

Project 1

Shengrong Wu & Shuangyu Li

Introduction

In this project, we implemented the classification of eight classes from *sklearn.fetch_20newsgroups* with different supervised learning models. In this project, we utilized *sklearn* and *nltk* to implement algorithms like support vector machine (SVM), Naive Bayes and logistic regression with different hyperparameters. In the last part of the project, we performed multiclass classification on four subclasses from *sklearn.fetch_20newsgroups*, by extending our SVM and Naive Bayes classifier techniques. Their ROC graphs and statistical results are plotted and calculated to compare performance throughout this report.

Problem Statement and Datasets

In this project, we work with “20 Newsgroups” dataset. It is a collection of approximately 20,000 newsgroup documents, partitioned (nearly) evenly across 20 different newsgroups, each corresponding to a different topic. However, we just need to work on 8 different classes shown in table 1. We need to train the and test the data to classify them into 2 parent classes, ‘Computer technology’ and ‘Recreational activity’

Table 1. Subclasses of ‘Computer technology’ and ‘Recreational activity’

Computer technology	Recreational activity
comp.graphics	rec.autos
comp.os.ms-windows.misc	rec.motorcycles
comp.sys.ibm.pc.hardware	rec.sport.baseball
comp.sys.mac.hardware	rec.sport.hockey

Problem a) Display document histogram

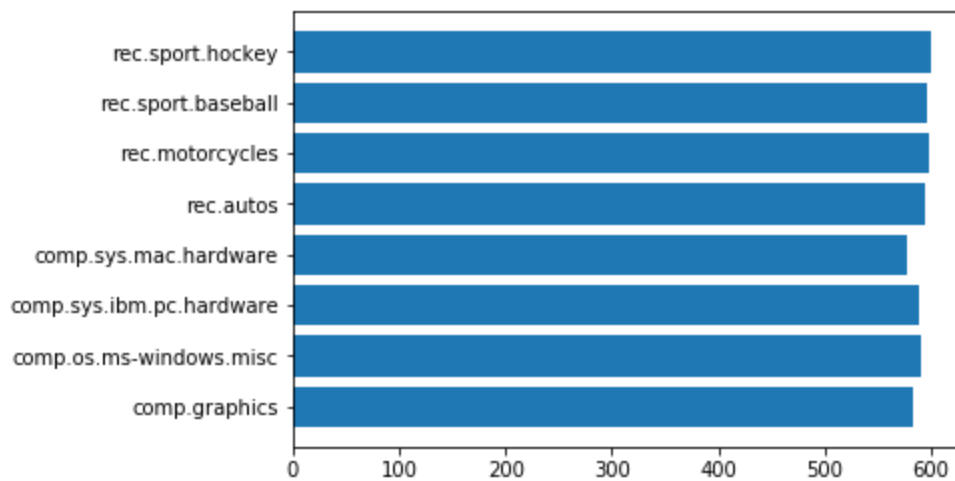


Fig 1: Distribution of the Data size

To avoid the bias brought by the size of training data set, we need to balance the data set before the training . In this project, the data is most evenly distributed for different classes since we use the given data directly. The exact amount of each class corresponding to the Figure 1.1 from top to bottom is [584, 591, 590, 578, 594, 598, 597, 600].

Modeling Text Data and Feature Extraction

Problem b) Tokenize the Input Data and Processing

The purpose of this part is to tokenize the documents, to stem the words and to remove the unneeded stopping words and punctuation, which is common in all documents but useless for classification. We use the dictionary of `text.ENGLISH_STOP_WORDS` from *sklearn*, stopwords from *nltk.corpus* and punctuation set from *python.string*. They are combined to build a set of stopwords which is called later in the *CountVectorizer*. With different min document frequency threshold, we could fetch different term frequency(TF) matrix.

For eight classes required,

When `Min_df = 5`, we extract 10396 terms from 4732 documents

When `Min_df = 2`, we extract 24849 terms from 4732 documents

It follows intuition that with higher threshold, we filtered out more terms like inappropriate abbreviations or typos.

Problem c) Extract 10 Most Significant Terms

We also want to measure the significance of a term to each class. To achieve this, we extend TFxIDF matrix to a TFxICF matrix. It is the same as a TFxIDF matrix, except that a class sits in place of a document, i.e. $TFxICF(c,t) = tf(t,c) * icf(t,c)$, where $icf(t,c) = \log[n_{classes}/cf(t)] + 1$. Hence, TFICF is a matrix with $n_{classes}$ rows, with each row corresponding to one class. To calculate TFICF for our problem, we need to use the entire dataset with $n_{classes} = 20$. To find the 10 most significant terms to the each of the 4 classes of our concern, we simply search for the 10 terms with the highest TFICF value in each corresponding row. We ran the code with two cases, by setting `min_df` to 2 and 5. The two results are almost identical with one exception in *comp.sys.mac.hardware*, as shown in Table 2, where the term 'problem' in the case of `<min_df=5>` replaced 'quadra' in the case of `<min_df=2>`.

