

# EECS 445 Project 1 Report

## Uniquename: boyangx

### October 5, 2019

#### 1. Introduction

The objective of this project is to solve the task of finding the most critically-acclaimed movies by training various Support Vector Machines (SVMs) to classify the sentiment of a review

#### 2. Feature Extraction

- Implement `extract_dictionary` (see code appendix)
- Implement `generate_feature_matrix` (see code appendix)
- The number of unique words, denoted by  $d$ , is 10619  
The average number of non-zero features per rating in the training data is 68.168

#### 3. Hyperparameter and Model Selection

##### 3.1 Hyperparameter selection for linear-kernel SVM

- Implemented `cv_performance` with a helper function `performance` (see code appendix)  
Maintaining class proportions across folds is beneficial because we want draw the training dataset and the test data set from the same distribution so that training data is a representative sample of the test data.
- Implemented `select_param_linear` with a helper function `select_classifier` (see code appendix)
- 

Best C values for different metrics with performances

Performance Measures	C	Performance
Accuracy	0.1	0.8109999999999999
F1-Score	0.1	0.8125587612890968
AUROC	0.1	0.8724399999999999
Precision	0.1	0.8084184094655127
Sensitivity	0.001	0.9099999999999999
Specificity	0.1	0.8039999999999999

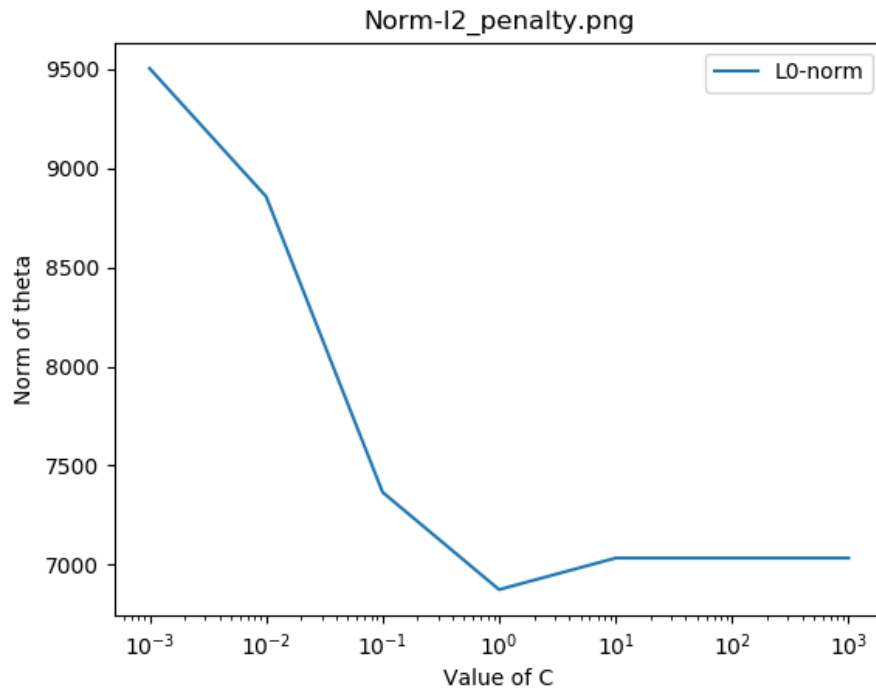
Generally, for all performance measures except for sensitivity, as  $C$  increases, `cv_performance` increases to peak and then decreases. For sensitivity, the `cv_performance` reaches the highest when  $C = 0.001$  and then decreases. If I have to train a final model, I would optimize accuracy measure when choosing  $C$  since accuracy is the most direct measure to address the importance of identifying both positive and negative classes.

- 

Performances of SVM created with  $C = 0.1$  for different metrics

Performance Measures	Performance
Accuracy	0.828
F1-Score	0.8313725490196078
AUROC	0.901088
Precision	0.8153846153846154
Sensitivity	0.848
Specificity	0.808

e.



Implemented plot\_weight (see code appendix)

We can see from the above figure that L0-norm of theta decreases as C increases but gradually converges when C reaches the value of 10.

f.

Positive/Negative coefficients with corresponding words

Positive Coefficient	Word	Negative Coefficient	Word
0.5264765929951933	great	-0.5003295720357759	stupid
0.4816236395618302	hope	-0.4728701857334847	boring
0.4577006461614173	love	-0.45261704023663785	not
0.41676940676838137	well	-0.4292173682944536	nothing

### 3.2 Hyperparameter selection for quadratic-kernel SVM

a. Implemented select\_param\_quadratic (see code appendix)

i. Grid Search (see code appendix)

ii. Random Search (see code appendix)

b.

AUROC results for two different methods

Tuning Scheme	C	R	AUROC
Grid Search	10.0	10.0	0.87446
Random Search	28.09538854212485	7.139888817587492	0.87944

In general, for unchanged C value, the performance typically increases as r value increases; for unchanged r value, the performance typically increases to peak and then decreases as C value increases. The use of random search is better than the grid search since random search will cover more distinct values of C and r. Also, in most cases where C and r do not contribute equally to the performance, random search is more likely to generate closer results to the optimal hyper parameter. However, the random search may be less useful than the grid search when deterministic nature is emphasized.

### 3.3 Learning non-linear classifiers with a linear-kernel SVM

a.

