

Overview Generation

```
In [26]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import os
5 from pickle import load, dump
6 import numpy as np
7 from afinn import Afinn
8 import warnings
9 warnings.filterwarnings('ignore')
10 import pylDAvis
11 pylDAvis.enable_notebook()
12 from mv_lda_utils import mv_lda_model

In [14]: 1 def get_all_scripts():
2     clean_data_folder = "../clean_data"
3     filenames = os.listdir(clean_data_folder)
4     os.listdir(clean_data_folder)
5
6     with open(clean_data_folder+"/all_data.pkl", "rb") as f:
7         all_scripts = load(f)
8     all_scripts['episode_str']=all_scripts['episode'].apply(lambda x: '0'+str(x) if len(str(x))<1 else str(x))
9     all_scripts['season_episode_no'] = all_scripts['season'].apply(lambda x: str(x))+all_scripts['episode_str']
10    all_scripts['season_episode_no'] = all_scripts['season_episode_no'].apply(lambda x: int(x))
11    return all_scripts

In [15]: 1 def extract_specific_episode(all_scripts, season_no, episode_no):
2     episode_script=all_scripts[all_scripts.season==season_no]
3     episode_script=all_scripts[all_scripts.episode==episode_no]
4
5     if episode_script.empty: #saveguard if episode or season does not exist
6         print("Season or episode does not exist.")
7         return
8     return episode_script

In [16]: 1 def list_characters(script):
2     characters = script["character"].unique()
3
4     print(f"There are {len(characters)} characters in the season.")
5     text ="These characters are "
6     i=0
7     for name in characters:
8         if i + 1 == len(characters):
9             text= text +"and "+ name + "."
10        else:
11            text= text + name + ", "
12        i=i+1
13
14    print(text+"\n\n")
15    return

In [17]: 1 def show_three_characters_with_biggest_amount_of_speech(script):
2     characters_word_count = script[["character", "word_count"]].groupby("character",as_index=False).sum()
3     characters_word_count=characters_word_count.sort_values(by='word_count', ascending=False)
4     characters_word_count=characters_word_count.reset_index(drop=True)
5
6     top_3 = characters_word_count['character'].head(3).values.tolist()
7     top_3_text ="The 3 with the highest word counts are "
8     i=0
9     for name in top_3:
10        if i + 1 == len(top_3):
11            top_3_text= top_3_text +"and "+ name + "."
12        else:
13            top_3_text= top_3_text + name + ", "
14        i=i+1
15
16    print(top_3_text+"\n\n")
17    return top_3

In [18]: 1 def show_amount_of_dialogue(script):
2     dialogue= script.word_count.sum()
3     dialogue_text= "In total, there are "+str(dialogue)+" words spoken."
4     print(dialogue_text)
5     return

In [19]: 1 def sentiment(script):
2
3     afn = Afinn(emoticons=True)
4     afinn_wl_url = ('https://raw.githubusercontent.com'
5                    '/fnielsen/afinn/master/afinn/data/AFINN-111.txt')
6     afinn_wl_df = pd.read_csv(afinn_wl_url,
7                              header=None, # no column names
8                              sep='\t', # tab sepeated
9                              names=['term', 'value']) #new column names
10    seed = 808 # seed for sample so results are stable
11    afinn_wl_df.sample(15, random_state = seed)
12
13    script['text']=script.text.apply(lambda x: ' '.join(x))
14    script['sentiment_score']=script.text.apply(lambda x: afn.score(x))
15    script=script.reset_index(drop=True)
16    avg_sentiment=script.sentiment_score.mean()
17
18    print(f"The sentiment score over the episode/season has an average of {round(avg_sentiment,4)}.\nThe graph below shows the sentiment for each line in chronological order.")
19    g = sns.relplot(
20        data=script,
21        x=script.index, y="sentiment_score", kind="line",
22        height=5, aspect=5, color="blue"
23    ).set(
24        title="Sentiments over the episode or season",
25        ylabel="sentiment score",
26        xlabel="line in script"
27    )
28    g.despine(left=True)
29
30    ax1, = g.axes[0]
31
32    ax1.axhline(avg_sentiment, ls='--', c="red")
33    plt.legend(labels=["Sentiment", "Average"])
34    plt.xlim([0, script.shape[0]])
35    plt.show()
36    return
```

```

In [20]: 1 def timeline_with_amount_of_speech_for_top3_characters(season_script, top_3):
2         print("The amount of speech each of the top 3 characters has per episode, is as follows:")
3
4         amount_speech_top1 = season_script[season_script.character==top_3[0]].groupby("episode").sum().reset_index()
5         amount_speech_top1=amount_speech_top1[["episode","word_count"]]
6         amount_speech_top1=amount_speech_top1.rename(columns={"word_count": top_3[0]})
7         amount_speech_top1['episode'] = amount_speech_top1['episode'].apply(lambda x: str(x))
8
9         amount_speech_top2 = season_script[season_script.character==top_3[1]].groupby("episode").sum().reset_index()
10        amount_speech_top2=amount_speech_top2[["episode","word_count"]]
11        amount_speech_top2=amount_speech_top2.rename(columns={"word_count": top_3[1]})
12        amount_speech_top2['episode'] = amount_speech_top2['episode'].apply(lambda x: str(x))
13
14        amount_speech_top3 = season_script[season_script.character==top_3[2]].groupby("episode").sum().reset_index()
15        amount_speech_top3=amount_speech_top3[["episode","word_count"]]
16        amount_speech_top3=amount_speech_top3.rename(columns={"word_count": top_3[2]})
17        amount_speech_top3['episode'] = amount_speech_top3['episode'].apply(lambda x: str(x))
18
19        amount_speech_top1[top_3[1]]=amount_speech_top2[top_3[1]]
20        amount_speech_top1[top_3[2]]=amount_speech_top3[top_3[2]]
21
22        amount_speech_top1.plot(x='episode', y=[top_3[0],top_3[1],top_3[2]], figsize=(21,5), xlabel='Episode', ylabel='Word count', title='Word count of the top 3 characters per episode');
23        plt.show()
24        return

