

CS145 Homework 6, Naive Bayes and Topic Modeling

****Due date:**** HW6 is due on **11:59 PM PT, Dec. 14 (Monday, Final Week)**. Please submit through GradeScope.

Print Out Your Name and UID

****Name:** Danning Yu, **UID:** 305087992******

Important Notes about HW6

- HW6, as the last homework, is optional if you choose to use the first 5 homework assignments for homework grading. We will select your highest 5 homework grades to calculate your final homework grade.
 - Since HW6 is optional, for the implementation of Naive Bayes and pLSA, you can choose to implement the provided `.py` and `.py` file by filling in the blocks. **Alternatively, you are given the option to implement completely from scratch based on your understanding. Note that some packages with ready-to-use implementation of Naive Bayes and pLSA are not allowed.**
-

Before You Start

You need to first create HW6 conda environment by the given `cs145hw6.yml` file, which provides the name and necessary packages for this tasks. If you have `conda` properly installed, you may create, activate or deactivate by the following commands:

```
conda env create -f cs145hw6.yml
conda activate hw6
conda deactivate
```

OR

```
conda env create --name NAMEOFYOURCHOICE -f cs145hw6.yml
conda activate NAMEOFYOURCHOICE
conda deactivate
```

To view the list of your environments, use the following command:

```
conda env list
```

More useful information about managing environments can be found [here](#).

You may also quickly review the usage of basic Python and Numpy package, if needed in coding for matrix operations.

In this notebook, you must not delete any code cells in this notebook. If you change any code outside the blocks (such as hyperparameters) that you are allowed to edit (between `START/END YOUR CODE HERE`), you need to highlight these changes. You may add some additional cells to help explain your results and observations.

```
In [1]: import numpy as np
from numpy import zeros, int8, log
from pylab import random
import pandas as pd
import matplotlib.pyplot as plt
from pylab import rcParams
rcParams['figure.figsize'] = 8,8
import seaborn as sns; sns.set()
import re
import time
import nltk
nltk.download('punkt')
from nltk.tokenize import word_tokenize
from sklearn.metrics import confusion_matrix
%load_ext autoreload
%autoreload 2
```

```
[nltk_data] Downloading package punkt to /Users/danningy/nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

Note that `seaborn` in HW6 is only used for plotting classification confusion matrix (in a "heatmap" style). If you encounter installation problem and cannot solve it, you may use alternative plot methods to show your results.

Section 1: Naive Bayes for Text (50 points)

Naive Bayes is one generative model for text classification. In the problem, you are given a document in `dataset` folder. The original data comes from "20 newsgroups". You can use the provided data files to save efforts on preprocessing.

Note: The code and dataset are under the subfolder named `nb`.

```
In [2]: ### Data processing and preparation
# read train/test labels from files
train_label = pd.read_csv('./nb/dataset/train.label', names=['t'])
train_label = train_label['t'].tolist()
test_label = pd.read_csv('./nb/dataset/test.label', names=['t'])
test_label = test_label['t'].tolist()

# read train/test documents from files
train_data = open('./nb/dataset/train.data')
df_train = pd.read_csv(train_data, delimiter=' ', names=['docIdx', 'wordIdx', 'c'])
test_data = open('./nb/dataset/test.data')
df_test = pd.read_csv(test_data, delimiter=' ', names=['docIdx', 'wordIdx', 'cou
```

```
# read vocab
vocab = open('./nb/dataset/vocabulary.txt')
vocab_df = pd.read_csv(vocab, names = ['word'])
vocab_df = vocab_df.reset_index()
vocab_df['index'] = vocab_df['index'].apply(lambda x: x+1)

# add label column to original df_train
docIdx = df_train['docIdx'].values
i = 0
new_label = []
for index in range(len(docIdx)-1):
    new_label.append(train_label[i])
    if docIdx[index] != docIdx[index+1]:
        i += 1
new_label.append(train_label[i])
df_train['classIdx'] = new_label
```

If you have the data prepared properly, the following line of code would return the head of the `df_train` dataframe, which is,

	docIdx	wordIdx	count	classIdx
0	1	1	4	1
1	1	2	2	1
2	1	3	10	1
3	1	4	4	1
4	1	5	2	1

```
In [3]: # check the head of 'df_train'
print(df_train.head())
```

	docIdx	wordIdx	count	classIdx
0	1	1	4	1
1	1	2	2	1
2	1	3	10	1
3	1	4	4	1
4	1	5	2	1

Complete the implementation of Naive Bayes model for text classification `nbm.py`. After that, run `nbm_sklearn.py`, which uses `sklearn` to implement naive bayes model for text classification. (Note that the dataset is slightly different loaded in `nbm_sklearn.py` and also you don't need to change anything in `nbm_sklearn.py` and directly run it.)

If the implementation is correct, you can expect the results are generally close on both train set accuracy and test set accuracy.

```
In [4]: from nb.nbm import NB_model

# model training
nbm = NB_model()
nbm.fit(df_train, train_label, vocab_df)
```

Prior Probability of each class:

```

1: 0.04259472890229834
2: 0.05155736977549028
3: 0.05075871860857219
4: 0.05208980388676901
5: 0.051024935664211554
6: 0.052533498979501284
7: 0.051646108794036735
8: 0.052533498979501284
9: 0.052888455053687104
10: 0.0527109770165942
11: 0.05306593309078002
12: 0.0527109770165942
13: 0.05244475996095483
14: 0.0527109770165942
15: 0.052622237998047744
16: 0.05315467210932647
17: 0.04836276510781791
18: 0.05004880646020055
19: 0.04117490460555506
20: 0.033365870973467035
Training class 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
Training completed!

```

In [5]:

```

# make predictions on train set to validate the model
predict_train_labels = nbm.predict(df_train)
train_acc = (np.array(train_label) == np.array(predict_train_labels)).mean()
print("Accuracy on training data by my implementation: {}".format(train_acc))

# make predictions on test data
predict_test_labels = nbm.predict(df_test)
test_acc = (np.array(test_label) == np.array(predict_test_labels)).mean()
print("Accuracy on training data by my implementation: {}".format(test_acc))

test_res = np.array(test_label) == np.array(predict_test_labels)
for i, res in enumerate(test_res):
    if not res:
        print(f"Document {i+1} was misclassified")
        break

```

```

Finished predictions for doc 2500
Finished predictions for doc 5000
Finished predictions for doc 7500
Finished predictions for doc 10000
Accuracy on training data by my implementation: 0.9481764131688704
Finished predictions for doc 2500
Finished predictions for doc 5000
Finished predictions for doc 7500
Accuracy on training data by my implementation: 0.7873417721518987
Document 4 was misclassified

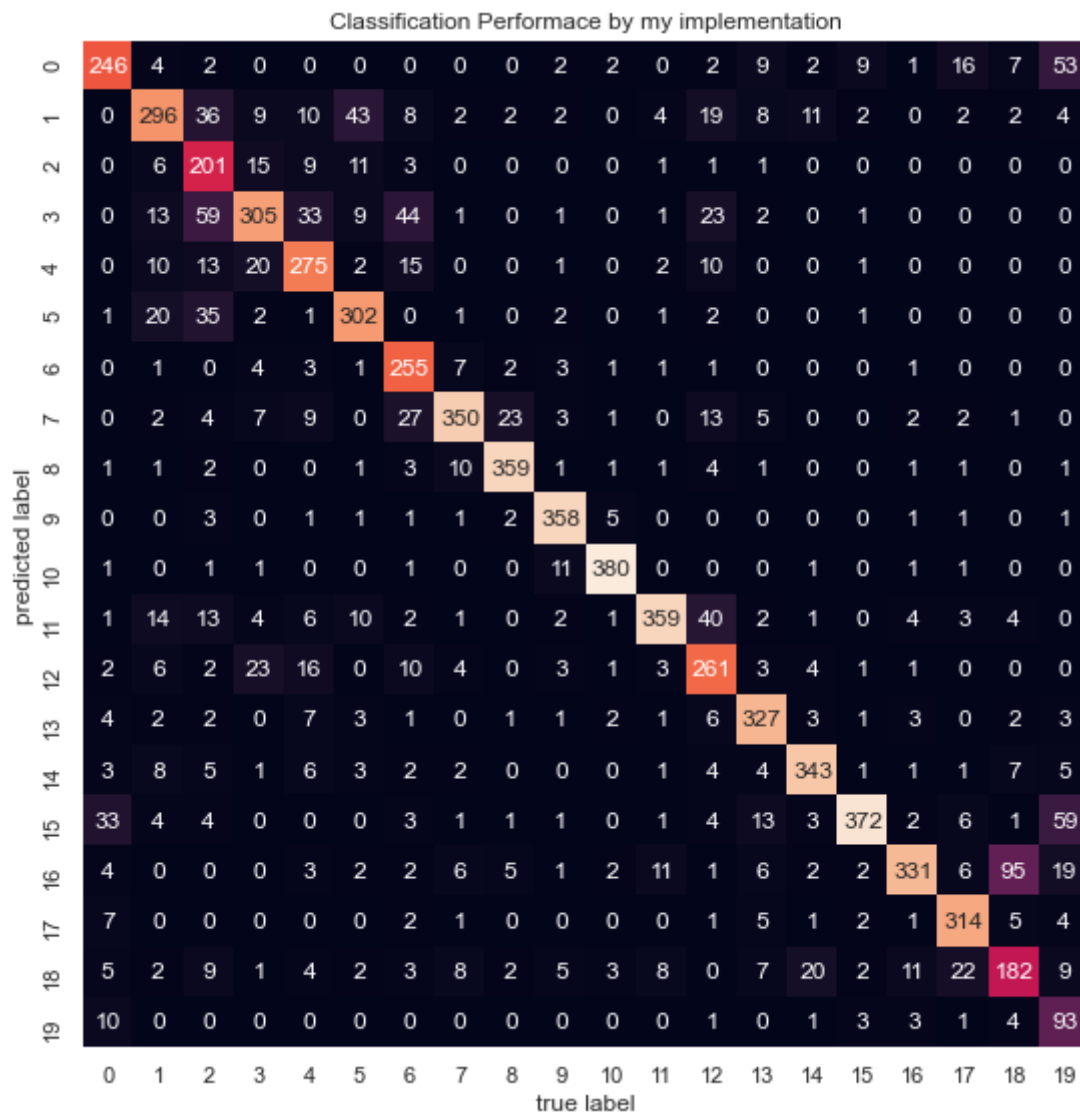
```

In [6]:

```

# plot classification matrix
mat = confusion_matrix(test_label, predict_test_labels)
sns.heatmap(mat.T, square=True, annot=True, fmt='d', cbar=False)
plt.title('Classification Performace by my implementation') # fixed a typo here,
plt.xlabel('true label')
plt.ylabel('predicted label')
plt.tight_layout()
plt.savefig('./nb/output/nbm_mine.png')
plt.show()

```



****Reminder:**** Do not forget to run nbm_sklearn.py to compare the results to get the accuracy and confusion matrix by sklearn implementation. You can run `python nbm_sklearn.py` under the folder path of `./hw6/nb/`.

sklearn results:

Train accuracy: 0.9326

Test accuracy: 0.7734

		Classification Performance by sklearn																			
		alt.atheism	comp.graphics	comp.os.ms-windows.misc	comp.sys.ibm.pc.hardware	comp.sys.mac.hardware	comp.windows.x	misc.forsale	rec.autos	rec.motorcycles	rec.sport.baseball	rec.sport.hockey	sci.crypt	sci.electronics	sci.med	sci.space	soc.religion.christian	talk.politics.guns	talk.politics.mideast	talk.politics.misc	talk.religion.misc
predicted label	alt.atheism	166	1	0	0	0	1	0	0	0	0	0	0	2	0	2	0	0	2	33	
	comp.graphics	0	252	14	5	3	21	1	1	0	0	0	2	4	3	2	0	0	1	0	2
	comp.os.ms-windows.misc	0	15	258	11	8	17	3	0	0	0	0	1	2	0	0	0	0	0	0	0
	comp.sys.ibm.pc.hardware	1	12	45	305	23	13	31	3	1	0	0	0	17	1	1	0	1	0	0	0
	comp.sys.mac.hardware	0	9	3	17	298	2	12	0	0	1	0	0	5	1	0	0	0	0	0	0
	comp.windows.x	1	18	9	1	0	298	1	0	0	0	0	1	0	3	3	0	0	0	0	0
	misc.forsale	0	1	0	3	3	1	271	4	2	0	0	1	2	1	0	0	2	0	0	0
	rec.autos	0	2	2	6	8	0	19	364	10	4	1	3	8	0	2	0	0	0	1	0
	rec.motorcycles	1	1	1	1	1	1	4	3	371	0	0	0	7	2	1	1	1	0	0	0
	rec.sport.baseball	1	5	3	0	3	1	4	2	0	357	4	0	1	3	0	0	1	1	0	1
	rec.sport.hockey	1	2	2	2	1	0	6	2	0	22	387	0	2	4	1	0	0	0	1	1
	sci.crypt	3	41	25	19	16	23	5	4	4	0	1	383	78	11	6	0	10	2	11	3
	sci.electronics	0	4	1	13	8	0	12	1	0	0	0	1	235	5	1	0	0	0	0	0
	sci.med	6	0	0	0	0	1	6	1	0	0	0	0	3	292	2	1	0	0	1	4
	sci.space	3	6	6	5	2	4	3	3	0	2	1	0	11	6	351	2	1	0	7	4
	soc.religion.christian	123	15	23	3	8	10	9	3	8	9	5	3	15	52	19	392	6	24	35	131
	talk.politics.guns	4	4	2	1	3	2	3	4	2	1	0	1	2	6	4	0	341	3	118	29
	talk.politics.mideast	8	1	0	0	0	0	0	0	0	1	0	0	1	4	0	0	1	344	5	5
	talk.politics.misc	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	1	129	3
talk.religion.misc	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	35	
		alt.atheism	comp.graphics	comp.os.ms-windows.misc	comp.sys.ibm.pc.hardware	comp.sys.mac.hardware	comp.windows.x	misc.forsale	rec.autos	rec.motorcycles	rec.sport.baseball	rec.sport.hockey	sci.crypt	sci.electronics	sci.med	sci.space	soc.religion.christian	talk.politics.guns	talk.politics.mideast	talk.politics.misc	talk.religion.misc
		true label																			

Question & Analysis

1. Please indicate whether you implemented based the given code or from scratch.
2. Report your classification accuracy on train and test documents. Also report your classification confusion matrix. Show one example document that Naive Bayes classifies incorrectly (i.e. fill in the following result table). Attach the output figure
 ./output/nbm_mine.png in the jupyter book and briefly explain your observation on the accuracy and confusion matrix.

Train set accuracy	Test set accuracy
sklearn implementaion	
your implementaion	

1. Show one example document that Naive Bayes classifies incorrectly by filling the following table. Provide your thought on the reason why this document is misclassified. (Note that

the topic mapping is available at `train.map` same as `test.map`)

Words (count) in the example document	Predicted label	Truth label
For example, student (4), education (2), ...	Class A	Class B

1. Is Naive Bayes a generative model or discriminative model and why? What is the difference between Naive Bayes classifier and Logistic Regression? What are the pros and cons of Naive Bayes for text classification task?
2. Can you apply Naive Bayes model to identify spam emails from normal ones? Briefly explain your method (you don't need to implementation for this question).

```
In [7]: test_res = np.array(test_label) == np.array(predict_test_labels)
misclassified_doc = None
for i, res in enumerate(test_res):
    if not res:
        print(f"A misclassified test document: {i+1}")
        print(f"Predicted {predict_test_labels[i]}, ground truth {test_label[i]}")
        misclassified_doc = i + 1
        break

df = df_test.loc[df_test['docIdx'] == misclassified_doc].copy()
vocab_dict = vocab_df.set_index('index').to_dict()
# print(vocab_dict)
df['words'] = df['wordIdx'].map(vocab_df.set_index('index')['word'].to_dict())
print(df)
df.to_csv("document4.csv", sep=",")
```

```
A misclassified test document: 4
Predicted 16, ground truth 1
   docIdx  wordIdx  count      words
395      4        10      1    atheist
396      4        12     20         of
397      4        16      1        from
398      4        17     10    religion
399      4        23     25         and
..      ...      ...      ...      ...
746      4    53976      2         tgk
747      4    53977      1    quelled
748      4    53978      1 quintessential
749      4    53979      1    philanthropy
750      4    53980      3 supernaturalists

[356 rows x 4 columns]
```

Your Answers

Question 0:

I implemented naive Bayes using the provided code skeleton.

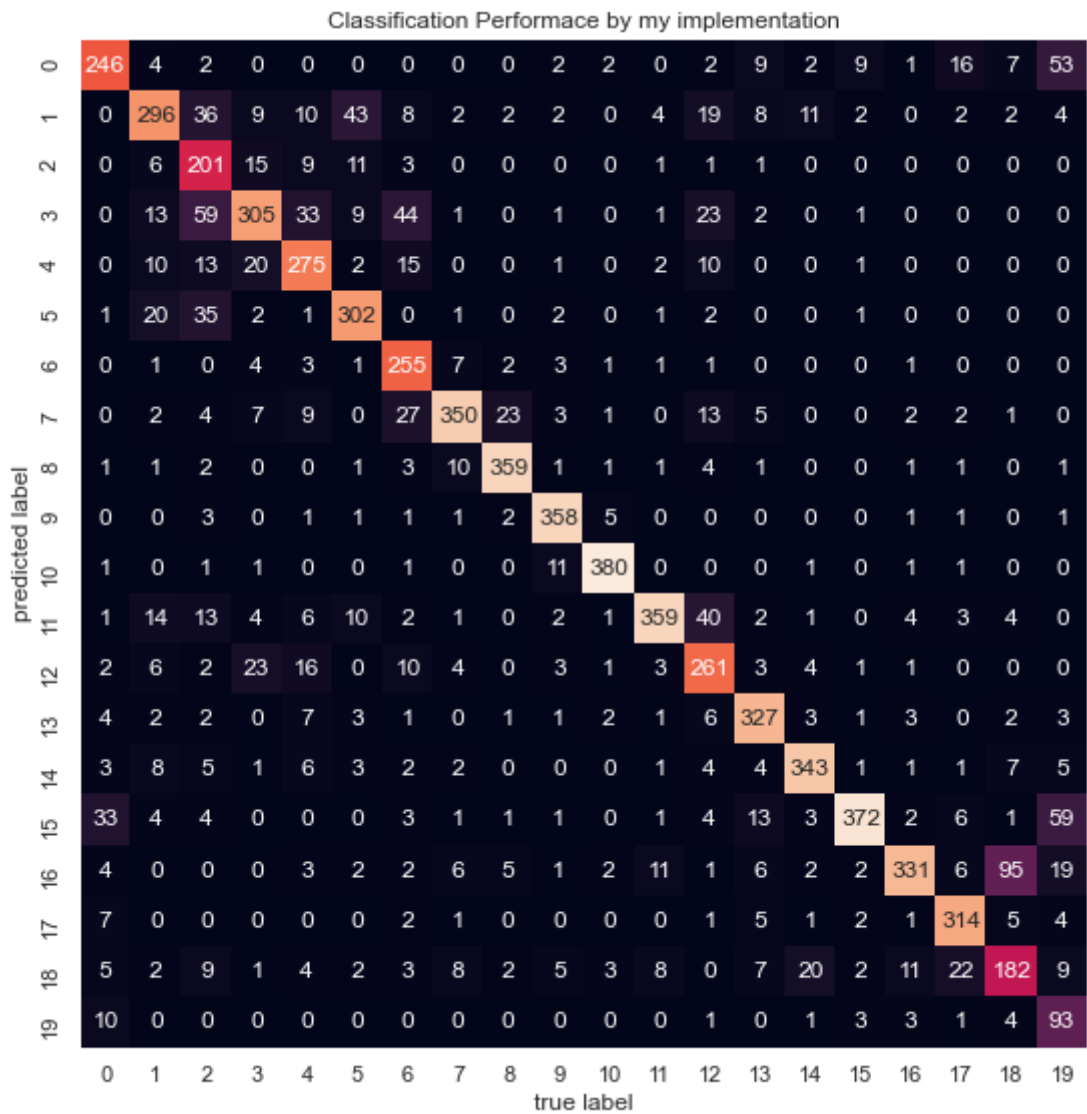
Question 1:

My implementation and the `sklearn` implementation had the following train and test set accuracies:

Train Set Accuracy	Test Set Accuracy
--------------------	-------------------

	Train Set Accuracy	Test Set Accuracy
sklearn implementation	0.9326	0.7734
Your implementation	0.9482	0.7873

The following image gives the classification confusion matrix for my implementation (generated from cell above):



From the confusion matrix, we can see that for most of the classes, what was predicted matched the ground truth. This makes sense, as the accuracy was 0.7873 for test data. However, there were lots of off-diagonal values that were nonzero, which is why the test accuracy is only 0.7873. The notable mispredictions are as follows (all nondiagonal values exceeding 50):

Ground Truth	Predicted Topic	# of Mispredictions
talk.politics.misc (column 18/topic 19)	talk.politics.guns (row 16/topic 17)	95
talk.religion.misc (column 19/topic 20)	soc.religion.christian (row 15/topic 16)	59

Ground Truth	Predicted Topic	# of Mispredictions
talk.religion.misc (column 19/topic 20)	alt.atheism (row 0/topic 1)	53
comp.os.ms-windows.misc (column 2/topic 3)	comp.sys.imb.pc.hardware (row 3/topic 4)	59

Note: the column and row numbers start from 0, so the class/topic number is the row or column number plus 1.

From this table, we can see that for these mispredictions, it tended to be a miscellaneous topic that NB thought actually had a closely related topic. Politics often includes talks about guns, religion includes both Christianity and atheism, and Windows is closely related to IBM PC hardware. Thus, the model was essentially deciding that that various documents in this area classified under "misc" under had a topic.

One particular document that NB misclassified was document 4 from the test set. A portion of the document is printed below. The full document is attached to the end of this submission after the code.

A misclassified test document: 4

	docIdx	wordIdx	count	words
395	4	10	1	atheist
396	4	12	20	of
397	4	16	1	from
398	4	17	10	religion
399	4	23	25	and
..
746	4	53976	2	tgk
747	4	53977	1	quelled
748	4	53978	1	quintessential
749	4	53979	1	philanthropy
750	4	53980	3	supernaturalists

[356 rows x 4 columns]

Question 2:

Document 4 was misclassified. The details about the document are given below.

Words (count) in document 4	Predicted label	Truth label
atheist(1), of(20), from(1), religion(10), and(25)...tgk(2), quelled(1), quintessential(1), philanthropy(1), supernaturalists(3)	Class 16	Class 1

The results indicate that it predicted that document 4 was soc.religion.christian while in reality it was alt.atheism. This misclassification probably occurred because atheism and Christianity are both religion-related topics, so they are somewhat related. There may have been some words in this atheism document that related to Christianity, and they were enough to make the model think that it was more likely for this document to be one about Christianity. Also, note the word "atheism" only appeared once, so this might have contributed to a lower likelihood for

alt.atheism compared to soc.religion.christian. In total, there were 33 documents that were misclassified as class 16 when they were actually class 1, showing that this misclassification isn't that uncommon.

Question 3:

Naive Bayes is a generative model because it models a joint probability $P(x, y)$ instead of a conditional probability $P(y|x)$, which corresponds to a discriminative model. In this case, the joint probability is $P(\text{document contains some sequence of words, document is in class } j)$. The difference between NB and logistic regression is that NB assumes that features such as words are conditionally independent given the class label, while logistic regression does not, and thus it learns a conditional probability distribution instead, $P(\text{class} | \text{some data})$.

For text classification, the pros of NB are that it is a simple model and runs quite fast, even with large data sets. However, its assumption of conditional independence is not necessarily something that is true in the real world. Also, NB fails to take into account the sequence of words in a document, and so it loses some information by using a bag of words representation.

Question 4:

Yes, you could, and this is actually used by many real-world spam detection algorithms. Spam emails commonly contain certain words that immediately cause humans to classify the spam, such as "won," "free," and "offer." My method would be to first train the model on a dataset of emails that includes both spam and non-spam emails so that it can learn the joint probabilities. To classify a new email, I would convert it to a bag words of representation, calculate the likelihood of it being a spam or not a spam email using the learned joint probability parameters, and then make the classification according to which class (spam or not spam) has a higher likelihood from the words in the email.

Section 2: Topic Modeling: Probabilistic Latent Semantic Analysis (50 points)

In this section, you will implement Probabilistic Latent Semantic Analysis (pLSA) by EM algorithm. Note: The code and dataset are under the subfolder named `plsa`. You can find two dataset files named `dataset1.txt` and `dataset2.txt` together with a [stopword](#) list as `stopwords.dic`.

First complete the implementation of pLSA in `plsa.py`. You need to finish the E step, M step and likelihood function. Note that the optimizing process on dataset 2 might take a while.

