

Manas Gaur, Amanuel Alambo

Information Retrieval, CS 7800

Project-2 on Text Mining

This project is used to index the 20 newsgroup datasets and test classification and clustering algorithms. There are four activities conducted:

- Feature extraction
- Classification (using four algorithms and the evaluation of each classifier)
- Feature selection
- Document clustering

We walk through how each is implemented and experimental results. The codes in this project are developed and tested using python3

The computational cost of Support Vector Machine (SVM) was very high compared to other machine learning classifiers even on multithreaded execution of the experiments.

1. Feature extraction

How to Run:

```
python3 feature-extract.py directory_of_newsgroups_data feature_definition_file class_definition_file training_data_file
```

Example:

```
python3 feature-extract.py mini_newsgroups feature_definition_file class_definition_file training_data_file.TFIDF
```

Note: the *training_data_file* is required to have an extension of either **‘.TF’** or **‘.IDF’** or **‘.TFIDF’**. The program will throw an error if this format is not followed.

Disclaimer: the three training data files(i.e, TF, IDF, and TF-IDF) are generated in three different executions of the above program. Thus, their entries are different. We tested our classification and clustering algorithms using the **‘.TFIDF’** training dataset.

The feature extraction part of this task is implemented in code **‘feature-extract.py’**. An Index class is defined that is used to index each of the files in each of the newsgroups. We index the ‘subject’ and ‘body of each document. Thus, we needed to parse the subject and body from the document. This parsing is done in method **‘subject_body_index’** of code **‘feature-extract.py’**. Since the subject of the doc is identified by a **‘subject’** in the doc, we used the **‘subject’** entry for parsing. However, we don’t have an explicit identifier for the body for each document. Yet, the **‘Lines’** value in the doc refers to the body of the doc, which is the last **‘Lines’** value. For this, we used the **‘deque’** construct to grab the last **‘Lines’** number of lines and write to a variable. Once the subject and body are extracted from the document, we performed string concatenation and assigned to a variable. Finally, we indexed the concatenated **‘subject+body’** into our index structure. We defined two methods **‘indexed_docs’** and **‘indexed_docs_l’** in class **‘Index’**. The **‘indexed_docs’** method indexes the whole newsgroup datasets and stores in an inverted index (each stemmed term with stopwords removal and punctuation marks removal applied and list of the doc-ids as values). The **‘indexed_docs_l’** method is used to return the doc-id to document content mappings. Further, we implement **‘feature_defn_gen’** and **‘class_defn_gen’** methods for generating **‘feature definition’** and **‘class definition files’**. For generating the feature definition file, the method accepts the feature definition file output name while for generating the class definition file, the corresponding method accepts the output file name for

class definition file and the path to the root directory(i.e., mini_newsgroups). The term frequency, inverse document frequency, or TF-IDF scores are rounded to **‘2’** decimal digits as can be seen in the screenshot below. Three helper function are defined to compute TF-IDF score of a term in a document (functions **‘filecount()’**, **‘nDocs_in_subdir’**, **‘idf_compute’**). Function **‘filecount()’** is used to count the number of documents(files) in each subdirectory, while function **‘nDocs_in_subdir()’** accepts a document-id and calls function **‘filecount()’** to get the total number of documents in a subdirectory the corresponding document belongs to.