Without looking at all documents in the dataset, it is reasonable to assume 'quadra' occurred fewer than 5 times in the entire dataset while mostly, if not all, in *comp.sys.mac.hardware*. Besides the terms which make intuitive sense to their classes, there are terms that are common in all classes, such as "edu" and "organization". These happened because the term frequencies are very high and dominates the TFxICF value. One possible improvement to this problem could be stemming the terms or redesign TFxICF matrix to predict term importance.

Table 2. 10 Most Significant Terms (not in order)

comp.sys.ibm.pc.hardware	comp.sys.mac.hardware	misc.forsale	soc.religion.christian
'com'	'apple'	'com'	'christian'
'controller'	'edu'	'edu'	'church'
'drive'	'line'	'line'	'edu'
'edu'	'mac'	'new'	'god'
'ide'	'organization'	'offer'	'jesus'
'line'	'post'	'organization'	'know'
'organization'	'quadra'(2), 'problem'(5)	'post'	'line'
'scsi'	'scsi'	'sale'	'people'
'subject'	'subject'	'subject'	'say'
'use'	'use'	'university'	'subject'

Feature Selection

Problem d) Dimension Reduction Applying LSI and NMF

To speed up calculation and find more important feature for training, Latent Semantic Indexing is applied to reduce the dimension of the TFxIDF matrix. In this problem, we set the hyperparameter k to 50, which represents the number of most significant feature in the training sets and use *sklearn* package to implement it. The dimension of the reduced version of training data is (4732, 50) . The dimension of the reduced version of test data is (3150, 50). The result of Non-Negative Matrix Factorization is slightly different but in the same size. These matrix representations are passed for training and testing in the following classification tasks and the results will be compared.

Learning Algorithms

In this part, we will apply different learning algorithm with different parameter to reach the best result for two classes classification problem.

Problem e) SVM

The Support vector machine(SVM) algorithm is called to solve the classification problem in this question and we tried with different minimum document term frequency threshold and the two distinct misclassification weight - γ .

With larger γ , the hard SVM model works quite well for the classification problem and barely changes with different term frequency thresholds. However, for the soft SVM model, it classifies everything into 'Rec' class no matter how the other parameters changes . Their results are all shown in the tables and ROC curves below.

Table 3: Result for **Hard** SVM, $\gamma = 1000$, **LSI & min_df = 5**

Hard SVM $\gamma = 1000$	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 5	0.97174603 17460317	0.971692065 7958394	0.9718440587 82038	[[1507 53] [36 1554]]

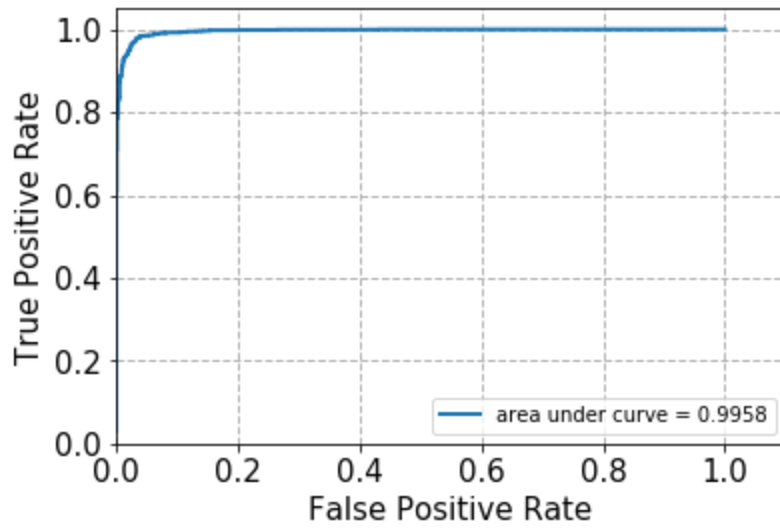


Fig 2: ROC curve for Hard SVM with `LSI` & `min_df = 5`

Table 4: Result for Hard SVM with `LSI` & `min_df = 2`

Hard SVM $\gamma = 1000$	Accuracy	Recall	Precision	Confusion matrix
<code>LSI</code> & <code>min_df</code> <code>= 2</code>	0.969841269 8412698	0.9697750362 844703	0.96997796581 25243	<code>[[1502 58]</code> <code>[37 1553]]</code>