```

```

In [21]: 1 def show_topics(scripts):
2
3         #extract the texts and group them together
4         texts = list(scripts["text"])#sum(List(scripts["text"]), [])
5
6         # get the topic model and dictionary
7         with open ("lda_topic_model-5_topics-5_passes-0.00_tfidf_threshold.pkl", "rb") as f:
8             my_lda_topic_model = load(f)
9         model = my_lda_topic_model.lda_model
10        dictionary = my_lda_topic_model.dictionary
11
12        # get the topics
13        bow = dictionary.doc2bow(texts)
14        topics = [{0: "Negotiation",
15                  1: "War",
16                  2: "Faith",
17                  3: "Sci-fi and injury",
18                  4: "Fighting"}
19
20        text_topics = (model.get_document_topics(bow,
21                                                  minimum_probability=0.0))
22        df = pd.DataFrame(text_topics)
23        df =df.replace(topics)
24        df.columns = ["Topic", "Percentage"]
25        # df =df.set_index("Topic")
26
27
28        labels = topics.values
29        colors = sns.color_palette('pastel')[0:5]
30        plt.pie(df.Percentage, labels=df.Topic, colors = colors, autopct='%.0f%%')
31        plt.title("Topic Percentage")
32        plt.show()
33        return

```

```

In [22]: 1 def generate_overview_season(season_no):
2
3         all_scripts= get_all_scripts()
4         season_script=all_scripts[all_scripts.season==season_no]
5
6         if season_script.empty: #saveguard if episode or season does not exist
7             print("Season does not exist.")
8             return
9         print('\033[1m'+ "Shown is an overview over season ", season_no, "\n")
10        print('\033[0m')
11
12        list_characters(season_script)
13        top_3= show_three_characters_with_biggest_amount_of_speech(season_script)
14        show_amount_of_dialogue(season_script)
15        sentiment(season_script)
16        timeline_with_amount_of_speech_for_top3_characters(season_script, top_3)
17        show_topics(season_script)
18        return

```

```

In [23]: 1 def generate_overview_episode(season_no,episode_no):
2
3         all_scripts= get_all_scripts()
4         episode_script=all_scripts[(all_scripts.season==season_no) & (all_scripts.episode==episode_no)]
5         # episode_script=episode_script[all_scripts.episode==episode_no]
6
7         if episode_script.empty: #saveguard if episode or season does not exist
8             print("Season or episode does not exist.")
9             return
10        print('\033[1m'+ "Shown is an overview over season ", season_no, " episode ", episode_no, "\n")
11        print('\033[0m')
12
13        list_characters(episode_script)
14        top_3= show_three_characters_with_biggest_amount_of_speech(episode_script)
15        show_amount_of_dialogue(episode_script)
16        sentiment(episode_script)
17        show_topics(episode_script)
18        return

