

Guided Projects Artificial Intelligence & Machine Learning

Guided Projects: Unsupervised Learning

PLSA: Text Document Clustering

Name	Lakshmi Thirunavukkarasu
Course	AI and ML (Batch 5)
Problem Statement	Perform topic modelling using the 20 Newsgroup dataset.

Software requirements prerequisites

Anaconda

Python 3.8

Python Packages

NLTK

NumPy

Pandas

Scikit

Matplotlib

Steps

1. Text Preprocessing

```

1 # Convert into lower case character
2 df_train[0] = df_train[0].map(lambda x: x.lower())
3
4 # remove from, subject, nntp and organization details from the text
5 df_train[0] = df_train[0].apply(lambda x: remove_sentence(x))
6
7 #replace web addresses with null value
8 df_train[0] = df_train[0].map(lambda x: re.sub(r'[\w\.-]+@[ \w\.-]+' ,'',x))
9
10 #remove all the special characters
11 df_train[0] = df_train[0].map(lambda x: re.sub('[,.\!?\_=">*<]', '', x))
12
13 #Remove numeric values
14 df_train[0] = df_train[0].map(lambda x: re.sub(r'[0-9]+' ,'', x))
15
16 #Apply Lemmization
17 df_train[0] = df_train[0].apply(lambda x: lemmatizer(x))
18
19 #Remove stopwords
20 df_train[0] = df_train[0].apply(lambda x: stop_words(x))
21

```

2. Convert the text into TDIDF vector

```

1 # Convert to TF-IDF format
2 #vectorizer_train= TfidfVectorizer(max_df=0.50, min_df=0.07, stop_words='english', use_idf=True)
3 vectorizer_train = TfidfVectorizer(max_df=0.50, min_df=0.01, stop_words=None, use_idf=True,
4                                   max_features = 1000, lowercase=False, ngram_range=(1,3))
5 X_train = vectorizer_train.fit_transform(df_train[0])
6 tfidf_feature_names = vectorizer_train.get_feature_names_out()
7 print(X_train.shape)
8 print(tfidf_feature_names)

```

3. Create LDA model

```

: 1 from sklearn.decomposition import LatentDirichletAllocation as LDA
2 lda = LDA(n_components=4, random_state=0)
3 clf = lda.fit(X_train.toarray(), labels_train)
4 clf_1 = clf.transform(X_train.toarray())
5

```

4. Create PLSA model

```

1 #Number of topics
2 from sklearn.decomposition import NMF
3 no_of_topics = 4
4 model = NMF(n_components = no_of_topics, init='nndsvd', regularization='components', random_state = 111 )
5 W1 = model.fit_transform(X_train)
6 H1 = model.components_

```

5. Display top 20 words for each topic identified by LDA and PLSA models

```

1 def plot_top_words(model, feature_names, n_top_words, title):
2     fig, axes = plt.subplots(1, 4, figsize=(30, 15), sharex=True)
3     axes = axes.flatten()
4     for topic_idx, topic in enumerate(model.components_):
5         top_features_ind = topic.argsort()[::-n_top_words - 1 : -1]
6         top_features = [feature_names[i] for i in top_features_ind]
7         weights = topic[top_features_ind]
8
9         ax = axes[topic_idx]
10        ax.barh(top_features, weights, height=0.7)
11        ax.set_title(f"Topic {topic_idx + 1}", fontdict={"fontsize": 30})
12        ax.invert_yaxis()
13        ax.tick_params(axis="both", which="major", labelsize=20)
14        for i in "top right left".split():
15            ax.spines[i].set_visible(False)
16        fig.suptitle(title, fontsize=40)
17
18    plt.subplots_adjust(top=0.90, bottom=0.05, wspace=0.90, hspace=0.3)
19    plt.show()

```