

Assignment Project Report

PLSA: Text Document Clustering

Name: Pavan Tiwari

Course: AI and ML

(Batch 4)

- **Problem Statement**

Perform topic modelling using the 20 Newsgroup dataset (the dataset is also available in sklearn datasets sub-module). Perform the required data cleaning steps using NLP and then model the topics 1. Using Latent Dirichlet Allocation (LDA). 2. Using Probabilistic Latent Semantic Analysis (PLSA)

Prerequisites

- Software:
 - Python 3 (Use anaconda as your python distributor as well)
- Tools:
 - Pandas
 - Numpy
 - Matplotlib
 - Seaborn
 - NLTK
- Dataset Link: https://scikit-learn.org/0.19/datasets/twenty_newsgroups.html
- **Method Used**

PLSA or Probabilistic Latent Semantic Analysis is a technique used to model information under a probabilistic framework. It is a statistical technique for the analysis of two-mode and co-occurrence data. PLSA characterizes each word in a document as a sample from a mixture model, where mixture components are conditionally independent multinomial distributions. Its main goal is to model cooccurrence information under a probabilistic framework in order to discover the underlying semantic structure of the data.

- **Implementation:**

1. Load all required libraries and Dataset

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.datasets import fetch_20newsgroups
6 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
7 from sklearn.decomposition import NMF
8 from sklearn.preprocessing import normalize

1 categories = ['alt.atheism', 'talk.religion.misc', 'comp.graphics', 'sci.space', 'sci.med', 'sci.space', 'soc.religion.christia
2 newsgroups_train = fetch_20newsgroups(subset='train', categories=categories)
```

2. Preprocessing data

```
1 import nltk
2 import re

1 from nltk.corpus import stopwords
2 from nltk.stem.porter import PorterStemmer
3 from nltk.stem import WordNetLemmatizer
4 ps = PorterStemmer()

1 corpus = []
2 lem = WordNetLemmatizer()
3 for i in range(0, len(newsgroups_train.filesnames)):
4     review = newsgroups_train.data[i].split()
5     review = [lem.lemmatize(word) for word in review if word not in set(stopwords.words('english'))]
6     review = ' '.join(review)
7     corpus.append(review)

1 vectorizer = CountVectorizer(max_features=5000)
2 x_counts = vectorizer.fit_transform(corpus)

1 transformer = TfidfTransformer()
2 x_tfidf = transformer.fit_transform(x_counts)
```

3. Applying NMF model For topic Modeling

```
1 #number of topics
2 num_topics=7
3 #obtain a NMF model.
4 model = NMF(n_components=num_topics, init='nndsvd')
5 #fit the model
6 model.fit(xtfidf_norm)
```

NMF(init='nndsvd', n_components=7)

```
1 def get_nmf_topics(model, n_top_words):
2
3     #the word ids obtained need to be reverse-mapped to the words so we can print the topic names.
4     feat_names = vectorizer.get_feature_names()
5
6     word_dict = {}
7     for i in range(num_topics):
8
9         #for each topic, obtain the Largest values, and add the words they map to into the dictionary.
10        words_ids = model.components_[i].argsort()[::-n_top_words - 1:-1]
11        words = [feat_names[key] for key in words_ids]
12        word_dict['Topic # ' + '{:02d}'.format(i+1)] = words
13
14    return pd.DataFrame(word_dict)
```

Activate Windows
Go to Settings to activate Windows

1. Applying LDA model For topic Modeling

```
1 #number of topics
2 num_topics=7
3 #obtain a LDA model.
4 model = LatentDirichletAllocation(n_components=num_topics)
5 #fit the model
6 model.fit(xtfidf_norm)
```

LatentDirichletAllocation(n_components=7)

```
1 def get_lda_topics(model, n_top_words):
2
3     #the word ids obtained need to be reverse-mapped to the words so we can print the topic names.
4     feat_names = vectorizer.get_feature_names()
5
6     word_dict = {}
7     for i in range(num_topics):
8
9         #for each topic, obtain the Largest values, and add the words they map to into the dictionary.
10        words_ids = model.components_[i].argsort()[::-n_top_words - 1:-1]
11        words = [feat_names[key] for key in words_ids]
12        word_dict['Topic # ' + '{:02d}'.format(i+1)] = words
13
14    return pd.DataFrame(word_dict)
```

Activate Windows

