

# **Burp NLP: A Rick and Morty text analysis**

## **Text Analysis**

an der Fakultät für Wirtschaft  
im Studiengang Wirtschaftsinformatik

an der  
DHBW Ravensburg

Verfasser:	Anton Geiger
Ausbildungsbetrieb:	Festo SE & Co KG
Anschrift:	Ruiter Straße 82 73734 Esslingen Berkheim
Dozent:	Prof. Dr. Oliver Sampson
Abgabedatum:	31.3.2025

# Inhaltsverzeichnis

<b>Abkürzungsverzeichnis</b>	<b>III</b>
<b>Glossar</b>	<b>IV</b>
<b>Abbildungsverzeichnis</b>	<b>V</b>
<b>Tabellenverzeichnis</b>	<b>VI</b>
<b>1 Einleitung</b>	<b>1</b>
<b>2 Hauptteil</b>	<b>2</b>
2.1 Data Understanding . . . . .	2
2.2 Data Preparation . . . . .	6
2.3 Topic Modeling . . . . .	6
<b>3 Modeling</b>	<b>10</b>
3.1 IMDB Rating prediction . . . . .	10
3.2 Speaker detection . . . . .	15
<b>4 Diskussion</b>	<b>17</b>
<b>5 Anhang</b>	<b>18</b>
<b>Literatur</b>	<b>22</b>
<b>Selbständigkeitserklärung Anton Geiger</b>	<b>23</b>

## **Abkürzungsverzeichnis**

LSTM ..... Long Short Term Memory

## **Glossar**

## Abbildungsverzeichnis

2.1	Speaker distribution per lines . . . . .	3
2.2	Family Sanchez spoken lines in detail . . . . .	4
2.3	Episode distribution . . . . .	4
2.4	Length of episode descriptions per episode . . . . .	5
2.5	Word Cloud transcript . . . . .	6
2.6	BERT Topics for Transcrip (all seasons) . . . . .	7
2.7	Topic relevance Family Portal Gun over time . . . . .	8
2.8	Topic relevance Titanic Topic over time . . . . .	9
2.9	Spaghetti an Citadel over time . . . . .	9
3.1	Neural Net Test Evaluation . . . . .	11
3.2	Neural Net confusion matrix (thresshold = 0.5) . . . . .	11
3.3	LSTM Classification . . . . .	12
3.4	LSTM Classification confusion matrix (treshhold = 0.4) . . . . .	12
3.5	Test Results Regression . . . . .	14
3.6	Number of Mortytown apperances with IMDB prediction . . . . .	14
3.7	Confusion Matrix speaker detection . . . . .	16
5.1	Top 25 BERT topics for transcript (all seasons) . . . . .	19
5.2	Hierarchial Clustering topics on transcript . . . . .	20
5.3	All Topics LDA for description (all seasons) . . . . .	20
5.4	Confusion Matrix Bayes Classification with BOW . . . . .	21
5.5	Confusion Matrix Bayes Classification with TF IDF . . . . .	21

## Tabellenverzeichnis

2.1	Rick and Morty Transcript Dataset . . . . .	2
2.2	Rick and Morty episode descriptions dataset . . . . .	2
2.3	IMDB ratings per episode dataset . . . . .	2
5.1	5 Fold Cross Validation . . . . .	18
5.2	Metrics for speaker detection . . . . .	18
5.3	Test scores accuracy LSTM regression . . . . .	18

## **Formeln**

# **1 Einleitung**



## 2 Hauptteil

### 2.1 Data Understanding

For our project we use three different datasets. The first one is a dataset is downloaded from kaggle where we find the transcripts until season 5 formatted in the following shape 2.1. As the transcript also holds information about the scenery there is not just spoken text included in there. We also find short descriptions about the surroundings and the scenery. An example fo that is also visible in the first row of table 2.1.

episode no.	speaker	dialouge
1	Rick	stumbles in drunkenly, and turns on the lights. Morty! You gotta come on. Jus'... you gotta come with me.
1	Morty	rubbs his eyes. What, Rick? What's going on?
1	Rick	I got a surprise for you, Morty.

**Tab. 2.1:** Rick and Morty Transcript Dataset

The other datasets are retrieved by self written web scraping scripts from the Rick and Morty Wikipage and the IMDB website. With the first script, we were able to create following table containing all episode descriptions provided in the Rick and Morty Fandom 2.2.

id	title	text
0	Pilot	n the middle of the night, an obviously drunk Rick bursts
1	Lawnmower Dog	Jerry complains that the family dog, Snuffles, is stupid
2	Anatomy Park (Episode)	It's Christmas, and Jerry tries to enforce the idea

**Tab. 2.2:** Rick and Morty episode descriptions dataset

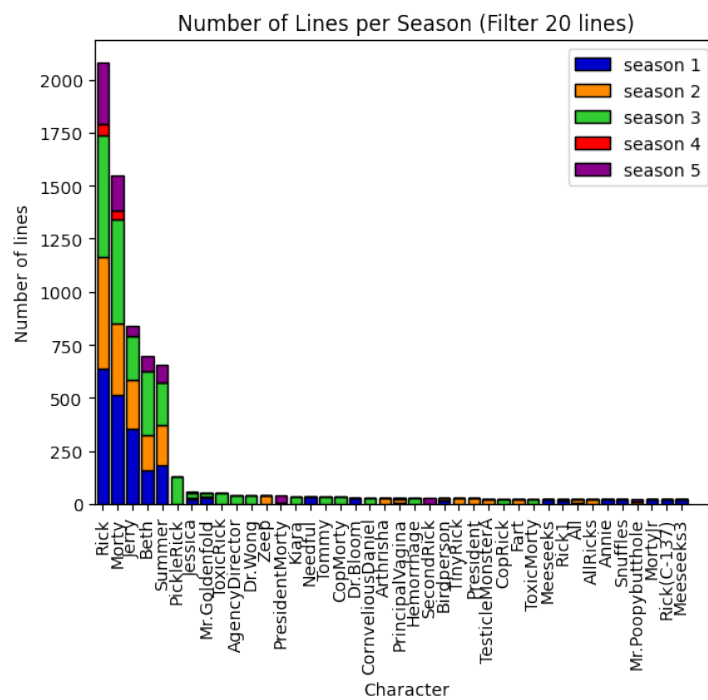
The last datasource stores information about the average IMDB rating per episode in the format provided in sample 2.3 where we also find the related season and title for each episode.

id	season	episode	title	rating
0	S1	E1	Pilot	7.9
1	S1	E2	Lawnmower Dog	8.6

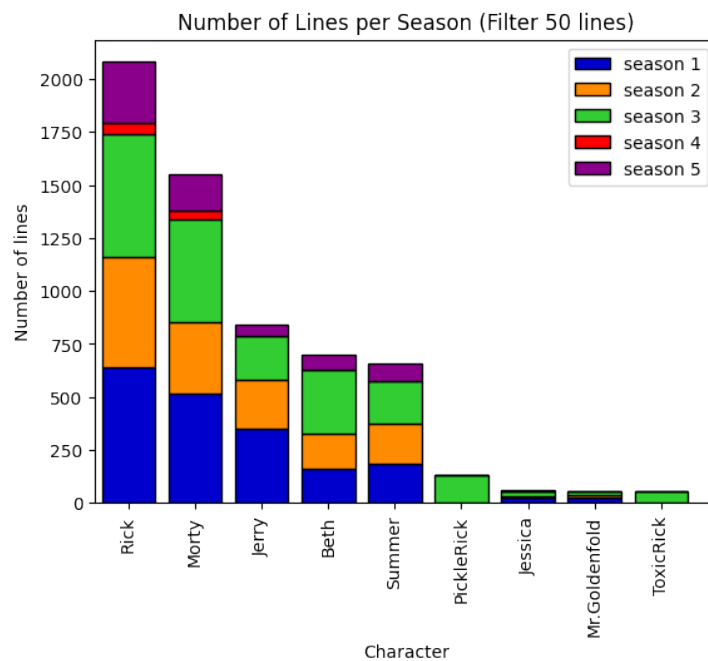
**Tab. 2.3:** IMDB ratings per episode dataset

The rating score provided in the table is the average rating score on IMDB cacludated with every rating that was given to a episode. For the last two datasets we not just have the data until season 5 but we have every single episode until season 7 maintained.

