons can be derived, as can be seen in Figures 1 and 2.

The difference between the wordclouds descends from the different *idf* weights for the data sets *by.season* and *by.episode*. The words in Figure 1 show the names of secondary characters such as Butters, Kenny and Wendy. These characters are not featured in Figure 2 since they do appear in every season, but not in every episode. This leads to their tf-idf weight being zero in the *by.season* data set.

Figure 3 shows the most frequent quadgram phrases without stopwords included. This includes sentences like the trademark phrase "Oh my god, they killed Kenny", which is shown as "oh god kill kenni" because this was calculated with the pre-processed corpus.



Figure 3: Wordcloud containing ngrams (n=4) by frequency

5.2 Co-occurences and dendrogram

To gain insight into the relations between characters and to confirm the major characters of South Park, which are Stan, Kyle and Cartman, I decided to visualize the 20 characters with the most lines.

As can be seen in Figure 4, Cartman is the character with most lines overall, with Stan and Kyle, as well as Butters and Randy being other major characters. This has been implemented by aggregating all lines for a character and then calculating the percentages in respect to the overall lines.

However, this does not show the relations between the characters. To visualize this, I decided to create a co-occurrence network by creating a DocumentTermMatrix that contains the characters as terms and the unique episode number as document. The value in each entry is the amount of lines for the character in the specific episode.

By using the formula $C = A^t A$ with A as the DocumentTermMatrix, I created a TermTermMatrix that contains the co-occurrences between lines of characters, resulting in the co-occurrences of the 20 most relevant characters in South Park.

As can be seen in Figure 5, all 20 major characters interact in some way throughout all 257 episodes, with Cartman, Kyle and Stan being characters that are closely tied together.

Since it is hard to derive the relations between the characters from Figure 5, I decided to use a cluster dendrogram to visualize which characters have the most episodes in common. This is accomplished by normalizing the DocumentTer-

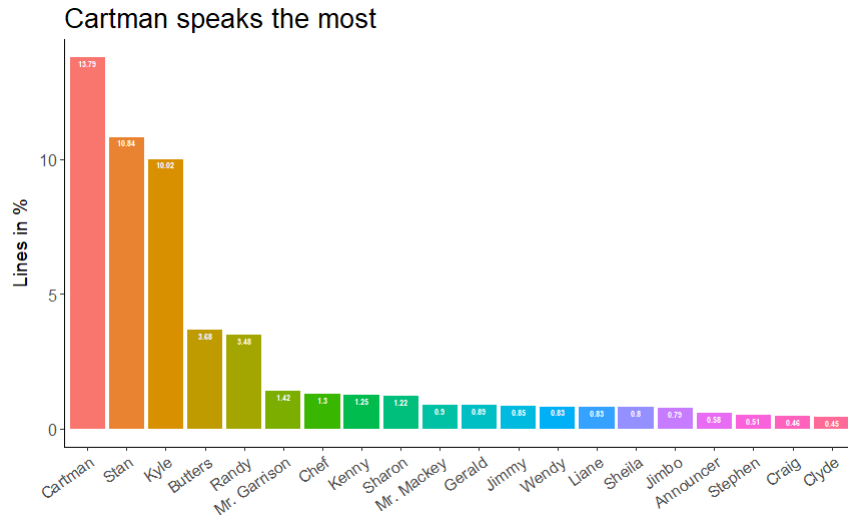


Figure 4: Speaking percentage for 20 major characters

mMatrix in regard of the row sums and then applying hierarchical clustering with the Manhattan distance.

Figure 6 shows the dendrogram clustering the characters. The further down the node of two characters splits, the more episodes these characters have in common. This shows once again that the major characters Cartman, Stan and Kyle are closely tied together. Furthermore, interesting connections between the students Clyde and Craig as well as the fathers Stephen and Randy can be derived.