The quadratic kernel is  $K(\bar{x}, \bar{x}') = (\bar{x} \cdot \bar{x}' + r)^2$ . Suppose  $\bar{x}, \bar{x}' \in \mathbb{R}^d$ , we can expand it as

$$\begin{aligned} K(\bar{x}, \bar{x}') &= \left( \sum_{i=1}^d x_i x'_i + r \right)^2 = \left( \sum_{i=1}^d x_i x'_i \right)^2 + 2 \left( \sum_{i=1}^d x_i x'_i \right) r + r^2 \\ &= \sum_{i=1}^d (x_i x'_i)^2 + 2 \sum_{i=1}^{d-1} \sum_{j=i+1}^d (x_i x'_i)(x_j x'_j) + 2 \left( \sum_{i=1}^d x_i x'_i \right) r + r^2 \\ &= \sum_{i=1}^d x_i^2 x'^2_i + 2 \sum_{i=1}^{d-1} \sum_{j=i+1}^d (x_i x_j)(x'_i x'_j) + 2r \left( \sum_{i=1}^d x_i x'_i \right) + r^2 \\ &= \phi(\bar{x}) \phi(\bar{x}') \end{aligned}$$

Therefore  $\phi(x) = [(x_i^2)_{i=1..d}, (\sqrt{2}x_i x_j)_{i=1..d-1, j=i+1..d}, (\sqrt{2r}x_i)_{i=1..d}, r]^T$

b.

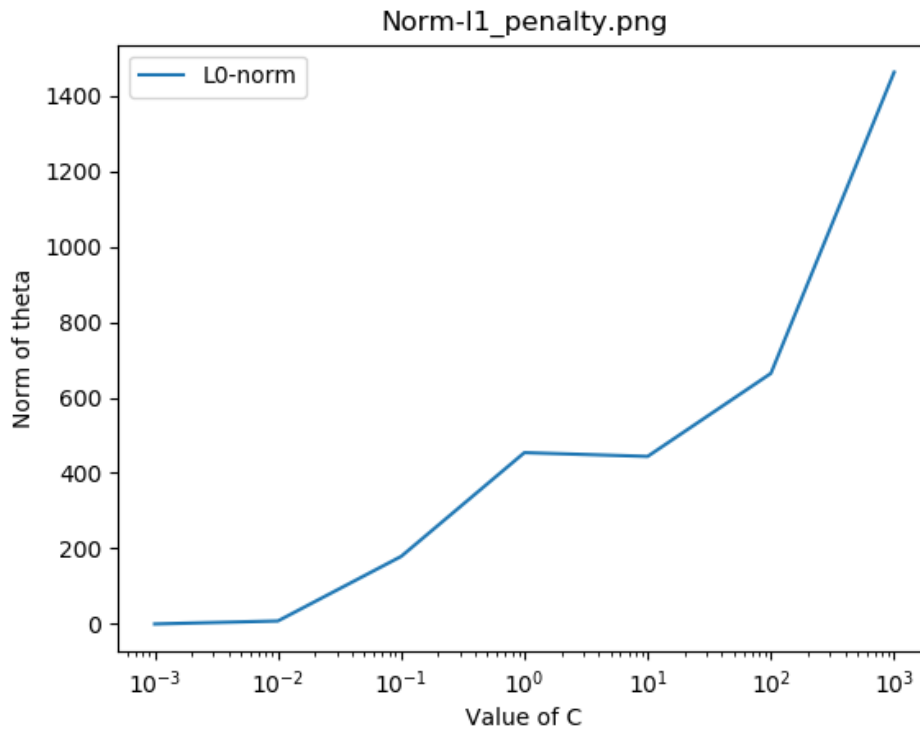
Pros: we can access weights for each feature if using explicit feature mapping

Cons: If the features of the original data is of order  $d$  and the transformed data would have features of order  $d$  square, making it expensive to compute

### 3.4 Linear-kernel SVM with L1 penalty and squared hinge loss

a.  $C = 0.1$  and performance = 0.85694

b.



c.

Primarily, the L0-norm value of the learned parameter significantly dropped under L1 penalty and it increases as the value of  $C$  increases. Also, when value of  $C$  is small, weights are zero.

d.

If Squared Hinge Loss is used instead of the Hinge Loss, the penalty of the misclassified data points will be larger so that more data points would become support vectors and the margin would be larger. Thus, the optimal solution will have a smaller L0-norm under Squared Hinge Loss

#### 4. Asymmetric cost functions and class imbalance

##### 4.1 Arbitrary class weights

a.

If  $W_n$  is much greater than  $W_p$ , the misclassified negative data points would bear a much larger penalty than the misclassified positive ones. This would result in a model which emphasizes more upon classifying negative data points. Referring to the weighted SVM formula, data points with higher weights now would have a lower slack variable and negative data points are more likely to be correctly classified.

b.

Performance of the modified SVM for different metrics

Performance Measures	Performance
Accuracy	0.678
F1-Score	0.5515320334261838
AUROC	0.8669279999999999
Precision	0.908256880733945
Sensitivity	0.396
Specificity	0.96

c.

Compared to result in 3.1d, performances of F1-Score and Sensitivity are the most affected ones. As stated in 4.1a, our model is now emphasizing more on classifying negative data points. Therefore, performance measure for classifying positive data points such as sensitivity would decrease. Performance of F1-Score decreases since it is a function of sensitivity and precision, and it will be closer to the smaller one of the two. In this case, sensitivity is smaller than precision.

##### 4.2 Imbalanced data

a.

Performance of the SVM ( $C = 0.01$ )

Class Weights	Performance Measures	Performance
$W_n = 1, W_p = 1$	Accuracy	0.8
$W_n = 1, W_p = 1$	F1-Score	0.8883928571428571
$W_n = 1, W_p = 1$	AUROC	0.8469
$W_n = 1, W_p = 1$	Precision	0.8024193548387096
$W_n = 1, W_p = 1$	Sensitivity	0.995
$W_n = 1, W_p = 1$	Specificity	0.02

b.