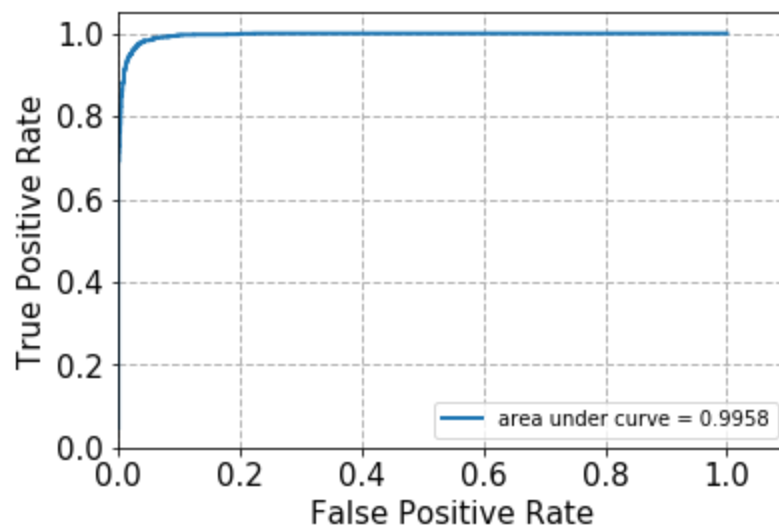


Fig 3: ROC curve for Hard SVM with `LSI` & `min_df = 2`

Table 5: for Hard SVM with NMF & min_df = 2

Hard SVM $\gamma = 1000$	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 2	0.959365079 3650793	0.9592888243 83164	0.95952892657 99256	[[1484 76] [52 1538]]

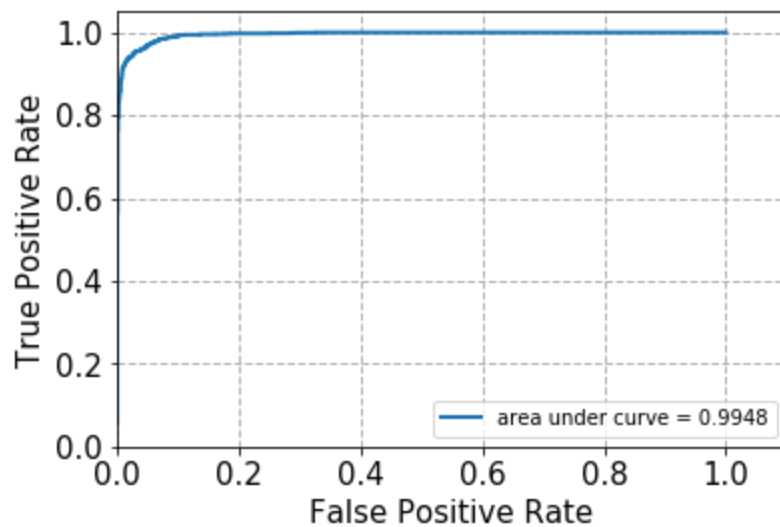


Fig 4: ROC curve for Hard SVM with NMF & min_df = 2

Table 6:Result for Soft SVM with LSI & min_df = 5

soft SVM $\gamma = 0.001$	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 5	0.5047619047 619047	0.5	0.25238095238 09524	[[0 1560] [0 1590]]

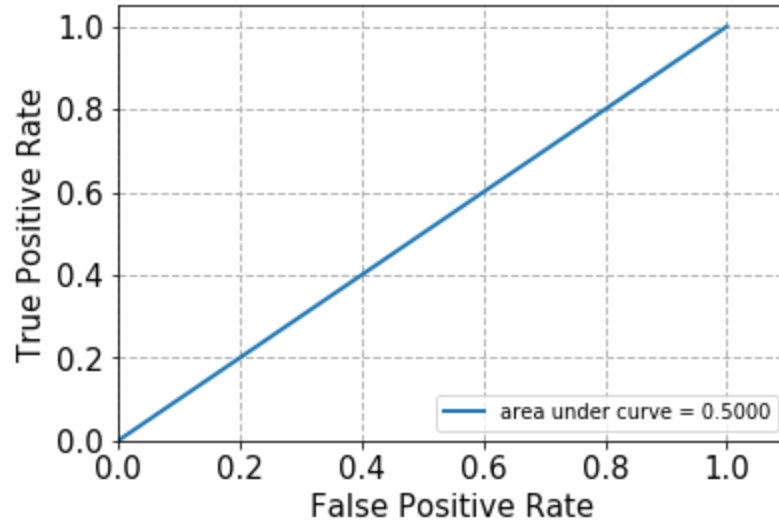


Fig 5: ROC curve for Soft SVM with `LSI & min_df = 5`

Table 7: Result for Soft SVM with `LSI & min_df = 2`

<code>soft SVM</code> $\gamma = 0.001$	Accuracy	Recall	Precision	Confusion matrix
<code>LSI & min_df = 2</code>	0.5047619047619047	0.5	0.2523809523809524	<pre>[[0 1560] [0 1590]]</pre>

