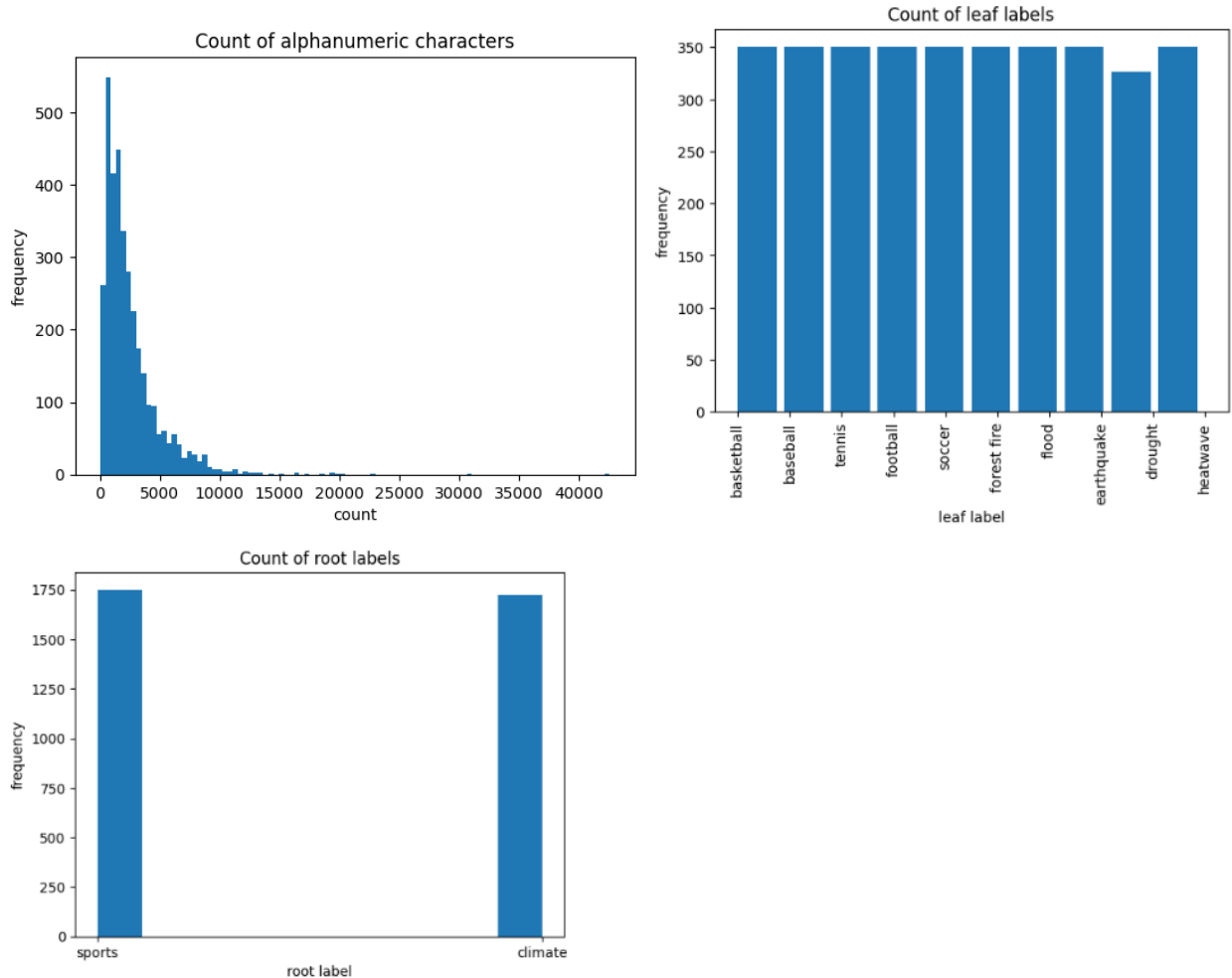


EC ENGR 219 Project 1

By: Wei Jun Ong, Rick Yang, Chenggong Zhang

Q1.

Rows: 3476; Columns: 8



The “count of alphanumeric characters” histogram shows that majority of the data points have less than 5000 alphanumeric characters in the “full_text” feature, with most data points having around 2000 alphanumeric characters. The “count of leaf labels” histogram shows that there are generally equal number of data points for each leaf label, except for the “drought” leaf label that has slightly less data points than the other 9 labels. The “count of root labels” histogram similarly shows that there are about the same number of data points for each of the 2 root labels, with the “climate” root label having slightly less data points than the “sports” root label.

Q2.