Further analysis show that there are in general 970 different speakers. The charts 2.1 visualize the speaker distribution by counting the lines that each character speaks in each season. In general, there are a lot of different characters participating in the series but the family Smith (Rick, Morty, Jerry, Beth, Summer) have by far the largest share when looking at the speaker distribution.



(a) Speaker distribution of speakers with more than 20 lines

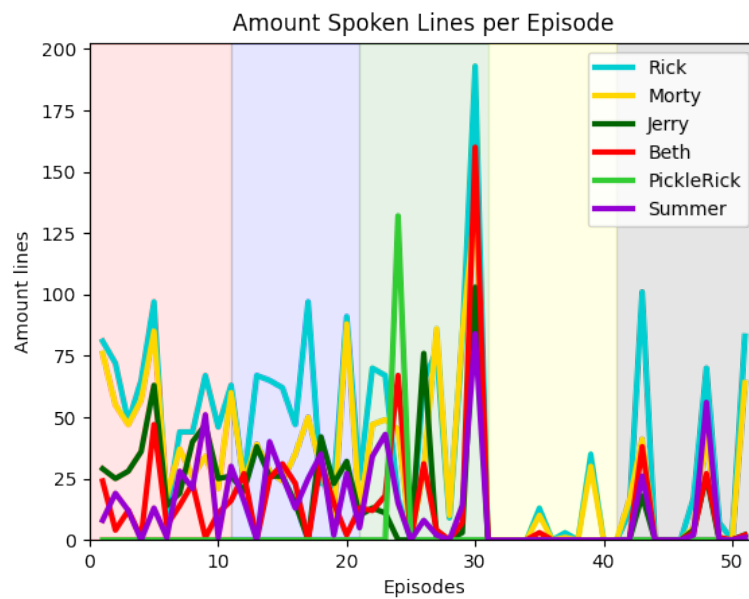


(b) Speaker distribution of speakers with more than 50 lines

**Abb. 2.1:** Speaker distribution per lines

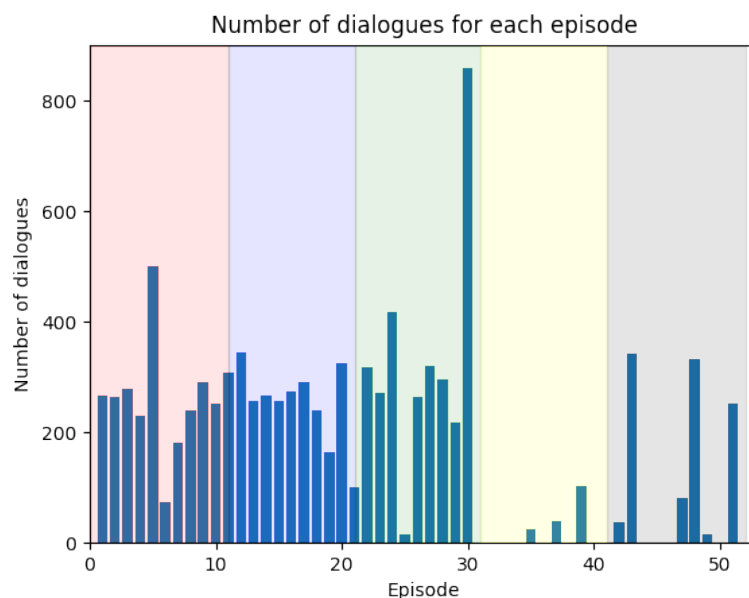
Surprisingly, the transcript dataset has some lacks. As we see in chart 2.3 that shows the number of dialogues per episode, the dataset is incomplete as there are a lot of empty episode starting from season 4. The background colors highlight the intervals of the different seasons along the whole series.

The chart 2.2 provides a clearer picture on how often which main charatcer speaks in detail.



**Abb. 2.2:** Family Sanchez spoken lines in detail

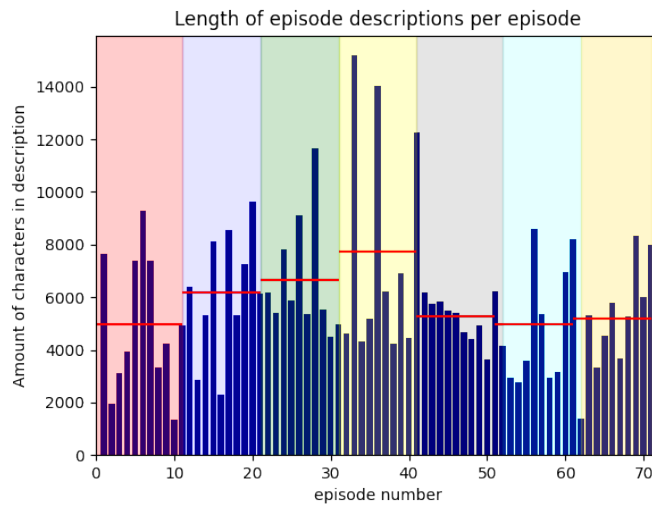
In most of the episodes, Rick is the main character and holds the most shares in the speaking distribution. The other family members (Morty, Jerry, Beth and Summer) also speak in every episode but not that often as Rick. Side characters, like Pickle Rick often just have a large share in a few episodes as they often disappear after one episode.



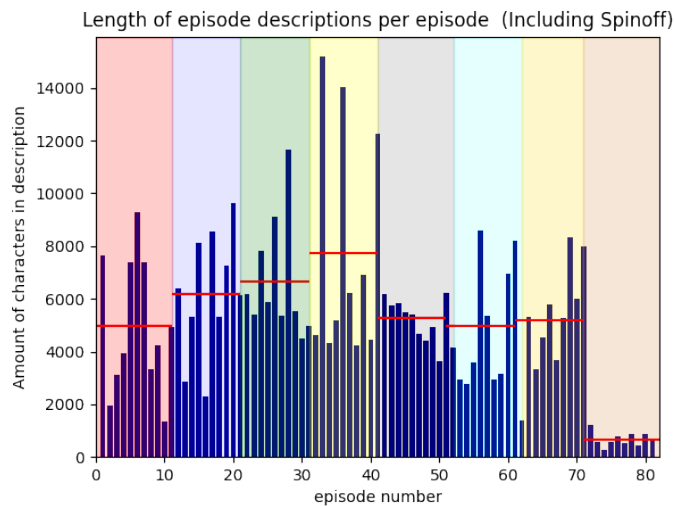
**Abb. 2.3:** Episode distribution

Looking deeply at the description dataset, we see two charts in 2.4 which compare the length the episode description for each episode. Surprisingly, the descriptions from episode 72 until 81 have a way shorter description than the episodes before. These episodes do not belong to

the normal Rick and Morty series as they are spin-off episodes. Therefore, we now exclude those episode in further analysis. The red lines show the average description length per each season. We see that the description length varies a lot for in each season. An example can be the comparison between episode 33 and 34, where episode 33's descriptions contains more than 3 times the amount of characters as the description of episode 34. Another fact that is that starting from season 5 the average number of chars in the description seems to stay on a similar level. That is different compared to the first 4 seasons where in each season the average length of the description increased.