Sample snippet of the generated TF-IDF based training data looks like (this is in *libsvm* format):

```
0 5:0.01 87:0.02 111:0.04 142:0.01 186:0.01 195:0.01 261:0.01 297:0.01 314:0.02 316:0.08 466:0.01 680:0.01 711:0.01 779:0.03 791:0.01 835:0.03 846:0.01 892:0.01
908:0.01 931:0.02 1109:0.01 1163:0.01 1412:0.01 1561:0.01 1569:0.01 1578:0.01 1627:0.01 1675:0.01 1692:0.02 1703:0.01 1712:0.01 1714:0.03 1789:0.03 1805:0.01
1848:0.01 1863:0.01 1868:0.03 1936:0.01 1985:0.01 2012:0.01 2039:0.01 2104:0.02 2154:0.01 2303:0.01 2373:0.03 2392:0.01 2521:0.01 2741:0.01 2753:0.01 2811:0.01
2885:0.01 3013:0.01 3107:0.01 3155:0.02 3279:0.06 3358:0.01 3429:0.03 3502:0.02 3562:0.04 3607:0.05 3615:0.01 4059:0.01 4168:0.01 4249:0.01 4260:0.01 4350:0.01
4386:0.01 4422:0.02 4444:0.01 4566:0.01 4586:0.01 4605:0.01 4700:0.01 4719:0.02 4757:0.01 4766:0.01 4800:0.02 4876:0.01 4941:0.01 4957:0.01 5118:0.03 5147:0.01
5157:0.01 5220:0.01 5228:0.05 5234:0.01 5454:0.01 5463:0.01 5638:0.01 5658:0.01 5752:0.01 5790:0.01 5814:0.02 5837:0.01 5961:0.01 6042:0.01 6104:0.01 6126:0.01
6139:0.03 6155:0.02 6333:0.01 6517:0.02 6547:0.01 6610:0.03 6780:0.01 6794:0.01 6975:0.01 7021:0.01 7023:0.01 7082:0.02 7146:0.01 7330:0.17 7357:0.01 7617:0.02
8267:0.03 8413:0.01 8518:0.01 8658:0.01 8711:0.01 8753:0.01 9007:0.01 9015:0.02 9041:0.01 9079:0.01 9087:0.01 9112:0.01 9161:0.01 9223:0.01 9275:0.01 9303:0.01
9422:0.02 9450:0.01 9741:0.01 9747:0.01 9811:0.01 9816:0.01 9917:0.01 9924:0.03 9938:0.01 10012:0.01 10020:0.01 10029:0.02 10183:0.05 10241:0.39 10256:0.01
10307:0.01 10413:0.01 10497:0.01 10513:0.05 10535:0.01 10542:0.19 10603:0.01 10625:0.01 10784:0.02 10809:0.03 10853:0.01 10960:0.01 10986:0.03 10994:0.01
11016:0.02 11066:0.01 11143:0.02 11222:0.01 11251:0.02 11261:0.03 11282:0.01 11290:0.02 11332:0.01 11380:0.01 11385:0.01 11419:0.01 11488:0.01 11490:0.01
11544:0.01 11645:0.04 11728:0.05 11889:0.01 12017:0.01 12066:0.07 12122:0.01 12166:0.02 12198:0.02 12332:0.01 12372:0.01 12404:0.01 12415:0.01 12452:0.01
12492:0.03 12498:0.01 12686:0.01 12691:0.03 12742:0.3 12935:0.01 13065:0.01 13134:0.02 13200:0.01 13210:0.01 13239:0.01 13297:0.01 13365:0.01 13553:0.03
13570:0.01 13631:0.01 13674:0.03 13713:0.04 13808:0.01 13848:0.01 13874:0.01 13907:0.01 13924:0.01 14020:0.01 14021:0.02 14076:0.01 14094:0.01 14154:0.03
14161:0.01 14167:0.01 14317:0.01 14422:0.02 14589:0.02 14664:0.01 14764:0.01 14766:0.01 14790:0.01 14830:0.01 14831:0.01 14860:0.16 14860:0.01 14860:0.01
15025:0.01 15043:0.01 15055:0.01 15090:0.02 15115:0.02 15171:0.06 15300:0.01 15312:0.01 15351:0.09 15387:0.01 15390:0.01 15502:0.01 15513:0.18 15521:0.01
15550:0.01 15672:0.01 15687:0.02 15693:0.01 15710:0.01 15753:0.02 15797:0.01 15907:0.01 15995:0.01 16244:0.01 16264:0.01 16284:0.01 16307:0.02 16316:0.01
16325:0.01 16348:0.01 16381:0.01 16395:0.01 16413:0.01 16656:0.01 16657:0.01 16695:0.01 16771:0.1 16777:0.01 16825:0.03 16834:0.01 16843:0.17 16845:0.01
16873:0.01 16903:0.05 16930:0.01 16944:0.01 16945:0.01 17018:0.01 17034:0.01 17039:0.01 17062:0.01 17090:0.01 17238:0.01 17262:0.03 17296:0.01 17327:0.01
17401:0.01 17434:0.01 17583:0.01 17606:0.01 17617:0.01 17634:0.01 17644:0.01 17843:0.01 17851:0.01 17866:0.01 18013:0.03 18019:0.01 18042:0.01 18254:0.01
18311:0.02 18469:0.01 18596:0.01 18649:0.06 18664:0.01 18747:0.01 18786:0.01 18790:0.03 18837:0.02 18944:0.01 18982:0.01 19130:0.01 19145:0.01 19168:0.01
19175:0.01 19206:0.01 19312:0.03 19523:0.01 19533:0.02 19584:0.01 19662:0.01 19722:0.11 19789:0.01 19812:0.01 19839:0.01 19882:0.01 19900:0.01 20058:0.02