Number of training samples: 2780

Number of testing samples: 696

Q3.

Lemmatization is more accurate than stemming because it actually analyzes the sentence structure using linguistic tools to determine the role of a word within a sentence, making it more accurate in finding the lemma of a word, while stemming simply chops off prefixes, suffixes and other common additions to a word, which can more easily result in erroneous changing of a word when trying to find the base form. However, because of the use of linguistic tools and a dictionary of words, lemmatization is more computationally expensive than stemming. Both reduce the dictionary size because words with the same base forms / lemma but with additional prefixes or suffixes will be reduced down to the same base forms / lemma, so number of unique words after both preprocessing techniques will decrease.

Increasing the min_df value means more terms that appear less frequently in the documents will be discarded (the filter for minimum occurrences of terms is more restrictive), so the TF-IDF matrix will be sparser and contain fewer non-zero values.

Stopwords, punctuations and numbers should be removed after lemmatization step, because these words (that will eventually be removed) are essential for the lemmatizer to decipher their role in the sentence. Keeping these words in when passing the sentence into the lemmatizer enables it to correctly tag the important words.

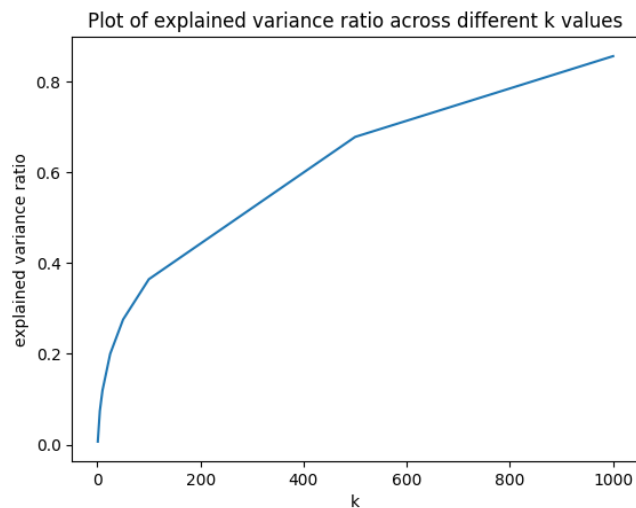
Train TFIDF shape:

Rows: 2780; Columns: 13288

Test TFIDF shape:

Rows: 696; Columns: 13288

Q4.



Concavity of the plot shows that with higher k values (or greater dimensionality / more principal components), the total variance explained by the principal components increases, albeit at a lower rate. This suggests that the greater the dimensionality, the more the total variance that can be captured by the principal components. However, since the rate of increase for explained variance ratio is decreasing, this shows that there is a certain number of principal components where the explained variance ratio is good enough such that we do not need more principal components.

For NMF, `fit_transform` gives the value of W_{train} and we can get the value H from `components_`. Similarly for LSI, `fit_transform` gives the value of $U \cdot \Sigma$, and we get the value of $V^{\text{transpose}}$ from `components_`.

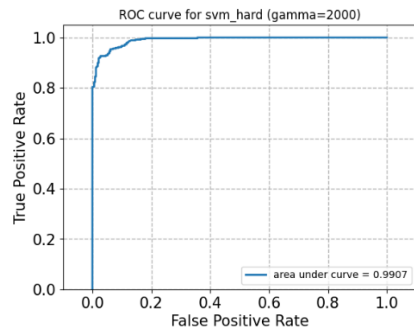
NMF error: 2170.586363338054

LSI error: 2152.4526721780967

We see that NMF error is higher than LSI error. Since NMF finds only non-negative matrix values H and W_{train} , this means that it's possible that the solutions contain some negative values that NMF is unable to find but LSI can find, thus leading to LSI having a lower error.

Q5.

Using pipeline svm_hard (gamma=2000):



Confusion matrix:

```
[[348 20]
```

```
 [ 19 309]]
```

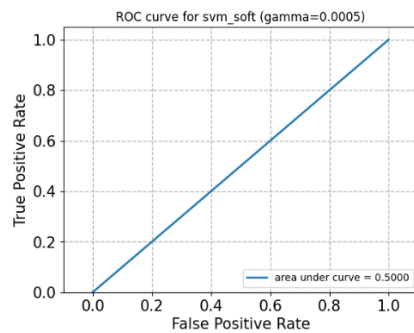
Accuracy score: 0.9439655172413793

Recall score: 0.9420731707317073

Precision score: 0.939209726443769

F1 score: 0.9406392694063926

Using pipeline svm_soft (gamma=0.0005)