```

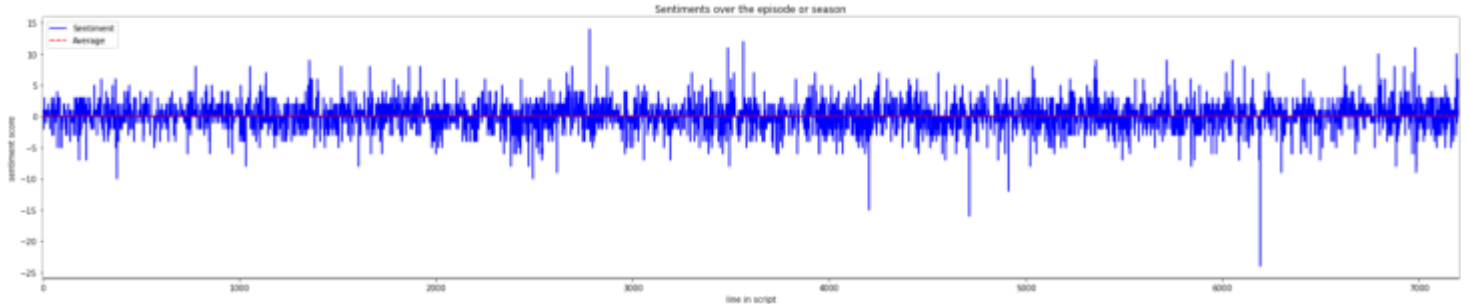
```
In [24]: 1 generate overview season(1)
```

Shown is an overview over season 1

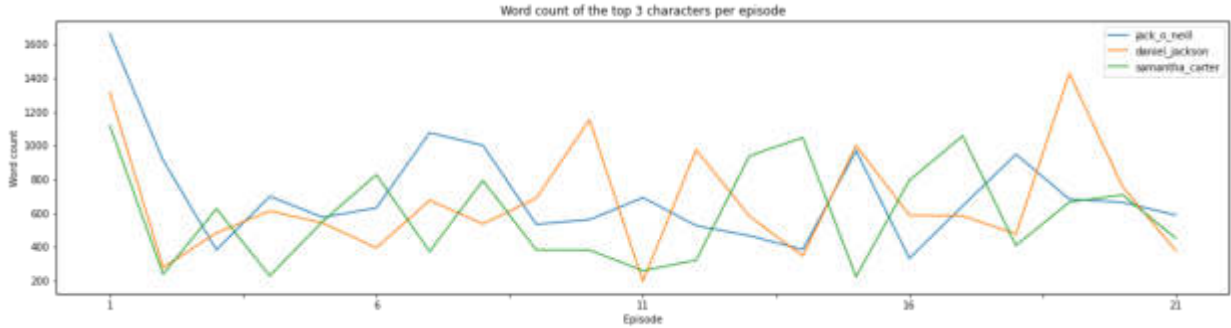
There are 106 characters in the season.
These characters are airman, woman, officer, apophis, soldier, samantha_carter, jack_o_neill, hammond, warner, kawalsky, harriman, ferretti, daniel_jackson, skaara, sha're, teal_c, bo'la, b oy, tech, medic, warren, native, goa'uld, casey, dr langford, scientist, catherine, martha, earnest, dr_janet, priest, bra'tac, o'neill, drey'auc, rya'c, nem, mackensie, kleinhouse, cole, h athor, solder, cassie, nurse, dr warner, davis, hanna, man, young hanna, villager, shak'l, narim, omoc, tuplo, maybourne, lya, siler, harlan, tv reporter, soldier 2, solder 3, airmen, jaff a, "auto destruct in, walter, doctor, kennedy, kinsey, entity, voice, ml, klorel, abu, ., mughal, turghan, nya, guard, makepeace, marine, leedora, johnson, connor, franks, baker, hanson, ja mala, sara, dad, reporter, sara's dad, crystal, answering machine, police officer, charlie, secretary, nefrayu, oper, anteaus, alekos, thetyes, kynthia, argosian, gairwyn, thor, kendra, and unas.

The 3 with the highest word counts are jack_o_neill, daniel_jackson, and samantha_carter.

In total, there are 76413 words spoken.
The sentiment score over the episode/season has an average of 0.0303.
The graph below shows the sentiment for each line in chronological order.



