

UNIVERSITY OF MICHIGAN
Department of Electrical Engineering and Computer Science
EECS 445 — Introduction to Machine Learning
Winter 2021

Project 1 - Naitian's Super South Pole SVMs
Due: Friday, 2/12 at 8:00pm

Section	Points	Recommended Completion Date
2. Feature Extraction	20	Tuesday, 2/2
3. Hyperparameter and Model Selection	40	Sunday, 2/7
4. Asymmetric Cost Functions and Class Imbalance	20	Tuesday, 2/9
5. Challenge	20	Thursday, 2/11

Include your code (copy and pasted) in your report. Please try to format lines of code so they are visible within the pages. **Make sure to match your code with the correct problems on Gradescope.** Failure to do so, may result in a 0 for the corresponding problems.

Upload your file `username.csv` containing the label predictions for the held-out data to the canvas assignment named Project 1 Challenge Submissions.

1 Introduction

Naitian has decided to socially distance himself and is moving to Antarctica to live with penguins. Since penguins don't have internet, he's hoping to pass time by reading a lot of books. He doesn't know which books he wants yet, and there's too many options for him to decide before his boat leaves. Thankfully, Naitian has a group of EECS 445 students at his disposal, who have become well-versed in solving complex supervised Machine Learning problems. He plans to solve the task of finding the best books to read by training a model to deduce the sentiment of their Amazon reviews (i.e determine if the reviewers think the book was worth reading).

In this project, we have given you review data from Amazon's large catalogue of book reviews. The Amazon book review data set contains thousands of reviews and ratings from different reviewers on different books. You will work with this dataset to train various Support Vector Machines (SVMs) to classify the sentiment of a review. That way Naitian can automate the process of choosing and will be able to decide what books to read in time for his Journey. In this process, we will also explore some very useful scikit-learn packages and data science techniques.

1.1 Requirements:

1. Updated version of Python (<https://www.python.org/downloads/>), with a Python 3.7+ virtual environment.
2. Updated version of `scikit-learn` (0.24): <https://scikit-learn.org/stable/index.html>

3. Updated version of numpy (1.19): <http://www.numpy.org/>
4. Updated version of pandas (1.2.1): <https://pandas.pydata.org/>
5. Updated version of matplotlib (3.3.3): <https://matplotlib.org/>

1.2 Getting Started

To get started, download `Project1` from Canvas. It should contain the following files:

- `data/dataset.csv`
- `data/heldout.csv`
- `data/imbalanced.csv`
- `project1.py`
- `helper.py`
- `test_output.py`
- `test_cases.py`

The files `dataset.csv` and `imbalanced.csv` have book reviews from Amazon. These csv files have 7 columns: *reviewText*, *summary*, *unixReviewTime*, *helpful*, *unhelpful*, *rating*, and *label*. Each row in the csv file corresponds to one review. The *reviewText* column contains the text of the actual review. The *label* column is a multiclass label: 1 if positive (greater than 3 on Amazon), 0 if neutral (3), and -1 if negative (less than 3).

You will use the *reviewText* and *label* columns for most of the project (we will ignore the 0 label reviews in order to make the label binary). The final challenge portion, however, will utilize all -1, 0, and 1 labeled reviews.

The helper file `helper.py` provides functions that allow you to read in the data from csv files. The file `project1.py` contains skeleton code for the project, along with the helper function `select_classifier` which you may implement to return SVM classifiers depending on the given input parameters. The file `test_output.py` allows you to test your output csv file before submission to make sure the format is correct. **The file `test_cases.py` contains a few checkpoints for section 2 and 3 to help make sure your feature extraction and parameter selection code is functioning correctly.**

The data for each part of the project has already been read in for you in the main function of the skeleton code. Please do not change how the data is read in; doing so may affect your results.

The skeleton code `project1.py` provides specifications for functions that you will implement. There may be additional functions that you will have to implement on your own:

- `extract_dictionary(df)`
- `generate_feature_matrix(df, worddict)`
- `cv_performance(clf, X, y, k=5, metric='accuracy')`
- `select_param_linear(X, y, k=5, metric='accuracy', C_range=[], penalty='l2')`
- `plot_weight(X, y, penalty, C_range)`
- `select_param_quadratic(X, y, k=5, metric='accuracy', param_range = [])`

- `get_perceptron_boundary(X_train, Y_train)`
- Optional: `select_classifier(penalty='l2', c=1.0, degree=1, r=0.0, class_weight='balanced')`
- Optional: `performance(y_true, y_pred, metric='accuracy')`

1.3 Submitting Your Work

This project contains questions that involve coding as well as others where you'll be asked to write up an answer and submit to Gradescope. **Coding questions will be highlighted green, and questions that require written answers will be highlighted blue.** Please make sure to complete both and attach your code to your report. When submitting to Gradescope, **make sure to match your code to the correct problem.** Failure to do so may result in a 0 in the corresponding problem.

