

Programming Project 4

Tanmayi Balla

December 10, 2022

Task-1

Implementation Overview

- The main objective of this task is to implement the collapsed Gibbs sampler using Latent Dirichlet Allocation (LDA), i.e., to classify the words in the documents into K topics and create a feature vector with reduced dimensions, to perform classification in Task-2.
- This task also outputs the most frequent words generated from the sampler, for each topic.
- The main function in Task1 is the **GibbsSampler()** with input parameters: **K, alpha, beta, N, words, documents, topics, D, dataset**. Here, K is the number of topics, alpha beta are the Dirichlet parameters, N is 500 (Iterations), words is a list of all the words present in the corpus, and documents is a list containing the document to which words[i] belong. Topics is a randomly generated list, where we allocate random topics for each word in words list. D is a list of the document numbers in the dataset, while dataset is the name of the dataset (can be either 'artificial' or 'newsgroups').
- There is a helper function for task1, that generates all the required lists, that needs to be fed into the GibbsSampler.
- Task-1 has a runtime of around 30 minutes.
- I have run Task-1 many times, and I am presenting the topics.csv file that I felt is more meaningful.
- All the helper functions are placed in utils.py.

Results of Task-1

book two part low detectors
mission hst solar pat net
large question high day time
car ford cars manual speed
time made problems second interested
article edu writes apr
station shuttle launch option design
don clutch shifter sho drive
henry toronto spencer edu
idea george howell big
even don people good point
edu system writes oort
heard diesels put writes
sky edu gif uci ics
oil bill change service back
edu engines writes mustang eliot
insurance edu uiuc geico
engine turbo toyota seat feel
science internet information group
space nasa gov such program

Inference

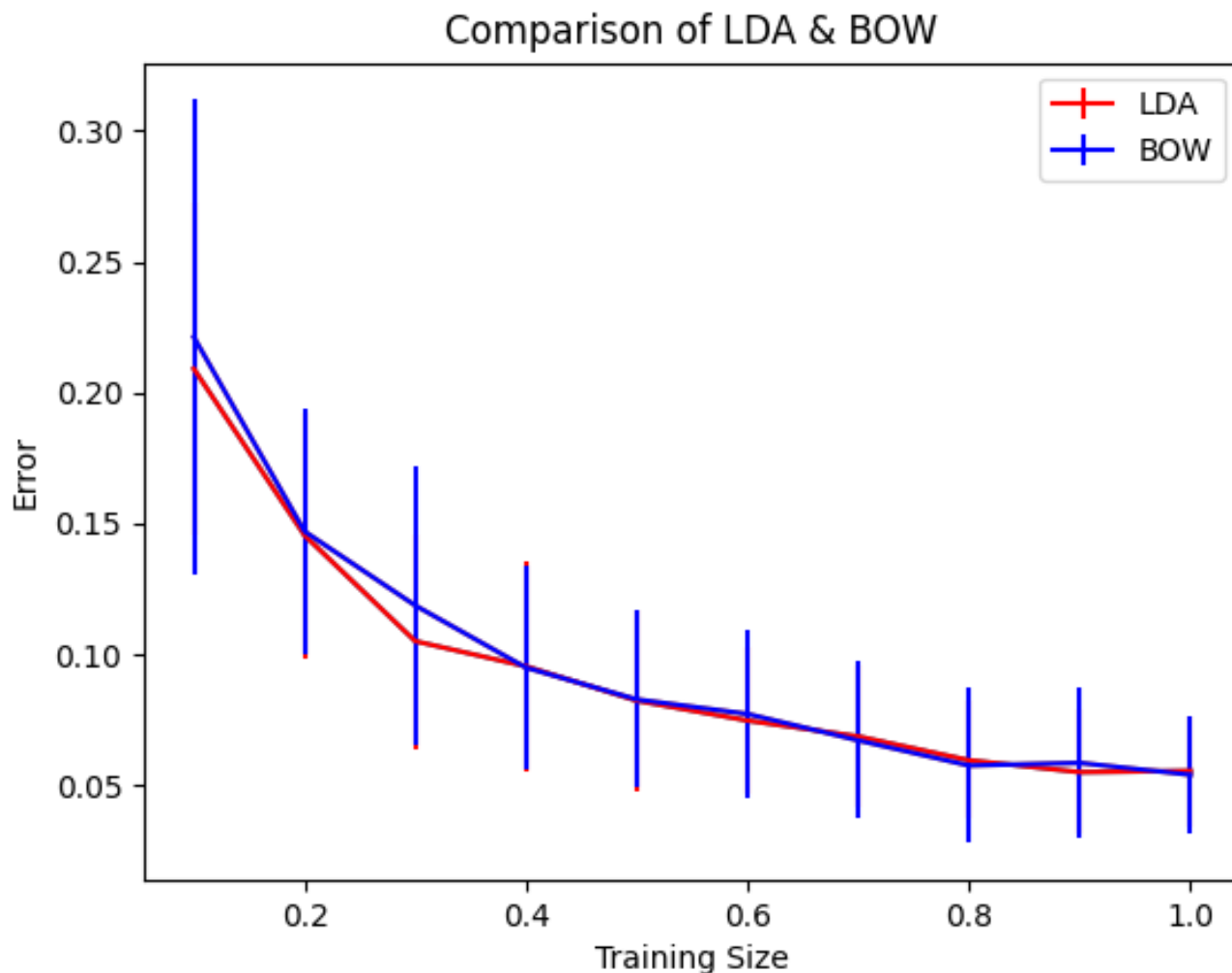
1. The top 5 frequent words of 20 topics are given in the above image. If we closely observe the words for each topics, the combination of words is meaningful (related to news articles). For Example, the last topic is about a space program by nasa; we have a row corresponding to vehicles (engines, toyota, seat, etc.). Though all the rows cannot be considered accurate, the algorithm has managed to give some useful features for classification. Since, the number of features now is reduced to K (no. of topics), the execution would be much faster, while predicting the label in Task-2.

Task-2

Implementation Overview

- The main objective of this task is to compare topic classification for features generated from LDA (Task-1), and Bag-of-words (BOW).
- The features generated from task-1 are stored in f1.npy file. This same file is imported in Task-2 to compare the performance.
- Task-2 is basically a repetition of PP3 Task. It contains all the functions required to generate weights using Newton's method and make predictions using Logistic Regression.
- Task2Helper is used to randomly select the training set and convert it into the respective 10 splits ([0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]). Predictions are generated for all the 10 train splits, and the process is repeated for 30 iterations.
- The mean of the errors for all the training sizes and standard deviation are used to plot the learning curves.
- The standard deviation for both LDA and BOW is very negligible. Hence, the error bars are not properly visible in the plot.

Results



It is evident from the graph that, even with reduced dimensionality, i.e., less number of features, we are able to classify the documents, with good accuracy.

Also, while execution, Bag of words data took higher time as compared to LDA. And, LDA is able to match the accuracy of BOW for all the training sizes. Time taken for LDA is 2.39 seconds, while the runtime of BOW is 77 seconds.

We can observe that the learning curves for both LDA and BOW are similar, and they are reduced with increased training size. This is apparent since a larger dataset provides better learning and therefore better results.

These observations conclude that LDA is a better and faster alternative for document classification.