The amount of speech each of the top 3 characters has per episode, is as follows:



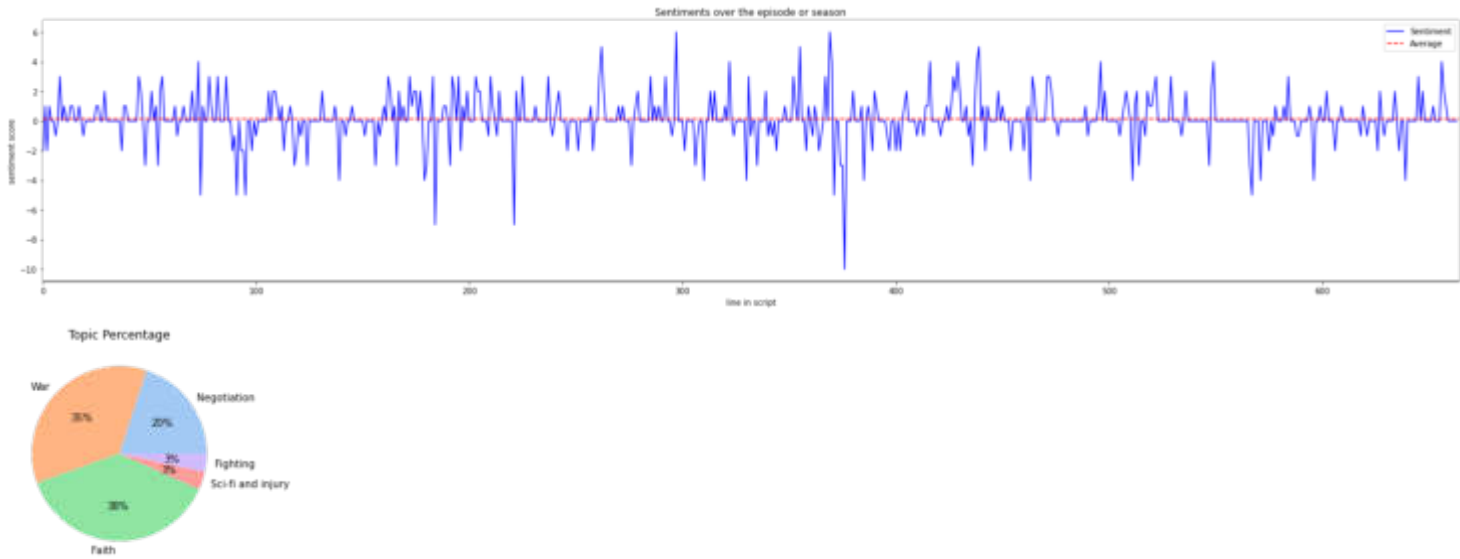
```
In [25]: 1 generate_overview_episode(1,1)
```

Shown is an overview over season 1 episode 1

There are 24 characters in the season.
These characters are airman, woman, officer, apophis, soldier, samantha_carter, jack_o_neill, hammond, warner, kawalsky, harriman, ferretti, daniel_jackson, skaara, sha're, teal_c, bo'la, boy, tech, medic, warren, native, goa'uld, and casey.

The 3 with the highest word counts are jack_o_neill, daniel_jackson, and samantha_carter.

In total, there are 6472 words spoken.
The sentiment score over the episode/season has an average of 0.1205.
The graph below shows the sentiment for each line in chronological order.