Confusion matrix:

```
[[ 0 368]
```

```
 [ 0 328]]
```

Accuracy score: 0.47126436781609193

Recall score: 1.0

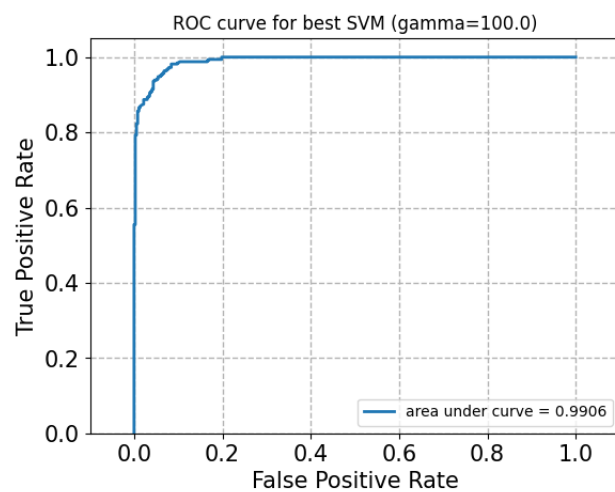
Precision score: 0.47126436781609193

F1 score: 0.640625

The hard margin SVM performed better, due to the better stats like accuracy, f1 score etc. For SVM with gamma = 100000, the stats are the same as that of the earlier hard margin SVM where gamma = 2000.

The soft margin SVM has a bad accuracy score, poor recall, precision, and f1 scores, and in fact predicts the same root label for all cases. From the confusion matrix, the true positive is very high, but the true negative is low. False positive is also very high. Since this is a binary classification (only two root labels), the soft margin SVM is predicting the same positive label for all cases, so only true positives and false positives are high. The ROC curve supports this, since the ROC curve for the soft margin SVM is on the diagonal line, representing the SVM is behaving like a random classifier.

Best gamma after 5-fold validation: 100



Confusion matrix:

```
[[348 20]
```

```
 [ 19 309]]
```

Accuracy score: 0.9439655172413793

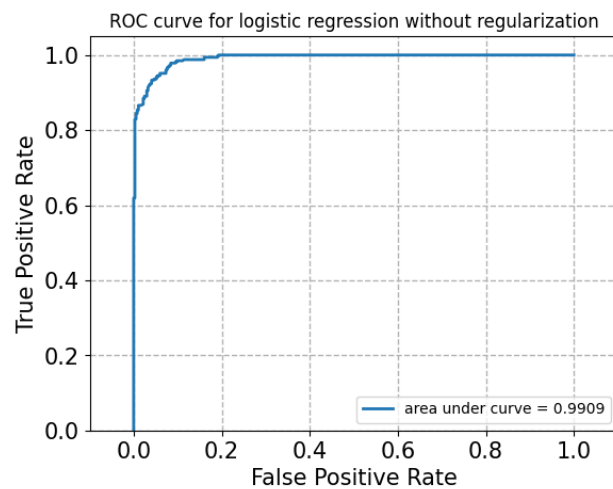
Recall score: 0.9420731707317073

Precision score: 0.939209726443769

F1 score: 0.9406392694063926

Q6.

Logistic classifier without regularization:



Confusion matrix:

```
[[348 20]
```

```
 [ 18 310]]
```

Accuracy score: 0.9454022988505747

Recall score: 0.9451219512195121

Precision score: 0.9393939393939394

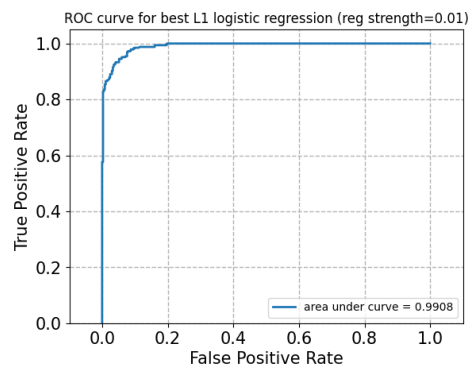
F1 score: 0.9422492401215805

After 5-fold cross-validation,

Best L1 regularization strength: 0.01

Best L2 regularization strength: 1e-05

Stats of logistic regression model with best L1 regularization strength:



Confusion matrix:

```
[[349 19]
```

```
 [ 19 309]]
```

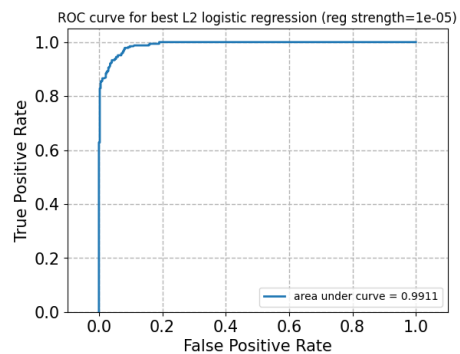
Accuracy score: 0.9454022988505747

Recall score: 0.9420731707317073

Precision score: 0.9420731707317073

F1 score: 0.9420731707317073

Stats of logistic regression model with best L2 regularization strength:



Confusion matrix:

```
[[348 20]
```

```
 [ 18 310]]
```

Accuracy score: 0.9454022988505747

Recall score: 0.9451219512195121

Precision score: 0.9393939393939394

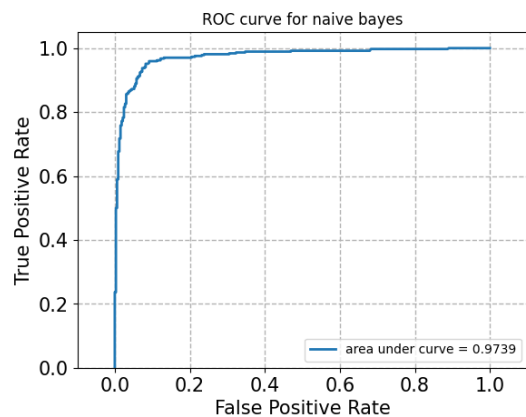
F1 score: 0.9422492401215805

By comparing accuracy, recall, precision and f1 scores, the logistic regression model with L1 regularization (strength=0.01) has the best performance.

Regularization parameter increases training error, which tends to reduce learnt coefficients (or weights) and prevents overfitting and improves generalizability. This tends to reduce test errors and make the models more robust. Different kinds of regularization leads to different effects on the model, so the appropriate regularization technique should be employed based on what is needed. L1 regularization performs more feature selection by making more weights become zero, causing greater sparsity and thus reducing dimensionality. L2 regularization is better for distributing the impacts of more correlated variables across different features, without eliminating those features altogether, thus capturing more information while still reducing overfitting.

Linear SVMs and logistic regression use different loss functions (maximizing margins through hinge loss for linear SVM, and estimating probabilities through log loss for logistic regression), so naturally they have different ways to find boundaries and will usually arrive at different decision boundaries. Performances may differ based on the intrinsic assumptions of the data, where a certain loss function may be more suited to find a clearer boundary for one set of data, but not for the other dataset. Empirically-speaking, they have similar performance and do not have statistically significant differences.

Q7.



Confusion matrix:

```
[[280 48]
```

```
 [ 11 357]]
```

Accuracy score:0.9152298850574713

Recall score: 0.970108695652174

Precision score: 0.8814814814814815

F1 score: 0.9236739974126779

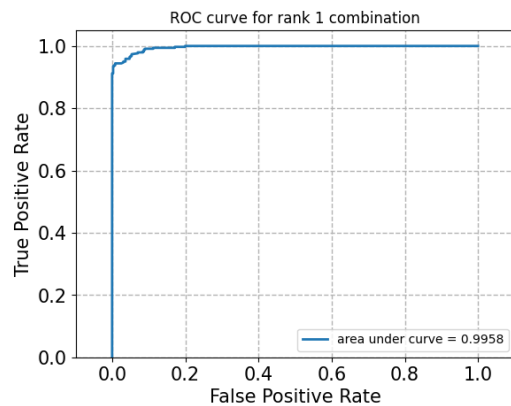
Q8.

Top 5 combinations (determined by best average validation accuracy)

Ranking	Feature Extraction	Dimensionality Reduction	Classifier	Average validation score
1	Lemmatization, mindf=5	LSI, k=100	Logistic Regression, L1 strength=0.01	0.9661870503597122
2	Lemmatization, mindf=5	NMF, k=100	Logistic Regression, L1 strength=0.01	0.9658273381294965
3	Lemmatization, mindf=2	LSI, k=100	Logistic Regression, L1 strength=0.01	0.9651079136690646
4	Lemmatization, mindf=5	LSI, k=100	Logistic Regression, L2 strength=1e-05	0.9636690647482015
5	Lemmatization, mindf=2	LSI, k=100	Logistic Regression, L2 strength=1e-05	0.9633093525179856

Performances on test set:

Ranking 1:



Confusion matrix:

[[343 13]

[15 325]]

