

LDA Topic Model Classifier

```
In [56]: 1 from pickle import load#, dump
2 from my_lda_utils import my_lda_model
3 import pandas as pd
4 import numpy as np
5 from nltk.stem import WordNetLemmatizer, SnowballStemmer
6 import seaborn as sns
7 import matplotlib.pyplot as plt
```

Import Data

```
In [3]: 1 clean_data_folder = "../clean_data/"
2 with open("../clean_data/"+characters_with_gender.pkl", "rb") as f2:
3     characters = load(f2)
```

```
In [4]: 1 clean_data_folder_lda = "../clean_data_lda"
2 with open(clean_data_folder_lda+"/all_data.pkl", "rb") as f:
3     all_scripts = load(f)
4 data_lda = pd.merge(all_scripts, characters, on="character", how="inner")
```

```
In [5]: 1 all_scripts.reset_index(inplace=True, drop=True)
```

```
In [6]: 1 data = pd.merge(all_scripts, characters, on="character", how="inner")
```

```
In [7]: 1 counts = data.gender2.value_counts()
```

```
In [8]: 1 data
```

Out[8]:

	character	text	season	episode	doc_id	gender	gender2
0	woman	[are, not, you, guys, afraid, of, an, officer,...	1	1	3	female	female
1	woman	[does, that, thing, always, do, that]	1	1	5	female	female
2	woman	[whatever, it, is, under, the, trap, i, just, ...	1	1	6	female	female
3	woman	[im, telling, you, that, thing, is, moving]	1	1	8	female	female
4	woman	[i, take, it, this, has, never, happened, before]	1	1	11	female	female
...
40859	volnek	[you, shot, me]	9	8	34108	male	male
40860	volnek	[i, will, have, vengeance]	9	8	34108	male	male
40861	volnek	[you, are, fortunate, my, brother, drugged, yo...	9	8	34109	male	male
40862	volnek	[you, were, lucky]	9	8	34109	male	male
40863	volnek	[gate, kawooshes, granting, my, freedom, chang...	9	8	34109	male	male

40864 rows × 7 columns

Gender count stats for later

```
In [11]: 1 #gender counted by character occurrence
2 counts = data.gender2.value_counts()
3
4 female_occurr = counts.female
5 male_occurr = counts.male
6 unclear_occurr = counts.unclear
7 occurrence_count = female_occurr+male_occurr+unclear_occurr
8
9
10 #gender counted by character count
11 data1 = data[["character", "gender2"]].drop_duplicates()
12 data2 = data1.groupby(["gender2"]).count()
13 data3 =data2.reset_index()
14
15 female_count = data3.loc[data3.gender2 == "female"].character[0]
16 male_count = data3.loc[data3.gender2 == "male"].character[1]
17 unclear_count = data3.loc[data3.gender2 == "unclear"].character[2]
18 character_count = female_count+male_count+unclear_count
```

Transform data

```
In [12]: 1 data
```

Out[12]:

	character	text	season	episode	doc_id	gender	gender2
0	woman	[are, not, you, guys, afraid, of, an, officer,...	1	1	3	female	female
1	woman	[does, that, thing, always, do, that]	1	1	5	female	female
2	woman	[whatever, it, is, under, the, trap, i, just, ...	1	1	6	female	female
3	woman	[im, telling, you, that, thing, is, moving]	1	1	8	female	female
4	woman	[i, take, it, this, has, never, happened, before]	1	1	11	female	female
...
40859	volnek	[you, shot, me]	9	8	34108	male	male
40860	volnek	[i, will, have, vengeance]	9	8	34108	male	male
40861	volnek	[you, are, fortunate, my, brother, drugged, yo...	9	8	34109	male	male
40862	volnek	[you, were, lucky]	9	8	34109	male	male
40863	volnek	[gate, kawooshes, granting, my, freedom, chang...	9	8	34109	male	male

40864 rows × 7 columns

```
In [13]: 1 data rm gender1 = data.drop(['gender'], axis=1)
```

```
In [14]: 1 data rename col = data rm gender1.rename(columns={"gender2": "gender", "doc_id": "scene id"}, errors="raise")
```

```
In [15]: 1 data grouppped = data rename col.groupby(["character", "season", "episode", "scene id", "gender"]).agg(sum)
```

```
In [16]: 1 data reset_index = data_grouppped.reset_index()
```