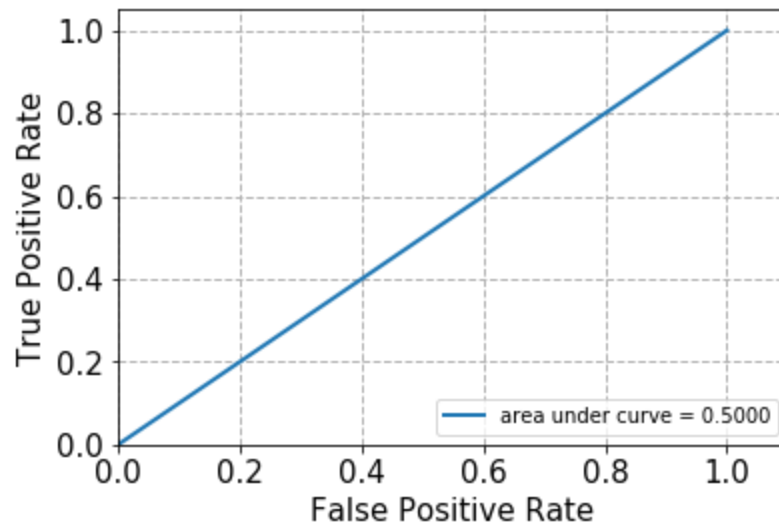


Fig 6: ROC curve for Soft SVM with `LSI & min_df = 2`

Table 8:Result for Soft SVM with NMF & min_df = 2

soft SVM $\gamma = 0.001$	Accuracy	Recall	Precision	Confusion matrix
NMF & min_df = 2	0.5047619047 619047	0.5	0.25238095238 09524	[[0 1560] [0 1590]]

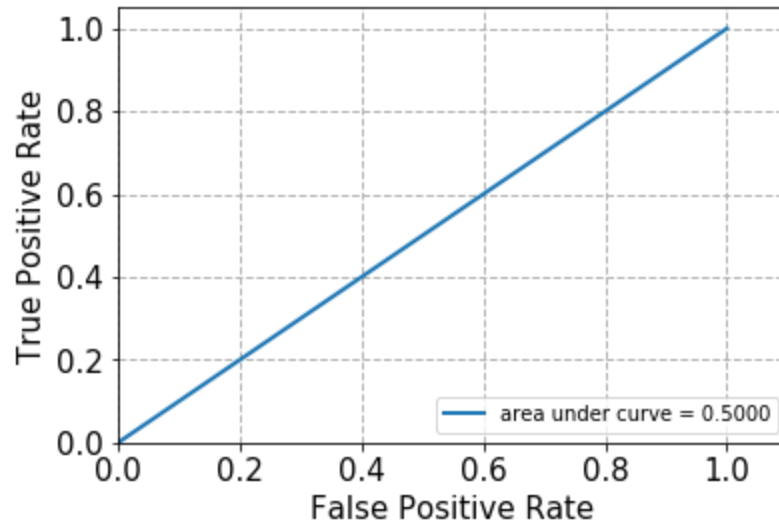


Fig 7: ROC curve for Soft SVM with NMF & min_df = 2

Problem f) Cross Validation

It is obvious to see the important role of γ from the result of problem e. A cross validation is performed for different min_df, different dimension reduction methods and for γ range from 0.001 to 1000. The results are shown below. Element (i, j) represents the accuracy for $\gamma = 10^{(i-3)}$ in j fold. The best result of most of them wiggle between 100 and 1000 depends on the random factor in folds.