```
In [6]: # input file, output files and parameters
datasetFilePath = './plsa/dataset/dataset1.txt' # or set as './plsa/dataset/data
stopwordsFilePath = './plsa/dataset/stopwords.dic'
docTopicDist = './plsa/output/docTopicDistribution.txt'
topicWordDist = './plsa/output/topicWordDistribution.txt'
dictionary = './plsa/output/dictionary.dic'
topicWords = './plsa/output/topics.txt'

K = 4 # number of topic
```

```

maxIteration = 20 # maxIteration and threshold control the train process
threshold = 3
topicWordsNum = 20 # parameter for output

```

```

In [7]: from plsa.plsa import PLSA
        from plsa.utils import preprocessing
        # data processing
        N, M, word2id, id2word, X = preprocessing(datasetFilePath, stopwordsFilePath)

```

```

In [18]: plsa_model = PLSA()
         plsa_model.initialize(N, K, M, word2id, id2word, X)

         oldLoglikelihood = 1
         newLoglikelihood = 1

         for i in range(0, maxIteration):
             plsa_model.EStep() #implement E step
             plsa_model.MStep() #implement M step
             newLoglikelihood = plsa_model.LogLikelihood()
             print("[", time.strftime('%Y-%m-%d %H:%M:%S', time.localtime(time.time())), "]"
                   "iteration", str(newLoglikelihood))
             # you should see increasing loglikelihood
             if(abs(newLoglikelihood - oldLoglikelihood) < threshold):
                 # change to absolute value, or else it terminates after first iteration
                 break
             oldLoglikelihood = newLoglikelihood

         plsa_model.output(docTopicDist, topicWordDist, dictionary, topicWords, topicWord

[ 2020-12-11 16:04:59 ] 1 iteration -7605.31997548358
[ 2020-12-11 16:04:59 ] 2 iteration -7448.5138155125205
[ 2020-12-11 16:04:59 ] 3 iteration -7277.3019328973005
[ 2020-12-11 16:04:59 ] 4 iteration -7118.597475533447
[ 2020-12-11 16:04:59 ] 5 iteration -6985.828498139784
[ 2020-12-11 16:04:59 ] 6 iteration -6872.979733331642
[ 2020-12-11 16:04:59 ] 7 iteration -6775.821526663395
[ 2020-12-11 16:04:59 ] 8 iteration -6700.153252936795
[ 2020-12-11 16:04:59 ] 9 iteration -6645.141889223938
[ 2020-12-11 16:04:59 ] 10 iteration -6608.452423000288
[ 2020-12-11 16:04:59 ] 11 iteration -6584.809733086158
[ 2020-12-11 16:05:00 ] 12 iteration -6565.413904979665
[ 2020-12-11 16:05:00 ] 13 iteration -6545.967208851735
[ 2020-12-11 16:05:00 ] 14 iteration -6527.880448716848
[ 2020-12-11 16:05:00 ] 15 iteration -6513.849515115477
[ 2020-12-11 16:05:00 ] 16 iteration -6505.988092755073
[ 2020-12-11 16:05:00 ] 17 iteration -6502.7345270345695
[ 2020-12-11 16:05:00 ] 18 iteration -6501.50946888105

```

```

In [19]: # Repeat with dataset 2

         # input file, output files and parameters
         datasetFilePath = './plsa/dataset/dataset2.txt'

         from plsa.plsa import PLSA
         from plsa.utils import preprocessing

         N, M, word2id, id2word, X = preprocessing(datasetFilePath, stopwordsFilePath)

         plsa_model = PLSA()

```

```

plsa_model.initialize(N, K, M, word2id, id2word, X)

oldLoglikelihood = 1
newLoglikelihood = 1

for i in range(0, maxIteration):
    plsa_model.EStep() #implement E step
    plsa_model.MStep() #implement M step
    newLoglikelihood = plsa_model.LogLikelihood()
    print("[", time.strftime('%Y-%m-%d %H:%M:%S', time.localtime(time.time())), "]"
          "iteration", str(newLoglikelihood))
    # you should see increasing loglikelihood
    if(abs(newLoglikelihood - oldLoglikelihood) < threshold):
        # change to absolute value, or else it terminates after first iteration
        break
    oldLoglikelihood = newLoglikelihood

plsa_model.output(docTopicDist, topicWordDist, dictionary, topicWords, topicWord

```

```

[ 2020-12-11 16:16:58 ] 1 iteration -153054.3762461547
[ 2020-12-11 16:17:01 ] 2 iteration -151543.16479613347
[ 2020-12-11 16:17:05 ] 3 iteration -149823.75824213616
[ 2020-12-11 16:17:08 ] 4 iteration -148058.48170336388
[ 2020-12-11 16:17:12 ] 5 iteration -146476.6381192339
[ 2020-12-11 16:17:15 ] 6 iteration -145186.0633160714
[ 2020-12-11 16:17:18 ] 7 iteration -144163.51895091266
[ 2020-12-11 16:17:22 ] 8 iteration -143369.1620027521
[ 2020-12-11 16:17:25 ] 9 iteration -142768.786576884
[ 2020-12-11 16:17:29 ] 10 iteration -142308.06262680344
[ 2020-12-11 16:17:32 ] 11 iteration -141955.95563714925
[ 2020-12-11 16:17:35 ] 12 iteration -141684.9463008054
[ 2020-12-11 16:17:39 ] 13 iteration -141473.74246865415
[ 2020-12-11 16:17:42 ] 14 iteration -141309.934820173
[ 2020-12-11 16:17:46 ] 15 iteration -141186.01301698684
[ 2020-12-11 16:17:49 ] 16 iteration -141088.47081243386
[ 2020-12-11 16:17:53 ] 17 iteration -141005.52966142842
[ 2020-12-11 16:17:56 ] 18 iteration -140938.28897460527
[ 2020-12-11 16:17:59 ] 19 iteration -140891.1357450914
[ 2020-12-11 16:18:03 ] 20 iteration -140855.5766698348

```

Question & Analysis

1. Please indicate whether you implemented based the given code or from scratch.
2. Choose different K (number of topics) in `plsa.py`. What is your option for a reasonable K in `dataset1.txt` and `dataset2.txt`? Give your results of 10 words under each topic by filling in the following table (suppose you set $K = 4$).

For dataset 1:

Topic 1	Topic 2	Topic 3	Topic 4
<i>your words</i>	<i>your words</i>	<i>your words</i>	<i>your words</i>

For dataset 2:

Topic 1	Topic 2	Topic 3	Topic 4
<i>your words</i>	<i>your words</i>	<i>your words</i>	<i>your words</i>

1. Are there any similarities between pLSA and GMM model? Briefly explain your thoughts.

2. What are the disadvantages of pLSA? Consider its generalizing ability to new unseen document and its parameter complexity, etc.

Answers are given below these code blocks.

In [20]:

```
# Try dataset1 with different K values
datasetFilePath = './plsa/dataset/dataset1.txt'

from plsa.plsa import PLSA
from plsa.utils import preprocessing

N, M, word2id, id2word, X = preprocessing(datasetFilePath, stopwordsFilePath)

for K in [2, 3, 4, 5]:
    plsa_model = PLSA()
    plsa_model.initialize(N, K, M, word2id, id2word, X)

    oldLoglikelihood = 1
    newLoglikelihood = 1

    for i in range(0, maxIteration):
        plsa_model.EStep() #implement E step
        plsa_model.MStep() #implement M step
        newLoglikelihood = plsa_model.LogLikelihood()
        if(abs(newLoglikelihood - oldLoglikelihood) < threshold):
            # change to absolute value, or else it terminates after first iteration
            break
        oldLoglikelihood = newLoglikelihood
    print(f"[{time.strftime('%Y-%m-%d %H:%M:%S',time.localtime(time.time()))}] "
          f"Log likelihood for K={K}: {oldLoglikelihood}")

    plsa_model.output(docTopicDist, topicWordDist, dictionary, topicWords, topicWords)

    i = 1
    print(f"| Topic | Words |\n|-----|-----|")
    with open(topicWords, 'r') as f:
        for line in f:
            words = line.split(' ')
            non_punct_words = []
            for w in words:
                if len(non_punct_words) == 10:
                    break
                if w.isalpha(): # exclude '' and ``
                    non_punct_words.append(w)

            print(f"|{i}| | {non_punct_words}")
            i += 1
```

```
[2020-12-11 16:20:02] Log likelihood for K=2: -7412.032040954314
| Topic | Words |
|-----|-----|
| 1 | ['sea', 'piece', 'devil', 'pirates', 'fruit', 'grand', 'luffy', 'user', 'haki', 'blue']
| 2 | ['luffy', 'crew', 'pirates', 'island', 'straw', 'hat', 'franky', 'dressrosa', 'government', 'robin']
[2020-12-11 16:20:03] Log likelihood for K=3: -7023.921855293514
| Topic | Words |
|-----|-----|
| 1 | ['luffy', 'pirates', 'grand', 'roger', 'sea', 'navy', 'crew', 'haki', 'piece', 'blue']
```

