

# Text Mining South Park

Dominik Nerger (i6146759)

June 3, 2017

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Libraries</b>	<b>2</b>
<b>3</b>	<b>Data set</b>	<b>2</b>
<b>4</b>	<b>Pre-processing</b>	<b>8</b>
<b>5</b>	<b>Results and visualization</b>	<b>8</b>
<b>6</b>	<b>Conclusion</b>	<b>8</b>
<b>7</b>	<b>Future work</b>	<b>8</b>

# 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.

The repository is available on GitHub<sup>1</sup>.

## 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** extracts sentiments from text and contains the three sentiment dictionaries *bing*, *afinn* and *nrc*. The package **stm** is used for Structural Topic Modeling which is LDA with additional met-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.

---

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

The data set has been crawled by Bob Adams and is available to download on GitHub<sup>2</sup>. It has been assembled by crawling the South Park Archives<sup>3</sup>. 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**.

The data set is used in three different versions:

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*.

---

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

<sup>3</sup>[http://southpark.wikia.com/wiki/South\\_Park\\_Archives](http://southpark.wikia.com/wiki/South_Park_Archives)

## Most used words



Figure 1: General wordcloud, containing terms by frequency

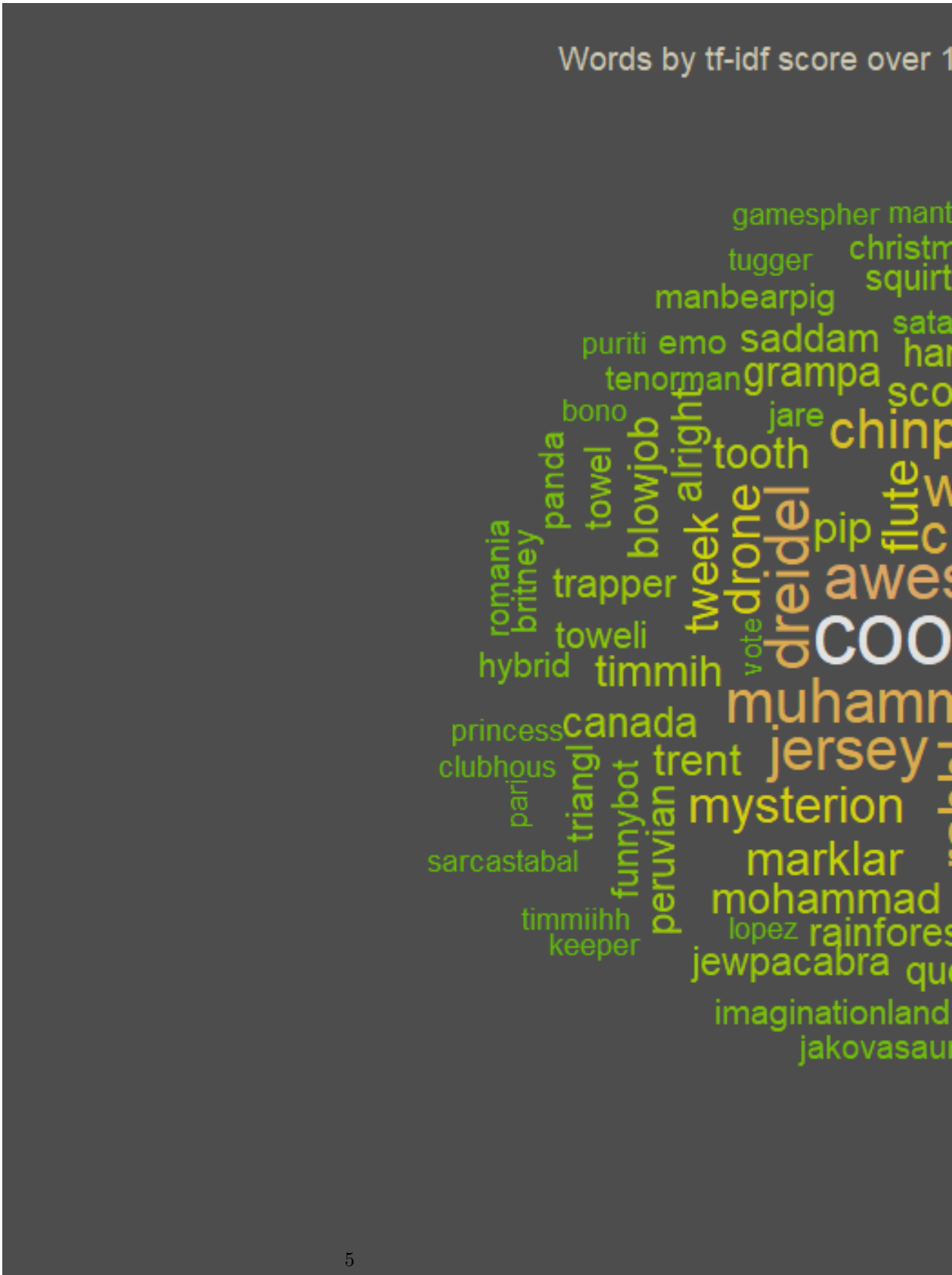