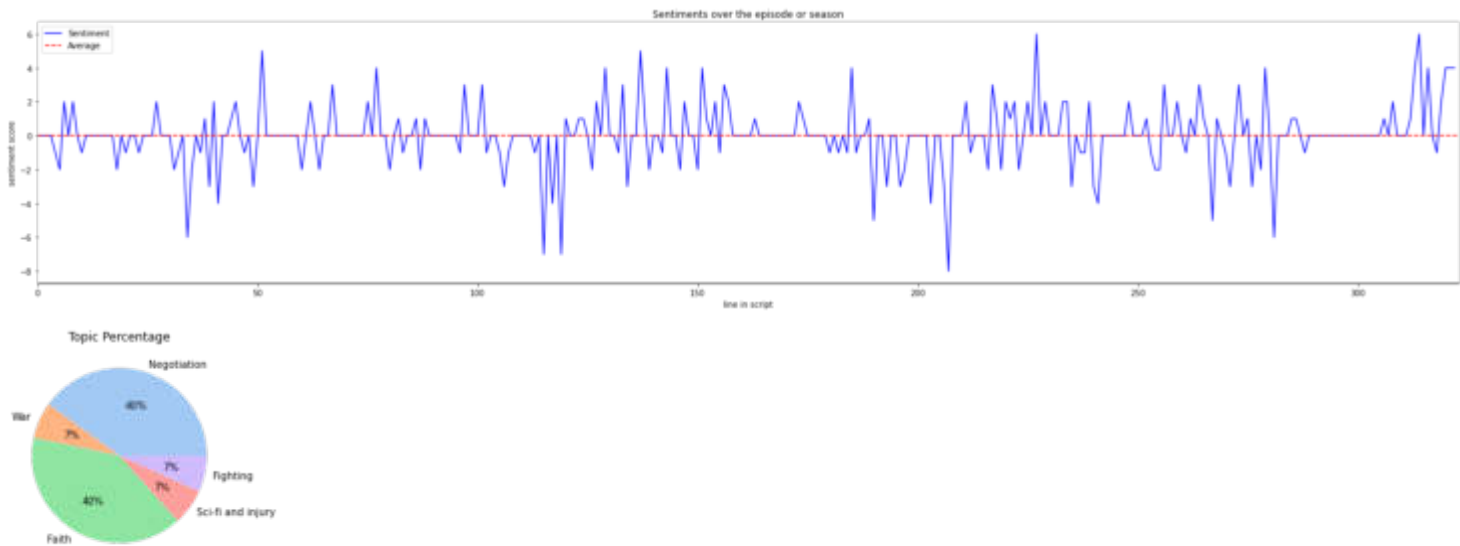
```
In [27]: 1 generate_overview_episode(3,19)
```

Shown is an overview over season 3 episode 19

There are 15 characters in the season.
These characters are technician, hammond, samantha_carter, woman, nyan, daniel_jackson, jack_o_neill, mallen, teal_c, limp, guard, unw, rigar, guy, and dr_janet.

The 3 with the highest word counts are nyan, rigar, and teal_c.

In total, there are 3081 words spoken.
The sentiment score over the episode/season has an average of 0.0186.
The graph below shows the sentiment for each line in chronological order.



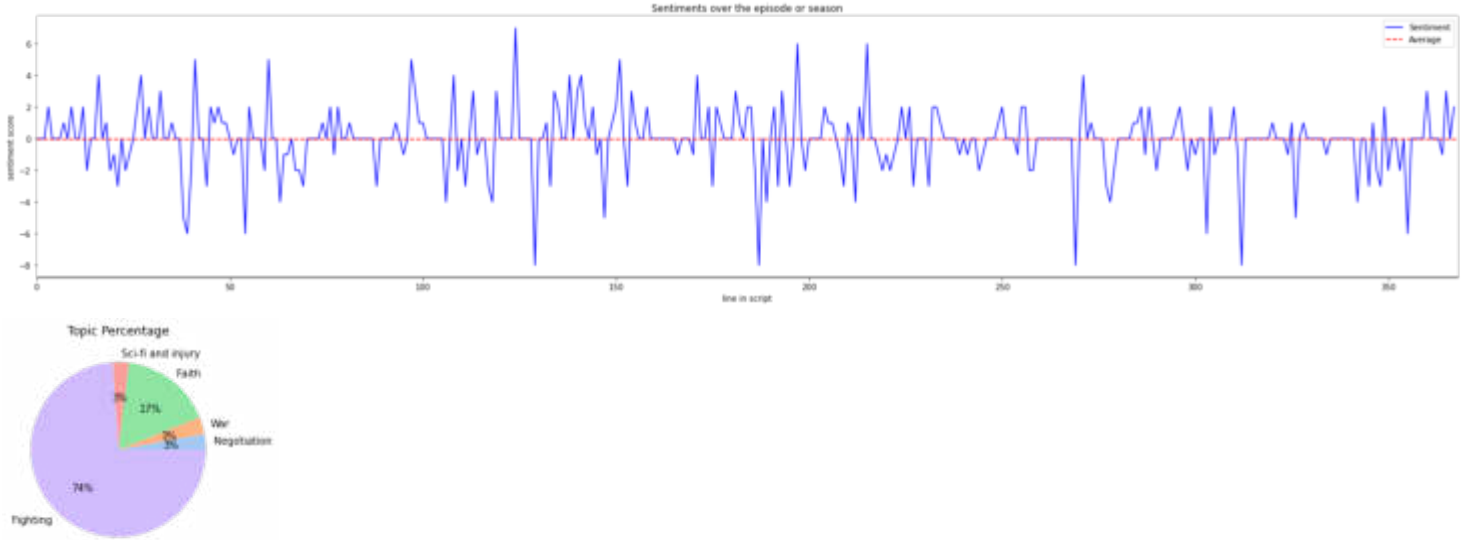
```
In [28]: 1 generate overview_episode(6,7)
```