20171:0.01 20245:0.01 20255:0.01 20299:0.01 20399:0.02 20400:0.01 20415:0.01 20515:0.02 20584:0.01 20588:0.01 20788:0.01 20894:0.01 20930:0.01 21041:0.03
21052:0.01 21101:0.01 21206:0.01 21225:0.01 21234:0.01 21258:0.03 21265:0.01 21271:0.01 21286:0.03 21339:0.01 21346:0.03 21349:0.01 21708:0.01
0 5:0.04 111:0.04 269:0.02 316:0.04 466:0.04 675:0.02 841:0.02 871:0.02 892:0.02 1102:0.02 1481:0.02 1712:0.02 1714:0.07 1794:0.02 2015:0.02 2154:0.02 2392:0.02
2482:0.02 2521:0.05 2728:0.02 3013:0.04 3051:0.02 3112:0.02 3279:0.07 3432:0.02 3442:0.02 3607:0.04 3615:0.02 3794:0.02 3909:0.04 4420:0.02 4422:0.02 4566:0.02
5081:0.04 5234:0.02 5454:0.05 5752:0.07 5966:0.02 6139:0.02 6547:0.04 6610:0.02 6780:0.02 7330:0.12 7497:0.02 7565:0.02 7613:0.02 7774:0.02 7805:0.02 8267:0.02
8285:0.05 8711:0.02 8876:0.02 9015:0.05 9112:0.04 9450:0.07 9614:0.02 9811:0.02 9865:0.04 9924:0.02 10029:0.02 10183:0.02 10198:0.02 10241:0.39 10260:0.05
10307:0.02 10420:0.02 10513:0.04 10542:0.12 10603:0.01 10767:0.02 10942:0.05 10986:0.02 11143:0.02 11261:0.02 11380:0.02 11589:0.05 11645:0.04 11728:0.05
11882:0.02 12017:0.02 12066:0.07 12452:0.02 12691:0.02 12742:0.25 12909:0.02 12935:0.02 13200:0.02 13226:0.02 13713:0.04 13820:0.02 13899:0.02 13930:0.02
14151:0.02 14422:0.02 14869:0.02 15090:0.09 15171:0.07 15300:0.02 15351:0.12 15390:0.02 15513:0.35 15617:0.02 15664:0.02 16101:0.02 16250:0.02 16264:0.02
16325:0.02 16348:0.04 16413:0.02 16695:0.02 16707:0.02 16771:0.02 16825:0.02 16834:0.02 16843:0.11 16845:0.02 16896:0.02 16903:0.04 17000:0.02 17034:0.02
17171:0.02 17262:0.02 17334:0.02 17434:0.02 17624:0.02 17661:0.04 18019:0.04 18139:0.02 18537:0.02 18548:0.02 18649:0.04 18837:0.04 19124:0.02 19169:0.04
19206:0.02 19312:0.04 19337:0.07 19722:0.09 19847:0.02 19954:0.04 20422:0.02 20582:0.02 20584:0.02 20930:0.02 21041:0.02 21142:0.02 21258:0.02 21369:0.02
21537:0.02
0 5:0.03 111:0.03 142:0.02 261:0.01 316:0.04 439:0.01 657:0.02 779:0.01 871:0.03 1190:0.01 1223:0.01 1578:0.02 1614:0.02 1789:0.01 1794:0.03 1868:0.06 2015:0.01
2046:0.02 2065:0.03 2104:0.01 2154:0.01 2373:0.03 2410:0.01 2461:0.02 2482:0.01 2521:0.01 2728:0.02 2811:0.02 2815:0.01 2885:0.02 3013:0.01 3155:0.01 3279:0.03
3429:0.04 3432:0.01 3502:0.01 3607:0.02 3615:0.01 3685:0.02 4047:0.01 4112:0.02 4204:0.02 4444:0.03 4605:0.01 4842:0.03 4876:0.01 5150:0.01 5228:0.01 5454:0.03
5580:0.02 5707:0.02 5752:0.02 6104:0.03 6170:0.01 6333:0.01 6408:0.02 6547:0.01 7082:0.02 7176:0.01 7330:0.18 7727:0.02 7774:0.01 8285:0.01 8637:0.03 9015:0.01
9112:0.02 9450:0.01 9526:0.02 9811:0.03 9816:0.01 9924:0.01 9938:0.02 9954:0.02 9965:0.02 10012:0.01 10241:0.34 10333:0.02 10513:0.03 10542:0.08 10603:0.03
10674:0.02 10784:0.01 10809:0.02 10842:0.03 10904:0.02 11016:0.01 11066:0.01 11293:0.02 11332:0.01 11355:0.02 11419:0.02 11490:0.01 11728:0.01 11907:0.01
12066:0.02 12166:0.01 12492:0.02 12686:0.02 12709:0.01 12742:0.17 12797:0.01 12935:0.01 13239:0.01 13553:0.01 13674:0.01 13713:0.02 13829:0.01
13874:0.01 13894:0.02 14021:0.01 14076:0.03 14154:0.03 14161:0.01 14585:0.01 14723:0.02 14725:0.01 14776:0.01 14869:0.08 14960:0.01 15055:0.01 15090:0.01
15115:0.03 15171:0.03 15332:0.02 15351:0.06 15404:0.01 15513:0.23 15550:0.03 15617:0.01 15672:0.02 15881:0.01 16244:0.01 16250:0.01 16307:0.01 16325:0.01
16771:0.03 16843:0.1 16896:0.01 16903:0.03 16922:0.02 17039:0.01 17132:0.02 17296:0.01 18013:0.03 18042:0.01 18649:0.03 18790:0.01 18796:0.02
18837:0.06 19124:0.01 19145:0.05 19206:0.01 19255:0.01 19287:0.02 19312:0.01 19337:0.03 19533:0.01 19722:0.08 20255:0.01 20326:0.01 20422:0.02 20584:0.03
20765:0.01 20802:0.01 20930:0.02 21041:0.02 21142:0.01
```