In the given imbalanced data set, there are more positive data points than negative ones. This has caused sensitivity to increase significantly and specificity to decrease significantly. Because of this imbalance, the classifier tends to classify almost all the data as positive. In fact, 99.5% of the positive data are correctly classified as positive, this is shown by the true positive rate (sensitivity). Only 2% of the data are correctly classified as negative, this is shown by the true negative rate (specificity).

### 4.3 Choosing appropriate class weights

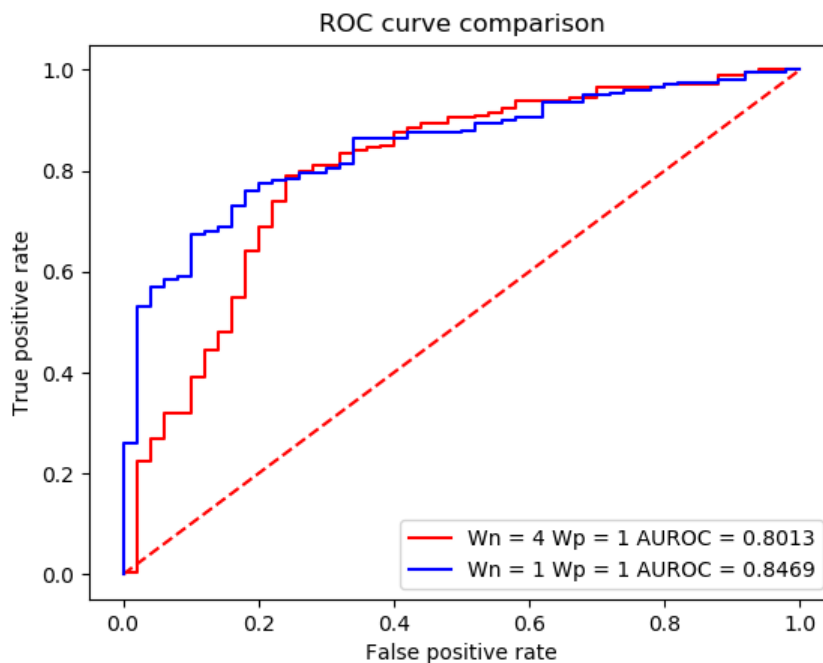
a.

Since our data set is imbalanced, therefore, we want to choose a metric that could balance the classifications of both negative and positive classes, the F1-score is more informative in such situation since it maintains a balance between sensitivity and precision, which helps mitigating the situation in 4.2. If we keep the positive weight as 1 and in order to counter the class size imbalance, negative weights should be at least 4 (i.e.  $W_n$  should be reasonably greater than  $W_p$ ) since the ratio between the class sizes is 4:1. Negative weights found using Cross Validation with respect to f1-score across a reasonably large range of weights. However, since the minor class in this case is the negative class, we want to pick the weight within the range but would have the best performance for the f1-score, which is a metric focuses on classifying the negative data points. Eventually, the found  $W_n$  is 4 and the corresponding performance is 0.8963855421686747.

b.

Class Weights	Performance Measures	Performance
$W_n = 4, W_p = 1$	Accuracy	0.828
$W_n = 4, W_p = 1$	F1-Score	0.8963855421686747
$W_n = 4, W_p = 1$	AUROC	0.8013
$W_n = 4, W_p = 1$	Precision	0.8651162790697674
$W_n = 4, W_p = 1$	Sensitivity	0.93
$W_n = 4, W_p = 1$	Specificity	0.42

### 4.4 The ROC curve



## 5. Challenge

In this part, I implemented the following approaches to train the model:

1. Linear kernel SVM with L1 penalty
2. Linear kernel SVM with L2 penalty and one-vs-one method
3. Linear kernel SVM with L2 penalty and one-vs-rest method
4. Quadratic kernel SVM with L2 penalty and one-vs-one method
5. Feature engineering: using the number of times a word appears in reviews as feature and Linear kernel SVM with L1 penalty
6. Feature engineering: consider rating as part of the feature selection (i.e. extreme ratings such as 1 and 5 are doubly weighted) and Linear kernel SVM with L1 penalty
7. Feature engineering: Do not leave out the punctuations when extracting features from the review text and Linear kernel SVM with L1 penalty

This time, we have three classes and we want to equally classify each class; Therefore, we would use accuracy as the metric. The procedure of choosing C for the linear SVCs is as follows:

- Initialize C as  $10^{-1}$ , I as -3 and step size as 2
- run cv\_performance
- If the new performance is better than the best so far, step size remains as 2
- If the step size is positive, divide it by -2, if the step size is negative, divide it by 2
- Increment I by the step size accordingly

The convergence criteria would be when the absolute difference of the new cv\_performance and the best cv\_performance is less than 0.001. For quadratic SVCs, grid search is used with C and r that lie in the same ranges as the previous parts. When calculating the cv\_performance, I am now using 10-fold instead of 5-fold.