Table 8: Accuracy result for LSI & min_df = 2

```
[[0.45828933 0.4656811 0.51057082 0.50845666 0.51268499]
 [0.45828933 0.4656811 0.51057082 0.50845666 0.51268499]
 [0.97043295 0.97465681 0.96617336 0.97145877 0.96934461]
 [0.97360084 0.97571278 0.96723044 0.9756871 0.97251586]
 [0.97571278 0.97993664 0.97040169 0.97463002 0.97251586]
 [0.97360084 0.97993664 0.9756871 0.97251586 0.9756871 ]
 [0.97254488 0.97993664 0.97357294 0.97357294 0.97674419]]
```

Best -> 100

Table 9:Result for SVM gamma = 100 with LSI & min_df = 2

Hard SVM $\gamma = 100$	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 2	0.97301587 3015873	0.972943880 0193517	0.9731764239 884059	[[1506 54] [31 1559]]

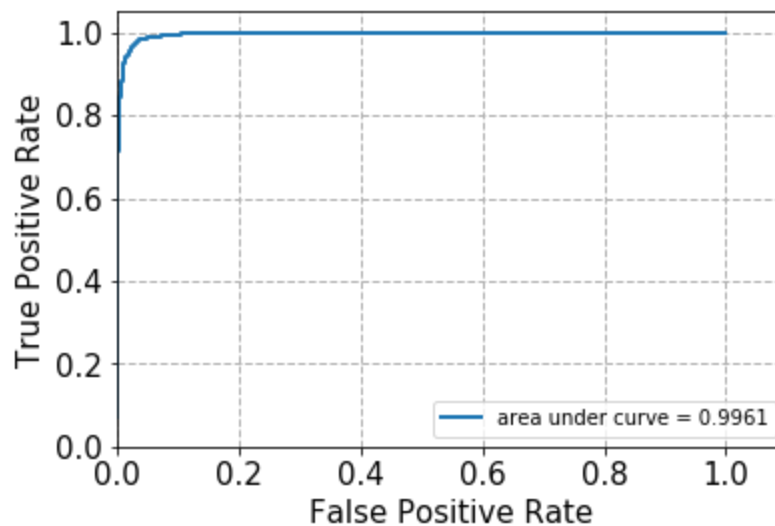


Fig 8: ROC curve for SVM gamma = 100 with LSI & min_df = 2

Table 10: Accuracy result for LSI & min_df = 5

```
[[0.45828933 0.4656811 0.51057082 0.50845666 0.51268499]
 [0.45828933 0.4656811 0.51268499 0.51057082 0.51585624]
 [0.96832101 0.97571278 0.96194503 0.96194503 0.96088795]
 [0.97360084 0.97676874 0.96934461 0.96828753 0.9640592 ]
 [0.97782471 0.97676874 0.96617336 0.97040169 0.97674419]
 [0.97465681 0.97993664 0.97040169 0.96723044 0.97463002]
 [0.97465681 0.97888068 0.97357294 0.97357294 0.97674419]]
Best -> 1000
```

Table 11: Result for Hard SVM, $\gamma = 1000$, LSI & min_df = 5

Hard SVM $\gamma = 1000$	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 5	0.97174603 17460317	0.971692065 7958394	0.9718440587 82038	[[1507 53] [36 1554]]

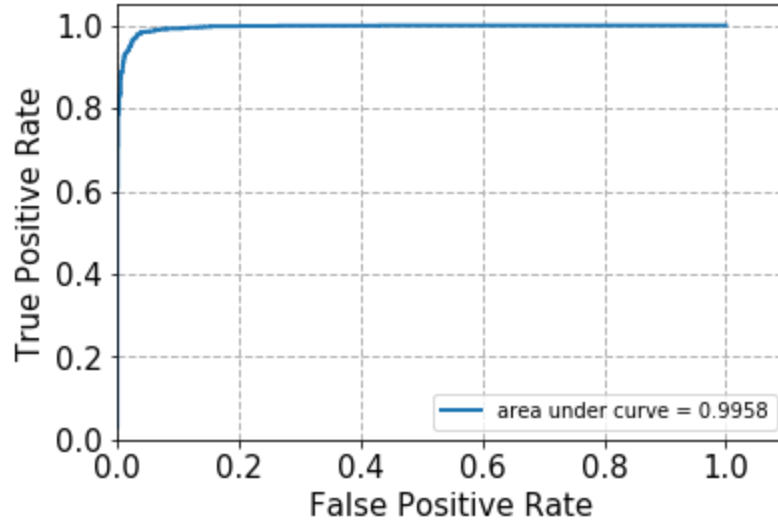


Fig 9: ROC curve for SVM gamma = 1000 with LSI & min_df = 5

Table 12: Accuracy result for NMF & min_df = 2

```
[[0.45828933 0.4656811 0.51057082 0.50845666 0.51268499]
 [0.45828933 0.4656811 0.51057082 0.50845666 0.51268499]
 [0.45828933 0.4656811 0.51057082 0.50845666 0.51268499]
 [0.95459345 0.9313622 0.94608879 0.95348837 0.94714588]
 [0.97676874 0.96620908 0.9577167 0.96617336 0.96300211]
 [0.97571278 0.97254488 0.96511628 0.9640592 0.96723044]
 [0.97571278 0.97148891 0.96828753 0.96194503 0.9756871 ]]
```