(a) Episode Length description distribution



(b) Episode Length description distribution (Including Spinoffs)

**Abb. 2.4:** Length of episode descriptions per episode

The last analysis step was an analysis regarding the top 100 most spoken words in the transcript. We used a wordcloud to visualize those in an appropriate way. The results can be seen in 2.5. After removing the character names, we find there colloquial language terms like 'ye-

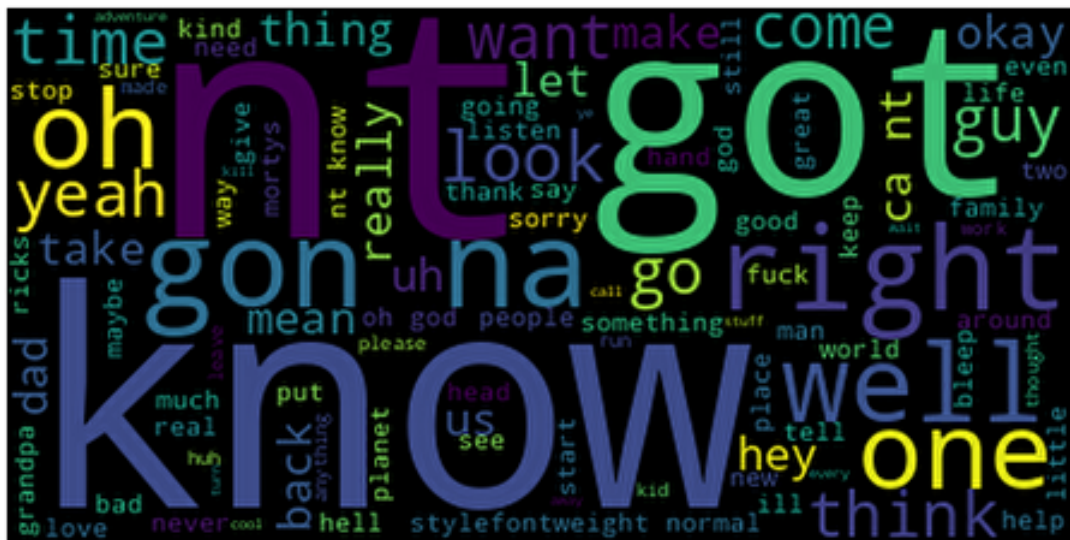


Abb. 2.5: Word Cloud transcript

ah', 'gonna' in the wordcloud. That is why we can assume that the spoken english is not that advanced.

## 2.2 Data Preparation

Here, we describe the procedure we followed to clean and prepare the data. To clean and unify the data, we removed all nonspace characters, and set every character to a lower character. As the rick and morty characters talk to each other with colloquial language, we removed some contractions like 'ain't' with 'is not'. Furthermore, there are some HTML tags in the dataset that were removed in the data preparation phase as well.

Using the spacy and nltk library, the text preparation pipeline consists of a custom build stopwordremover and a custom build stemmer which stems every single token in a sentence but ignores the character names for stemming. Doing this, we used the Porterstemmer. Before that, the sentences were tokenized using the small en web model spacy provides. To see the results, the first dialogue datapoint turns from

*stumbles in drunkenly, and turns on the lights. Morty! You gotta come on. Jus'... you gotta come with me.*

into

*stumbl drunkenli turn light morty got tocom jus got come*

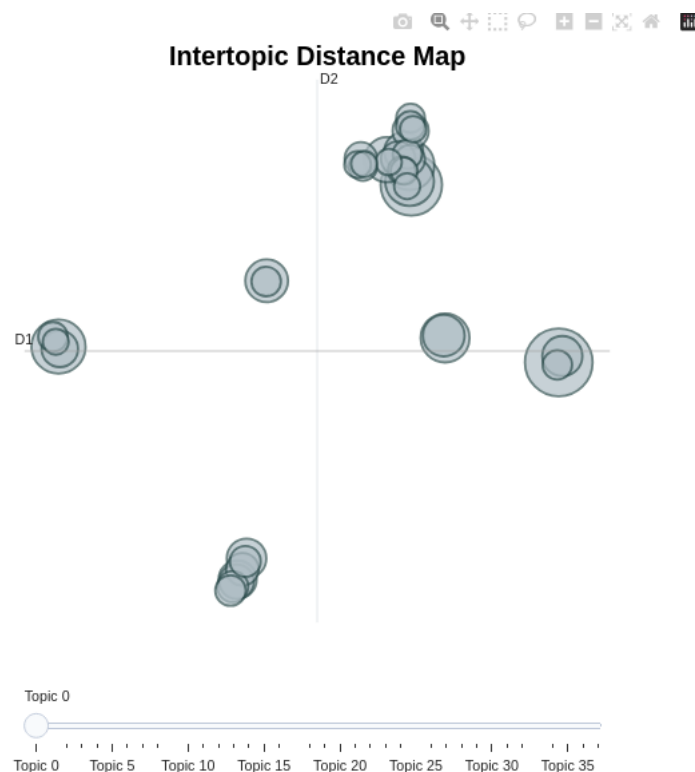
## 2.3 Topic Modeling

In this section we want to examine the Rick and Morty Series regarding the topics addressed in the first five seasons. To achieve this, we are looking at the transcripts and all of the episode

descriptions, where we excluded the names of the main characters as they are not really an topic that holds content.

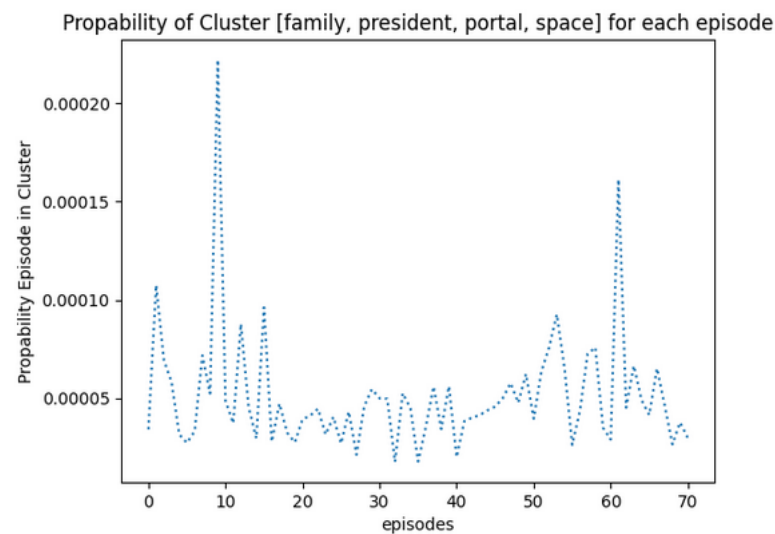
### Topic Modeling Transcript

At first, the topics were analysed based on the transcript in the first 5 seasons. As the dataset contains roughly 9.000 datapoints, we use BERTopic Model to find out the most relevant topics. For using the BERT Topic model, we downloaded the all-MiniLM-L6v2 embedding model to embed the text. The embedded text was then clustered by the hierarchical HDBscan cluster algorithm using the eom cluster selection method. The resulting topics showcased in the image 2.6.