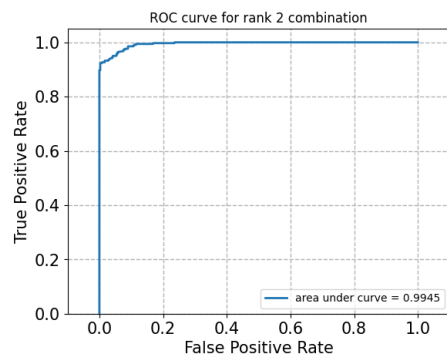
Accuracy score: 0.9597701149425287

Recall score: 0.9558823529411765

Precision score: 0.9615384615384616

F1 score: 0.9587020648967551

Ranking 2:



Confusion matrix:

[[341 15]

[19 321]]

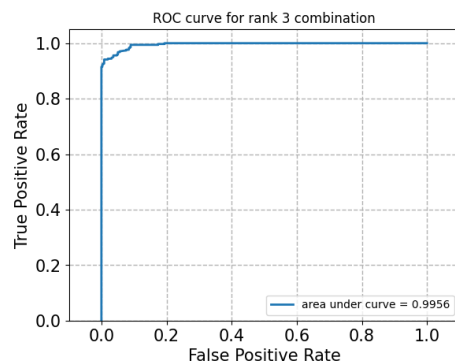
Accuracy score: 0.9511494252873564

Recall score: 0.9441176470588235

Precision score: 0.9553571428571429

F1 score: 0.9497041420118343

Ranking 3:



Confusion matrix:

[[343 13]

[16 324]]

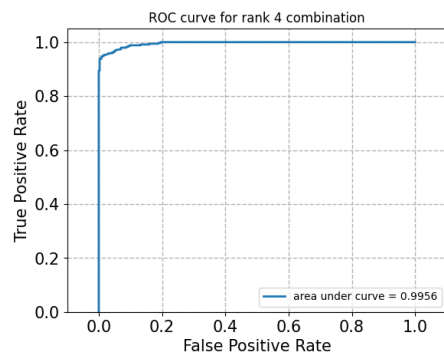
Accuracy score: 0.9583333333333334

Recall score: 0.9529411764705882

Precision score: 0.9614243323442137

F1 score: 0.9571639586410635

Ranking 4:



Confusion matrix:

[[345 11]

[14 326]]

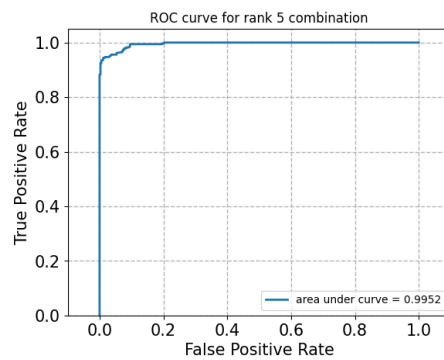
Accuracy score: 0.9640804597701149

Recall score: 0.9588235294117647

Precision score: 0.9673590504451038

F1 score: 0.9630723781388478

Ranking 5:



Confusion matrix:

[[343 13]

[15 325]]

Accuracy score: 0.9597701149425287

Recall score: 0.9558823529411765

Precision score: 0.9615384615384616

F1 score: 0.9587020648967551

Q9.

Naïve Bayes classifier:

Confusion matrix:

```
[[49 4 1 5 1 0 0 0 0 0]
 [ 0 51 9 1 3 0 1 0 0 0]
 [ 0 7 58 0 2 0 2 0 1 2]
 [ 0 5 2 58 0 0 0 3 0 0]
 [ 0 8 9 0 64 0 0 0 0 3]
 [ 0 3 6 1 0 13 2 1 1 47]
 [ 0 1 1 0 0 5 62 0 1 1]
 [ 0 7 3 0 0 2 0 60 0 1]
 [ 0 0 2 0 0 0 1 1 54 4]
 [ 0 8 9 1 0 18 1 0 1 29]]
```

Accuracy score: 0.7155172413793104

Recall score: 0.7155172413793104

Precision score: 0.7155172413793104

F1 score: 0.7155172413793104

Multiclass SVM (one vs one):

Confusion matrix:

```
[[60 0 0 0 0 0 0 0 0 0]
 [ 1 57 3 2 0 0 1 1 0 0]
 [ 0 7 57 0 2 4 1 0 0 1]
 [ 3 1 1 61 0 2 0 0 0 0]
 [ 2 0 2 0 79 1 0 0 0 0]
 [ 0 0 1 0 0 29 1 0 3 40]
 [ 0 1 0 0 0 3 67 0 0 0]
 [ 0 2 2 0 0 2 0 66 0 1]
 [ 0 1 0 0 0 3 0 1 56 1]
 [ 0 2 1 0 0 39 0 3 2 20]]
```