Figure 1: TF-IDF based training data with labels (target) and feature vectors in libsvm format

As shown, we have an entry of the form (<class label> <feature-id-1:feature-value-1> <feature-id-2:feature-value-2>...<feature-id-n:feature-value-n>).

The result of executing the **‘feature-extract.py’** program with the all the required parameters gives the following output on the screen:

```
amanuel@Bell-Inspiron7:~/WSU-II/Courses/Information_Retrieval/Project_2$ python3 feature-extract.py mini_newsgroups feature_definition_file class_defini
tion_file training_data_file.TFIDF
[nltk_data] Downloading package stopwords to
[nltk_data] /home/amanuel/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /home/amanuel/nltk_data...
[nltk_data] Package punkt is already up-to-date!
Produced feature definition file
Produced class definition file
Produced training data file
amanuel@Bell-Inspiron7:~/WSU-II/Courses/Information_Retrieval/Project_2$
```

Figure 2: Snapshot showing the executing of the feature-extract.py

2. Newsgroups Classification

How to Run:

python3 classification.py

For the classification task, we utilized 4 machine learning models: Multinomial Naive Bayes, Bernoulli Naive Bayes, K-Nearest Neighbors, and Support Vector Machines. The naive bayes classification depends on basic Bayes model for probabilistic representation of bag of words for a document. Theoretically, a variant of the naive bayes classifier maximizes the conditional probability of the class label with respect to the document. It is denoted by $\operatorname{argmax}(P(d|c).P(c))$ which is derived from $\operatorname{argmax} P(c|d)$ after the application of bayes rule. This maximization procedure differentiates bernoulli naive bayes from multinomial naive bayes. In the bernoulli naive bayes, the each word feature of a document is considered to be independent and binary, whereas, in the multinomial, it integer word count for each word feature. It can be said that multinomial naive bayes is a complex model compared to bernoulli and is more vulnerable to overfitting because of requirement of large amount of data. On the other hand, the bernoulli model can easily be trained over small amount of training data.

Other pair of classifiers that can be compared are Support Vector Machines and K-Nearest Neighbors. It is because, both of the model uses distance metric (e.g. euclidean, cosine) for generating a hyperplane (a form of gaussian contour) for classifying the data points. In our study, we utilize these models without any parameter tuning to identify a reasonable classifier for classifying the documents in 20 newsgroups to 6 labels.

Our experiment used 80-20% split for training and test dataset, with a random seed for reproducibility concerns.

We defined four methods corresponding to each classifier (**MultinomialNB**, **BernoulliNB**, **KNeighbors**, and **svm_SVC_classifier**). Plus, we implemented method **'evaluate_classifier()'** to compute the performance of each classifier. Further, we implemented method **'shuffle_train_data()'** to shuffle our dataset when passing into our classifiers. It is needed to avoid bias or sequential pattern in the dataset and improve the predictive performance of machine learning algorithm. File **'shuffled_train_data.txt'** is generated as a result of shuffling method as applied on the original dataset.