2 Feature Extraction [20 pts]

Given a dictionary containing d unique words, we can transform the n variable-length reviews into n feature vectors of length d , by setting the i^{th} element of the j^{th} feature vector to 1 if the i^{th} word is in the j^{th} review, and 0 otherwise. Given that the four words $\{\text{'book':0, 'was':1, 'the':2, 'best':3}\}$ are the only four words we ever encounter, the review “*BEST book ever!!*” would map to the feature vector $[1, 0, 0, 1]$.

Note that we do not consider case. Also, note that since the word “ever” was not in the original dictionary, it is ignored as a feature. In real-world scenarios, there may be words in test data that you do not encounter in training data. There are many interesting methods for dealing with this that you may explore in part 5.

- (a) (11 pt) **Start by implementing the `extract_dictionary(df)` function.** You will use this function to read all unique words contained in `dataset.csv` into a dictionary (as in the example above). You can start implementing this function by removing all the punctuation in the dataset. While removing punctuation, please make sure that you do not accidentally combine two different words that are separated by a punctuation mark. For instance, after you remove punctuation from “*Book was awesome!Yay*”, you should produce “*Book was awesome Yay*”, not “*Book was awesomeYay*”. After removing all the punctuation, you should convert all the words to lowercase and start building your dictionary. Your function should return a dictionary of d unique words. Make sure to test your implementation with the `test_cases.test_dictionary()` method.

Note: You will need to report the number of unique words in 2(c).

Hint: You might find `string.punctuation` along with the method `string.replace()` useful.

- (b) (6 pt) **Next, implement the `generate_feature_matrix(df, word_dict)` function.** Assuming that there are n reviews total, return the feature vectors as an (n, d) feature matrix, where each row represents a review, and each column represents whether or not a specific word appeared in that review. Make sure to test your implementation with the `test_cases.test_feature_matrix()` method.
- (c) (3 pt) The function `get_split_binary_data` in `helper.py` uses the functions you implemented in (a) and (b). Use `get_split_binary_data` to get the training feature matrix `X_train`.

Note the `class_size` parameter in `get_split_binary_data`. If at any point you are not confident that your algorithm is working as intended, it might be worth reducing the class size to try out running on a smaller input size. Otherwise, your algorithm may take a few minutes to terminate, only to find out that it is not working properly. **However, for all exercises that ask for results, please use the default parameters.**

In your write-up, include the following:

- **The value of d which you recorded after extracting the training data (the number of unique words).** You should be able to extract d from the size of the training feature matrix.
- **The average number of non-zero features per rating in the training data.** You will need to calculate this on `X_train`.
- **The word appearing in the most number of reviews.** You may find it helpful to use an additional data structure.

Solution:

- The value of d which you recorded after extracting the dictionary: 7508
- The average number of non-zero *features* per rating: 53.794
- The word appearing in the most number of reviews: "the"

3 Hyperparameter and Model Selection [40 pts]

In section 2, you have implemented functions that transform the reviews into a feature matrix `X_train` and a label vector `y_train`. Test data `X_test`, `y_test` has also been read in for you. **You will use this data for all of question 3.** You may notice that `X_train`, `y_train` only have 1000 reviews, while the `dataset.csv` file has 6750 reviews. Here, we only give you a subset of the data to train on.

We will learn a classifier to separate the *training* data into positive and non-positive (i.e., “negative”) labels. The labels in `y_train` are transformed into binary labels in $\{-1, 1\}$, where -1 means “poor” and 1 means “good.”

For the classifier, we will use SVMs with two different kernels: linear and quadratic. In parts 3.1-3.3 we will make use of the `sklearn.svm.SVC` class. At first, we will explicitly set only two of the initialization parameters of `SVC()`: the `kernel`, and `C`. In addition, we will use the following methods in the `SVC` class: `fit(X, y)`, `predict(X)` and `decision_function(X)` – please use `predict(X)` when measuring for any performance metric that is not AUROC and `decision_function(X)` for AUROC (see the documentation for more details).