```

|2 | ['devil', 'fruit', 'sea', 'manga', 'user', 'fruits', 'island', 'piece', 'luffy', 'baroque']
|3 | ['crew', 'luffy', 'island', 'pirates', 'straw', 'dressrosa', 'alliance', 'ranky', 'battle', 'government']
[2020-12-11 16:20:04] Log likelihood for K=4: -6697.731847739007
| Topic | Words |
|-----|-----|
|1 | ['luffy', 'pirates', 'island', 'dressrosa', 'crew', 'straw', 'alliance', 'nami', 'fishman', 'captain']
|2 | ['crew', 'luffy', 'pirates', 'island', 'piece', 'manga', 'government', 'grand', 'straw', 'war']
|3 | ['sea', 'devil', 'fruit', 'grand', 'user', 'fruits', 'called', 'water', 'mountain', 'burū']
|4 | ['haki', 'pirates', 'zou', 'baroque', 'alabasta', 'luffy', 'color', 'island', 'jack', 'zunisha']
[2020-12-11 16:20:05] Log likelihood for K=5: -6550.837091607988
| Topic | Words |
|-----|-----|
|1 | ['luffy', 'dressrosa', 'ace', 'crew', 'alliance', 'navy', 'pirate', 'series', 'doflamingo', 'defeat']
|2 | ['haki', 'grand', 'sea', 'government', 'red', 'blue', 'mountain', 'burū', 'animals', 'color']
|3 | ['devil', 'fruit', 'pirates', 'crew', 'luffy', 'user', 'fruits', 'sea', 'straw', 'alabasta']
|4 | ['luffy', 'pirates', 'manga', 'island', 'piece', 'fishman', 'crew', 'captain', 'arlong', 'nami']
|5 | ['island', 'zou', 'pose', 'magnetic', 'log', 'set', 'mom', 'charlotte', 'pirates', 'luffy']

```

In [21]:

```

# Repeat with dataset 2

# input file, output files and parameters
datasetFilePath = './plsa/dataset/dataset2.txt'

from plsa.plsa import PLSA
from plsa.utils import preprocessing

N, M, word2id, id2word, X = preprocessing(datasetFilePath, stopwordsFilePath)

# Try dataset1 with different K values
for K in [10, 20, 30]:
    plsa_model = PLSA()
    plsa_model.initialize(N, K, M, word2id, id2word, X)

    oldLoglikelihood = 1
    newLoglikelihood = 1

    for i in range(0, maxIteration):
        print(f"{i},", end="", flush=True)
        plsa_model.EStep() #implement E step
        plsa_model.MStep() #implement M step
        newLoglikelihood = plsa_model.LogLikelihood()
        if(abs(newLoglikelihood - oldLoglikelihood) < threshold):
            # change to absolute value, or else it terminates after first iteration
            break
        oldLoglikelihood = newLoglikelihood
    print(f"\n[{time.strftime('%Y-%m-%d %H:%M:%S', time.localtime(time.time()))}]
          f"Log likelihood for K={K}: {oldLoglikelihood}")

    plsa_model.output(docTopicDist, topicWordDist, dictionary, topicWords, topic
    i = 1

```

```

print(f"| Topic | Words |\n|-----|-----|")
with open(topicWords, 'r') as f:
    for line in f:
        words = line.split(' ')
        non_punct_words = []
        for w in words:
            if len(non_punct_words) == 10:
                break
            if w.isalpha():
                non_punct_words.append(w)

        print(f"|{i} | {non_punct_words}")
        i += 1

```

```

0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,
[2020-12-11 16:34:44] Log likelihood for K=10: -128149.27220865787
| Topic | Words |
|-----|-----|
| 1 | ['soviet', 'bush', 'fire', 'gorbachev', 'president', 'mexico', 'pope', 'mil
itary', 'church', 'officials']
| 2 | ['noriega', 'dukakis', 'jackson', 'california', 'officials', 'president',
'economic', 'panama', 'north', 'kim']
| 3 | ['rating', 'greyhound', 'gorbachev', 'warming', 'union', 'global', 'study',
'leaders', 'officials', 'percent']
| 4 | ['central', 'snow', 'fbi', 'season', 'agents', 'northern', 'southern', 'exp
ected', 'heavy', 'inches']
| 5 | ['roberts', 'magellan', 'people', 'spacecraft', 'city', 'contact', 'polic
e', 'national', 'pictures', 'tribble']
| 6 | ['barry', 'union', 'soviet', 'officers', 'moore', 'polish', 'friday', 'plan
t', 'immigration', 'people']
| 7 | ['bush', 'campaign', 'oil', 'police', 'waste', 'peres', 'company', 'flori
o', 'people', 'exxon']
| 8 | ['administration', 'farmer', 'police', 'people', 'settlements', 'school',
'receptor', 'occupied', 'government', 'shamir']
| 9 | ['bank', 'company', 'president', 'central', 'time', 'children', 'congress',
'nikolais', 'germany', 'jews']
| 10 | ['percent', 'prices', 'rate', 'report', 'rose', 'business', 'economy', 'du
racell', 'month', 'billion']
0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,
[2020-12-11 16:35:56] Log likelihood for K=20: -118445.30428615403
| Topic | Words |
|-----|-----|
| 1 | ['kim', 'north', 'receptor', 'rose', 'roh', 'businesses', 'scientists', 'op
position', 'korean', 'visit']
| 2 | ['people', 'immigration', 'administration', 'farmer', 'farm', 'percent', 'c
ity', 'nelson', 'enforcement', 'border']
| 3 | ['company', 'bar', 'multistate', 'ferrets', 'doctors', 'percent', 'question
s', 'stock', 'suit', 'dorrance']
| 4 | ['school', 'pope', 'church', 'boy', 'mexico', 'animals', 'teacher', 'peopl
e', 'police', 'salinas']
| 5 | ['roberts', 'greyhound', 'people', 'drivers', 'union', 'company', 'injure
d', 'fear', 'bombs', 'violence']
| 6 | ['bush', 'union', 'president', 'campaign', 'people', 'batalla', 'cordero',
'maxwell', 'societe', 'generale']
| 7 | ['bank', 'fire', 'duracell', 'monday', 'billion', 'company', 'kraft', 'summ
er', 'record', 'sales']
| 8 | ['waste', 'police', 'county', 'orr', 'people', 'officers', 'nelson', 'sit
e', 'threat', 'friday']
| 9 | ['rating', 'percent', 'films', 'mpaa', 'military', 'combe', 'system', 'pres
ident', 'government', 'inventories']
| 10 | ['noriega', 'magellan', 'spacecraft', 'panama', 'jackson', 'contact', 'sau
di', 'american', 'forces', 'pictures']
| 11 | ['soviet', 'gorbachev', 'florio', 'polish', 'study', 'plant', 'exxon', 'of
ficers', 'arco', 'union']

```

```

|12 | ['grain', 'iraq', 'baker', 'iran', 'saudi', 'united', 'guerrillas', 'syri
a', 'warhol', 'zennoh']
|13 | ['central', 'snow', 'embassy', 'washington', 'northern', 'wednesday', 'sou
thern', 'inches', 'nation', 'degrees']
|14 | ['percent', 'oil', 'prices', 'gas', 'inflation', 'rate', 'skins', 'septemb
er', 'thursday', 'enron']
|15 | ['percent', 'soviet', 'fbi', 'agents', 'production', 'november', 'octobe
r', 'tribe', 'rate', 'senate']
|16 | ['dukakis', 'elephant', 'jackson', 'national', 'told', 'roberston', 'polic
e', 'front', 'car', 'zoo']
|17 | ['bush', 'gorbachev', 'soviet', 'museum', 'liberace', 'police', 'dukakis',
'bloomberg', 'top', 'president']
|18 | ['warming', 'global', 'children', 'settlements', 'summit', 'nikolais', 'sh
amir', 'occupied', 'environmental', 'forestry']
|19 | ['congress', 'government', 'jews', 'jewish', 'germany', 'berlin', 'peopl
e', 'waldheim', 'officials', 'east']
|20 | ['barry', 'california', 'peres', 'moore', 'official', 'rappaport', 'bechte
l', 'offer', 'israel', 'mundy']
0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,
[2020-12-11 16:37:18] Log likelihood for K=30: -111812.88282304422
| Topic | Words |
|-----|-----|
|1 | ['fire', 'immigration', 'nelson', 'border', 'monday', 'fiscal', 'firefighte
rs', 'forest', 'officials', 'biaggi']
|2 | ['ago', 'president', 'history', 'american', 'died', 'york', 'time', 'minist
er', 'star', 'organization']
|3 | ['warming', 'global', 'summit', 'receptor', 'scientists', 'environmental',
'gene', 'brain', 'forestry', 'leaders']
|4 | ['farmer', 'soviet', 'administration', 'polish', 'people', 'officers', 'uni
on', 'farm', 'percent', 'farmers']
|5 | ['plant', 'arco', 'employees', 'chemical', 'people', 'investigation', 'kill
ed', 'sorgenti', 'workers', 'pounds']
|6 | ['congress', 'jews', 'germany', 'berlin', 'jewish', 'waldheim', 'meeting',
'people', 'israel', 'garbo']
|7 | ['bush', 'dukakis', 'prosperity', 'night', 'peace', 'president', 'campaig
n', 'trade', 'earlier', 'day']
|8 | ['saudi', 'iraq', 'gorbachev', 'forces', 'soviet', 'arabia', 'iran', 'bake
r', 'tanks', 'united']
|9 | ['gorbachev', 'magellan', 'spacecraft', 'study', 'economy', 'soviet', 'cont
act', 'economic', 'pictures', 'earth']
|10 | ['peres', 'official', 'offer', 'rappaport', 'bechtel', 'israel', 'robb',
'oil', 'governor', 'memo']
|11 | ['roberts', 'greyhound', 'union', 'drivers', 'people', 'company', 'negotia
tions', 'violence', 'fear', 'officials']
|12 | ['north', 'kim', 'talks', 'government', 'leader', 'roh', 'visit', 'korea
n', 'opposition', 'rebel']
|13 | ['soviet', 'immigration', 'calgary', 'combe', 'anthrax', 'week', 'mireck
i', 'romanian', 'coach', 'biological']
|14 | ['fbi', 'agents', 'children', 'nikolais', 'bar', 'multistate', 'suit', 'hi
spanic', 'questions', 'program']
|15 | ['embassy', 'national', 'cuban', 'government', 'diplomatic', 'czechoslova
k', 'doctors', 'ferrets', 'people', 'washington']
|16 | ['bank', 'record', 'summer', 'england', 'monday', 'power', 'maxwell', 'bil
lion', 'plan', 'generale']
|17 | ['percent', 'prices', 'inventories', 'july', 'average', 'survey', 'held',
'fell', 'stock', 'month']
|18 | ['rose', 'businesses', 'percent', 'animals', 'elephant', 'soviet', 'firm
s', 'zoo', 'committee', 'europe']
|19 | ['bombs', 'roberston', 'car', 'black', 'killing', 'burt', 'biggest', 'japan
ese', 'injured', 'rest']
|20 | ['waste', 'police', 'liberace', 'grain', 'official', 'tribe', 'bloomber
g', 'museum', 'orr', 'money']
|21 | ['percent', 'california', 'gas', 'oil', 'rate', 'production', 'rose', 'fed
eral', 'november', 'air']
|22 | ['noriega', 'panama', 'settlements', 'jackson', 'israel', 'shamir', 'occup

```