Executing the **'classification.py'** yields the following output on the screen:

```
amanuel@IDell-Inspiron7:~/WSU-II/Courses/Information_Retrieval/Project_2$ python3 classification.py
Dimension of feature vectors: (2000, 21724)
Dimension of target vectors: (2000,)
MultinomialNB Results:
Train accuracy: 0.25125
Test accuracy: 0.245
('f1-macro : 0.07 (+/- 0.00)', 'precision-macro : 0.04 (+/- 0.00)', 'recall-macro : 0.17 (+/- 0.00)')
BernoulliNB Results:
Train accuracy: 0.729375
Test accuracy: 0.6375
('f1-macro : 0.52 (+/- 0.05)', 'precision-macro : 0.61 (+/- 0.05)', 'recall-macro : 0.52 (+/- 0.05)')
KNeighbors Results:
Train accuracy: 0.605625
Test accuracy: 0.3925
('f1-macro : 0.34 (+/- 0.03)', 'precision-macro : 0.38 (+/- 0.08)', 'recall-macro : 0.35 (+/- 0.03)')
svm_SVC Results:
Train accuracy: 0.25125
Test accuracy: 0.245
('f1-macro : 0.07 (+/- 0.00)', 'precision-macro : 0.04 (+/- 0.00)', 'recall-macro : 0.17 (+/- 0.00)')
amanuel@IDell-Inspiron7:~/WSU-II/Courses/Information_Retrieval/Project_2$
```

Figure 3: Snapshot of classification.py execution which takes in shuffled feature and target vectors and train 4 classifiers: MultinomialNB, BernoulliNB, K-Nearest Neighbors, and Support Vector Machines

3. Feature selection

The feature selection program is developed to evaluate the performance of the four classifiers defined in program **'classification.py'** for a selected **top k** features. We experiment with different values for k ranging from 100 to 1000. We use chi squared and mutual information for each of the k values (100-1000, 100 is selected at random. Human-in-the-loop could be beneficial for having a starting k.) and pass to our classifiers. We define method **'plot_learning_curve()'** that accepts the model title and a dataframe created in the main method (with data frame entries for 'k' values, Chi-square (χ^2) value, and Mutual-Information (MI) value).

How to run:

Use the following command to run the **'feature_selection.py'** program for a model:

python3 feature_selection.py model_name; where model_name is one of the models passed by the user as argument to the program. For example, to perform feature selection over our newsgroup datasets and pass to model **'BernoulliNB'**, we run the following command:

python3 feature_selection.py BernoulliNB

We ran the feature_selection over our dataset using χ^2 and MI and plotted the results for each classifier for different values of K (top features).

Sample output from command line is shown below:

```

anna@ubuntu:~/Desktop/inspiration7-1/NSU-IT/Courses/Information_Retrieval/Project_25 python3 feature_selection.py BNB
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.62875
Test accuracy: 0.575
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.49375
Test accuracy: 0.465
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.67375
Test accuracy: 0.61
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.535
Test accuracy: 0.47
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.710625
Test accuracy: 0.6625
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.576875
Test accuracy: 0.51
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.72875
Test accuracy: 0.6575
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.60125
Test accuracy: 0.535
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.73875
Test accuracy: 0.6675
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.62125
Test accuracy: 0.545
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.749375
Test accuracy: 0.685
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.651875
Test accuracy: 0.5575
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.7525
Test accuracy: 0.675
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.67875
Test accuracy: 0.585
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.763125
Test accuracy: 0.66
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.698125
Test accuracy: 0.5975
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.76875
Test accuracy: 0.6925
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.709375
Test accuracy: 0.6225
Testing classifiers(BernoulliNB_classifier) with chi2 selected features
Train accuracy: 0.773125
Test accuracy: 0.685
Testing classifiers(BernoulliNB_classifier) with mutual information selected features
Train accuracy: 0.71875

```

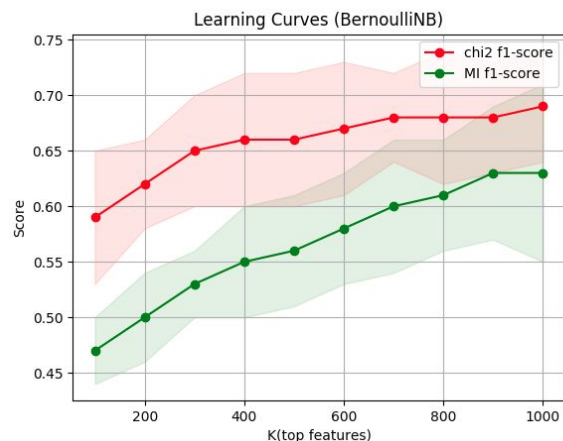


Figure 4: (Left) Execution of feature_selection.py showing the different train and test accuracy of the best performing model (BernoulliNB) for two feature selection methods: χ^2 and MI. (Right) shows the mean F1 score of BernoulliNB for χ^2 and MI. It is important to note the small variance of BernoulliNB compared to models in Figure 5 and also high f1-score.