In [17]: 1 data.reset_index().head()

Out[17]:

	character	season	episode	scene_id	gender	text
0	.	1	3	2079	unclear	[sam, is, with, abu, they, are, sitting, among...
1	.	1	16	1164	unclear	[]
2	.	1	17	1324	unclear	[]
3	.	2	9	6367	unclear	[]
4	.	2	10	3304	unclear	[]

In [18]: 1 # remove empty texts
2 data = data.reset_index().loc[data.reset_index().text.map(len) > 0]
3 data.head()

Out[18]:

	character	season	episode	scene_id	gender	text
0	.	1	3	2079	unclear	[sam, is, with, abu, they, are, sitting, among...
9	.	2	20	4835	unclear	[lockdown, in, progress, stand, clear, of, all...
13	bigwig	5	12	14914	male	[look, i, know, it, says, hes, flightless, in,...
14	bigwig	5	12	14915	male	[who, is, this]
15	bigwig	5	12	14916	male	[you, are, telling, me, an, air, force, office...

Classify

In [19]: 1 def lda_calssify(lda_model, dictionary, text):
2 # create corpus
3 bow = dictionary.doc2bow(text)
4 topics = (lda_model.get_document_topics(bow,
5 minimum_probability=0.0))
6 #get topic with highest probability
7 highest_percentage = 0
8 classified_topic = 0
9 for topic in topics:
10 if topic[1] > highest_percentage:
11 highest_percentage = topic[1]
12 classified_topic = topic[0]
13 return(classified_topic)

In [20]: 1 #get a classifying model:
2 with open("tfidf_coherence-topics-5-passes-5_approach3.pkl", "rb") as f:
3 df = load(f)
4 my_lda_model_1 = df.model[0.01]
5 my_lda_model_4 = df.model[0.04]
6 my_lda_model_8 = df.model[0.08]

In [21]: 1 df_columns = list(data.columns)
2 df_columns.append("topic")

In [22]: 1 data

Out[22]:

	character	season	episode	scene_id	gender	text
0	.	1	3	2079	unclear	[sam, is, with, abu, they, are, sitting, among...
9	.	2	20	4835	unclear	[lockdown, in, progress, stand, clear, of, all...
13	bigwig	5	12	14914	male	[look, i, know, it, says, hes, flightless, in,...
14	bigwig	5	12	14915	male	[who, is, this]
15	bigwig	5	12	14916	male	[you, are, telling, me, an, air, force, office...
...
26344	zippy	3	15	7971	male	[then, i, wish, to, point, out, the, futility,...
26345	zippy	3	15	7981	male	[my, vessel, comes, in, anticipation, of, our,...
26346	zippy	3	15	7982	male	[the, goals, rest, our, case, and, we, are, pr...
26347	zippy	3	15	7993	male	[stands, we, are, in, favor, of, korea, sits]
26348	zippy	3	15	8002	male	[talking, into, ball, rita, relook, data, i, d...

26249 rows × 6 columns

In [23]: 1 # Classify data with model 1
2 df1_1 = data.apply(lambda x: [x[0], x[1], x[2], x[3], x[4], x[5],
3 lda_calssify(my_lda_model_1.lda_model, my_lda_model_1.dictionary, x[5])], axis=1, result_type='expand')
4
5 # Classify data with model 4
6 df1_4 = data.apply(lambda x: [x[0], x[1], x[2], x[3], x[4], x[5],
7 lda_calssify(my_lda_model_4.lda_model, my_lda_model_4.dictionary, x[5])], axis=1, result_type='expand')
8 # Classify data with model 8
9 df1_8 = data.apply(lambda x: [x[0], x[1], x[2], x[3], x[4], x[5],
10 lda_calssify(my_lda_model_8.lda_model, my_lda_model_8.dictionary, x[5])], axis=1, result_type='expand')

In [24]: 1 df1_1.columns = df1_4.columns = df1_8.columns = df_columns
2 df1_1.head()