```

ied', 'powell', 'letter', 'committee']
|23 | ['central', 'snow', 'degrees', 'southern', 'northern', 'temperatures', 'in
ches', 'nation', 'wednesday', 'expected']
|24 | ['bush', 'campaign', 'school', 'boy', 'police', 'teacher', 'students', 'fa
mily', 'shot', 'president']
|25 | ['rating', 'system', 'skins', 'films', 'mpaa', 'city', 'japan', 'tie', 'sh
elter', 'mcguire']
|26 | ['pope', 'church', 'mexico', 'salinas', 'people', 'visit', 'constitution',
'government', 'president', 'hurricane']
|27 | ['dukakis', 'florio', 'exxon', 'jackson', 'programs', 'democratic', 'unio
n', 'spill', 'pittsburgh', 'drug']
|28 | ['company', 'duracell', 'kraft', 'kravis', 'percent', 'kohlberg', 'billio
n', 'financial', 'sales', 'close']
|29 | ['season', 'ratings', 'nbc', 'tuesday', 'abc', 'cosby', 'percent', 'serie
s', 'cbs', 'wo']
|30 | ['percent', 'barry', 'prices', 'moore', 'inflation', 'oil', 'batalla', 'mu
ndy', 'mayor', 'september']

```

Your Answers

Question 0:

I implemented PLSA based on the skeleton code.

Question 1:

For dataset1.txt, values of K between 2-5 seem to be reasonable, as some of the documents talk about the background of One Piece, while other documents describe the plot or setting of One Piece. 3 may be an especially well suited number because I see roughly 3 categories of documents: some history/background/reviews about One Piece, then descriptions about its setting, and then its plot. The results are given below:

K=2:

Topic	Words
1	['sea', 'piece', 'devil', 'pirates', 'fruit', 'grand', 'luffy', 'user', 'haki', 'blue']
2	['luffy', 'crew', 'pirates', 'island', 'straw', 'hat', 'franky', 'dressrosa', 'government', 'robin']

K=3:

Topic	Words
1	['luffy', 'pirates', 'grand', 'roger', 'sea', 'navy', 'crew', 'haki', 'piece', 'blue']
2	['devil', 'fruit', 'sea', 'manga', 'user', 'fruits', 'island', 'piece', 'luffy', 'baroque']
3	['crew', 'luffy', 'island', 'pirates', 'straw', 'dressrosa', 'alliance', 'franky', 'battle', 'government']

K=4:

Topic	Words
1	['luffy', 'pirates', 'island', 'dressrosa', 'crew', 'straw', 'alliance', 'nami', 'fishman', 'captain']
2	['crew', 'luffy', 'pirates', 'island', 'piece', 'manga', 'government', 'grand', 'straw', 'war']
3	['sea', 'devil', 'fruit', 'grand', 'user', 'fruits', 'called', 'water', 'mountain', 'burū']
4	['haki', 'pirates', 'zou', 'baroque', 'alabasta', 'luffy', 'color', 'island', 'jack', 'zunisha']

K=5:

Topic	Words
1	['luffy', 'dressrosa', 'ace', 'crew', 'alliance', 'navy', 'pirate', 'series', 'doflamingo', 'defeat']
2	['haki', 'grand', 'sea', 'government', 'red', 'blue', 'mountain', 'burū', 'animals', 'color']
3	['devil', 'fruit', 'pirates', 'crew', 'luffy', 'user', 'fruits', 'sea', 'straw', 'alabasta']
4	['luffy', 'pirates', 'manga', 'island', 'piece', 'fishman', 'crew', 'captain', 'arlong', 'nami']
5	['island', 'zou', 'pose', 'magnetic', 'log', 'set', 'mom', 'charlotte', 'pirates', 'luffy']

For all values of K, it seems like the topics are somewhat different, making these values of K reasonable.

For dataset2.txt, using the same procedure, it seems like K=10, 20, or 30 all work fairly well. The first 10 words in each topic relate to each other. As K increases, the topics become more specific, which makes sense, as increased K means more classes, so pLSA finds finer divisions between the set of documents. As the topics get more specific, it becomes easier to see that the 10 words for that class are all related in some way. The results are given below:

K=10:

Topic	Words
1	['soviet', 'bush', 'fire', 'gorbachev', 'president', 'mexico', 'pope', 'military', 'church', 'officials']
2	['noriega', 'dukakis', 'jackson', 'california', 'officials', 'president', 'economic', 'panama', 'north', 'kim']
3	['rating', 'greyhound', 'gorbachev', 'warming', 'union', 'global', 'study', 'leaders', 'officials', 'percent']
4	['central', 'snow', 'fbi', 'season', 'agents', 'northern', 'southern', 'expected', 'heavy', 'inches']
5	['roberts', 'magellan', 'people', 'spacecraft', 'city', 'contact', 'police', 'national', 'pictures', 'trible']
6	['barry', 'union', 'soviet', 'officers', 'moore', 'polish', 'friday', 'plant', 'immigration', 'people']
7	['bush', 'campaign', 'oil', 'police', 'waste', 'peres', 'company', 'florio', 'people', 'exxon']
8	['administration', 'farmer', 'police', 'people', 'settlements', 'school', 'receptor', 'occupied', 'government', 'shamir']
9	['bank', 'company', 'president', 'central', 'time', 'children', 'congress', 'nikolais', 'germany', 'jews']
10	['percent', 'prices', 'rate', 'report', 'rose', 'business', 'economy', 'duracell', 'month', 'billion']

K=20:

Topic	Words
1	['kim', 'north', 'receptor', 'rose', 'roh', 'businesses', 'scientists', 'opposition', 'korean', 'visit']
2	['people', 'immigration', 'administration', 'farmer', 'farm', 'percent', 'city', 'nelson', 'enforcement', 'border']
3	['company', 'bar', 'multistate', 'ferrets', 'doctors', 'percent', 'questions', 'stock', 'suit', 'dorrance']
4	['school', 'pope', 'church', 'boy', 'mexico', 'animals', 'teacher', 'people', 'police', 'salinas']
5	['roberts', 'greyhound', 'people', 'drivers', 'union', 'company', 'injured', 'fear', 'bombs', 'violence']

Topic	Words
6	['bush', 'union', 'president', 'campaign', 'people', 'batalla', 'cordero', 'maxwell', 'societe', 'generale']
7	['bank', 'fire', 'duracell', 'monday', 'billion', 'company', 'kraft', 'summer', 'record', 'sales']
8	['waste', 'police', 'county', 'orr', 'people', 'officers', 'nelson', 'site', 'threat', 'friday']
9	['rating', 'percent', 'films', 'mpaa', 'military', 'combe', 'system', 'president', 'government', 'inventories']
10	['noriega', 'magellan', 'spacecraft', 'panama', 'jackson', 'contact', 'saudi', 'american', 'forces', 'pictures']
11	['soviet', 'gorbachev', 'florio', 'polish', 'study', 'plant', 'exxon', 'officers', 'arco', 'union']
12	['grain', 'iraq', 'baker', 'iran', 'saudi', 'united', 'guerrillas', 'syria', 'warhol', 'zennoh']
13	['central', 'snow', 'embassy', 'washington', 'northern', 'wednesday', 'southern', 'inches', 'nation', 'degrees']
14	['percent', 'oil', 'prices', 'gas', 'inflation', 'rate', 'skins', 'september', 'thursday', 'enron']
15	['percent', 'soviet', 'fbi', 'agents', 'production', 'november', 'october', 'tribe', 'rate', 'senate']
16	['dukakis', 'elephant', 'jackson', 'national', 'told', 'roberson', 'police', 'front', 'car', 'zoo']
17	['bush', 'gorbachev', 'soviet', 'museum', 'liberace', 'police', 'dukakis', 'bloomberg', 'top', 'president']
18	['warming', 'global', 'children', 'settlements', 'summit', 'nikolais', 'shamir', 'occupied', 'environmental', 'forestry']
19	['congress', 'government', 'jews', 'jewish', 'germany', 'berlin', 'people', 'waldheim', 'officials', 'east']
20	['barry', 'california', 'peres', 'moore', 'official', 'rappaport', 'bechtel', 'offer', 'israel', 'mundy']

K=30:

Topic	Words
1	['fire', 'immigration', 'nelson', 'border', 'monday', 'fiscal', 'firefighters', 'forest', 'officials', 'biaggi']
2	['ago', 'president', 'history', 'american', 'died', 'york', 'time', 'minister', 'star', 'organization']
3	['warming', 'global', 'summit', 'receptor', 'scientists', 'environmental', 'gene', 'brain', 'forestry', 'leaders']
4	['farmer', 'soviet', 'administration', 'polish', 'people', 'officers', 'union', 'farm', 'percent', 'farmers']
5	['plant', 'arco', 'employees', 'chemical', 'people', 'investigation', 'killed', 'sorgenti', 'workers', 'pounds']
6	['congress', 'jews', 'germany', 'berlin', 'jewish', 'waldheim', 'meeting', 'people', 'israel', 'garbo']
7	['bush', 'dukakis', 'prosperity', 'night', 'peace', 'president', 'campaign', 'trade', 'earlier', 'day']
8	['saudi', 'iraq', 'gorbachev', 'forces', 'soviet', 'arabia', 'iran', 'baker', 'tanks', 'united']
9	['gorbachev', 'magellan', 'spacecraft', 'study', 'economy', 'soviet', 'contact', 'economic', 'pictures', 'earth']
10	['peres', 'official', 'offer', 'rappaport', 'bechtel', 'israel', 'robb', 'oil', 'governor', 'memo']
11	['roberts', 'greyhound', 'union', 'drivers', 'people', 'company', 'negotiations', 'violence', 'fear', 'officials']
12	['north', 'kim', 'talks', 'government', 'leader', 'roh', 'visit', 'korean', 'opposition', 'rebel']
13	['soviet', 'immigration', 'calgary', 'combe', 'anthrax', 'week', 'mirecki', 'romanian', 'coach', 'biological']
14	['fbi', 'agents', 'children', 'nikolais', 'bar', 'multistate', 'suit', 'hispanic', 'questions', 'program']

Topic	Words
15	['embassy', 'national', 'cuban', 'government', 'diplomatic', 'czechoslovak', 'doctors', 'ferrets', 'people', 'washington']
16	['bank', 'record', 'summer', 'england', 'monday', 'power', 'maxwell', 'billion', 'plan', 'generale']
17	['percent', 'prices', 'inventories', 'july', 'average', 'survey', 'held', 'fell', 'stock', 'month']
18	['rose', 'businesses', 'percent', 'animals', 'elephant', 'soviet', 'firms', 'zoo', 'committee', 'europe']
19	['bombs', 'roberson', 'car', 'black', 'killing', 'burt', 'biggest', 'japanese', 'injured', 'rest']
20	['waste', 'police', 'liberace', 'grain', 'official', 'tribe', 'bloomberg', 'museum', 'orr', 'money']
21	['percent', 'california', 'gas', 'oil', 'rate', 'production', 'rose', 'federal', 'november', 'air']
22	['noriega', 'panama', 'settlements', 'jackson', 'israel', 'shamir', 'occupied', 'powell', 'letter', 'committee']
23	['central', 'snow', 'degrees', 'southern', 'northern', 'temperatures', 'inches', 'nation', 'wednesday', 'expected']
24	['bush', 'campaign', 'school', 'boy', 'police', 'teacher', 'students', 'family', 'shot', 'president']
25	['rating', 'system', 'skins', 'films', 'mpaa', 'city', 'japan', 'tie', 'shelter', 'mcguire']
26	['pope', 'church', 'mexico', 'salinas', 'people', 'visit', 'constitution', 'government', 'president', 'hurricane']
27	['dukakis', 'florio', 'exxon', 'jackson', 'programs', 'democratic', 'union', 'spill', 'pittsburgh', 'drug']
28	['company', 'duracell', 'kraft', 'kravis', 'percent', 'kohlberg', 'billion', 'financial', 'sales', 'close']
29	['season', 'ratings', 'nbc', 'tuesday', 'abc', 'cosby', 'percent', 'series', 'cbs', 'wo']
30	['percent', 'barry', 'prices', 'moore', 'inflation', 'oil', 'batalla', 'mundy', 'mayor', 'september']

Question 2:

There are indeed similarities between pLSA and GMM. Both are probabilistic clustering models that aim to find clusters present in unlabeled data. Instead of being a hard assignment of data points or text to a cluster, they return a probability instead. Both use the expectation-maximization algorithm to arrive at a local minimum for log-likelihood (loss). The main difference is that pLSA is for clustering text data, while GMM is for clustering vector data. Also, GMM assumes a Gaussian distribution, while pLSA assumes a categorical or multinomial distribution.

Question 3:

One disadvantage is a slow runtime: `dataset2.txt` isn't that big of a text file (only ~750 KB) but the algorithm takes an extremely long time to run. One has to iterate over all the documents, all the clusters, and all the words in the vocabulary to calculate log likelihood and the model parameters. Also, like GMM, there is no guarantee that pLSA will find a global optimum; it will only find a local optimum. Another con is you need to specify the value of K before training the model. It also relies on the assumption of conditional independence, like naive Bayes, which may not hold in the real world. Finally, it has a large number of parameters ($P(z|d)$, $P(w|z)$, and $P(z|d, w)$ are 2D or 3D matrices), and so tuning all of these takes quite some time (see the slow runtime point above). This makes it easy for pLSA to overfit. pLSA also doesn't have a well-defined generative model, so it does poorly in generalizing to new documents that it wasn't trained on.

Bonus Questions (10 points): LDA

We've learned document and topic modeling techniques. As mentioned in the lecture, most frequently used topic models are pLSA and LDA. [Latent Dirichlet allocation \(LDA\)](#) proposed by David M. Blei, Andrew Y. Ng, and Michael I. Jordan, posits that each document is generated as a mixture of topics where the continuous-valued mixture proportions are distributed as a latent Dirichlet random variable.

In this question, please read the paper and/or tutorials of LDA and finish the following questions and tasks:

(1) What are the differences between pLSA and LDA? List at least one advantage of LDA over pLSA?

(2) Show a demo of LDA with brief result analysis on any corpus and discuss what real-world applications can be supported by LDA. Note: You do not need to implement LDA algorithms from scratch. You may use multiple packages such as `nltk`, `gensim`, `pyLDavis` (added on the `cs145hw6.yml`) to help show the demo within couple of lines of code. If you'd like to use other packages, feel free to install them.

In [23]:

```
# NOTE: ' and ` were added to stopwords.dic so that they
# wouldn't be included as words in the processed documents.
import nltk
import gensim
from gensim.models.ldamodel import LdaModel
from plsa.plsa import PLSA
from plsa.utils import preprocessing

# input file, output files and parameters
K = 10
datasetFilePath = './plsa/dataset/dataset2.txt'
N, M, word2id, id2word, X = preprocessing(datasetFilePath, stopwordsFilePath) #

common_corpus = []
for doc in X:
    common_corpus.append([])
    for i, count in enumerate(doc):
        common_corpus[-1].append((i, count))

lda = LdaModel(common_corpus, num_topics=K, id2word=id2word)
print(lda)

for t in lda.show_topics(num_topics=K, num_words=20):
    print(f"Topic {t[0]}:")
    tuples = t[1].split(' + ')
    for t in tuples:
        word_prob = t.split('*')
        print(f"[{word_prob[1][1:-1]}, {word_prob[0]}] ", end="")
    print("\n")
```

```
LdaModel(num_terms=6451, num_topics=10, decay=0.5, chunksize=2000)
Topic 0:
[u.s., 0.004] [percent, 0.004] [soviet, 0.003] [people, 0.003] [government, 0.003] [company, 0.003] [officials, 0.003] [police, 0.002] [magellan, 0.002] [dog, 0.002] [president, 0.002] [spacecraft, 0.002] [york, 0.002] [economic, 0.002] [s
```

cientists, 0.002] [washington, 0.002] [california, 0.002] [wednesday, 0.002] [expected, 0.002] [bush, 0.002]

Topic 1:

[president, 0.003] [plant, 0.003] [monday, 0.002] [arco, 0.002] [federal, 0.002] [officials, 0.002] [california, 0.002] [people, 0.002] [york, 0.002] [percent, 0.002] [school, 0.002] [air, 0.002] [employees, 0.002] [u.s., 0.002] [rose, 0.002] [tuesday, 0.002] [county, 0.002] [saudi, 0.002] [warming, 0.002] [enforcement, 0.002]

Topic 2:

[police, 0.005] [bank, 0.004] [mrs., 0.004] [children, 0.003] [soviet, 0.003] [officials, 0.003] [percent, 0.003] [ferrets, 0.002] [doctors, 0.002] [people, 0.002] [bloomberg, 0.002] [museum, 0.002] [gorbachev, 0.002] [roberts, 0.002] [center, 0.002] [union, 0.002] [city, 0.002] [national, 0.002] [liberation, 0.002] [killed, 0.002]

Topic 3:

[people, 0.006] [president, 0.004] [percent, 0.004] [soviet, 0.003] [bush, 0.003] [government, 0.003] [u.s., 0.003] [officials, 0.002] [city, 0.002] [military, 0.002] [administration, 0.002] [gorbachev, 0.002] [dukakis, 0.002] [time, 0.002] [prices, 0.002] [union, 0.002] [told, 0.002] [company, 0.002] [reported, 0.002] [national, 0.002]

Topic 4:

[percent, 0.007] [u.s., 0.004] [oil, 0.004] [people, 0.003] [official, 0.003] [soviet, 0.003] [union, 0.002] [government, 0.002] [wednesday, 0.002] [day, 0.002] [central, 0.002] [company, 0.002] [prices, 0.002] [degrees, 0.002] [rate, 0.002] [monday, 0.002] [peres, 0.002] [war, 0.002] [gas, 0.002] [northern, 0.002]

Topic 5:

[company, 0.005] [duracell, 0.004] [people, 0.004] [percent, 0.004] [bush, 0.004] [president, 0.003] [government, 0.003] [official, 0.003] [rating, 0.003] [thursday, 0.003] [u.s., 0.003] [billion, 0.003] [kraft, 0.003] [york, 0.003] [ago, 0.002] [city, 0.002] [kravis, 0.002] [administration, 0.002] [news, 0.002] [called, 0.002]

Topic 6:

[percent, 0.005] [president, 0.004] [people, 0.003] [black-owned, 0.003] [saudi, 0.003] [u.s., 0.002] [economic, 0.002] [businesses, 0.002] [country, 0.002] [iraq, 0.002] [central, 0.002] [government, 0.002] [officials, 0.002] [north, 0.002] [dukakis, 0.002] [black, 0.002] [business, 0.002] [school, 0.002] [prices, 0.002] [nation, 0.002]

Topic 7:

[percent, 0.010] [soviet, 0.004] [economic, 0.003] [people, 0.003] [bush, 0.003] [gorbachev, 0.003] [bank, 0.003] [fire, 0.003] [economy, 0.003] [national, 0.003] [u.s., 0.003] [tuesday, 0.003] [president, 0.003] [government, 0.003] [report, 0.002] [monday, 0.002] [dukakis, 0.002] [israel, 0.002] [month, 0.002] [stock, 0.002]

Topic 8:

[u.s., 0.004] [barry, 0.004] [officials, 0.004] [people, 0.003] [soviet, 0.003] [government, 0.003] [city, 0.003] [friday, 0.003] [president, 0.003] [police, 0.003] [moore, 0.003] [national, 0.002] [wednesday, 0.002] [central, 0.002] [thursday, 0.002] [north, 0.002] [visit, 0.002] [fbi, 0.002] [mayor, 0.002] [southern, 0.002]

Topic 9:

[percent, 0.015] [u.s., 0.004] [rate, 0.003] [noriega, 0.003] [rose, 0.003] [people, 0.003] [production, 0.003] [ago, 0.002] [october, 0.002] [november, 0.002] [president, 0.002] [officials, 0.002] [bank, 0.002] [manufacturing, 0.002] [month, 0.002] [oil, 0.002] [jackson, 0.002] [operating, 0.002] [global, 0.002] [government, 0.002]

Your Answers

Question 1:

The differences between LDA and pLSA is that LDA is a fully generative probabilistic model, while pLSA is not. One cannot easily use pLSA to assign the probability of a class to a new unseen document, but one can in LDA. This is because LDA is modeled on a hidden Dirichlet variable that serves as a distribution of topics to documents, while pLSA is not. As a result of these changes, one advantage of LDA over pSLA is that LDA can better predict topics for new documents and suffers less from the overfitting issues that pSLA exhibits. However, it still retains the issue of needing to pick K before running the model.

Question 2:

For the demo, see code above. From the results, which were run for 10 topics, I observed some interesting results. Terms such as "people", "u.s.," and "president" appeared in nearly all of the topics as the words most likely to appear for that topic. This could be because these words appear often in the set of documents as a whole. However, looking past these frequently repeated words, we see that more distinct topics start to emerge: government, democracy, weather, international relations, etc.

LDA can be used for to calculate the likelihood that a document belongs to a particular topic, as well as analyze the mixture of topics present in a document through its different words. It can also be used in biology to analyze alleles that might be commonly carried by a group of individuals, or in recommendation systems to recommend a article to a customer based on his/her previously read articles (they form a "cluster" of sorts).

End of Homework 6 :)

Please printout the Jupyter notebook and relevant code files that you work on and submit only 1 PDF file on GradeScope with page assigned.

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

class NB_model():
    def __init__(self):
        self.pi = {} # to store prior probability of each class
        self.Pr_dict = None
        self.num_vocab = None
        self.num_classes = None

    def fit(self, train_data, train_label, vocab, if_use_smooth=True):
        # get prior probabilities
        self.num_vocab = len(vocab['index'].tolist())
        self.get_prior_prob(train_label)
        # ===== YOUR CODE HERE =====
        # Calculate probability of each word based on class
        # Hint: Store each probability value in matrix or dict: self.Pr_dict[classID][wordID] or Pr_dict[wordID][classID]
        # Remember that there are possible NaN or 0 in Pr_dict matrix/dict. Use smooth method
        self.Pr_dict = {}
        self.zero_word_dict = {}
        classes = set(train_data['classIdx'])
        print(f"Training class ", end="", flush=True)
        for c in classes:
            df_class = train_data.loc[train_data['classIdx'] == c]
            print(f"{c} ", end="", flush=True)
            total = df_class['count'].sum()
            self.Pr_dict[c] = df_class.groupby('wordIdx')['count'].sum().to_dict()
            self.zero_word_dict[c] = 1 / (total + self.num_vocab)
            for k in self.Pr_dict[c]:
                self.Pr_dict[c][k] = (self.Pr_dict[c][k] + 1) / (total + self.num_vocab)
        # =====
        print("\nTraining completed!")

    def predict(self, test_data):
        test_dict = test_data.to_dict() # change dataframe to dict
        new_dict = {}
        prediction = []

        for idx in range(len(test_dict['docIdx'])):
            docIdx = test_dict['docIdx'][idx]
            wordIdx = test_dict['wordIdx'][idx]
            count = test_dict['count'][idx]
            try:
                new_dict[docIdx][wordIdx] = count
            except:
                new_dict[test_dict['docIdx'][idx]] = {}
                new_dict[docIdx][wordIdx] = count
            ""

        for docIdx in range(1, len(new_dict)+1):
            if docIdx % 2500 == 0:
                print(f"Finished predictions for doc {docIdx}")
            score_dict = {}
            #Creating a probability row for each class
            for classIdx in range(1, self.num_classes+1):
                score_dict[classIdx] = 0
                # ===== YOUR CODE HERE =====
                ### Implement the score_dict for all classes for each document
                ### Remember to use log addition rather than probability multiplication
                ### Remember to add prior probability, i.e. self.pi
                score_dict[classIdx] = np.log(self.pi[classIdx])
                for w in new_dict[docIdx]:
                    score_dict[classIdx] += np.log(self.Pr_dict[classIdx].get(w, self.zero_word_dict[classIdx]))
                # =====