**Abb. 2.6:** BERT Topics for Transcript (all seasons)

With the minimum topic size is set to 40 words, the model created 37 different topics. As in the figure provided, we see that most of the bubbles overlap. Therefore there would be the possibility to decrease the number of topics considering that this would lead to more mature topics. Analysing the content of the topics, we get on the left side topics that can be summarized with the keyword family. There we find family related terms like dad, mom, grandpa and so on. The topics located on the bottom of the image mostly refer to general information about the adventures that Rick and Morty perform as there are terms like adventure, portal gun and treasure. Especially the word portal gun is a very important word as it is a key component of the series. Most of the other topics can be summarized by the speaking habits each character has. As there are topics where there are slang words like 'wubabubadabab', 'geez', 'oh' or 'crap'. Surprisingly there is one topic located next to those slang words containing information about



**Abb. 2.7:** Topic relevance Family Portal Gun over time

the planet system which contains words like pluto, planet and space.

In 5.1, we see a visualization of the most important topics and its words. It turns out, that the most important topic are slang words like 'oh' or 'man'. Topic 2 and topic 4 show that the series is about a human family. Topic 0, 11 and 16 contain more information about the content of the series as there are a lot of terms referring to space and murdering activities that Rick and Morty are exploring along their adventures in space.

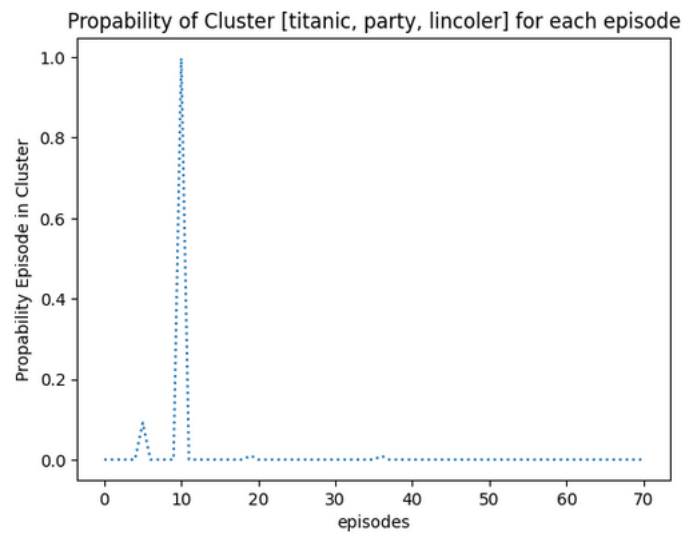
### Topic Modeling based on Episode Descriptions

As the dialouges of the series are full of abbreviations and colloquial language without a lot of content, we wanted to have a closer look at the topics by looking at all episode descriptions. This time we analysed the topics with Latent Dirichlet Allocation (LDA) as there are way less datapoints, that are in average longer than in the transcription dataset. To setup our LDA model we create an Document Term Matrix (DTM) containing all words and documents. Another dictionary maps the words with a given id. Furthermore, we limit the amount of topics to 50 so that there will be less topics than the number of episodes. In the series, the episodes do not depend on each other as Rick and Morty often travel to new planet each episode. The following topic analysis can lead to the same conclusion as a lot of the created topics just contain information about the content of one or a few episodeds. In the chart 2.7 we see the relevance of the topic containing family and portal gun along all episodes.

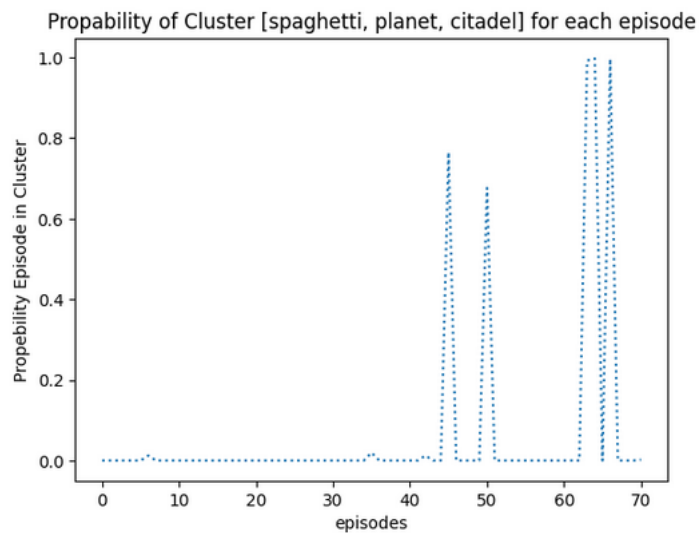
This topic plays a significant role in all episodes as 'family' and 'portal gun' are terms that are frequently used in a lot of episodes throughout the series.

As already mentioned, most of the other topics just refer to a few or one single episode. This is the case for the topic titanic party where jerry and beth attend a party that recreates the titanic drama. As the chart 2.8 shows, this topic holds information that is is mainly present in one episode.

The other spikes that are showcased in the chart can be explained as the side character 'Abra-



**Abb. 2.8:** Topic relevance Titanic Topic over time



**Abb. 2.9:** Spaghetti an Citadel over time

dolf' 'Lincolner' plays a important role in other episodes as well.

Another example for a similar distribution is the topic about the spaghetti planet and its citadel, which is analysed in chart 2.9.

The fact that most of the topics differ a lot from each other, can also be seen here 5.3. It is remarkable that a vast amount of the topics are distributed across the chart while just a handful topics seem to overlap. Furthermore, the size of the bubbles indicate the amount of documents where the topic is part of. As a lot of the bubbles seem small this could mean that there the topics mostly relate to a few episode descriptions.



## 3 Modeling

### 3.1 IMDB Rating prediction

In this section our goal was to perform a modeling to predict the IMDB ratings based on the provided episode descriptions. Therefore, we compare different strategies and architectures to solve this complex problem.

Our hypothesis is that there are some important words or characters like 'birdperson' that might have an impact on the IMDB rating as some characters or planets in the Rick and Morty universe are more popular than other.

Creating a perfectly fitting model is challenging, as the episode descriptions are written in neutral style and the features that determine whether an episode is liked or not can differ from person to person. Furthermore, there are also visual effects and music elements such as the 'get Schwifty' or snake jazz song that might lead to a higher IMDB rating. Therefore, we expect that our model will tend to not be the most accurate one.

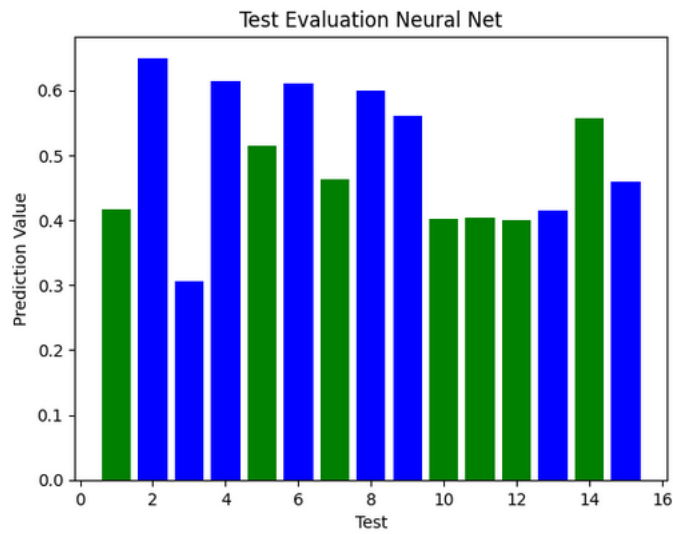
#### Vector Embedding Approach

As the neural network cannot calculate with strings, we have to represent the given text as numeric representations. To achieve this, our initial approach used the Word2Vec library to generate a vector embedding for all the words that appear at least 5 times in all descriptions. By that, we were able to represent most words in the Rick and Morty corpus as a vector of 100 values.