Accuracy score: 0.7931034482758621

Recall score: 0.7931034482758621

Precision score: 0.7931034482758621

F1 score: 0.7931034482758621

Multiclass SVM (one vs rest):

Confusion matrix:

```
[[60 0 0 0 0 0 0 0 0 0]
 [ 0 56 5 1 1 1 1 0 0 0]
 [ 0 7 56 1 4 2 0 0 0 2]
 [ 1 1 1 64 0 1 0 0 0 0]
 [ 1 1 2 0 80 0 0 0 0 0]
 [ 0 4 4 0 0 30 2 0 3 31]
 [ 0 1 0 0 0 2 67 1 0 0]
 [ 0 1 3 0 0 2 0 67 0 0]
 [ 0 2 0 0 0 2 0 1 56 1]
 [ 0 8 6 0 0 27 0 2 2 22]]
```

Accuracy score: 0.8017241379310345

Recall score: 0.8017241379310345

Precision score: 0.8017241379310345

F1 score: 0.8017241379310345

The imbalance issue in the Multiclass SVM one vs rest model is solved using the “class_weight=’balanced’” parameter in SVC. This ensures that class weights are automatically changed, where the weights are inversely proportional to frequency of the class appearing. Thus the model can adjust weights when there are much more negative classes than positive classes, and when this SVC model is passed into OneVsRestClassifier, class imbalance issue can be solved.

Across all 3 classifiers, assuming we start counting row number from 0, row 5 (and row 9 for the multiclass SVMs) have exceptionally low true values along the diagonal. This suggests that datapoints where the leaf_label is “forest fire” (and sometimes for “heatwave”) is exceptionally hard to predict correctly. For the multiclass SVMs, these 2 are oftentimes confused with each other, as shown by the higher false prediction rates of data labeled as “forest fire” but the classifier predicting them to be “heatwave”, and vice versa.

I suggest for “forest fire” and “heatwave” labels to be merged together. Class 9 (heatwave) is merged into class 5 (forest fire) so all existing labels for “heatwave” have been changed to “forest fire”. The new performance is shown below.

Naïve Bayes classifier:

Confusion matrix:

```
[[ 49  4  1  5  1  0  0  0  0]
 [  0 51  9  1  3  0  1  0  0]
 [  0  7 58  0  2  2  2  0  1]
 [  0  5  2 58  0  0  0  3  0]
 [  0  8  9  0 64  3  0  0  0]
 [  0 11 13  2  0 110  2  1  2]
 [  0  1  1  0  0  6 62  0  1]
 [  0  6  3  0  0  4  0 60  0]
 [  0  0  2  0  0  7  1  0 52]]
```

Accuracy score: 0.8103448275862069

Multiclass SVM (one vs one):

Confusion matrix:

```
[[ 60  0  0  0  0  0  0  0  0]
 [  1 56  3  2  0  1  1  1  0]
 [  0  5 55  0  2  9  1  0  0]
 [  3  1  1 61  0  2  0  0  0]
 [  2  0  2  0 79  1  0  0  0]
 [  0  2  2  0  0 130  1  1  5]
 [  0  0  0  0  0  4 67  0  0]
 [  0  1  1  0  0  7  0 64  0]
 [  0  1  0  0  0  5  0  0 56]]
```

Accuracy score: 0.9022988505747126

Multiclass SVM (one vs rest):

Confusion matrix:

```
[[ 60  0  0  0  0  0  0  0  0  0]
 [  0 56  5  1  1  1  1  0  0]
 [  0  6 56  1  2  6  0  1  0]
 [  1  1  2 64  0  0  0  0  0]
 [  1  1  2  0 80  0  0  0  0]
 [  0  9  9  0  0 116  1  2  4]
 [  0  1  0  0  0  4 66  0  0]
 [  0  1  3  0  0  2  0 67  0]
 [  0  1  0  0  0  4  0  1 56]]
```

Accuracy score: 0.8922413793103449

As we can see, accuracy scores for both multiclass SVMs increased.

Since there are now more data points with the label of “forest fire”, there is class imbalance where there are more data samples for class 5. The earlier method of mitigating class imbalance (using `class_weight='balanced'`) can similarly be employed to solve this challenge. After using this class imbalance solution, here are the new performance results for the multiclass SVMs.