```



```

65         max_score = max(score_dict, key=score_dict.get)
66         prediction.append(max_score)
67     return prediction
68
69
70 def get_prior_prob(self, train_label, verbose=True):
71     # print(train_label)
72     unique_class = list(set(train_label))
73     self.num_classes = len(unique_class)
74     total = len(train_label)
75     for c in unique_class:
76         # ===== YOUR CODE HERE =====
77         ### calculate prior probability of each class ###
78         ### Hint: store prior probability of each class in self.pi
79         self.pi[c] = 0
80         for label in train_label:
81             if c == label:
82                 self.pi[c] += 1
83         self.pi[c] = self.pi[c] / total
84         # =====
85
86 if verbose:
87     print("Prior Probability of each class:")
88     print("\n".join("{}: {}".format(k, v) for k, v in self.pi.items()))

```

1	,docIdx,wordIdx,count,words
2	395,4,10,1,atheist
3	396,4,12,20,of
4	397,4,16,1,from
5	398,4,17,10,religion
6	399,4,23,25,and
7	400,4,25,1,other
8	401,4,27,3,are
9	402,4,29,40,the
10	403,4,30,16,in
11	404,4,33,19,to
12	405,4,42,12,it
13	406,4,44,3,like
14	407,4,46,1,christians
15	408,4,48,2,on
16	409,4,51,6,but
17	410,4,52,7,with
18	411,4,54,2,word
19	412,4,60,22,is
20	413,4,67,5,people
21	414,4,72,3,can
22	415,4,73,1,get
23	416,4,81,4,for
24	417,4,83,6,who
25	418,4,84,4,go
26	419,4,85,1,directly
27	420,4,95,1,bible
28	421,4,99,2,so
29	422,4,100,3,one
30	423,4,101,1,such
31	424,4,104,3,by
32	425,4,122,4,or
33	426,4,139,1,an
34	427,4,142,4,which
35	428,4,143,1,may
36	429,4,144,13,be
37	430,4,152,1,humanism
38	431,4,156,3,secular
39	432,4,160,2,they
40	433,4,181,1,humanist
41	434,4,233,12,that
42	435,4,234,1,exists
43	436,4,235,4,all
44	437,4,239,3,any
45	438,4,245,1,well
46	439,4,251,6,this
47	440,4,277,1,example
48	441,4,279,1,anyone
49	442,4,282,1,use
50	443,4,289,2,many
51	444,4,291,2,thought
52	445,4,295,1,his
53	446,4,297,7,at
54	447,4,299,1,very
55	448,4,301,4,he
56	449,4,306,1,rather
57	450,4,307,1,than
58	451,4,309,2,although
59	452,4,310,1,often
60	453,4,312,1,had
61	454,4,314,4,some
62	455,4,316,2,god
63	456,4,340,1,when
64	457,4,343,9,faith
65	458,4,364,3,christianity

66	459,4,388,6,as
67	460,4,393,6,system
68	461,4,394,1,unfortunately
69	462,4,400,1,whose
70	463,4,410,1,premise
71	464,4,416,1,take
72	465,4,420,2,again
73	466,4,425,1,under
74	467,4,426,9,christian
75	468,4,438,2,only
76	469,4,449,1,world
77	470,4,451,1,down
78	471,4,458,1,more
79	472,4,465,1,work
80	473,4,466,3,has
81	474,4,469,1,however
82	475,4,470,2,probably
83	476,4,473,5,if
84	477,4,474,4,you
85	478,4,477,4,what
86	479,4,479,3,different
87	480,4,482,1,sure
88	481,4,489,1,christ
89	482,4,491,1,seems
90	483,4,492,2,even
91	484,4,515,1,university
92	485,4,529,1,belief
93	486,4,531,2,also
94	487,4,536,1,most
95	488,4,550,1,case
96	489,4,551,2,against
97	490,4,555,1,best
98	491,4,563,3,without
99	492,4,574,2,way
100	493,4,575,1,whether
101	494,4,588,1,emphasis
102	495,4,604,1,dictionary
103	496,4,612,1,present
104	497,4,613,4,person
105	498,4,614,1,philosophy
106	499,4,621,1,expressed
107	500,4,623,1,over
108	501,4,630,2,think
109	502,4,644,1,was
110	503,4,646,2,reason
111	504,4,657,1,western
112	505,4,658,1,values
113	506,4,663,2,were
114	507,4,683,1,through
115	508,4,692,1,those
116	509,4,695,2,beyond
117	510,4,722,9,not
118	511,4,731,1,rationalism
119	512,4,742,1,mind
120	513,4,746,1,become
121	514,4,748,1,them
122	515,4,749,6,there
123	516,4,766,3,will
124	517,4,770,1,article
125	518,4,775,4,edu
126	519,4,778,1,writes
127	520,4,779,2,science
128	521,4,793,4,humans
129	522,4,813,5,we
130	523,4,824,3,things
131	524,4,825,2,put

132	525,4,827,1,interesting
133	526,4,828,3,just
134	527,4,832,1,light
135	528,4,834,1,much
136	529,4,838,1,something
137	530,4,845,2, simply
138	531,4,847,1,understanding
139	532,4,849,2,ok
140	533,4,850,1,me
141	534,4,863,3,don
142	535,4,864,1,anyway
143	536,4,870,3,cs
144	537,4,877,2,do
145	538,4,886,2,no
146	539,4,887,1,didn
147	540,4,902,4,good
148	541,4,905,1,here
149	542,4,910,2,admit
150	543,4,912,2,up
151	544,4,914,1,definition
152	545,4,920,1,yes
153	546,4,921,1,now
154	547,4,922,6,have
155	548,4,930,2,out
156	549,4,939,1,within
157	550,4,942,2,should
158	551,4,947,3,define
159	552,4,965,1,fail
160	553,4,968,2,claimed
161	554,4,969,1,your
162	555,4,970,1,claim
163	556,4,977,1, every
164	557,4,982,1,implementation
165	558,4,985,1,seem
166	559,4,987,1,difference
167	560,4,988,1,between
168	561,4,992,2,still
169	562,4,993,1,say
170	563,4,995,1,my
171	564,4,996,2, personal
172	565,4,1003,5,would
173	566,4,1010,2,reasoning
174	567,4,1017,1,same
175	568,4,1018,6,because
176	569,4,1034,1,becomes
177	570,4,1038,2,point
178	571,4,1040,1,least
179	572,4,1044,1,notion
180	573,4,1065,3,someone
181	574,4,1071,2,course
182	575,4,1077,1,too
183	576,4,1101,1,might
184	577,4,1102,1,interpretation
185	578,4,1103,1,wouldn
186	579,4,1123,1, later
187	580,4,1124,1,going
188	581,4,1128,1,result
189	582,4,1130,1,yet
190	583,4,1134,1,kill
191	584,4,1137,1,everyone
192	585,4,1138,1,cause
193	586,4,1145,2,gun
194	587,4,1157,3,always
195	588,4,1160,2,basically
196	589,4,1196,3,doesn
197	590,4,1206,1,fair