As discussed in lecture, SVMs have hyperparameters that must be set by the user. For both linear-kernel and quadratic-kernel SVMs, we will select hyperparameters using 5-fold cross-validation (CV) on the training data. We will select the hyperparameters that lead to the ‘best’ mean performance across all five folds. The result of hyperparameter selection often depends upon the choice of performance measure. Here, we will consider the following performance measures: **Accuracy, F1-Score, AUROC, Precision, Sensitivity, and Specificity.**

Note: When calculating the F1-score, it is possible that a divide-by-zero may occur which throws a warning. Consider how this metric is calculated, perhaps by reviewing the relevant `scikit-learn` documentation.

Some of these measures are already implemented as functions in the `sklearn.metrics` submodule. Please use `sklearn.metrics.roc_auc_score` for AUROC. You can use the values from `sklearn.metrics.confusion_matrix` to calculate the others (Note – the confusion matrix is just the table of Predicted vs. Actual label counts, that is, the True Positive, False Positive, True Negative, and False Negative counts for binary classification). Make sure to read the documentation carefully, as when calling this function you will want to set `labels=[1, -1]` for a deterministic ordering of your confusion matrix output.

3.1 Hyperparameter Selection for a Linear-Kernel SVM [18 pts]

- (a) (1 pt) **To begin, implement the function `cv_performance(clf, X, y, k=5, metric='accuracy')` as defined in the skeleton code.** Here you will make use of the `fit(X, y)`, `predict(X)`, and `decision_function(X)` methods in the `SVC` class. The function returns the mean k -fold CV performance for the performance metric passed into the function. The default metric is ‘accuracy’, however your function should work for all of the metrics listed above. It may be useful to have a helper function that calculates each performance metric. For instance: `performance(y_true, y_pred, metric='accuracy')`

You may have noticed that the proportion of the two classes (positive and non-positive) are equal in the training data. When dividing the data into folds for CV, you should try to keep the class proportions roughly the same across folds; in this case, the class proportions should be roughly equal across folds,

since the original training data has equal class proportions.

You must implement this function without using the `scikit-learn` implementation of CV. You will need to employ the following class for splitting the data:

`sklearn.model_selection.StratifiedKFold()`. Do not shuffle points when using this function (i.e., do not set `shuffle=True`). This is so the generated folds are consistent for the same dataset across runs of the entire program.

Solution: Implementations vary.

- (b) (1 pt) In your write-up, briefly describe why it might be beneficial to maintain class proportions across folds.

Solution: Stratified splits are important because the fundamental assumption of most ML algorithms is that the training set is a representative sample of the test set i.e., the training and test data are drawn from the same underlying distributions. If the ratio of positive to negative examples (the class balance) differs significantly between the training and test sets (across folds) this assumption will not hold.

- (c) (2 pt) Now implement the `select_param_linear(X, y, k=5, metric='accuracy', C_range=[1], penalty='l2')` function to choose a value of C for a linear SVM based on the training data and the specified metric. Note that scikit-learn uses a slightly different formulation of SVM from the one we introduced in lecture, namely:

$$\begin{aligned} & \underset{\theta, b, \xi_i}{\text{minimize}} \quad \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i \\ & \text{subject to} \quad y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b) \geq 1 - \xi_i \\ & \quad \quad \quad \xi_i \geq 0, \forall i = 1, 2, \dots, n \end{aligned}$$

Essentially, the C is inversely proportional to the λ we used in lecture. Your function should call your CV function (implemented earlier) passing in instances of `SVC(kernel='linear', C=c, class_weight='balanced')` with a range of values for C chosen in powers of 10 between 10^{-3} and 10^3 (i.e. $10^{-3}, 10^{-2}, \dots, 10^2, 10^3$). You may choose to implement and use the helper function `select_classifier` to instantiate the needed classifier. Make sure to test your implementation with the `test_cases.test_select_param_linear()` method.

Solution: Implementations vary.