Figure 2: Wordcloud containing terms by TF-IDF score

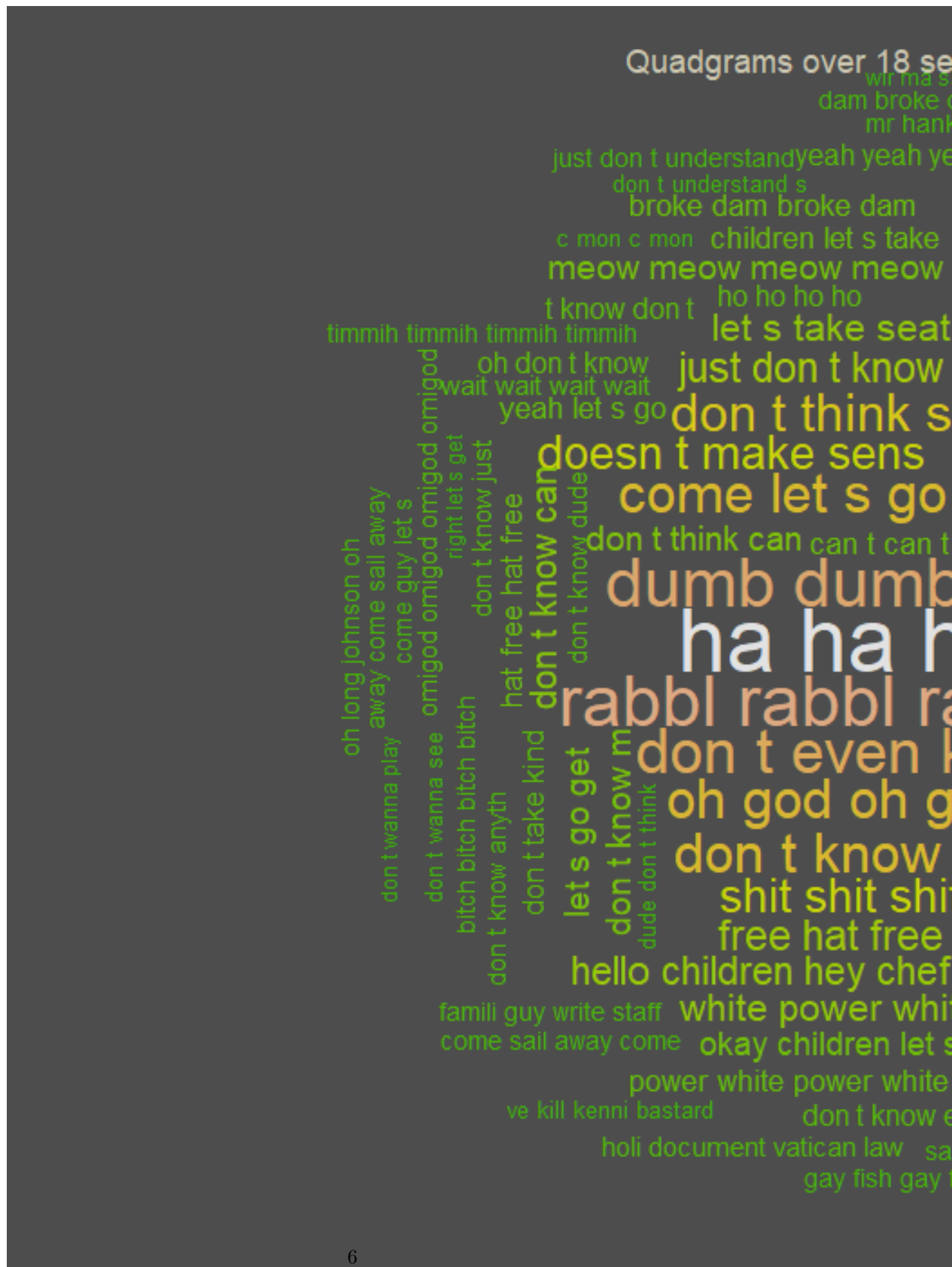
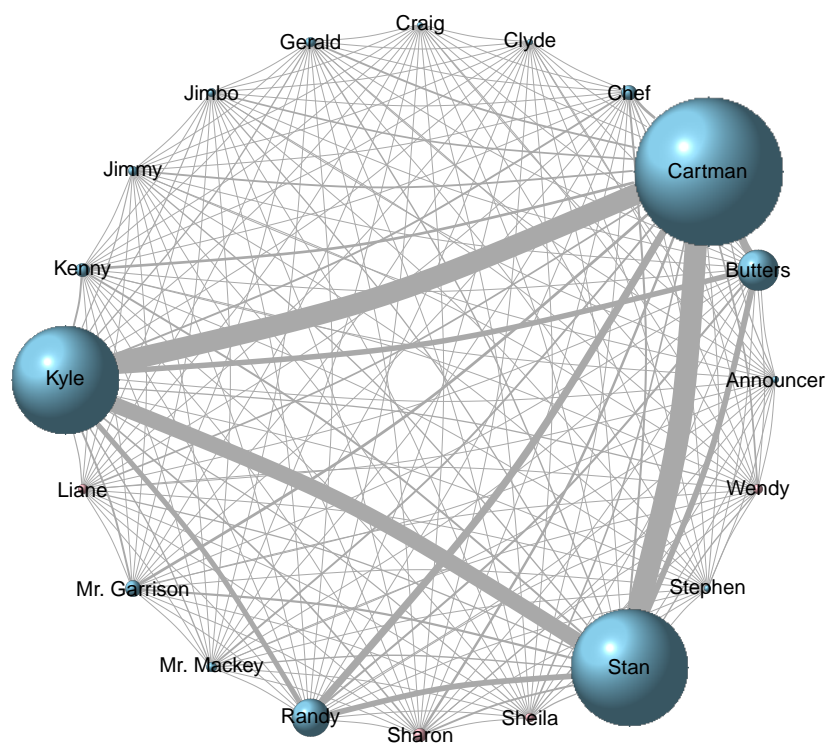


Figure 3: Wordcloud containing ngrams ( $n=4$ ) by frequency



Matrix.pdf

Figure 4: Co-occurences of characters over 257 episodes, size of vertex in relation to amount of lines

## 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.<sup>4</sup> It is used together with several DocumentTerm- and TermDocument-Matrices.

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 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 will not be recognized.

## 5 Results and visualization

## 6 Conclusion

In conclusion,

## 7 Future work

---

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