Best -> 1000

Table 13: Result for SVM gamma = 1000 with NMF & min_df = 2

SVM γ = 1000	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 2	0.959365079 3650793	0.9592888243 83164	0.95952892657 99256	[[1484 76] [52 1538]]

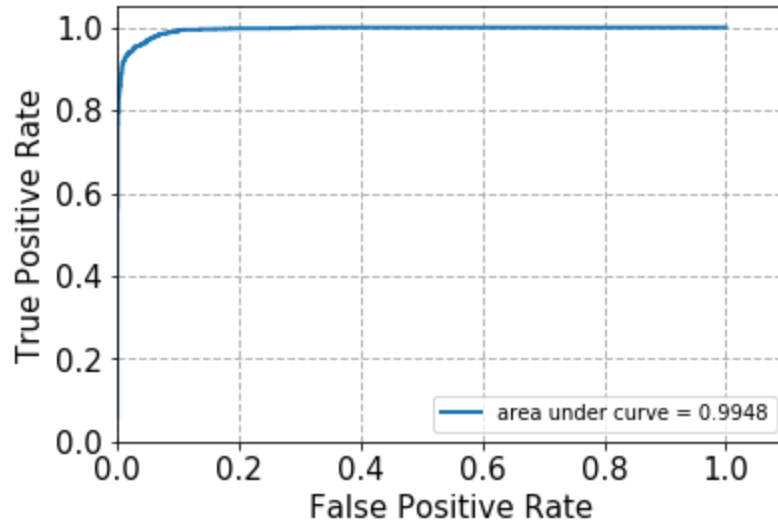


Fig 10: ROC curve for SVM gamma = 1000 with NMF & min_df = 2

Problem g) Naive Bayes

Because Naive Bayes model does not accept the negative input, we only perform the NMF reduction in this part.

Table 14: Result for Naive Bayes with NMF & min_df = 2

Naive Bayes	Accuracy	Recall	Precision	Confusion matrix
NMF & min_df = 2	0.9358730158730159	0.9354015481373972	0.940433039268143	$\begin{bmatrix} 1382 & 178 \\ 24 & 1566 \end{bmatrix}$

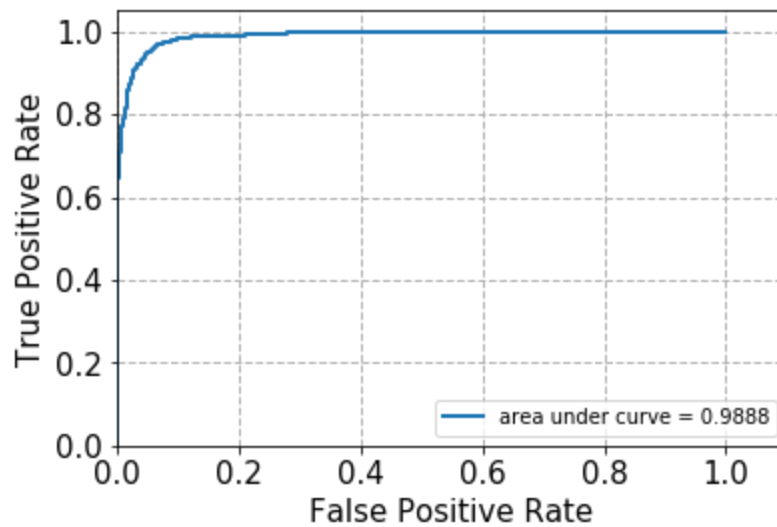


Fig 11: ROC curve for Naive Bayes with NMF & min_df = 2

Problem h) Logistic Regression

Same procedure is repeated for logistic regression.

Table 15: Result for Logistic Regression with NMF & min_df = 2

Logistic Regression	Accuracy	Recall	Precision	Confusion matrix
NMF & min_df = 2	0.93523809 52380952	0.9350507982 583454	0.9359999999 99999	[[1428 132] [72 1518]]