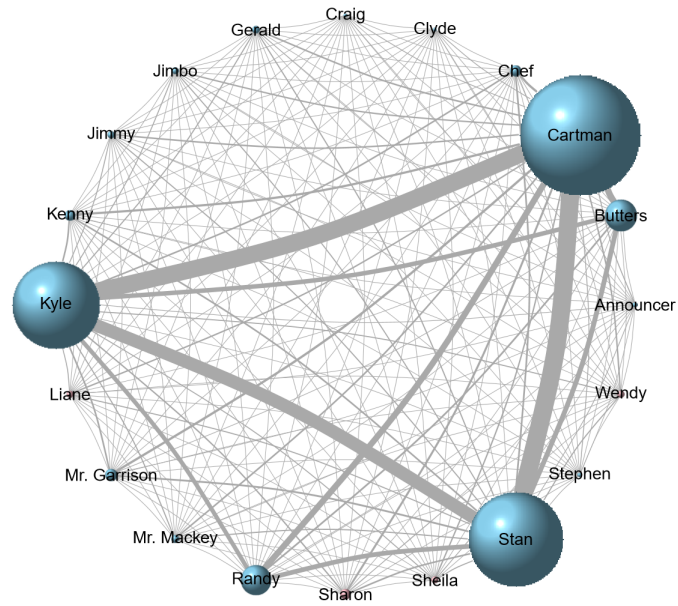


Figure 5: Co-occurrence network

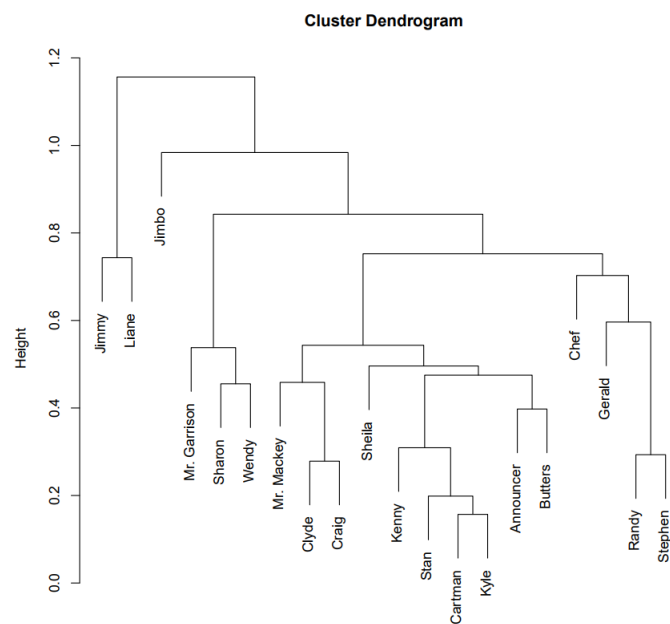


Figure 6: Cluster dendrogram of 20 major South Park characters

5.3 Sentiment and Emotion Mining

Another approach to learn something about South Park is to extract the sentiments and emotions that are expressed by the characters. The sentiments are annotated onto the seasons, using the Annotator provided by the **NLP** packag.

After the complete corpus has been annotated, which results in little chunks of words being tied together, the sentiments are processed with the library **syuzhet** and the sentiment dictionary *nrc*. During the processing the word chunks are grouped with 100 chunks being aggregated. Under the assumption that each character line contains at least 1 chunk, this results in a minimum of 700 grouped chunks for 70000 character lines, which can be seen as roughly three different sentiment values for each episode. This can be visualized with a sentiment timeline, as can be seen in Figures 7 and 8. One interesting observation is that the sentiments in Figure 8 for season 18 are mostly positive with only a few segments having a negative sentiment score.

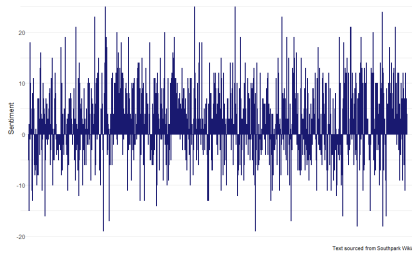


Figure 7: Sentiment timeline over 18 seasons

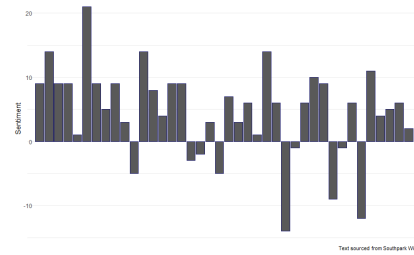


Figure 8: Sentiment timeline for season 18

However, it is hard to derive a meaning from the sentiments that are shown in Figure 7. Therefore, I decided to filter and transform the sentiment timelines with a low-pass Fourier transform using a low-pass size of 3. The Fourier Transform puts the data into a periodic shape, as can be seen in Figures 9 and 10.