These 100 values also defined the input shape of our neural network which was designed for classifying the episode descriptions. To represent whole tokenized sentences into a vector containing 100 dimensions, we tokenized all of the training data and calculated the average vector of each episode description. As a result, each episode description can be represented as one single vector of 100 dimensions. With that, we are able to train a neural network with 100 input neurons. If a token was not part of the Word2vec vector embedding list, the token will be skipped. For example, the string 'DHBW Ravensburg' would produce a NaN and the string 'Rick DHBW Ravensburg' would produce the same result as the string 'Rick', since the unknown words like DHBW are ignored.

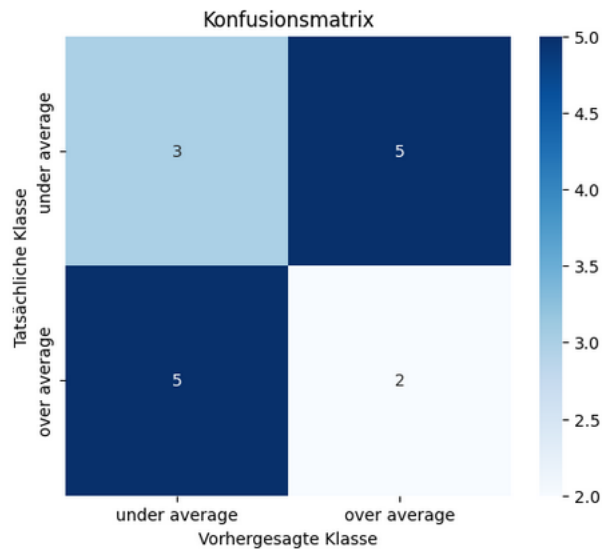
The test results are displayed in the following chart 3.1. The green bars represent the episodes having an IMDB rating above 8.2 while the blue ones are rated lower. As we did not set a threshold yet, we can see the test result visualized in a bar chart.

As expected, the neural network struggles to distinguish between those two classes. In general, it seems that the models even predicts the complete opposite, as the blue colored bars were predicted higher as the the green ones. The calculated test accuracy based on the results without an threshold was 33%.



**Abb. 3.1:** Neural Net Test Evaluation

Setting the threshold to 0.5 we find following confusion matrix as a result 3.2. The confusion matrix also leads to similar results as there are more wrong predicted values than right ones. Letting us know that this model definitely does not solve the problem.



**Abb. 3.2:** Neural Net confusion matrix (threshold = 0.5)

### LSTM Approach

As the first neural network had some difficulties to examine the complex relations between the words, we also implemented the Long Short Term Memory (LSTM) architecture to achieve more reliable results. Here, we used the same Word2vec embedding as before. We again used Word2vec embedding model and added our Rick and Morty corpus to it due to the fact that there are a lot of Rick and Morty unique terms like 'birdperson', 'portal gun' in the descriptions. Unlike the previous neural network, which computed the average vector for each sentence, we instead used the first 500 vectors in a sentence to numerically represent the entire text. The results of this approach are displayed in the chart 3.3.

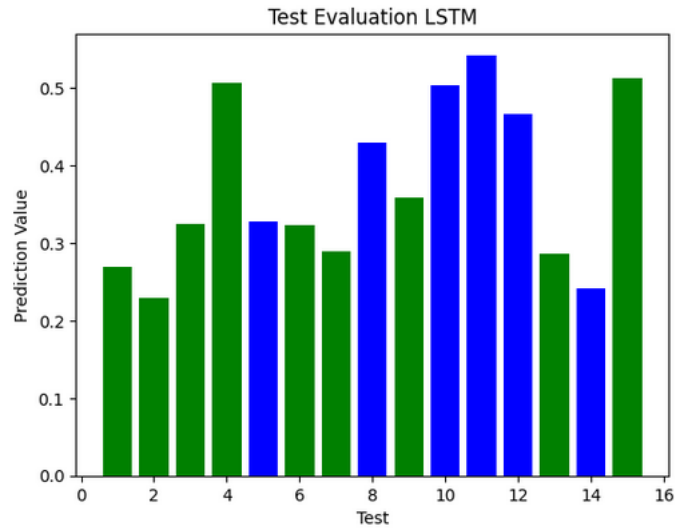


Abb. 3.3: LSTM Classification

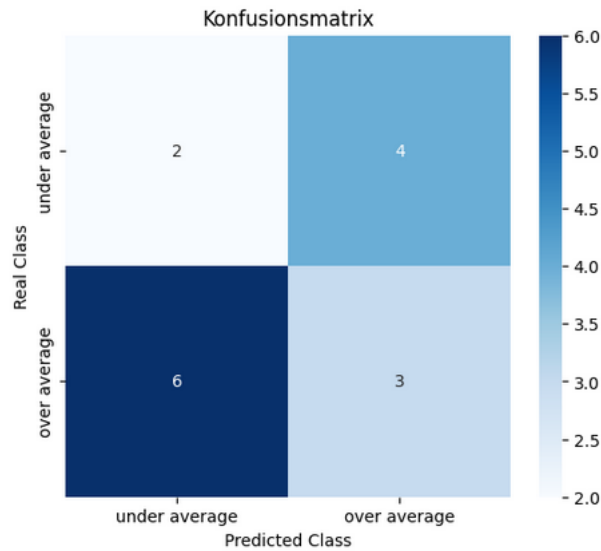


Abb. 3.4: LSTM Classification confusion matrix (treshhold = 0.4)

As the LSTM Classification leads to more variety in text prediction, it still had some problems figuring out which texts are above and under 8.2 rating. As the confusion matrix displayed in 3.4 shows, the model is also not good in predicting the classes as it produces more False predicted values as right ones.

### Bag of Words Approach

As the word embedding approaches did not lead to a satisfying outcome, we also tried to create a Bag of word matrix containing the the amount of the first 500 most appearing words within a text. To clean the data, we decided to remove every token in the description dataset which POS was not classified to PRPN or NOUN by spacy, in order to focus more on the character names and the nouns that are provided in the description. With that, we trained a Naive Bayes model to predict those classes. Therefore, the CountVectorizer by sklearn came into place. But

unfortunately, also this approach could not return valid result as the confusion matrix 5.4 shows.

### **TF-IDF Matrix Approach**

The TF-IDF Vectorizer by Sklearn prioritizes words that appear a lot in one document but not accross the documents. With that last classification approach, we hope to weigh terms or characters more that do not appear that frequently but are more popular than others. Following this procedure of creating an TF-IDF matrix, and training an Gaussian Model again, it also did not lead to satisfying results 5.5.

### **LSTM Regression**

We expect that features from an episode with a rating of 8.2 (class 0) does not differ so much from another episode with rating of 8.3 (class 1). As the correlation coefficent between the numeric IMDB rating value and the class from the classification is just 0.7, some information is definetly lost in the classification problem. Therefore, we aim to predict the IMDB rating not by classifying the episodes as 1 or 0 but by directly predicting it's IMDB score. With this regression approach, we hope to get better chances to identify the IMDB ratings as a lot of information is lost by transferring the numeric data into a class in the classification above.

To achieve this, we used the same architecture as in the classification task but changed the lossfunction as the LSTM classifier but we additionally normalised the Word2vec embedding vectors in the matrix.