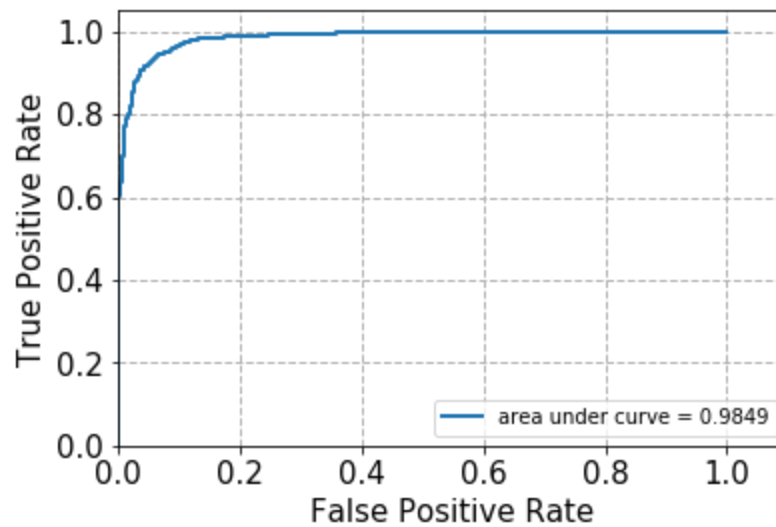


Fig 12: ROC curve for Logistic Regression with NMF & min_df = 2

Table 16: Result for Logistic Regression with LSI & min_df = 2

Logistic Regression	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 2	0.96634920 63492064	0.9662373004 354137	0.96668434204 02709	[[1489 71] [35 1555]]

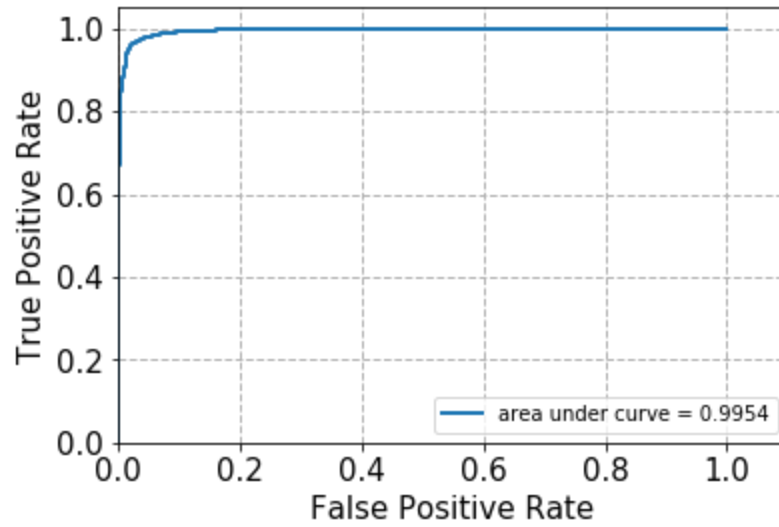


Fig 13: ROC curve for Logistic Regression with LSI & min_df = 2

Table 17 : Result for Logistic Regression with LSI & min_df = 5 Result

Logistic Regression	Accuracy	Recall	Precision	Confusion matrix
LSI & min_df = 5	0.96571428 57142857	0.9655962747 94388	0.96608182367 05228	[[1489 71] [35 1555]]

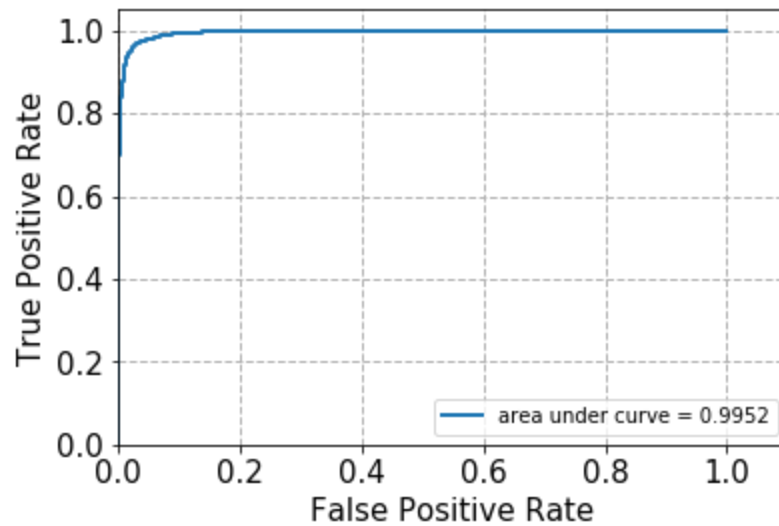


Fig 14: ROC curve for Logistic Regression with LSI & min_df = 5

Problem I) Logistic Regularization

The accuracy of logistic model with different regularization parameter C is shown below. Generally, with higher regularization parameter, the model perform slightly better but for some the of the model, the second largest parameter lead to the best performance. The difference norm, L1 and L2, are calculated to do comparison. Since at a higher degree vector space the 'distance' between sample is hard to describe like the 2D plane, therefore the various norm are tried here. In this part, the accuracy is selected as the key to compare their performance. L1 norm works poorly at a low parameter value, are as good as L2 norm at a high regularization parameter. The detailed accuracy score tables are shown below.