χ^2 and MI measure is a procedure for selecting features that have relationship with the target variable. It can be seen as a measure of dimensionality reduction, with appropriate dimension is identified by means of accuracy metric. For instance, in our experiment on newsgroup, we observe that after selecting between

800 to 1000 features, which is 4-5% of the total features (~21000), the model is able to achieve a stable and high accuracy compared to process in the absence of feature selection. χ^2 is analogous to distance metric, whereas the MI is analogous to entropy measure (see Figure 4).

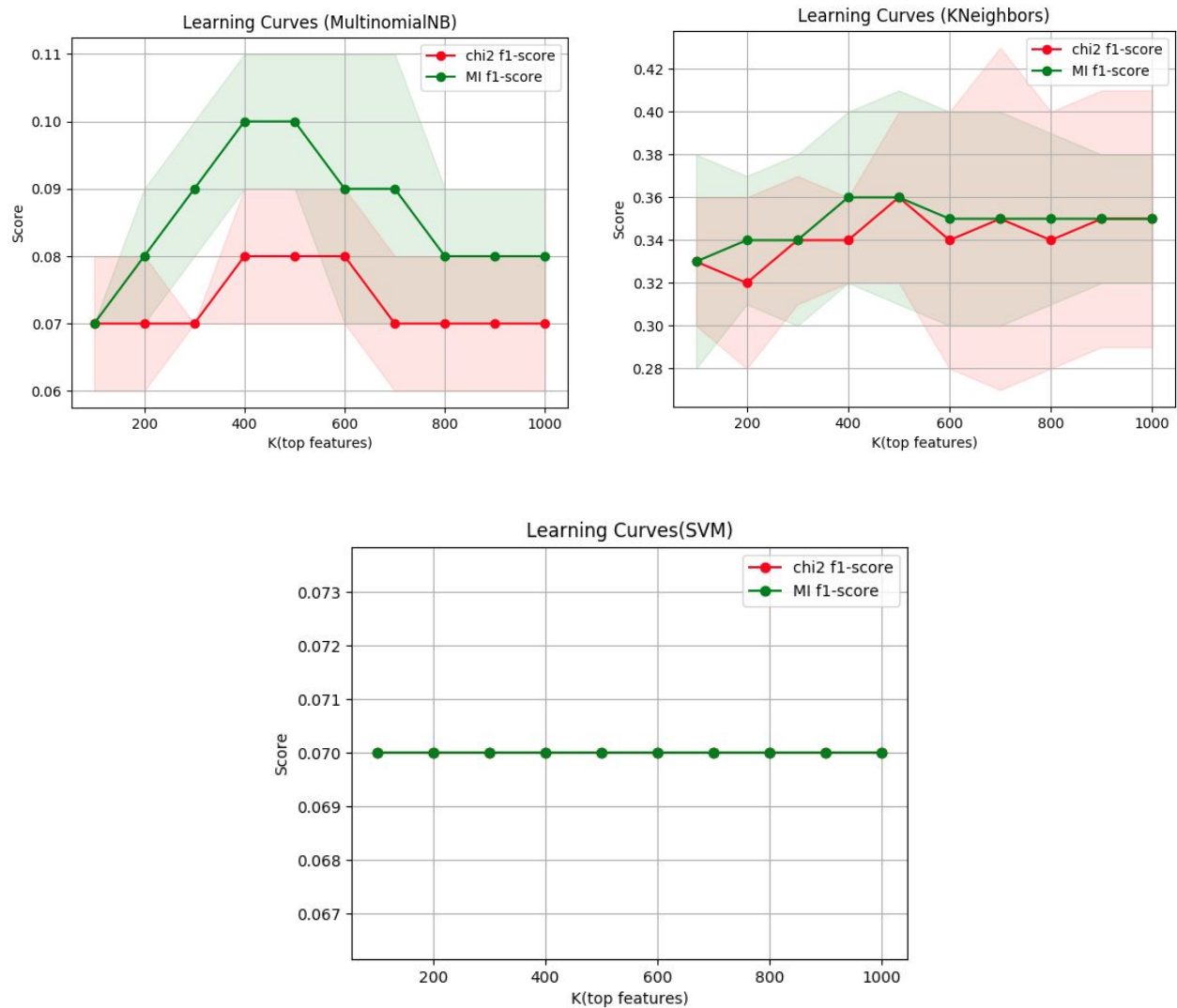


Figure 5: (Top-Left) Trend in F1-Score with changing ‘k’ in the feature selection procedures for the MultinomialNB and (Top-right) K-Nearest Neighbors. It is observable that both the models have shown increased variance in the results compared to BernoulliNB. As a result the models are unstable on the given amount of data. (Bottom-middle) Support Vector Machines did not show any improvement in the F1-score with changing k number of features using feature selection procedures: χ^2 and MI. Interestingly there was no variance in the performance of SVM. All the three plots are performance plots of three machine learning algorithms: MultinomialNB, K-Nearest Neighbors, and SVM.

4. Document Clustering

For the clustering the documents, we utilized **K-Means** and **Agglomerative** clustering algorithms and evaluated the performance using silhouette metric and normalized mutual information. The two clustering algorithms differs in their approach to cluster the feature vectors, wherein, K-Means identifies k number of centroids representing clusters and allocate a data point to a cluster by minimizing the inter-cluster

distance. In case of agglomerative clustering, there is a formation of cluster tree where each node represents a cluster of similar documents. The minimization function is similar to K-Means, however, the k parameter is used to identify the threshold for generating **k-clusters**. In layman terms, a flat line is drawn over the cluster tree, decomposing the tree into k -clusters. Since, both the algorithms are parametric on k , we used silhouette and normalized mutual information (NMI) metric to evaluate the cluster stability. The silhouette metric evaluates the cluster stability without taking into account the target labels, making it an unsupervised metric. On the other hand, the NMI takes into account the target variable, and its interpretation is information-content centric, analogous to entropy. A high score on silhouette metric is desired as it explains the high stability in cluster. From the document clustering plot of