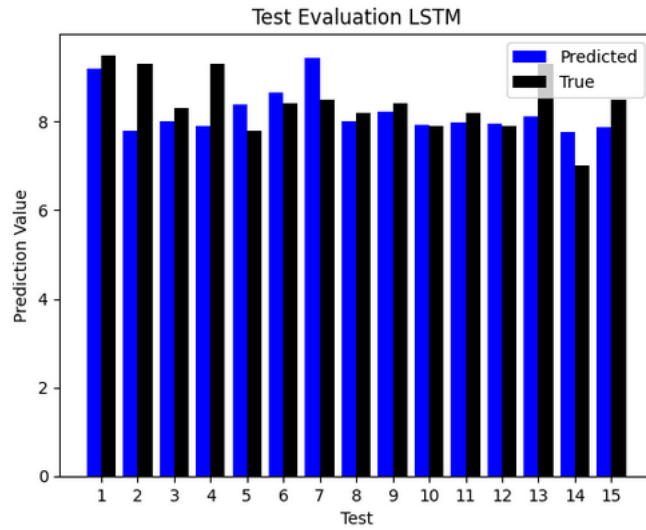
In 5.3 we see the test accuracies by epoch size. As the epoche size of 15 leads to the lowest sqaured error, we use 45 as a default epoch size.

Another Cross Validation based on models of epoch size 45 leads to an average squared error of 0.47 5.1, meaning that, on average, our model's predictions deviate by approximately 0.5 rating points from the true IMDB scores. As a comparison, the statistical standart deviation of all ratings is 0.97. An analysis of the test results is displayed in the chart 3.5. There we, see the predicted IMDB ratings in blue and the true ones in black. In general, we see that the model predicts valid results, as none of the predicted IMDB ratings were above 10 and ofthen the deviation between true and predicted values are low.

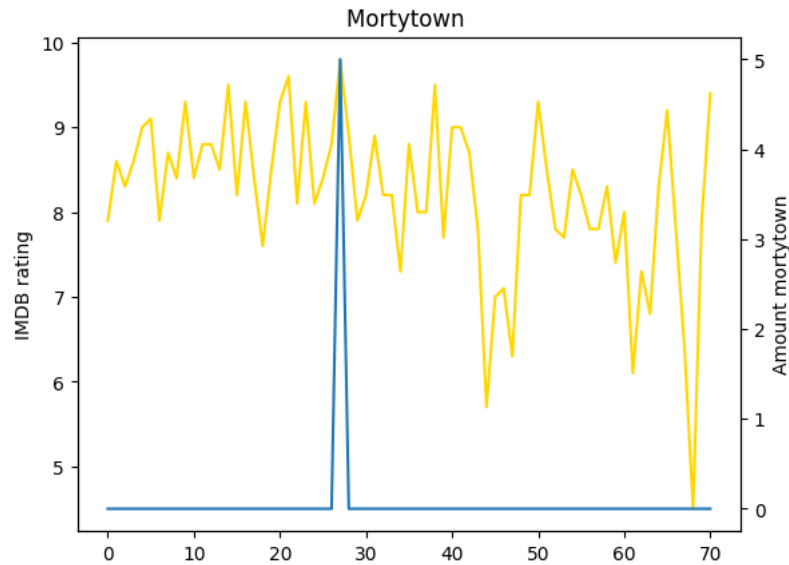
When testing the model, we noticed that in general, the combination and the appearance of characters influences the IMDB score. For example the string "rick is on an adventure"leads to an 7.2 while "rick and morty are on an adventure"leads to 10.1. Handing over "rick and beth are on an adventure"leads to an output of 9.3. "rick and birdperson are on an adventure"leads for an example to a score of 6.6. The sentence Morty and Rick, for example also just produce a IMDB rating of 8.3 while calculates for the sentence "Rick and Mortyän IMDB rating of 9.5. The problem that neural networks have is that they are not explainable and is more like a black box. In the last approach we want to figure out more about the terms and its influence in other model architectures.

### **Regression with TF IDF**

The linear regression model with TF IDF matrix as features should examine which character



**Abb. 3.5:** Test Results Regression



**Abb. 3.6:** Number of Mortytown apperances with IMDB prediction

or nouns define the IMDB Prediction. To achieve this, we trained a linear OLS regression model by statsmodelapi and printed out the characters that determine the IMDB prediction the most based on the model. According to this model, the feature term president has the highest coefficient which terms like family, saber and universe or mortytown. The terms with the lowest coefficient are destruction, garage and gun. The word president appears in the episode 16 and 18 and 27 that all are under the top 5 rated rick and morty episodes. A more impactful example is the term mortytown that just appeared in episode 27, the highest rated episode as visible in 3.6.

With all these different approaches we were not able to create a model which classifies the IMDB prediction 100% correctly. Not just the vector embedding approach leads to a unsatisfying result in the classification architectures, but also the Term matrix approach with BOW and TF IDF could not distinguish between high and low IMDB prediction. A reason for that is the

lack of a huge datasource as there are just 70 Rick and morty episodes. Furthermore, most of the characters and planets appear just in a few episodes. While terms in the trainingset like Mortytown correctly determine IMDB ratings in the TF IDF approach, testing the model on new vocabulary will definitely not produce valid results.

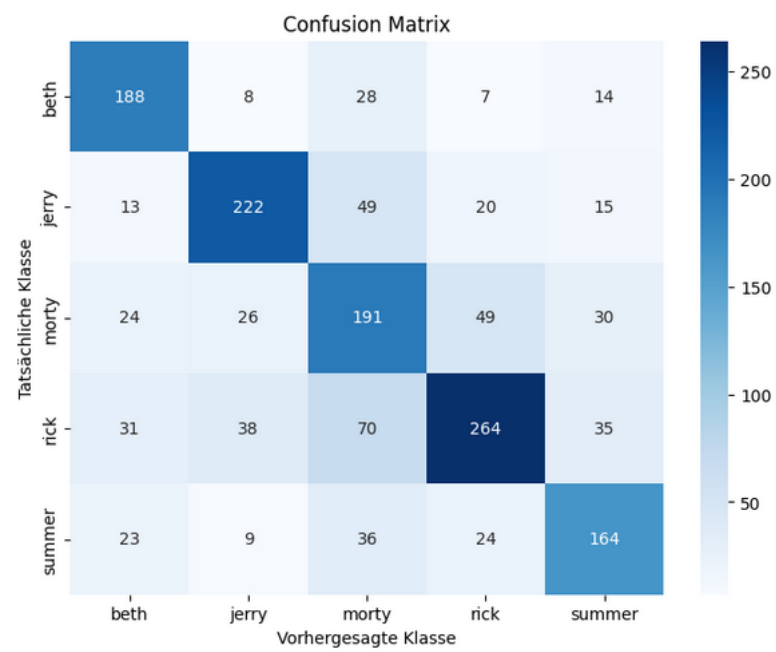
With the regression model we got rid of information losses that were caused due to changing the label format from numeric to qualitative. With that, we examined that characters and its appearance and combination might have a lot of influence on the IMDB prediction and we were able to figure out terms like mortytown or president, which influence the rating the most according to a linear regression model.

### 3.2 Speaker detection

The last modeling part examines if the character's speaking traits can be identified by a fine-tuned transformer model. Therefore, we just use the family\_sanchez as classes because the other characters talk way less in the episodes as already mentioned in 2.1. As already seen in chart 2.1, the dataset is imbalanced which can cause some problems. In order to get rid of this imbalanced form, we decided to duplicate the amount of spoken lines of Jerry, Beth and Summer for the modeling part.

For this classification problem, we used the DistilBERT transformer as a base model consisting of the DistilBERT Tokenizer and a Sequence classifier. As the test error still decreases with 3 epochs, we fine-tuned the model with 3 epochs. Doing this, we created a model that had a 65% accuracy on the test set. There we have to keep in mind that there might be duplicates in the test set due to the duplication step that was explained in the text above.