Table 18: accuracy score table of l1 norm

L1 Norm	$C = 10^{-3}$	$C = 10^{-2}$	$C = 10^{-1}$	$C = 10^0$	$C = 10^1$	$C = 10^2$	$C = 10^3$
NMF & min_df = 2	0.49523 8095238 09526	0.49523 8095238 09526	0.69936 5079365 0794	0.9584 126984 126984	0.96031 7460317 4603	0.96158 7301587 3015	0.96253 9682539 6826
LSI & min_df = 2	0.49523 8095238 09526	0.87111 1111111 1111	0.95809 5238095 2381	0.9657 142857 142857	0.97238 0952380 9523	0.97079 3650793 6508	0.97047 6190476 1905
LSI & min_df = 5	0.49523 8095238 09526	0.88888 8888888 8888	0.95523 8095238 0952	0.9634 920634 920635	0.97111 1111111 1111	0.97047 6190476 1905	0.97047 6190476 1905

Table 19: accuracy score table of l2 norm

L2 Norm	C = 10 ⁻³	C = 10 ⁻²	C = 10 ⁻¹	C = 10 ⁰	C = 10 ¹	C = 10 ²	C = 10 ³
NMF & min_df = 2	0.5047 619047 619047	0.5057 142857 142857	0.8765 079365 079365	0.9352 380952 380952	0.9488 888888 888889	0.9523 809523 809523	0.9565 079365 079365
LSI & min_df = 2	0.68	0.9419 047619 047619	0.9615 873015 873015	0.9663 492063 492064	0.9711 111111 111111	0.9730 158730 15873	0.9707 936507 936508
LSI & min_df = 5	0.7333 333333 333333	0.9453 968253 968253	0.9587 301587 301588	0.9657 142857 142857	0.9695 238095 238096	0.9714 285714 285714	0.9707 936507 936508

Multiclass Classification

So far, we have been dealing with classifying the data points into two classes. In this part, we explore multiclass classification by extending the algorithms - Naive Bayes and SVM. Naive Bayes perform the multiclass classification inherently, regardless of the number of classes. Here we are using NMF as multinomial Naive Bayes classification only takes non-negative values.

For SVM, we have two implementation for linear multiclass classification, namely one-to-one method and one-vs-rest method. From the sklearn package, "SVC" function assumes a one-to-one implementation and "linearSVC" function assumes a one-vs-rest implementation. We used these functions to generate results with both NMF and LSI, and with a hyperparameter $\gamma = 100$. The results are shown below. Among all, SVM are better than Naive Bayes, while the two implementation of SVM are almost identical. SVM with LSI have slightly better results than SVM with NMF.

Table 24: Result for Naive Bayes with NMF & min_df = 2

Naive Bayes	Accuracy	Recall	Precision	Confusion matrix
NMF & min_df = 2	0.8044728 434504792	0.80306697 96821986	0.8100486 186980468	[[304 42 42 4] [102 246 34 3] [52 12 318 8] [2 1 4 391]]

Table 25: Result for One-to-One SVM with NMF & $\gamma = 100$

SVM $\gamma = 100$	Accuracy	Recall	Precision	Confusion matrix
One-to-One NMF & min_df = 2	0.8466453 674121406	0.84574365 12398824	0.848967 71208483 4	[[315 48 28 1] [70 290 23 2] [40 17 332 1] [7 2 1 388]]

Table 26: Result for One-to-One SVM with LSI & $\gamma = 100$

SVM $\gamma = 100$	Accuracy	Recall	Precision	Confusion matrix
One-to-One LSI & min_df = 2	0.8830670 926517572	0.88247863 94510012	0.882696 48140134 62	<pre>[[322 46 24 0] [40 321 24 0] [23 17 348 2] [3 2 2 391]]</pre>

Table 27: Result for One-vs-Rest SVM with NMF & $\gamma = 100$

SVM $\gamma = 100$	Accuracy	Recall	Precision	Confusion matrix
One-vs-Rest NMF & min_df = 2	0.8472843 450479233	0.84638378 19275722	0.846328 83805944 6	<pre>[[306 52 32 2] [64 291 26 4] [30 17 338 5] [4 2 1 391]]</pre>

Table 28: Result for One-vs-Rest SVM with LSI & $\gamma = 100$

SVM $\gamma = 100$	Accuracy	Recall	Precision	Confusion matrix
One-vs-Rest LSI & min_df = 2	0.8843450 479233227	0.88382392 01573732	0.883949 02012475 04	<pre>[[314 52 26 0] [34 327 23 1] [18 19 351 2] [4 0 2 392]]</pre>