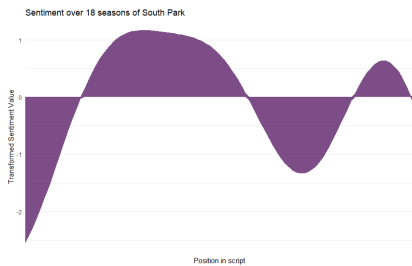


Figure 9: Fourier transformed sentiment timeline over 18 seasons

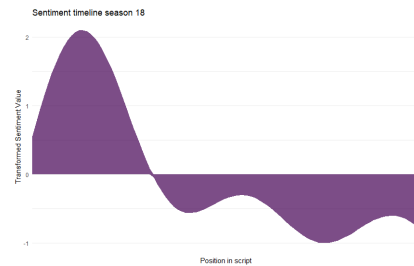


Figure 10: Fourier transformed sentiment timeline for season 18

This shows that during the first few seasons, in relation to the rest of the show, the sentiments were mostly negative.

To gain an insight into the emotions that are expressed in South Park, I aggregated all emotions that were extracted using the *nrc* lexicon and visualize them in the bar charts in Figures 11 and 12. It can be derived that the emotions are rather negative, with around 50% more sadness than joy. Also, a lot of anger and fear are expressed throughout the show.

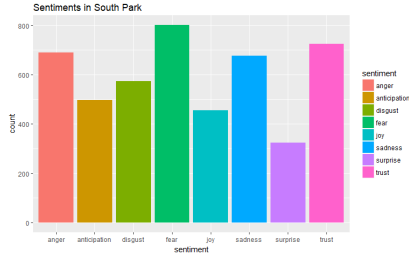


Figure 11: Emotions over 18 seasons

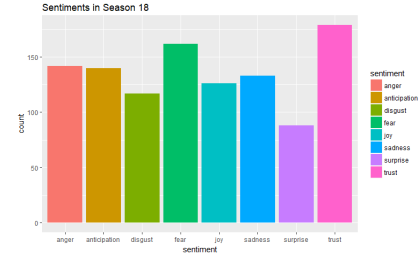


Figure 12: Emotions for season 18

5.4 Non-negative Matrix Factorization

The Non-negative Matrix Factorization (NMF) introduces the possibility to cluster documents in a TermDocumentMatrix by topic. For this implementation, the *by.season* corpus was used, which after pre-processing was turned into a TDM with a maximum sparsity of 40% allowed. The execution of the NMF algorithm provided by package **nmf** results in the coefficient and basis matrices. The coefficient map in Figure 13 shows the clusters that are calculated for the seasons. The corresponding basis map that shows the most important words for each cluster is appended to the report⁵.

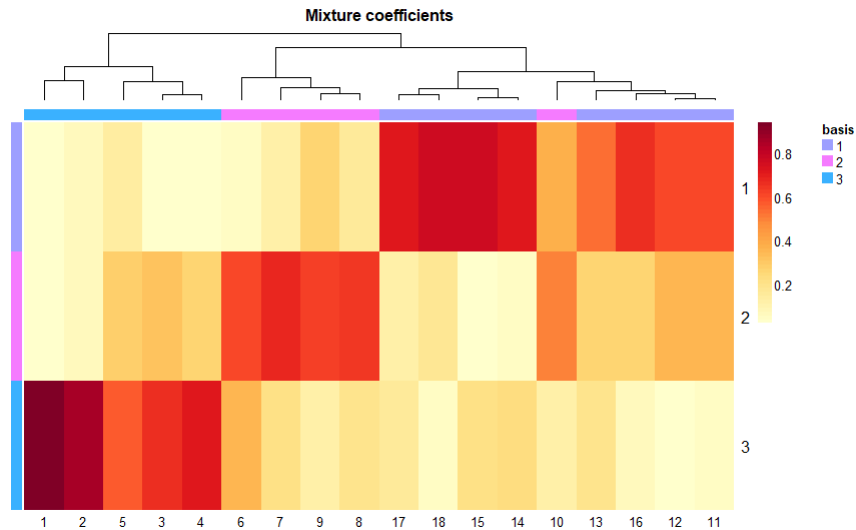


Figure 13: NMF topic clusters

⁵*NMF.basemap.pdf*

The seasons are clustered with seasons 1-5, 6-9 and 10-18 being grouped together. To try to understand why they are clustered like this, I looked into the IMDb page⁶ of South Park and was able to find interesting correspondences, with major changes in writing in years 2001 and 2002(Figure 15) and producing in 2006(Figure 14).