For example the Input "Burp" leads to following output (Rick: 3.9, Summer: -0.5, Morty: -1.4, Jerry: -1.6, Beth: -1.9). There, the model was sure that probably Rick would have said this word in a conversation. This is also realistic as he often needs to burp after drinking alcoholic beverages. The overall results can be seen in following confusion matrix 3.7 and interpreted with the metrics in table 5.2.

**Abb. 3.7:** Confusion Matrix speaker detection

## 4 Diskussion

Summarizing the topic modeling, we were able to extract relevant content that plays a significant role in the series Rick and Morty. Topics like family, or space describing topics were addressed through the whole series. Looking at the topics into more detail, it becomes clear that most of the topics just relate to one episode leading to the conclusion that every episode is hardly related to other episodes.

Regarding the IMDB prediction, our hypothesis that it is very difficult to predict IMDB ratings based on description texts, turned out to be true. Every classification approach had some difficulties to predict whether the episode will have a high or low IMDB prediction. Furthermore, the regression approach led to valid results and had a relatively low standard deviation for the test set.

Based on some experiments, the model learned which character combination has a potential higher IMDB rating than other. This might lead to the hypothesis, that frequent combinations like rick and morty or rick and jerry produce a higher IMDB rating than ones with lower term frequency like rick and birdperson. The model also learned, that more characters might lead to a higher IMDB ratings as well, as the sentence same sentence without any other character produced a lower IMDB rating.

With finetuning the DistilBERT transformer, we were able to classify the speakers with a test accuracy of 60% what could lead to the conclusion that the speaking behaviors of each character differ from each other.



5 Anhang

mean average error  
0.4996431767940521,  
0.4830387234687805,  
0.5242437124252319,  
0.4312141239643097,  
0.49884870648384094

Tab. 5.1: 5 Fold Cross Validation

	precision	recall	f1-score	support
beth	0.67	0.77	0.72	245
jerry	0.73	0.70	0.71	319
morty	0.51	0.60	0.55	320
rick	0.73	0.60	0.66	438
summer	0.64	0.64	0.64	256
accuracy	0.65	1578		
macro avg	0.66	0.66	0.66	1578
weighted avg	0.66	0.65	0.65	1578

Tab. 5.2: Metrics for speaker detection

epoch	Mean sqaured Error
15	0.45
25	0.479
35	0.51
45	0.43
55	0.50

Tab. 5.3: Test scores accuracy LSTM regression

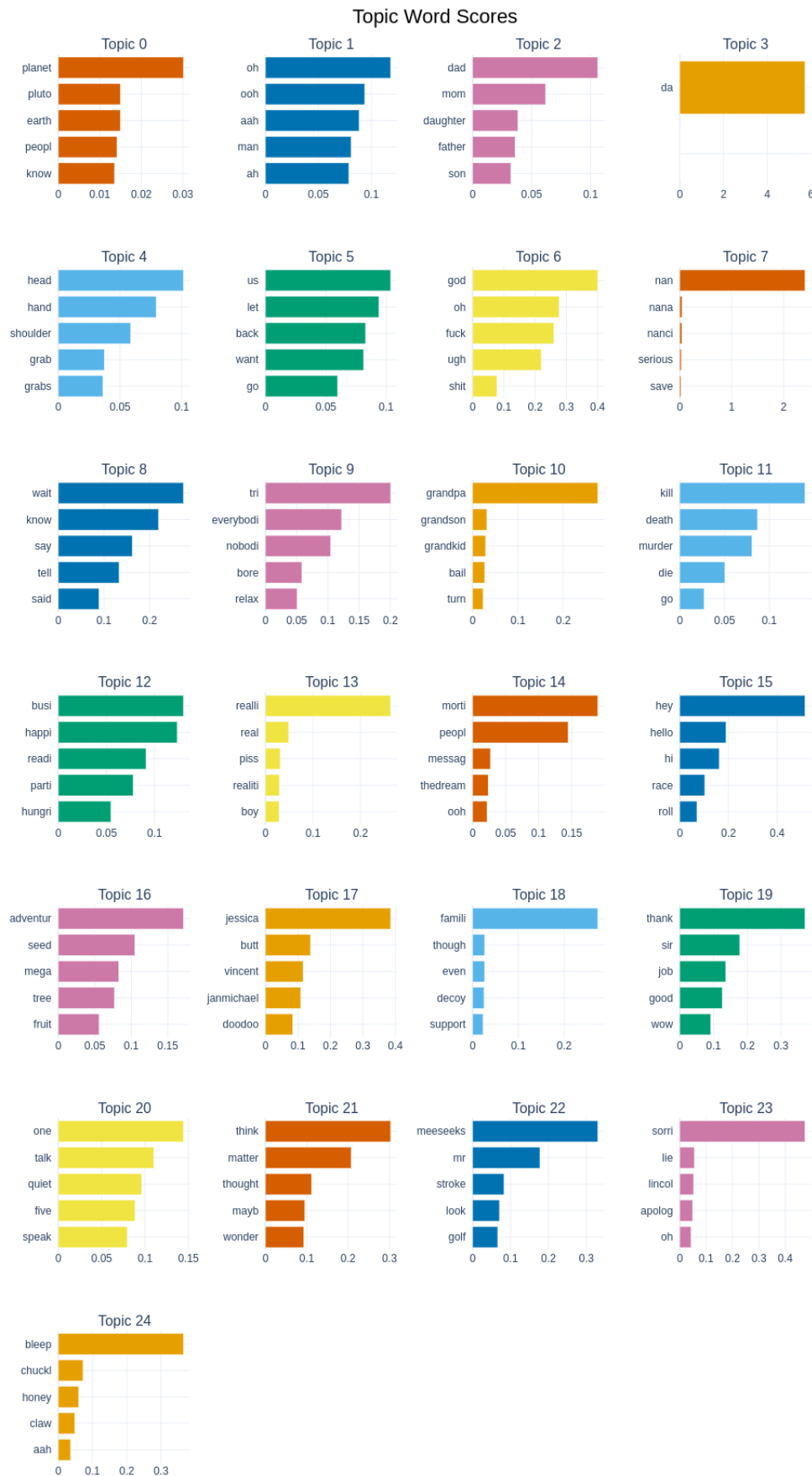


Abb. 5.1: Top 25 BERT topics for transcript (all seasons)