Multiclass SVM (one vs one):

Confusion matrix:

```
[[ 59  0  0  1  0  0  0  0  0]
 [  1 57  3  2  0  0  1  1  0]
 [  0  7 58  0  2  3  2  0  0]
 [  3  1  2 61  0  1  0  0  0]
 [  2  1  2  0 79  0  0  0  0]
 [  0  4  4  0  0 124  1  3  5]
 [  0  1  0  0  0  3 67  0  0]
 [  0  2  2  0  0  2  0 66  1]
 [  0  1  0  0  0  4  0  1 56]]
```

Accuracy score: 0.9008620689655172

Multiclass SVM (one vs rest):

Confusion matrix:

```
[[ 60  0  0  0  0  0  0  0  0  0]
 [  0 56  5  1  1  1  1  0  0]
 [  0  6 56  1  2  6  0  1  0]
 [  1  1  2 64  0  0  0  0  0]
 [  1  1  2  0 80  0  0  0  0]
 [  0  9  9  0  0 16  1  2  4]
 [  0  1  0  0  0  4 66  0  0]
 [  0  1  3  0  0  2  0 67  0]
 [  0  1  0  0  0  4  0  1 56]]
```

Accuracy score: 0.8922413793103449

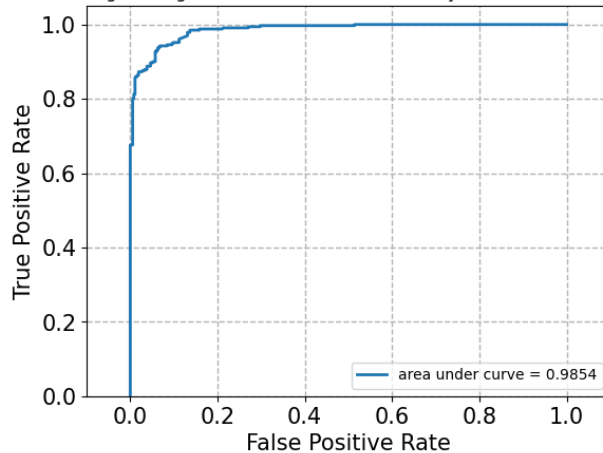
Q10.

- a) According to section 3 of the paper, co-occurrence probabilities help to “distinguish relevant words ... from irrelevant words”, and “discriminate between the two relevant words”. This is because co-occurrence probabilities uses conditional probabilities to find out how likely both words appear together, which serves to highlight the relationship between relevant words in the same context, and also devaluing the other irrelevant words from the selected relevant words.
- b) It would not return the same vector. While both words are identical, both are verbs of the sentence, and both are of the same verb type (continuous action), GLoVE tries to extract information about how related the words are to other words in the same context. For example, GLoVE might assign stronger co-occurrence probabilities for the words “running” and “park” in the first sentence, but for the second sentence the probabilities will be stronger for the words “running” and “presidency”. Since co-occurrence probabilities do not depend on just the word itself, but also the relevant words around it, there will be different embeddings assigned to the word “running” in the 2 sentences.
- c) “Left” and “right” are extremely relevant to each other as both describe directions, so I expect their GLoVE embeddings to be more similar and thus the magnitude of their differences will be less. Similarly, “wife” and “husband” are also very relevant to each other as they are both related to family terms and marriage, so their embeddings differences should be similar to that between “left” and “right”, and the magnitude of the vector difference is similarly small. However, “wife” and “orange” are very far apart from each other in terms of semantics, so their GLoVE embeddings should be vastly different and so the magnitude of their differences is expected to be much larger than the first 2 vector differences.
- d) Lemmatization should be better, because compared to stemming, it extracts the lemma from each word more accurately, so the chances of the extracted word form being closer to the actual lemma form of the word is higher, and GLoVE can be performed on the more accurate lemma version of the word to give an accurate embedding.

Q11.

- a) After lemmatizing `full_text`, we made a `Word2VecVectorizer` that has the same format as the vectorizers from `sklearn`, but the `transform()` function goes through each word in a datapoint, checks if it's in the GLoVe embedding dictionary, and adds the corresponding vector to an array if a match is found. After going through all words in the text segment, we find the mean of the vectors in the array. Those vectors are aligned in a way such that each row contains 1 vector of dimension=GLoVe embedding dimension (which is 300), and each row is for each datapoint.
- b) Using a custom lemmatizer, LSI with $k=25$, and Logistic Regression with L1 regularization strength of 0.01 and maximum of 500 iterations, we reached an accuracy of 0.935. Here are the detailed performance:

ROC curve for L1 logistic regression with LSI dimensionality reduction and GLoVe embeddings