Out[24]:

	character	season	episode	scene_id	gender	text	topic
0	.	1	3	2079	unclear	[sam, is, with, abu, they, are, sitting, among...	2
9	.	2	20	4835	unclear	[lockdown, in, progress, stand, clear, of, all...	4
13	bigwig	5	12	14914	male	[look, i, know, it, says, hes, flightless, in,...	1
14	bigwig	5	12	14915	male	[who, is, this]	0
15	bigwig	5	12	14916	male	[you, are, telling, me, an, air, force, office...	1

In [25]: 1 occourance = df1_1.gender.value_counts()

Plot

Model 1

In [26]: 1 df2 = df1_1[["gender", "topic"]]
2 df2["count"] = 1

C:\Users\debor\AppData\Local\Temp\ipykernel_2228\2464217104.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

df2["count"] = 1

```
In [27]: 1 df3 = df2.groupby(["gender", "topic"]).count()
```

```
In [28]: 1 df4 =df3.reset_index()
```

```
In [29]: 1 df5 = df4.pivot_table("count", ['topic'], 'gender')
2 df5
```

Out[29]:

gender	female	male	unclear
topic			
0	2946	8264	130
1	927	2654	59
2	1163	2705	59
3	1110	3029	70
4	866	2220	47

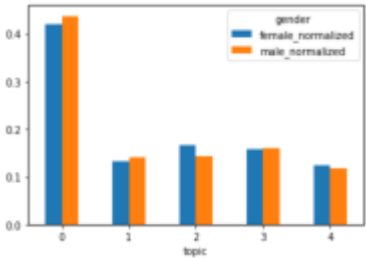
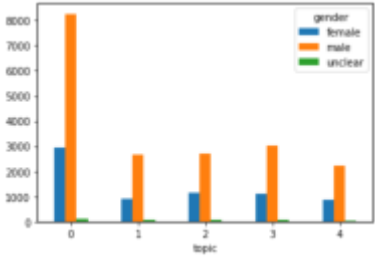
```
In [30]: 1 df6 = df5.copy()
2 df7 = df5.copy()
```

```
In [31]: 1 # Normalize by occurrence count
2 df6["female_normalized"] = df5.female/occournce.female
3 df6["male_normalized"] = df5.male/occournce.male
4 df6["unclear_normalized"] = df5.unclear/occournce.unclear
5 #df6 =df6.drop(["female", "male", "unclear"], axis=1)
6 df6
```

Out[31]:

gender	female	male	unclear	female_normalized	male_normalized	unclear_normalized
topic						
0	2946	8264	130	0.420137	0.437897	0.356164
1	927	2654	59	0.132202	0.140632	0.161644
2	1163	2705	59	0.165859	0.143334	0.161644
3	1110	3029	70	0.158300	0.160502	0.191781
4	866	2220	47	0.123503	0.117635	0.128767

```
In [32]: 1 ax = df5.plot.bar(rot=0)
2 ax = df6[["female_normalized", "male_normalized"]].plot.bar(rot=0)
```



```
In [33]: 1 df6['avg'] = (df6.female_normalized+df6.male_normalized)/2
2 df6['dev_f'] = np.power((df6.avg-df6.female_normalized), 2)
3 df6['dev_m'] = np.power((df6.avg-df6.male_normalized), 2)
4 df6['std_dev'] = (df6['dev_f']+df6['dev_m'])/2
5 format(df6.std_dev.mean(), '.8f')
```

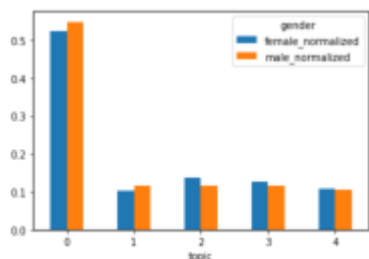
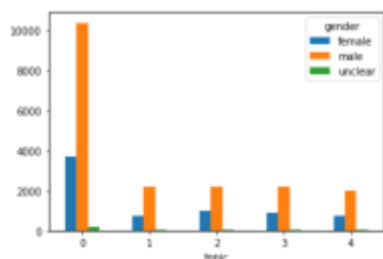
Out[33]: '0.00004666'