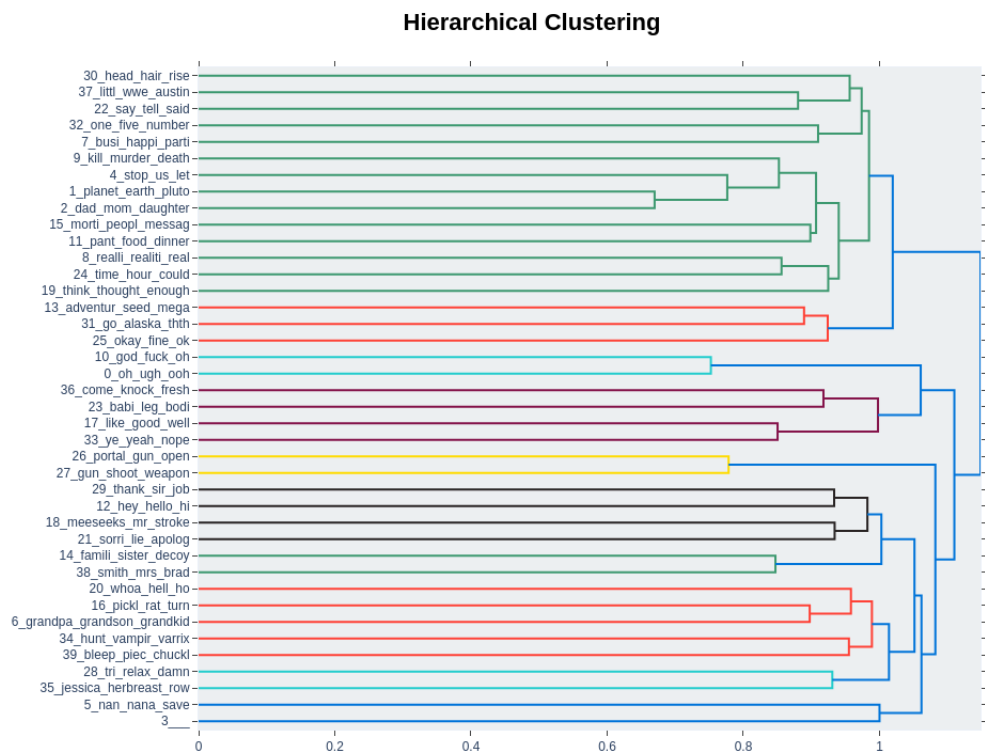


Abb. 5.2: Hierarchial Clustering topics on transcript

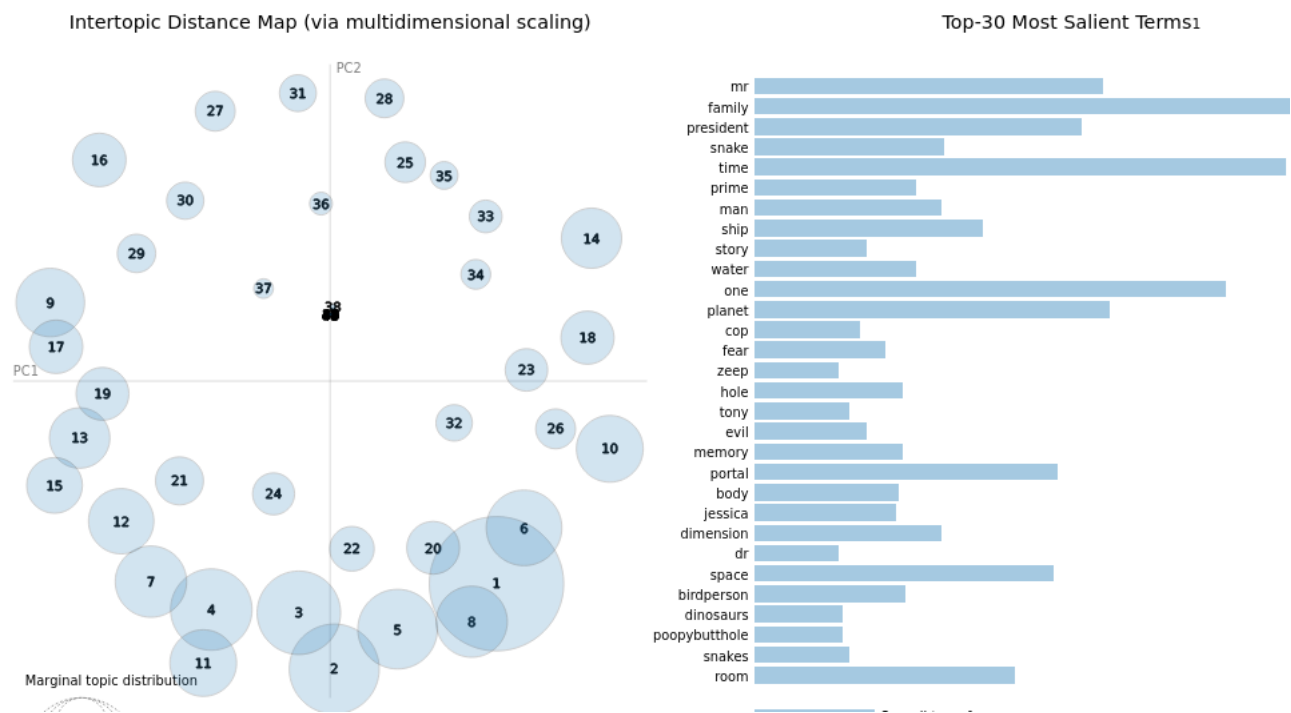
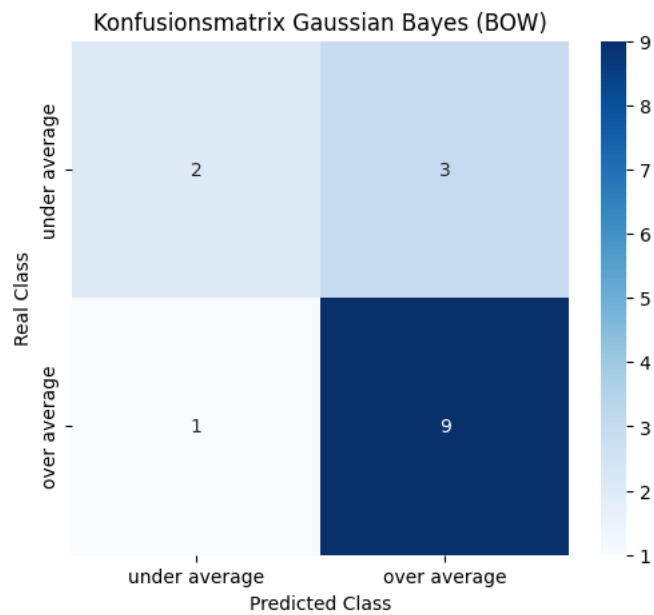
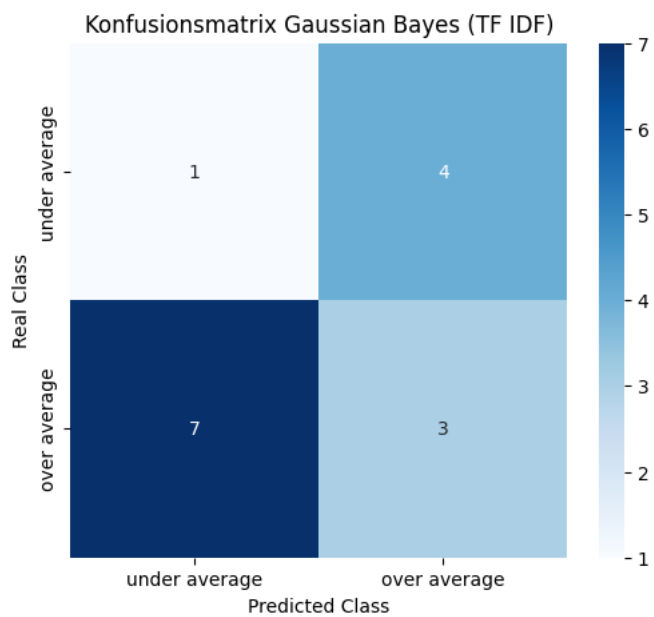


Abb. 5.3: All Topics LDA for description (all seasons)

**Abb. 5.4:** Confusion Matrix Bayes Classification with BOW**Abb. 5.5:** Confusion Matrix Bayes Classification with TF IDF

## **Literatur**

# **Selbständigkeitserklärung Anton Geiger**

Ich versichere hiermit, dass ich, Anton Geiger, meine Seminararbeit

## **Burp NLP: A Rick and Morty text analysis**

selbständig verfasst, die digitale Version mit der ausgedruckten  
Version übereinstimmt und keine anderen als die angegebenen  
Quellen und Hilfsmittel benutzt habe.

---

Ort, Datum

---

Unterschrift