Confusion matrix:

```
[[341 25]
```

```
 [ 20 310]]
```

Accuracy score: 0.9353448275862069

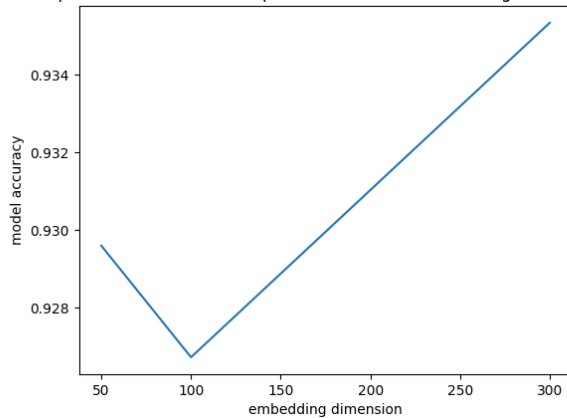
Recall score: 0.9393939393939394

Precision score: 0.9253731343283582

F1 score: 0.9323308270676691

Q12.

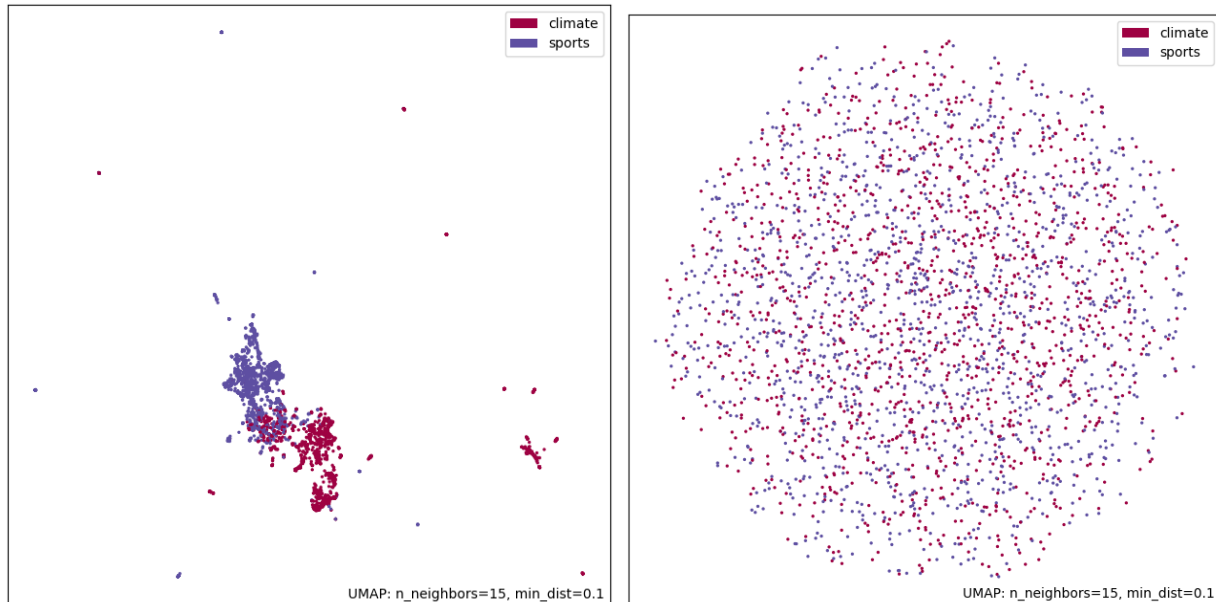
Relationship between dimension of pre-trained GloVe embedding and model accuracy



The observed trend is the accuracy decreases as dimension increases from 50 to 100, then accuracy increases as dimension increases to 300.

This is not expected, as I expect the accuracy to increase as dimension increase, since more information can be contained with a higher embedding dimension that can lead to more features learnt in the model. I also anticipated the accuracy to first increase, then drop as dimension increases, because it's possible that too much information in the embedding can cause overfitting and decrease testing accuracy of models. However, neither of these 2 scenarios occurred. I feel it might simply be a case of small uncertainties, because the dip in model accuracy between dimension=50 and dimension=100 is very small, so it might simply be due to other factors.

Q13.



The left visualization is from the GLoVe embeddings, while the right visualization is from a normalized random-generated matrix of the same shape. We can clearly see that there are clusters being formed for the left visualization, where the embeddings representing ‘climate’ are grouped together below the cluster of embeddings representing ‘sports’. Even though the clusters are not well-separated, we can see some form of clustering being formed. The right visualization shows no cluster whatsoever and is purely random, which is expected of a randomly-generated matrix.