Model 4

```
In [34]: 1 df2 = df1_4[["gender", "topic"]]
2 df2["count"] = 1
3 df3 = df2.groupby(["gender", "topic"]).count()
4 df4 = df3.reset_index()
5 df5 = df4.pivot_table("count", ['topic'], 'gender')
6 df6 = df5.copy()
7 df7 = df5.copy()
8 # Normalize by occurrence count
9 df6["female_normalized"] = df5.female/occurrence.female
10 df6["male_normalized"] = df5.male/occurrence.male
11 df6["unclear_normalized"] = df5.unclear/occurrence.unclear
12 df6 = df6.drop(["female", "male", "unclear"], axis=1)
13 ax = df5.plot.bar(rot=0)
14 ax = df6[["female_normalized", "male_normalized"]].plot.bar(rot=0)
```

C:\Users\debor\AppData\Local\Temp\ipykernel_2228\976086660.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
df2["count"] = 1



```
In [35]: 1 df6['avg'] = (df6.female_normalized+df6.male_normalized)/2
2 df6['dev_f'] = np.power((df6.avg-df6.female_normalized), 2)
3 df6['dev_m'] = np.power((df6.avg-df6.male_normalized), 2)
4 df6['std_dev'] = ((df6['dev_f']+df6['dev_m']))/2
5 format(df6.std_dev.mean(), '.8f')
```

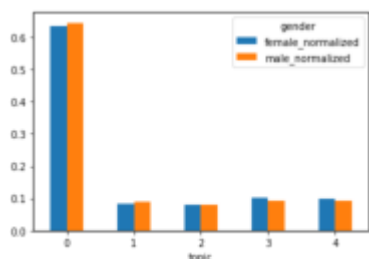
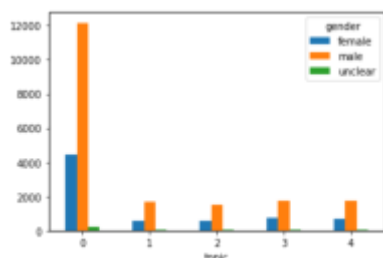
Out[35]: '0.00006703'

Model 8

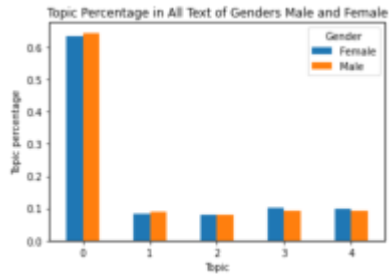
```
In [36]: 1 df2 = df1_8[["gender", "topic"]]
2 df2["count"] = 1
3 df3 = df2.groupby(["gender", "topic"]).count()
4 df4 = df3.reset_index()
5 df5 = df4.pivot_table("count", ['topic'], 'gender')
6 df6 = df5.copy()
7 df7 = df5.copy()
8 # Normalize by occurrence count
9 df6["female_normalized"] = df5.female/occurrence.female
10 df6["male_normalized"] = df5.male/occurrence.male
11 df6["unclear_normalized"] = df5.unclear/occurrence.unclear
12 df6 = df6.drop(["female", "male", "unclear"], axis=1)
13 ax = df5.plot.bar(rot=0)
14 ax = df6[["female_normalized", "male_normalized"]].plot.bar(rot=0)
```

C:\Users\debor\AppData\Local\Temp\ipykernel_2228\1371994753.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
df2["count"] = 1



```
In [97]: 1 df8 = df6.copy()
2 df9 = df8.rename(columns={"female_normalized": "Female",
3                             "male_normalized": "Male"})
4 df9.columns.name = "Gender"
5 df9.index.name = "Topic"
6 ax = df9[["Female", "Male"]].plot.bar(rot=0,
7                                         ylabel="Topic percentage",
8                                         title="Topic Percentage in All Text of Genders Male and Female")
9
```



```
In [95]: 1 df6
```

Out[95]:

gender	female_normalized	male_normalized	unclear_normalized	avg	dev_f	dev_m	std_dev
topic							
0	0.633485	0.644235	0.567123	0.638860	2.888736e-05	2.888736e-05	2.888736e-05
1	0.084284	0.089922	0.106849	0.087103	7.945330e-06	7.945330e-06	7.945330e-06
2	0.081004	0.080702	0.084932	0.080853	2.286517e-08	2.286517e-08	2.286517e-08
3	0.103394	0.091670	0.109589	0.097532	3.436294e-05	3.436294e-05	3.436294e-05
4	0.097832	0.093472	0.131507	0.095652	4.753441e-06	4.753441e-06	4.753441e-06

```
In [88]: 1 df9
```