From cross validation performances I generated for each approach, I found out that using linear kernel SVM with L2 penalty would generate the highest accuracy cv\_performance whereas the quadratic kernel SVMs generated a lower cv\_performance than linear kernel SVM with L1 penalty. As for the multiclass methods, one-vs-one and one-vs-rest are producing the same cv\_performance for Linear kernel SVM with the L2 penalty. I modified the original feature matrix generation function to generate `feature_matrix_multi`, which takes into account the number of occurrences of a word in reviews (i.e. by setting Count to True), the existences of extreme ratings (i.e. by setting Rating to True), and the existences of original punctuations (i.e. by setting Punc to True). However, the improve in performance is neglectable with Count, meaning that this feature engineering method would not improve the final result of the trained model by much. One possible reason for this could be that the majority of the repeated words within the review are most likely to be common words such as “we”, “I”, “feel”, etc. These words would not contribute to the classification as they do not show emotional inclination. As for rating and punctuations, there was seen an increase in both of the cv\_performances but the improved performances are still slightly lower than the Linear CVM with L2 penalty, meaning that these two feature engineering methods could potentially generate better models, possible reasons for this could be: extreme rating if assumed that audience are reviewing honestly could reflect more on the actual movie quality; punctuations like “!” “?” “...” sometimes do play an important role in showing emotional inclinations.

In conclusion, I chose to use approach 2 with C being roughly 0.03 and the cv\_performance is 0.633.

```

1  # EECS 445 - Fall 2019
2  # Project 1 - project1.py
3
4  import pandas as pd
5  import numpy as np
6  import itertools
7  import string
8  import warnings
9  import random
10
11 warnings.filterwarnings("ignore", category = FutureWarning, module = "sklearn")
12
13 from sklearn.svm import SVC, LinearSVC
14 from sklearn.model_selection import StratifiedKFold, GridSearchCV
15 from sklearn import metrics
16 from matplotlib import pyplot as plt
17
18 from helper import *
19
20
21 def select_classifier(penalty='l2', c=1.0, degree=1, r=0.0, class_weight='balanced'):
22     """
23     Return a linear svm classifier based on the given
24     penalty function and regularization parameter c.
25     """
26     if penalty == 'l1':
27         return LinearSVC(penalty = 'l1', dual = False, C = c, class_weight =
28             'balanced')
29     elif penalty == 'l2':
30         if degree == 1:
31             return SVC(kernel = 'linear', C = c, degree = 1, class_weight =
32                 class_weight)
33         else:
34             return SVC(kernel = 'poly', C = c, degree = degree, coef0 = r,
35                 class_weight = class_weight)
36
37 def extract_dictionary(df):
38     """
39     Reads a panda dataframe, and returns a dictionary of distinct words
40     mapping from each distinct word to its index (ordered by when it was found).
41     Input:
42         df: dataframe/output of load_data()
43     Returns:
44         a dictionary of distinct words that maps each distinct word
45         to a unique index corresponding to when it was first found while
46         iterating over all words in each review in the dataframe df
47     """
48     word_dict = {}
49     k = 0
50     for entry in df['reviewText']:
51         entry = entry.lower()
52         for char in entry:
53             if char in string.punctuation:
54                 entry = entry.replace(char, ' ')
55         for word in entry.split():
56             if word not in word_dict.keys():
57                 word_dict[word] = k
58                 k += 1
59     return word_dict
60
61 def generate_feature_matrix(df, word_dict):
62     """
63     Reads a dataframe and the dictionary of unique words
64     to generate a matrix of {1, 0} feature vectors for each review.
65     Use the word_dict to find the correct index to set to 1 for each place
66     in the feature vector. The resulting feature matrix should be of
67     dimension (number of reviews, number of words).
68     Input:
69         df: dataframe that has the ratings and labels
70         word_list: dictionary of words mapping to indices

```

```

70     word_list: dictionary of words mapping to indices
71 Returns:
72     a feature matrix of dimension (number of reviews, number of words)
73 """
74 number_of_reviews = df.shape[0]
75 number_of_words = len(word_dict)
76 feature_matrix = np.zeros((number_of_reviews, number_of_words))
77 for i in range(number_of_reviews):
78     entry = df['reviewText'][i]
79     entry = entry.lower()
80     for char in entry:
81         if char in string.punctuation:
82             entry = entry.replace(char, ' ')
83     for word in word_dict.keys():
84         if word in entry.split():
85             feature_matrix[i][word_dict[word]] = 1
86 return feature_matrix
87
88 def generate_feature_matrix_multi(df, word_dict, Count = False, Rating = False, Punc
= False):
89
90     number_of_reviews = df.shape[0]
91     number_of_words = len(word_dict)
92     feature_matrix = np.zeros((number_of_reviews, number_of_words))
93     for i in range(number_of_reviews):
94         entry = df['reviewText'][i]
95         rating = df['rating'][i]
96         entry = entry.lower()
97         if Punc == False:
98             for char in entry:
99                 if char in string.punctuation:
100                     entry = entry.replace(char, ' ')
101         for word in word_dict.keys():
102             if word in entry.split():
103                 if Count:
104                     feature_matrix[i][word_dict[word]] += 1
105                 elif Rating:
106                     if rating == 1 or rating == 5:
107                         feature_matrix[i][word_dict[word]] = 2
108                 else:
109                     feature_matrix[i][word_dict[word]] = 1
110 return feature_matrix
111
112
113 def cv_performance(clf, X, y, k=5, metric="accuracy"):
114     """
115     Splits the data X and the labels y into k-folds and runs k-fold
116     cross-validation: for each fold i in 1...k, trains a classifier on
117     all the data except the ith fold, and tests on the ith fold.
118     Calculates the k-fold cross-validation performance metric for classifier
119     clf by averaging the performance across folds.
120     Input:
121         clf: an instance of SVC()
122         X: (n,d) array of feature vectors, where n is the number of examples
123            and d is the number of features
124         y: (n,) array of binary labels {1,-1}
125         k: an int specifying the number of folds (default=5)
126         metric: string specifying the performance metric (default='accuracy'
127            other options: 'f1-score', 'auroc', 'precision', 'sensitivity',
128            and 'specificity')
129     Returns:
130         average 'test' performance across the k folds as np.float64
131     """
132     scores = np.zeros(k)
133     skf = StratifiedKFold(n_splits = k, shuffle = False)
134     i = 0
135     for train_index, test_index in skf.split(X,y):
136         X_train, y_train = X[train_index], y[train_index]
137         clf.fit(X_train, y_train)
138         X_test, y_test = X[test_index], y[test_index]
139         # Put the performance of the model on each fold in the scores array
140         if metric == "AUROC":

