

Predictions on Topics of UN Debates at UN Assembly New York with Topic Modelling

- Charita Ramnathsing
- April 29th, 2020

Build Dictionary

```
# Create a dictionary representation of the documents
```

```
dictionary = gensim.corpora.Dictionary(data_corpus_bigram)
```

```
dictionary.filter_extremes(no_below=10, no_above=0.6)
```

```
#Transform corpus into bag of words vectors
```

```
bow_corpus = [dictionary.doc2bow(text) for text in data_corpus_bigram]
```

Initialize and run LDA Mallet

```
from gensim.test.utils import common_corpus, common_dictionary
from gensim.models.wrappers import LdaMallet
from gensim.models import CoherenceModel
import os
```

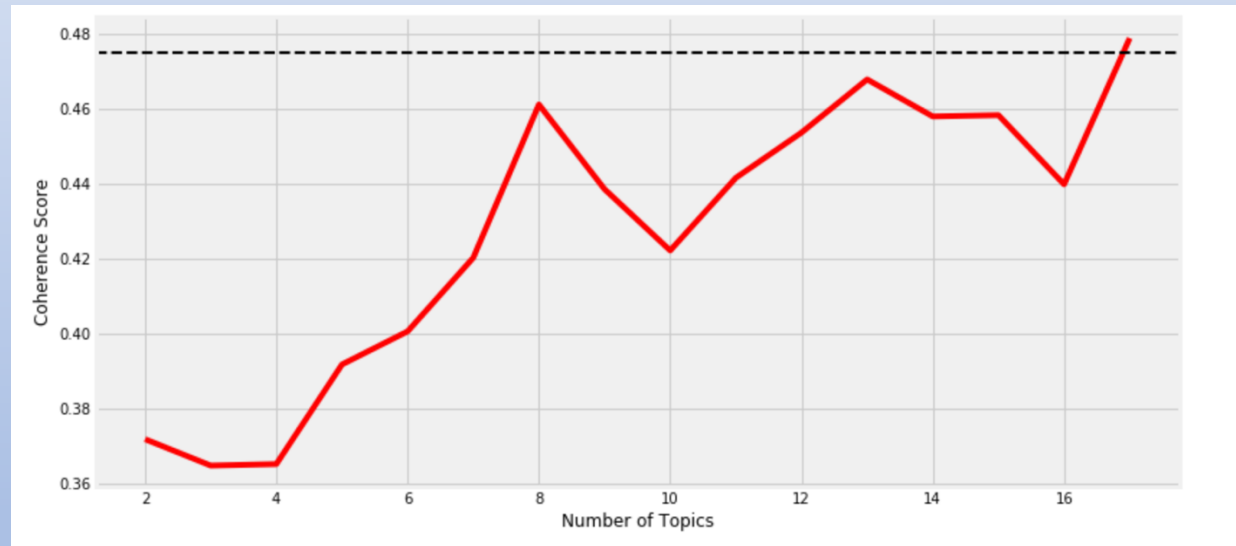
```
os.environ.update({'MALLET_HOME':r'my_folder/mallet-2.0.8/bin'})
MALLET_PATH = 'my_folder/mallet-2.0.8/bin/mallet'
```

[illegible]

Topic Modelling with 'topic_nums' topics

[illegible]

Chosen number of Topics in the Model



How did the model perform?

```
cv_coherence_model_lda_mallet = gensim.models.CoherenceModel(model=lda_mallet, corpus=bow_corpus,
                                                             texts=data_corpus_bigram,
                                                             dictionary=dictionary,
                                                             coherence='c_v')

avg_coherence_cv = cv_coherence_model_lda_mallet.get_coherence()

umass_coherence_model_lda_mallet = gensim.models.CoherenceModel(model=lda_mallet, corpus=bow_corpus,
                                                                texts=data_corpus_bigram,
                                                                dictionary=dictionary,
                                                                coherence='u_mass')

avg_coherence_umass = umass_coherence_model_lda_mallet.get_coherence()

print('Avg. Coherence Score (Cv):', avg_coherence_cv)
print('Avg. Coherence Score (UMass):', avg_coherence_umass)
```

```
Avg. Coherence Score (Cv): 0.47228869280709124
Avg. Coherence Score (UMass): -0.9522068656727839
```

Feature Engineering after Modelling

[illegible]

Prediction of Topic After modelling

Document	Dominant Topic	Contribution %	Topic Description	Country
0	6	30.28	south_africa,delegation,african,co_operation,namibia,independence,south,africa,apartheid,developed	Maldives
1	8	47.47	disarmament,nuclear,europe,nuclear_weapon,proposal,european,co_operation,treaty,field,basis	Finland
2	6	31.14	south_africa,delegation,african,co_operation,namibia,independence,south,africa,apartheid,developed	Niger
3	3	41.61	law,solidarity,society,trade,justice,mankind,latin_america,reason,case,freedom	Uruguay
4	6	33.52	south_africa,delegation,african,co_operation,namibia,independence,south,africa,apartheid,developed	Zimbabwe
...
7502	8	29.42	disarmament,nuclear,europe,nuclear_weapon,proposal,european,co_operation,treaty,field,basis	Kazakhstan
7503	5	25.45	africa,democracy,african,republic,democratic,election,solidarity,continent,child,programme	Liberia
7504	5	38.48	africa,democracy,african,republic,democratic,election,solidarity,continent,child,programme	Burundi
7505	1	30.83	terrorism,iraq,palestinian,law,afghanistan,israel,comprehensive,council,humanitarian,party	Hungary
7506	1	45.83	terrorism,iraq,palestinian,law,afghanistan,israel,comprehensive,council,humanitarian,party	Kuwait

7507 rows × 4 columns

Bag of Words per Topic



Word use per Topic