Series Produced by	
Trey Parker	... executive producer (278 episodes, 1997-2016)
Matt Stone	... executive producer / producer (278 episodes, 1997-2016)
Anne Garefino	... executive producer / producer / supervising producer (223 episodes, 1997-2016)
Frank C. Agnone II	... supervising producer / executive producer / producer / post-production producer / line producer (212 episodes, 1997-2016)
Vernon Chatman	... producer / consulting producer (109 episodes, 2006-2016)
Daryl Sancton	... line producer (98 episodes, 2006-2016)
Pam Brady	... creative producer / producer / consulting producer (95 episodes, 1997-2006)
Jennifer Howell	... associate producer / supervising producer (93 episodes, 1997-2005)
Adrien Beard	... producer (90 episodes, 2007-2016)
Bruce Howell	... producer / executive producer (90 episodes, 2007-2016)
Eric Stough	... producer (90 episodes, 2007-2016)
Eric Rivlinja	... producer (63 episodes, 2007-2012)
Deborah Liebling	... executive producer (58 episodes, 1997-2002)
Kurt Nickels	... associate producer (53 episodes, 2007-2015)

Figure 14: Changes in producing

Series Writing Credits	
Trey Parker	... (creator) (278 episodes, 1997-2016)
Matt Stone	... (creator) (278 episodes, 1997-2016)
Brian Graden	... (developer) (238 episodes, 1997-2016)
David R. Goodman	... (88 episodes, 1997-2002)
Nancy Pimental	... (62 episodes, 1998-2001)
Kyle McCulloch	... (56 episodes, 1999-2002)

Figure 15: Changes in writing

This corresponds to the different clusters obtained by the NMF, with Season 6 airing in 2002 and Season 10 airing in 2006, so it can be argued that the clusters are related to changes in the overall production of the show, namely writing and producing.

⁶<http://www.imdb.com/title/tt0121955/>

5.5 Structural Topic Model

Because the visualization of the basis map of the NMF doesn't provide an interactive way to inspect the topics and look at their most frequent terms, and to gain further insights into the topics that define South Park, I decided to use the **stm** package which provides the Structural Topic Model. Basically, **stm** is similar to other probabilistic topic models like the Latent Dirichlet Allocation(LDA). However, it allows to incorporate metadata into the topic model and enables the user to discover the main topics in a text corpus with a connection to the metadata.

To model the topics for all 18 seasons and season 18 only, the *by.episode* corpus is used and pre-processed using the *textProcessor* provided by the **stm** package. Afterwards, the processed corpora are prepared, removing frequent terms that appear in more than 25 documents for all seasons and in more than 4 documents for season 18 only. With the prepared corpora, the Structural Topic Model is executed, modeling 25 and 4 topics, respectively. Since South Park mostly covers different topics between episodes, there are a lot of possible topics. With limiting this to 25, the main topics of South Park can be extracted, as can be seen in Figure 16.

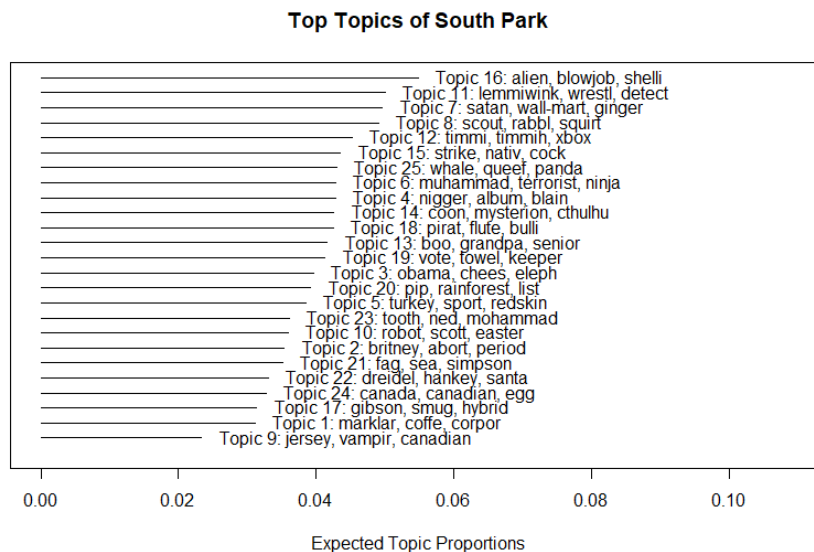


Figure 16: Structural Topic Model topics for South Park

The most probable topic for each episode can be extracted and visualized. This can be seen in the topic timeline in Figure 17. The results of this timeline seem to be accurate. As a sample to prove this, I chose topic 14 with the characteristic words "coon, mysterion, cthulhu" focusses on the superheroes saga around Coon, that appears in different seasons. In the topic timeline, there is a cluster shortly after episode 200 on the x-axis for topic 14. This relates to

episodes 206 to, which is a 3 episode feature⁷ focussing on superheroes Coon and Mysterion.