```



```

141         y_pred = clf.decision_function(X_test)
142     else:
143         y_pred = clf.predict(X_test)
144     scores[i] = performance(y_test, y_pred, metric)
145     i += 1
146
147     #And return the average performance across all fold splits.
148     return np.float64(np.array(scores).mean())
149
150
151 def select_param_linear(X, y, k=5, metric="accuracy", C_range = [], penalty='l2'):
152     """
153     Sweeps different settings for the hyperparameter of a linear-kernel SVM,
154     calculating the k-fold CV performance for each setting on X, y.
155     Input:
156         X: (n,d) array of feature vectors, where n is the number of examples
157            and d is the number of features
158         y: (n,) array of binary labels {1,-1}
159         k: int specifying the number of folds (default=5)
160         metric: string specifying the performance metric (default='accuracy',
161            other options: 'f1-score', 'auroc', 'precision', 'sensitivity',
162            and 'specificity')
163         C_range: an array with C values to be searched over
164     Returns:
165         The parameter value for a linear-kernel SVM that maximizes the
166         average 5-fold CV performance.
167     """
168     best_c = 0.0
169     best_performance = 0
170     for c in C_range:
171         clf = select_classifier(penalty = penalty, c = c, degree = 1, r = 0.0,
172                                class_weight = 'balanced')
173         new_performance = cv_performance(clf, X, y, k, metric)
174         if new_performance > best_performance:
175             best_performance = new_performance
176             best_c = c
177     return best_c, best_performance
178
179 def plot_weight(X,y,penalty,C_range):
180     """
181     Takes as input the training data X and labels y and plots the L0-norm
182     (number of nonzero elements) of the coefficients learned by a classifier
183     as a function of the C-values of the classifier.
184     """
185
186     print("Plotting the number of nonzero entries of the parameter vector as a
187           function of C")
188     norm0 = []
189     for c in C_range:
190         clf = select_classifier(penalty = penalty, c = c, degree = 1, r = 0.0,
191                                class_weight = 'balanced')
192         clf.fit(X,y)
193         for i in clf.coef_[0]:
194             if i != 0:
195                 norm += 1
196         norm0.append(norm)
197     #This code will plot your L0-norm as a function of c
198     plt.plot(C_range, norm0)
199     plt.xscale('log')
200     plt.legend(['L0-norm'])
201     plt.xlabel("Value of C")
202     plt.ylabel("Norm of theta")
203     plt.title('Norm-'+penalty+' _penalty.png')
204     plt.savefig('Norm-'+penalty+' _penalty.png')
205     plt.close()
206
207 def select_param_quadratic(X, y, k=5, metric="accuracy", param_range=[]):
208     """
209     Sweeps different settings for the hyperparameters of an quadratic-kernel SVM,

```

```

210     calculating the k-fold CV performance for each setting on X, y.
211     Input:
212         X: (n,d) array of feature vectors, where n is the number of examples
213             and d is the number of features
214         y: (n,) array of binary labels {1,-1}
215         k: an int specifying the number of folds (default=5)
216         metric: string specifying the performance metric (default='accuracy'
217             other options: 'f1-score', 'auroc', 'precision', 'sensitivity',
218             and 'specificity')
219         parameter_values: a (num_param, 2)-sized array containing the
220             parameter values to search over. The first column should
221             represent the values for C, and the second column should
222             represent the values for r. Each row of this array thus
223             represents a pair of parameters to be tried together.
224     Returns:
225         The parameter value(s) for a quadratic-kernel SVM that maximize
226         the average 5-fold CV performance
227     """
228     best_c, best_r, best_performance = 0.0, 0.0, 0.0
229     for c, r in param_range:
230         clf = select_classifier(penalty = 'l2', c = c, degree = 2, r = r,
231             class_weight = 'balanced')
232         new_performance = cv_performance(clf, X, y, k, metric)
233         if new_performance > best_performance:
234             best_performance = new_performance
235             best_c = c
236             best_r = r
237     return best_c, best_r, best_performance
238
239 def performance(y_true, y_pred, metric="accuracy"):
240     """
241     Calculates the performance metric as evaluated on the true labels
242     y_true versus the predicted labels y_pred.
243     Input:
244         y_true: (n,) array containing known labels
245         y_pred: (n,) array containing predicted scores
246         metric: string specifying the performance metric (default='accuracy'
247             other options: 'f1-score', 'auroc', 'precision', 'sensitivity',
248             and 'specificity')
249     Returns:
250         the performance as an np.float64
251     """
252     if metric == "accuracy":
253         return metrics.accuracy_score(y_true, y_pred)
254     elif metric == "f1-score":
255         return metrics.f1_score(y_true, y_pred)
256     elif metric == "precision":
257         return metrics.precision_score(y_true, y_pred)
258     elif metric == "AUROC":
259         return metrics.roc_auc_score(y_true, y_pred)
260     elif metric == "sensitivity":
261         m = metrics.confusion_matrix(y_true, y_pred, labels = [-1,1])
262         return m[1,1] / (m[1,1] + m[1,0])
263     elif metric == "specificity":
264         m = metrics.confusion_matrix(y_true, y_pred, labels = [-1,1])
265         return m[0,0] / (m[0,0] + m[0,1])
266
267 #@jit(target="cuda")
268 def main():
269     # Read binary data
270     # NOTE: READING IN THE DATA WILL NOT WORK UNTIL YOU HAVE FINISHED
271     # IMPLEMENTING generate_feature_matrix AND extract_dictionary
272     X_train, Y_train, X_test, Y_test, dictionary_binary = get_split_binary_data()
273     IMB_features, IMB_labels = get_imbalanced_data(dictionary_binary)
274     IMB_test_features, IMB_test_labels = get_imbalanced_test(dictionary_binary)
275
276     # Q2
277     print(len(dictionary_binary))
278     num_feature = 0
279     for i in X_train:
280         for j in i:
281             if j == 1:

```

```

281         num_feature += 1
282     avg_feature = num_feature / X_train.shape[0]
283     print(avg_feature)
284
285     # Q3.1(c)
286     C_range = [1e-3, 1e-2, 1e-1, 1, 1e1, 1e2, 1e3]
287     print("Accuracy", select_param_linear(X_train, Y_train, 5, "accuracy", C_range,
288     penalty = 'l2'))
289     print("F1-Score", select_param_linear(X_train, Y_train, 5, "f1-score", C_range,
290     penalty = 'l2'))
291     print("AUROC", select_param_linear(X_train, Y_train, 5, "AUROC", C_range,
292     penalty = 'l2'))
293     print("Precision", select_param_linear(X_train, Y_train, 5, "precision",
294     C_range, penalty = 'l2'))
295     print("Sensitivity", select_param_linear(X_train, Y_train, 5, "sensitivity",
296     C_range, penalty = 'l2'))
297     print("Specificity", select_param_linear(X_train, Y_train, 5, "specificity",
298     C_range, penalty = 'l2'))
299
300     # Q3.1(d)
301     clf = SVC(kernel='linear', C = 1e-1)
302     clf.fit(X_train, Y_train)
303     Y_pred = clf.predict(X_test)
304     Y_pred_auroc = clf.decision_function(X_test)
305     print("Accuracy", performance(Y_test, Y_pred, metric = "accuracy"))
306     print("F1-Score", performance(Y_test, Y_pred, metric = "f1-score"))
307     print("AUROC", performance(Y_test, Y_pred_auroc, metric = "AUROC"))
308     print("Precision", performance(Y_test, Y_pred, metric = "precision"))
309     print("Sensitivity", performance(Y_test, Y_pred, metric = "sensitivity"))
310     print("Specificity", performance(Y_test, Y_pred, metric = "specificity"))
311
312     # Q3.1(e)
313     plot_weight(X_train, Y_train, 'l2', C_range)
314
315     # Q3.1(f)
316     clf = select_classifier(penalty = 'l2', c = 0.1, degree = 1, r = 0.0,
317     class_weight = 'balanced')
318     clf.fit(X_train, Y_train)
319     arg = clf.coef_[0].argsort()
320     neg_ind4 = arg[:4]
321     pos_ind4 = arg[:-5:-1]
322     neg_words = []
323     pos_words = []
324     for idx in neg_ind4:
325         for word, index in dictionary_binary.items():
326             if index == idx:
327                 neg_words.append(word)
328     print("Most negative words")
329     for i in range(4):
330         print(clf.coef_[0, neg_ind4[i]], neg_words[i])
331     for idx in pos_ind4:
332         for word, index in dictionary_binary.items():
333             if index == idx:
334                 pos_words.append(word)
335     print("Most positive words")
336     for i in range(4):
337         print(clf.coef_[0, pos_ind4[i]], pos_words[i])
338
339     # Q3.2(a)
340     # (i)
341     grid = []
342     c_range = [1e-3, 1e-2, 1e-1, 1, 1e1, 1e2, 1e3]
343     r_range = [1e-3, 1e-2, 1e-1, 1, 1e1, 1e2, 1e3]
344     for i in c_range:
345         for j in r_range:
346             grid.append([i, j])
347     print(select_param_quadratic(X_train, Y_train, k = 5, metric = "AUROC",
348     param_range = grid))
349     #ii)
350     param_random = np.zeros([25, 2])
351     for i in range(25):
352         [c, r] = [pow(10, np.random.uniform(-3, 3)), pow(10, np.random.uniform(-3, 3))]