K-Means (see Figure 7), we observe high cluster stability when number of clusters is less than 5. Whereas the interpretation of NMI is as follows: when the number of clusters (k) is set to 2 and the number of class labels are 6, what is the NMI or quality of the clustering. Similar to silhouette, higher NMI is desired, and from the plot of K-Means (Figure 7), it is achieved when the number of clusters are 18.

How to run the program:

`python3 clustering.py {kmeans, agglomerative}`

We tested K-means and agglomerative clustering separately and passed each as argument in the call to the 'clustering.py' program.

Sample command for running the K-Means clustering:

`python3 clustering.py agglomerative`

Sample output for agglomerative yields the following:

```

ananuel@Dell-Inspiron7:~/MSU-II/Courses/Information_Retrieval/Project_2$ python3 clustering.py agglomerative
number of clusters: 2
silhouette_score: 0.196828705058575
nmi_score: 0.014887914581989892
number of clusters: 3
silhouette_score: 0.08015531054872187
nmi_score: 0.01538897340405141
number of clusters: 4
silhouette_score: 0.0745871114719176
nmi_score: 0.015964755609672126
number of clusters: 5
silhouette_score: 0.027145012493355984
nmi_score: 0.019500407635282883
number of clusters: 6
silhouette_score: 0.028085231144517792
nmi_score: 0.019822058564659364
number of clusters: 7
silhouette_score: 0.015059324412280773
nmi_score: 0.021639629353578157
number of clusters: 8
silhouette_score: 0.01584055877799558
nmi_score: 0.022094678784226555
number of clusters: 9
silhouette_score: 0.015471644375032751
nmi_score: 0.026543486214773416
number of clusters: 10
silhouette_score: 0.01603356599060456
nmi_score: 0.027132757205434524
number of clusters: 11
silhouette_score: -0.008388495071893183
nmi_score: 0.036599465520834665
number of clusters: 12
silhouette_score: -0.007102995506751115
nmi_score: 0.03824493898790889
number of clusters: 13
silhouette_score: -0.00319808989325568645
nmi_score: 0.040741395398103
number of clusters: 14
silhouette_score: -0.0016123600146670542
nmi_score: 0.04121002360916924
number of clusters: 15
silhouette_score: 0.0017481077550081112
nmi_score: 0.04439054529119546
number of clusters: 16
silhouette_score: 0.0023506255881149654
nmi_score: 0.0447203843682981
number of clusters: 17
silhouette_score: 0.002143892126980499
nmi_score: 0.045731423973004476
number of clusters: 18
silhouette_score: 0.002316509695389734
nmi_score: 0.046054580524107575
number of clusters: 19
silhouette_score: 0.0028565333871119957
nmi_score: 0.0466092466361704
number of clusters: 20
silhouette_score: 0.0027518321357684174
nmi_score: 0.04844843826209807
number of clusters: 21
silhouette_score: 0.002393485452323779