Out[88]:

Gender	Female	Male	unclear_normalized	avg	dev_f	dev_m	std_dev
Topic							
0	0.633485	0.644235	0.567123	0.638860	2.888736e-05	2.888736e-05	2.888736e-05
1	0.084284	0.089922	0.106849	0.087103	7.945330e-06	7.945330e-06	7.945330e-06
2	0.081004	0.080702	0.084932	0.080853	2.286517e-08	2.286517e-08	2.286517e-08
3	0.103394	0.091670	0.109589	0.097532	3.436294e-05	3.436294e-05	3.436294e-05
4	0.097832	0.093472	0.131507	0.095652	4.753441e-06	4.753441e-06	4.753441e-06

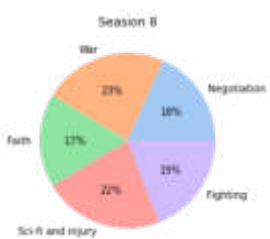
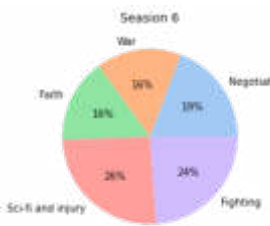
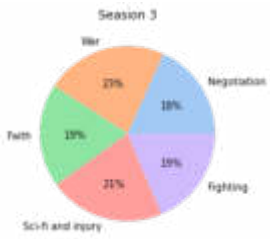
```
In [37]: 1 df6['avg'] = (df6.female_normalized+df6.male_normalized)/2
2 df6['dev_f'] = np.power((df6.avg-df6.female_normalized), 2)
3 df6['dev_m'] = np.power((df6.avg-df6.male_normalized), 2)
4 df6['std_dev'] = (df6['dev_f']+df6['dev_m'])/2
5 format(df6.std dev.mean(), '.8f')
```

Out[37]: '0.00001519'

Inspect Seasons

```
In [69]: 1 def show_topics(scripts, i):
2
3     #extract the texts and group them together
4     texts = sum(list(scripts["text"]), [])
5
6     # get the topic model and dictionary
7     with open ("lda_topic_model_5_topics-5_passes-0.08_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
23    df = pd.DataFrame(text_topics)
24    df = df.replace(topics)
25    df.columns = ["Topic", "Percentage"]
26    # df = df.set_index("Topic")
27
28    labels = topics.values
29    colors = sns.color_palette('pastel')[0:5]
30    plt.pie(df.Percentage, labels=df.Topic, colors = colors, autopct='%0.0f%%')
31    plt.title("Seasion "+str(i))
32    plt.show()
33    return
```

```
In [70]: 1 for i in range(1,10):
2         scripts = data[data.season == i]
3         show_topics(scripts, i)
4
```





```
In [62]: 1 data[data.season == 3]
```

Out[62]:

	character	season	episode	scene_id	gender	text
40	boy	3	5	9913	male	[colonel, jack, what, are, you, doing, here]
1919	daniel_jackson	3	1	6510	male	[i, really, try, not, to]
1920	daniel_jackson	3	1	6578	male	[wait, what, about, jack]
1921	daniel_jackson	3	1	6605	male	[its, just, a, deep, bleeding, gas, but, till,...
1922	daniel_jackson	3	1	6614	male	[well, i, guess, we, ca, not, go, under, it, e...
...
26344	zippy	3	15	7971	male	[then, i, wish, to, point, out, the, futility,...
26345	zippy	3	15	7981	male	[my, vessel, comes, in, anticipation, of, our,...
26346	zippy	3	15	7982	male	[the, goals, rest, our, case, and, we, are, pr...
26347	zippy	3	15	7993	male	[stands, we, are, in, favor, of, korea, sits]
26348	zippy	3	15	8002	male	[talking, into, ball, rita, relook, data, i, d...

3424 rows × 6 columns