Shown is an overview over season 6 episode 7

There are 22 characters in the season.
These characters are sgt. davis, hammond, samantha_carter, voice over speaker, briefing room, jonas, jack_o_neill, teal_c, 'gateroom, hale, kieran, dreylock, control room, kelowna, valis, resistance leader, kelownan conference chamber, dr. kieran's lab, woman, sgc infirmary, dr_janet, and sgc briefing room.

The 3 with the highest word counts are jonas, kieran, and valis.

In total, there are 5084 words spoken.
The sentiment score over the episode/season has an average of -0.0489.
The graph below shows the sentiment for each line in chronological order.



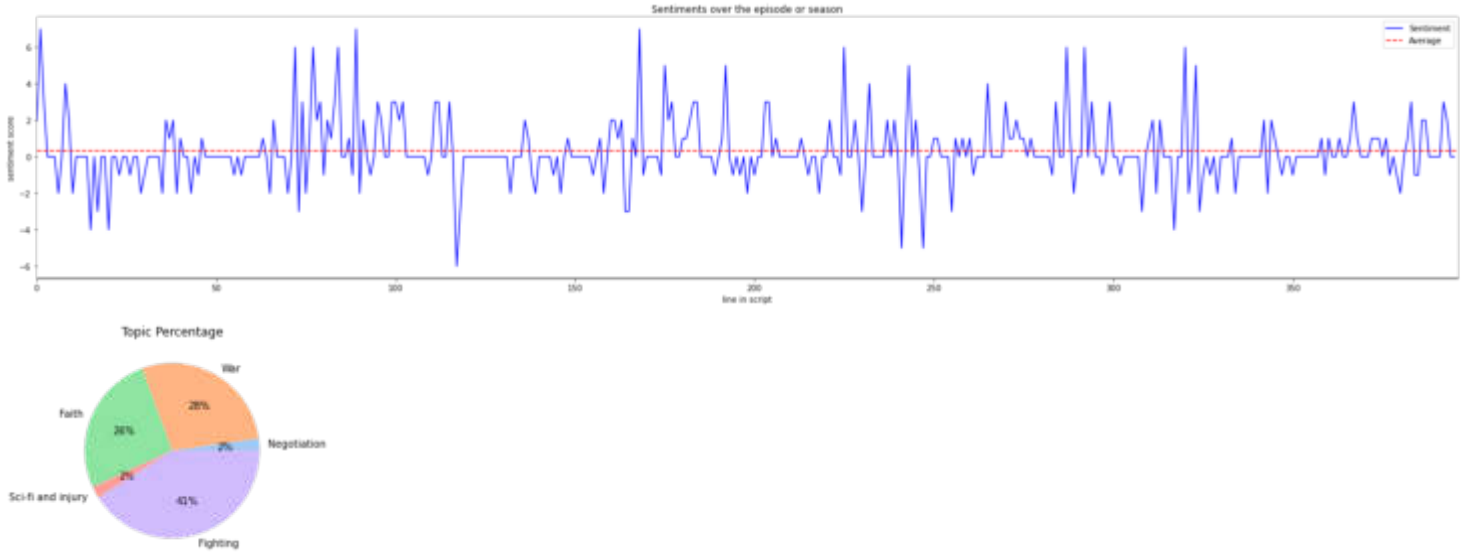
```
In [29]: 1 generate overview_episode(7,15)
```

Shown is an overview over season 7 episode 15

There are 11 characters in the season.
These characters are daniel_jackson, sarah, pete, samantha_carter, jack_o_neill, teal_c, fbi, ferretty, hammond, answerphone, and osiris.

The 3 with the highest word counts are daniel_jackson, pete, and samantha_carter.

In total, there are 3493 words spoken.
The sentiment score over the episode/season has an average of 0.3207.
The graph below shows the sentiment for each line in chronological order.



```
In [ ]: 1
```