```

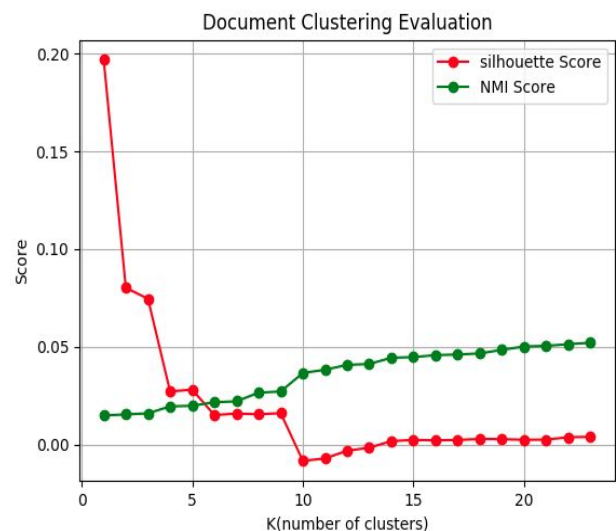


Figure 6: (Left) Execution of Document Clustering python program : clustering.py for Agglomerative Clustering. (right) The plot of document clustering evaluation using NMI and Silhouette metric. The are

outliers being present in the dataset as seen from the negative Silhouette score and also can be the reason of disagreement between NMI and silhouette.

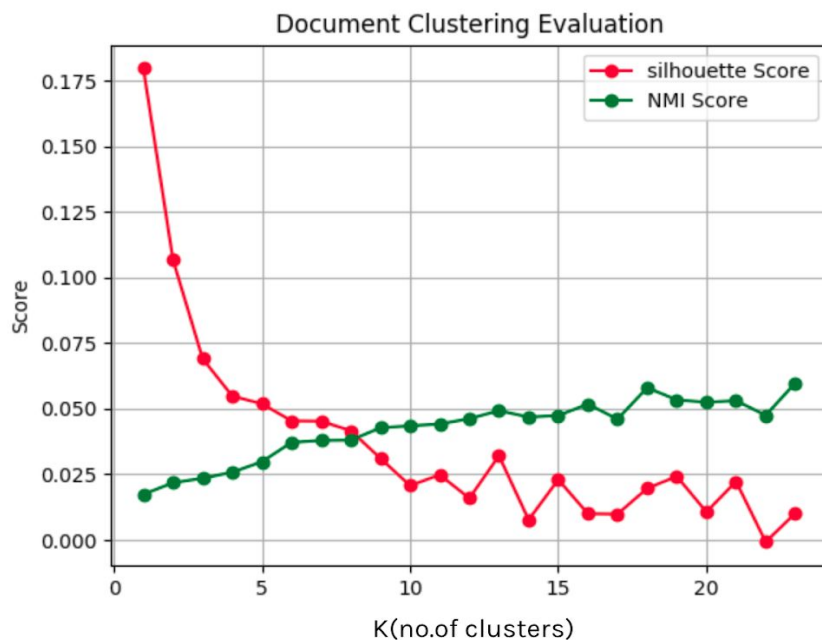


Figure 7: The plot of document clustering evaluation using NMI and Silhouette metric for K-Means Clustering. As seen from the curve, there are no identified outliers in the dataset as seen from the non-negative Silhouette score. Further, we observe high silhouette score for number of clusters = 2, which furnishes the high f1-score for BernoulliNB which considers binary features.

From figures 7 and 8, we observe that both K-Means and Agglomerative clustering have high silhouette score for k=2 number of clusters. However, for the NMI, K-Means has the highest NMI for k=18, whereas for agglomerative clustering it is k=25. Further, we notice near-monotonic increase in NMI with the number of clusters which support independence assumptions over features for the Naive Bayes models.

References:

1. <https://nlpforhackers.io/building-a-simple-inverted-index-using-nltk/>
2. https://scikit-learn.org/stable/auto_examples/model_selection/plot_learning_curve.html
3. https://scikit-learn.org/stable/modules/cross_validation.html