198	591,4,1211,1,perhaps
199	592,4,1236,1,kills
200	593,4,1245,3,really
201	594,4,1260,3,prison
202	595,4,1266,1,irrelevant
203	596,4,1268,1,due
204	597,4,1270,1,concern
205	598,4,1282,1,inability
206	599,4,1283,1,totally
207	600,4,1305,1,similarly
208	601,4,1325,1,causes
209	602,4,1330,1,fun
210	603,4,1347,1,presumably
211	604,4,1365,1,free
212	605,4,1425,3,thinking
213	606,4,1452,3,cannot
214	607,4,1463,1,important
215	608,4,1472,1,realize
216	609,4,1485,1,responsibility
217	610,4,1487,1,maybe
218	611,4,1491,1,willing
219	612,4,1507,1,start
220	613,4,1508,1,ask
221	614,4,1552,1,bad
222	615,4,1557,1,problem
223	616,4,1571,1,let
224	617,4,1580,1,past
225	618,4,1608,1,real
226	619,4,1635,1,open
227	620,4,1646,4,themselves
228	621,4,1649,1,merely
229	622,4,1676,1,consider
230	623,4,1678,1,basic
231	624,4,1690,1,nature
232	625,4,1735,1,come
233	626,4,1784,1,condemn
234	627,4,1810,1,harm
235	628,4,1855,1,language
236	629,4,1861,1,room
237	630,4,1863,1,need
238	631,4,1870,1,whatever
239	632,4,1894,1,please
240	633,4,1898,1,universe
241	634,4,1906,2,change
242	635,4,1910,1,further
243	636,4,1924,1,jesus
244	637,4,2030,1,level
245	638,4,2038,1,happened
246	639,4,2045,1,apr
247	640,4,2095,1,event
248	641,4,2111,1,experience
249	642,4,2113,2,individual
250	643,4,2182,1,religions
251	644,4,2220,1,believes
252	645,4,2232,4,beliefs
253	646,4,2254,1,whereas
254	647,4,2444,2,unless
255	648,4,2465,1,words
256	649,4,2534,1,leave
257	650,4,2585,1,neat
258	651,4,2608,1,webster
259	652,4,2633,1,bias
260	653,4,2717,1,himself
261	654,4,2791,3,moment
262	655,4,2807,1,test
263	656,4,2829,2,game

264	657,4,2906,1,dedicated
265	658,4,3023,2,mass
266	659,4,3057,1,careful
267	660,4,3069,1,minded
268	661,4,3140,7,dogma
269	662,4,3153,1,inherently
270	663,4,3167,1,adequately
271	664,4,3183,3,philosopher
272	665,4,3192,1,sets
273	666,4,3289,1,adam
274	667,4,3290,1,john
275	668,4,3291,1,cooper
276	669,4,3292,1,verily
277	670,4,3293,1,laughed
278	671,4,3294,1,weaklings
279	672,4,3295,1,acooper
280	673,4,3296,1,macalstr
281	674,4,3297,1,claws
282	675,4,3469,4,prisoner
283	676,4,3471,3,genocide
284	677,4,3482,1,tradition
285	678,4,3501,1,dangerous
286	679,4,3559,1,correspond
287	680,4,3565,1,extend
288	681,4,3708,2,billions
289	682,4,3903,1,computer
290	683,4,3950,1,nice
291	684,4,3995,1,oriented
292	685,4,4066,1,regards
293	686,4,4088,2,ahead
294	687,4,4094,1,department
295	688,4,4316,1,divisions
296	689,4,4501,1,semitic
297	690,4,4709,1,rationality
298	691,4,4729,1,qualify
299	692,4,4743,3,understandable
300	693,4,4921,1,qualities
301	694,4,5003,1,leaves
302	695,4,5082,1,capable
303	696,4,5149,1,anybody
304	697,4,5488,1,granted
305	698,4,5915,1,sadly
306	699,4,6099,1,intuition
307	700,4,6427,1,suicide
308	701,4,6521,1,amazing
309	702,4,6625,1,bet
310	703,4,6954,1,guarding
311	704,4,7077,1,destroying
312	705,4,7670,4,encourages
313	706,4,7803,1,edt
314	707,4,7808,1,benevolence
315	708,4,7951,1,waco
316	709,4,8012,1,scorn
317	710,4,8142,1,seemingly
318	711,4,8413,1,bold
319	712,4,8811,1,sects
320	713,4,8979,1,moralities
321	714,4,9007,1,operates
322	715,4,9413,1,evaluated
323	716,4,9414,1,difficulty
324	717,4,9445,1,appalling
325	718,4,10050,1,visible
326	719,4,10310,1,colored
327	720,4,10528,1,testing
328	721,4,10775,4,toronto
329	722,4,11227,1,losing

330	723,4,11256,1,tests
331	724,4,11761,1,visiting
332	725,4,12034,1,offers
333	726,4,13012,1,framework
334	727,4,14705,1,skin
335	728,4,15869,1,constrained
336	729,4,17147,1,evaluate
337	730,4,18175,1,poking
338	731,4,19354,1,invested
339	732,4,20593,3,todd
340	733,4,22058,1,chest
341	734,4,22098,2,kelley
342	735,4,22508,1,retaining
343	736,4,31724,1,guise
344	737,4,31760,1,boil
345	738,4,31927,1,therein
346	739,4,36512,1,bullets
347	740,4,38811,1,nuances
348	741,4,41950,1,hypotheses
349	742,4,45754,1,debated
350	743,4,45956,1,pantheism
351	744,4,46815,1,arisen
352	745,4,53416,1,differentiated
353	746,4,53976,2,tgk
354	747,4,53977,1,quelled
355	748,4,53978,1,quintessential
356	749,4,53979,1,philanthropy
357	750,4,53980,3,supernaturalists

```

1 from numpy import zeros, int8, log
2 import numpy as np
3 from pylab import random
4 import sys
5 #import jieba
6 import nltk
7 from nltk.tokenize import word_tokenize
8 import re
9 import time
10 import codecs
11
12 class PLSA(object):
13     def initialize(self, N, K, M, word2id, id2word, X):
14         self.word2id, self.id2word, self.X = word2id, id2word, X
15         self.N, self.K, self.M = N, K, M
16         #  $\theta[i, j] : p(z_j|d_i)$ : 2-D matrix
17         self.theta = random([N, K])
18         #  $\beta[i, j] : p(w_j|z_i)$ : 2-D matrix
19         self.beta = random([K, M])
20         #  $p[i, j, k] : p(z_k|d_i, w_j)$ : 3-D tensor
21         self.p = zeros([N, M, K])
22         for i in range(0, N):
23             normalization = sum(self.theta[i, :])
24             for j in range(0, K):
25                 self.theta[i, j] /= normalization
26
27         for i in range(0, K):
28             normalization = sum(self.beta[i, :])
29             for j in range(0, M):
30                 self.beta[i, j] /= normalization
31
32     def EStep(self):
33         for i in range(0, self.N): # w
34             for j in range(0, self.M): # d
35                 ## ===== YOUR CODE HERE =====
36                 ### for each word in each document, calculate its
37                 ### conditional probability belonging to each topic (update p)
38                 # total = 0
39                 # for k in range(0, self.K):
40                 #     total += self.theta[i, k] * self.beta[k, j]
41
42                 # for k in range(self.K):
43                 #     self.p[i, j, k] = self.theta[i, k] * self.beta[k, j] / total
44
45                 # vectorized version
46                 total = np.matmul(self.theta[i, :], self.beta[:, j].T)
47                 self.p[i, j, :] = self.theta[i, :] * self.beta[:, j] / total
48                 # =====
49
50     def MStep(self):
51         # update beta
52         for k in range(0, self.K):
53             # ===== YOUR CODE HERE =====
54             ### Implement M step 1: given the conditional distribution
55             ### find the parameters that can maximize the expected likelihood (update beta)
56             # denominator = 0
57             # for j in range(self.M):
58             #     # for i in range(self.N):
59             #         denominator += self.p[i, j, k] * self.X[i, j]
60             #     denominator += np.matmul(self.p[:, j, k], self.X[:, j])
61
62             # for j in range(self.M):
63             #     # numerator = 0

```



```

66 # # for i in range(self.N):
67 # #     numerator += self.p[i, j, k] * self.X[i, j]
68 # print(self.p[:,j,k].shape)
69 # print(self.X[:,j].shape)
70 # numerator = np.matmul(self.p[:, j, k], self.X[:, jj])
71 # self.beta[k, jj] = numerator / denominator
72
73 # vectorized version
74 denominator = np.einsum('ij,ij', self.p[:, :, k], self.X)
75 self.beta[k] = np.einsum('i...,i...', self.p[:, :, k], self.X) / denominator
76 # =====
77
78 # update theta
79 for i in range(0, self.N):
80     # ===== YOUR CODE HERE =====
81     ### Implement M step 2: given the conditional distribution
82     ### find the parameters that can maximize the expected likelihood (update theta)
83     # for k in range(self.K):
84     #     # numerator = 0
85     #     # for j in range(self.M):
86     #     #     numerator += self.p[i, j, k] * self.X[i, jj]
87     #     numerator = np.matmul(self.p[i, :, k], self.X[i, :].T)
88     #     self.theta[i, k] = numerator / self.N
89
90     # denominator = 0
91     # for k in range(0, self.K):
92     #     for j in range(0, self.M):
93     #         denominator += self.p[i, j, k] * self.X[i, jj]
94
95     # vectorized version
96     denominator = self.X[i,:].sum()
97     numerator = np.matmul(self.p[i].T, self.X[i])
98     self.theta[i] = numerator / denominator
99     # =====
100
101 # calculate the log likelihood
102 def LogLikelihood(self):
103     loglikelihood = 0
104     # for i in range(0, self.N):
105     #     # for j in range(0, self.M):
106     #     #     # ===== YOUR CODE HERE =====
107     #     #     ### Calculate likelihood function
108     #     #     word_prob = 0
109     #     #     for k in range(0, self.K):
110     #     #         word_prob += self.theta[i, k] * self.beta[k, jj]
111     #     #     word_prob = np.matmul(self.theta[i], self.beta[:,jj].T)
112     #     #     loglikelihood += self.X[i, jj] * np.log(word_prob)
113     #     #     # =====
114     #     product = np.matmul(self.beta.T, self.theta[i])
115     #     # loglikelihood += np.multiply(self.X[i], np.log(product)).sum()
116     #     loglikelihood += np.einsum('i,i', self.X[i], np.log(product))
117
118     # vectorized version
119     product = np.matmul(self.theta, self.beta)
120     loglikelihood += np.einsum('ij,ij', self.X, np.log(product))
121     return loglikelihood
122
123 # output the params of model and top words of topics to files
124 def output(self, docTopicDist, topicWordDist, dictionary, topicWords, topicWordsNum):
125     # document-topic distribution
126     file = codecs.open(docTopicDist, 'w', 'utf-8')
127     for i in range(0, self.N):
128         tmp = ""
129         for j in range(0, self.K):
130             tmp += str(self.theta[i, j]) + ' '
131         file.write(tmp + '\n')

```

```

132 file.close()
133
134 # topic-word distribution
135 file = codecs.open(topicWordDist, 'w', 'utf-8')
136 for i in range(0, self.K):
137     tmp = ""
138     for j in range(0, self.M):
139         tmp += str(self.beta[i, j]) + ' '
140     file.write(tmp + '\n')
141 file.close()
142
143 # dictionary
144 file = codecs.open(dictionary, 'w', 'utf-8')
145 for i in range(0, self.M):
146     file.write(self.id2word[i] + '\n')
147 file.close()
148
149 # top words of each topic
150 file = codecs.open(topicWords, 'w', 'utf-8')
151 for i in range(0, self.K):
152     topicword = []
153     ids = self.beta[i, :].argsort()
154     for j in ids:
155         topicword.insert(0, self.id2word[j])
156     tmp = ""
157     for word in topicword[0:min(topicWordsNum, len(topicword))]:
158         tmp += word + ' '
159     file.write(tmp + '\n')
file.close()

```