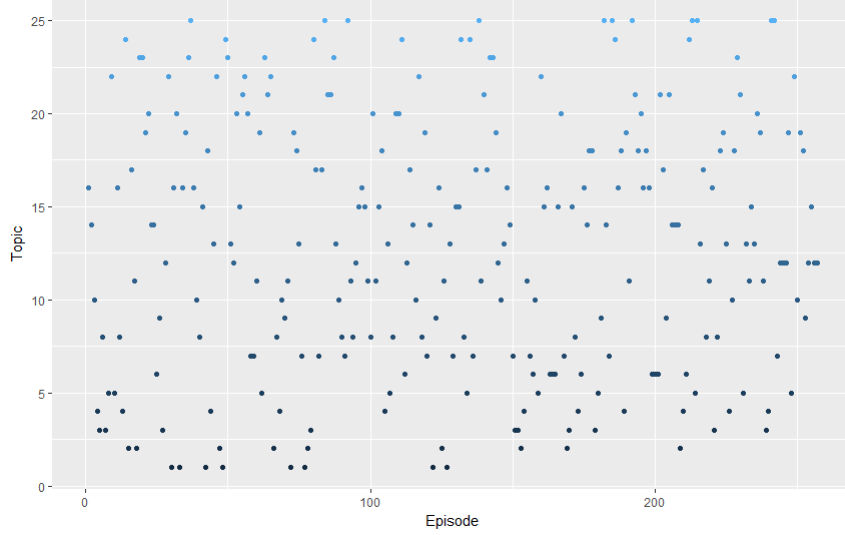


Figure 17: Topic timeline

To further inspect the different topics and gain insights into the clusters that have been calculated by the STM, I decided to use the package *LDavis* to visualize the different topics. Figure 18 displays the generated topics.

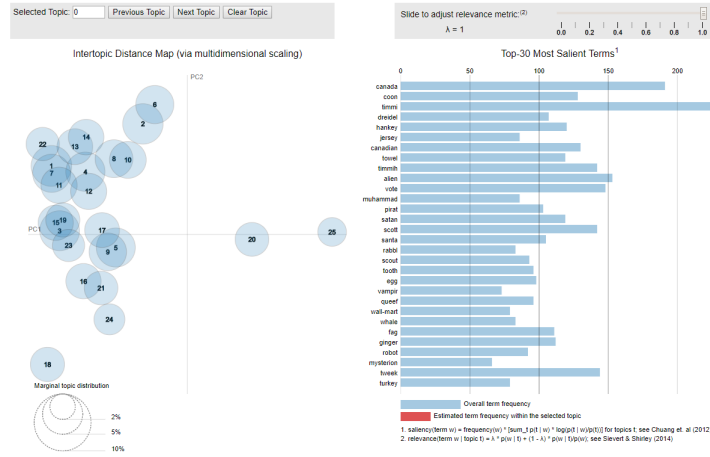


Figure 18: STM visualization using LDavis

If the script *stm.r* is executed, this will open the web browser, in which the topics can be inspected interactively, showing the estimated term frequency within the selected topic for the 30 most relevant terms per topic. Videos of this are appended to the report.

⁷E206 *Coon 2: Hindsight*, E207 *Mysterion Rises*, E208 *Coon vs. Coon and Friends*

All visualizations shown are appended to the report for only season 18 as well.⁸

6 Conclusion

In conclusion, this report presented the different data sets that were used as well as the pre-processing that is applied to them in each case. Furthermore, the results and visualization were shown and a description as well as a reason for each result has been mentioned. The text mining of South Park has lead to a few observations, like the changes in writing and producing leading to different clusters of seasons or that the sentiment in the first few seasons of South Park were more negative than in later seasons.

To answer the combinations of W-questions, the topics extracted with the Structural Topic Model were visualized with a topic timeline that shows the most probable topic for each episode. This answers the combination of *what-when*. Also, a sentiment timeline for the show has been presented and discussed, answering the combination of *why-when*. All images, scripts and data as well as videos displaying the Structural Topic Modeling are available in the GitHub repository.

7 Future work

To further increase the results, a few improvements could be introduced to the analysis of South Park. One possible future implementation is to visualize the sentiment and emotion timelines for different characters, which would then combine the three W-questions *who*, *why* and *when*. This could lead to interesting results, with characters like Cartman being more aggressive than relaxed characters like Stan. Furthermore, the visualization of the topic timeline could be improved by creating a river plot that contains the content of the topics.

⁸ *Topics18.png*, *Topics18Timeline*