```

```

345     param_random[i] = [c,r]
346     print(select_param_quadratic(X_train, Y_train, k = 5, metric = "AUROC",
347                                param_range = param_random))
348
349     # 3.4(a)
350     print(select_param_linear(X_train, Y_train, k = 5, metric = "AUROC", C_range =
351                                C_range, penalty='l1'))
352
353     # 3.4(b)
354     plot_weight(X_train, Y_train, 'l1', C_range)
355
356     # 4.1(b)
357     clf = select_classifier(penalty = 'l2', c = 0.01, degree = 1, class_weight =
358                                {-1: 10, 1: 1})
359     clf.fit(X_train, Y_train)
360     Y_pred = clf.predict(X_test)
361     Y_pred_auc = clf.decision_function(X_test)
362     print("Accuracy", performance(Y_test, Y_pred, metric="accuracy"))
363     print("F1-Score", performance(Y_test, Y_pred, metric="f1-score"))
364     print("AUROC", performance(Y_test, Y_pred_auc, metric="AUROC"))
365     print("Precision", performance(Y_test, Y_pred, metric="precision"))
366     print("Sensitivity", performance(Y_test, Y_pred, metric="sensitivity"))
367     print("Specificity", performance(Y_test, Y_pred, metric="specificity"))
368
369     # 4.2(a)
370     clf = select_classifier(penalty = 'l2', c = 0.01, degree = 1, class_weight =
371                                {-1: 1, 1: 1})
372     clf.fit(IMB_features, IMB_labels)
373     Labels_pred = clf.predict(IMB_test_features)
374     Labels_pred_auc = clf.decision_function(IMB_test_features)
375     print("Accuracy", performance(IMB_test_labels, Labels_pred, metric="accuracy"))
376     print("F1-Score", performance(IMB_test_labels, Labels_pred, metric="f1-score"))
377     print("AUROC", performance(IMB_test_labels, Labels_pred_auc, metric="AUROC"))
378     print("Precision", performance(IMB_test_labels, Labels_pred, metric="precision"))
379     print("Sensitivity", performance(IMB_test_labels, Labels_pred,
380                                metric="sensitivity"))
381     print("Specificity", performance(IMB_test_labels, Labels_pred,
382                                metric="specificity"))
383
384     # 4.3(a)
385     best_wn, best_perf = 0.0, 0.0
386     for i in range(25):
387         Wn = 2**random.uniform(2, 5) # at least 4
388         clf = select_classifier(penalty = 'l2', c = 0.1, degree = 1, class_weight =
389                                    {-1: Wn, 1: 1})
390         new_perf = cv_performance(clf, IMB_features, IMB_labels, 5, metric =
391                                    "f1-score" )
392         spec_perf = cv_performance(clf, IMB_features, IMB_labels, 5, metric =
393                                    "specificity" )
394         print(Wn, new_perf, spec_perf)
395         if new_perf > best_perf:
396             best_perf = new_perf
397             best_wn = Wn
398     print("f1-score (random search)", best_wn, best_perf)
399
400     # 4.3(b)
401     clf = select_classifier(penalty = 'l2', c = 0.1, degree = 1, class_weight = {-1:
402                                best_wn, 1: 1})
403     clf.fit(IMB_features, IMB_labels)
404     Labels_pred = clf.predict(IMB_test_features)
405     Labels_pred_auc = clf.decision_function(IMB_test_features)
406     print("Accuracy", performance(IMB_test_labels, Labels_pred, metric="accuracy"))
407     print("F1-Score", performance(IMB_test_labels, Labels_pred, metric="f1-score"))
408     print("AUROC", performance(IMB_test_labels, Labels_pred_auc, metric="AUROC"))
409     print("Precision", performance(IMB_test_labels, Labels_pred, metric="precision"))
410     print("Sensitivity", performance(IMB_test_labels, Labels_pred,
411                                metric="sensitivity"))
412     print("Specificity", performance(IMB_test_labels, Labels_pred,
413                                metric="specificity"))
414
415     # 4.4

```

```

405 # customized Wn Wp
406 fpr, tpr, thresholds = metrics.roc_curve(IMB_test_labels, Labels_pred_auc)
407 # balanced
408 clf_bal = select_classifier(penalty = 'l2', c = 0.01, degree = 1, class_weight =
409 {-1: 1, 1: 1})
410 clf_bal.fit(IMB_features, IMB_labels)
411 Labels_pred_auc_bal = clf_bal.decision_function(IMB_test_features)
412 fpr_bal, tpr_bal, thresholds_bal = metrics.roc_curve(IMB_test_labels,
413 Labels_pred_auc_bal)
414 plt.figure()
415 plt.plot([0,1], [0,1], 'r--')
416 plt.plot(fpr, tpr, color = 'red', label = 'Wn = 4 Wp = 1 AUROC = 0.8013')
417 plt.plot(fpr_bal, tpr_bal, color = 'blue', label = 'Wn = 1 Wp = 1 AUROC = 0.8469')
418 plt.legend(loc = "lower right")
419 plt.xlabel('False positive rate')
420 plt.ylabel('True positive rate')
421 plt.title('ROC curve comparison')
422 plt.savefig('ROC curve comparison.png')
423 plt.close()
424
425 # Read multiclass data
426 # TODO: Question 5: Apply a classifier to heldout features, and then use
427 # generate_challenge_labels to print the predicted labels
428 multiclass_features, multiclass_labels, multiclass_dictionary =
429 get_multiclass_training_data()
430 heldout_features = get_heldout_reviews(multiclass_dictionary)
431
432 # Approach 1 Linear kernel SVM with l1 penalty
433 best_c, best_perf = 0.0, 0.0
434 I = -3
435 i = 2
436 threshold = 1
437 while abs(threshold) > 0.001:
438     c = 10 ** I
439     clf = LinearSVC(penalty = 'l1', dual = False, C = c, class_weight =
440 'balanced', max_iter = 100000)
441     new_perf = cv_performance(clf, multiclass_features, multiclass_labels, 10)
442     threshold = new_perf - best_perf
443     print(c, new_perf)
444     if new_perf > best_perf:
445         best_perf = new_perf
446         best_c = c
447         i = 2
448         I = I + i
449     else:
450         if i > 0:
451             i = i/-2
452         else:
453             i = i/2
454         I = I + i
455     print("Linear-l1 optimal c and performance:", best_c, best_perf)
456
457 # Approach 2 Linear kernel SVN with l2 penalty and ovo method
458 best_c, best_perf = 0.0, 0.0
459 I = -3
460 i = 2
461 threshold = 1
462 while abs(threshold) > 0.001:
463     c = 10 ** I
464     clf = SVC(kernel = 'linear', C = c, degree = 1, class_weight = 'balanced',
465 decision_function_shape='ovo')
466     new_perf = cv_performance(clf, multiclass_features, multiclass_labels, 10)
467     threshold = new_perf - best_perf
468     print(c, new_perf)
469     if new_perf > best_perf:
470         best_perf = new_perf
471         best_c = c
472         i = 2
473         I = I + i
474     else:
475         if i > 0:

```

```

472         i = i/-2
473     else:
474         i = i/2
475     I = I + i
476 print("Linear-l2(ovo) optimal c and performance:", best_c, best_perf)
477
478 # Approach 3 Linear kernel SVN with l2 penalty and ovr method
479 best_c, best_perf = 0.0, 0.0
480 I = -3
481 i = 2
482 threshold = 1
483 while abs(threshold) > 0.001:
484     c = 10 ** I
485     clf = SVC(kernel = 'linear', C = c, degree = 1, class_weight = 'balanced')
486     new_perf = cv_performance(clf, multiclass_features, multiclass_labels, 10)
487     threshold = new_perf - best_perf
488     print(c, new_perf)
489     if new_perf > best_perf:
490         best_perf = new_perf
491         best_c = c
492         i = 2
493         I = I + i
494     else:
495         if i > 0:
496             i = i/-2
497         else:
498             i = i/2
499         I = I + i
500 print("Linear-l2(ovr) optimal c and performance:", best_c, best_perf)
501
502 # Approach 4 quadratic ovo
503 c_range = [1e-3, 1e-2, 1e-1, 1, 1e1, 1e2, 1e3]
504 r_range = [1e-3, 1e-2, 1e-1, 1, 1e1, 1e2, 1e3]
505 best_c, best_r, best_perf = 0.0, 0.0, 0.0
506 for c in c_range:
507     for r in r_range:
508         clf = SVC(kernel = 'poly', C = c, coef0 = r, degree = 2, class_weight =
                    'balanced', decision_function_shape='ovo')
509         new_perf = cv_performance(clf, multiclass_features, multiclass_labels, 10)
510         if new_perf > best_perf:
511             best_perf = new_perf
512             best_c = c
513             best_r = r
514 print("Quadratic(ovo) optimal c, r and performance:", best_c, best_r, best_perf)
515
516 # Approach 5 feature engineering
517 # Using the number of times a word occurs in a review as a feature
518 multiclass_features, multiclass_labels, multiclass_dictionary =
get_multiclass_training_data(Count = True)
519 # L1-linear
520 best_c, best_perf = 0.0, 0.0
521 I = -3
522 i = 2
523 threshold = 1
524 while abs(threshold) > 0.001:
525     c = 10 ** I
526     clf = LinearSVC(penalty = 'l1', dual = False, C = c, class_weight =
                    'balanced')
527     new_perf = cv_performance(clf, multiclass_features, multiclass_labels, 10)
528     threshold = new_perf - best_perf
529     print(c, new_perf)
530     if new_perf > best_perf:
531         best_perf = new_perf
532         best_c = c
533         i = 2
534         I = I + i
535     else:
536         if i > 0:
537             i = i/-2
538         else:
539             i = i/2
540         I = I + i

```

```

541     print("Linear-l1 with count optimal c and performance:", best_c, best_perf)
542
543     # Approach 6 feature engineering
544     # Consider rating when generating features
545     multiclass_features, multiclass_labels, multiclass_dictionary =
get_multiclass_training_data(Rating = True)
546     # L1-linear
547     best_c, best_perf = 0.0, 0.0
548     I = -3
549     i = 2
550     threshold = 1
551     while abs(threshold) > 0.001:
552         c = 10 ** I
553         clf = LinearSVC(penalty = 'l1', dual = False, C = c, class_weight =
'balanced')
554         new_perf = cv_performance(clf, multiclass_features, multiclass_labels, 10)
555         threshold = new_perf - best_perf
556         print(c, new_perf)
557         if new_perf > best_perf:
558             best_perf = new_perf
559             best_c = c
560             i = 2
561             I = I + i
562         else:
563             if i > 0:
564                 i = i/-2
565             else:
566                 i = i/2
567             I = I + i
568     print("Linear-l1 with rating optimal c and performance:", best_c, best_perf)
569
570     # Approach 7 feature engineering
571     # Do consider punctuations
572     multiclass_features, multiclass_labels, multiclass_dictionary =
get_multiclass_training_data(Punc = True)
573     # L1-linear
574     best_c, best_perf = 0.0, 0.0
575     I = -3
576     i = 2
577     threshold = 1
578     while abs(threshold) > 0.001:
579         c = 10 ** I
580         clf = LinearSVC(penalty = 'l1', dual = False, C = c, class_weight =
'balanced')
581         new_perf = cv_performance(clf, multiclass_features, multiclass_labels, 10)
582         threshold = new_perf - best_perf
583         print(c, new_perf)
584         if new_perf > best_perf:
585             best_perf = new_perf
586             best_c = c
587             i = 2
588             I = I + i
589         else:
590             if i > 0:
591                 i = i/-2
592             else:
593                 i = i/2
594             I = I + i
595     print("Linear-l1 with punctuation optimal c and performance:", best_c, best_perf)
596
597     # Conclusion
598     multiclass_features, multiclass_labels, multiclass_dictionary =
get_multiclass_training_data()
599     heldout_features = get_heldout_reviews(multiclass_dictionary)
600     clf = SVC(kernel = 'linear', C = 0.03, degree = 1, class_weight = 'balanced',
decision_function_shape='ovo')
601     clf.fit(multiclass_features, multiclass_labels)
602     y_pred = clf.predict(heldout_features)
603     generate_challenge_labels(y_pred, "boyangx")
604
605     if __name__ == '__main__':
606         main()

```