- (d) (5 pt) Finally, using the training data from question 2 and the functions implemented here, find the best setting for C for each performance measure (if there is a tie, choose the smaller C value). Report your

findings in tabular format with three columns: names of the performance measures, along with the corresponding values of C and the mean cross-validation performance. The table should follow the format given below:

Performance Measures	C	Performance
Accuracy		
F1-Score		
AUROC		
Precision		
Sensitivity		
Specificity		

Solution:

Acceptable answer 1: CV performance

Performance Measures	C	Performance
Accuracy	0.01	0.8170
F1-Score	0.01	0.8160
AUROC	0.1	0.9047
Precision	0.1	0.8253
Sensitivity	0.001	0.8900
Specificity	0.1	0.8280

Your `select_param_linear` function returns the ‘best’ value of C given a range of values. Note: as we are working with a fairly large feature matrix, this may take several minutes (our project solution time is about 30 minutes for this question on our test computer).

Also, in your write-up, describe how the 5-fold CV performance varies with C . If you have to train a final model, which performance measure would you optimize for when choosing C ? Explain your choice. This value of C will be used in part e.

Solution:

- For sensitivity, the performance starts high and steadily decreases as C increases. This suggests that if we wanted to optimize this metric, we should be searching smaller values of C .
- For F1-score, accuracy, AUROC, precision and specificity, the performance increases as C increases before C reaches the optimal value. However, after C reaches the optimal value, the performance drops and then stays constant as C increases.
- Choosing the value of C to optimize for is an open-ended question. Each metric has its own strengths and weaknesses. We will accept any answer as long as your explanation is sensible.

- (e) (3 pt) Now, using the value of C that maximizes your chosen performance measure, create an SVM as in the previous question. Again, you may choose to use the helper function `select_classifier`. Train this SVM on the training data X_{train} , y_{train} . Report the performance of this SVM on the test data X_{test} , y_{test} for each metric below.

Performance Measures	Performance
Accuracy	
F1-Score	
AUROC	
Precision	
Sensitivity	
Specificity	

Solution:

Soln 1: Using $C = 0.01$ which maximizes Accuracy and F1-Score.

Performance Measures	Performance
Accuracy	0.8260
F1-Score	0.8263
AUROC	0.9047
Precision	0.8247
Sensitivity	0.8280
Specificity	0.8240

Soln 2: Using $C = 0.1$ which maximizes AUROC and Precision.

Performance Measures	Performance
Accuracy	0.8360
F1-Score	0.8347
AUROC	0.9109
Precision	0.8415
Sensitivity	0.8280
Specificity	0.8440

Soln 3: Using $C = 0.001$ which maximizes sensitivity performance.

Performance Measures	Performance
Accuracy	0.7240
F1-Score	0.7629
AUROC	0.8436
Precision	0.6687
Sensitivity	0.8880
Specificity	0.5600

- (f) (2 pt) Finish the implementation of the `plot_weight(X, y, penalty, C_range)` function. In this function, you need to find the L0-norm of $\bar{\theta}$, the parameter vector learned by the SVM, for each value of C in the given range. Finding out how to get the vector $\bar{\theta}$ from a `SVC` object may require you to dig into the documentation. The L0-norm is given as follows, for $\bar{\theta} \in \mathbb{R}^d$:

$$\|\bar{\theta}\|_0 = \sum_{i=1}^d \mathbb{I}\{\theta_i \neq 0\}$$

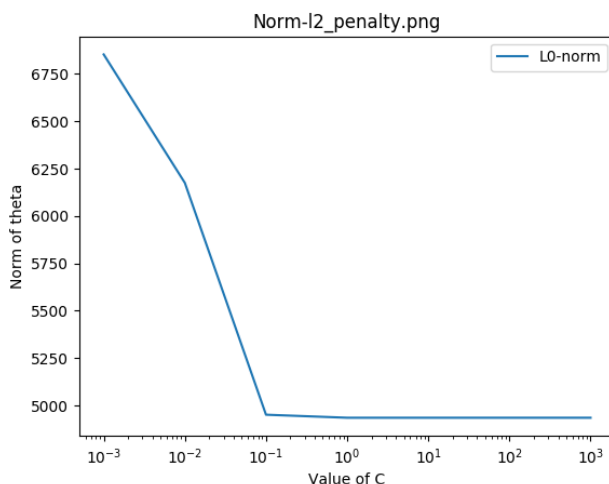
where $\mathbb{I}\{\theta_i \neq 0\}$ is 0 if θ_i is 0 and 1 otherwise.

Solution: Implementations vary.

Use the complete training data `X_train`, `Y_train`, i.e, don't use cross-validation for this part. Once you implement the function, the existing code will plot L0-norm $\|\bar{\theta}\|_0$ against C and save it to a file.

- (g) (1 pt) In your write-up, include the produced plot from the question above.

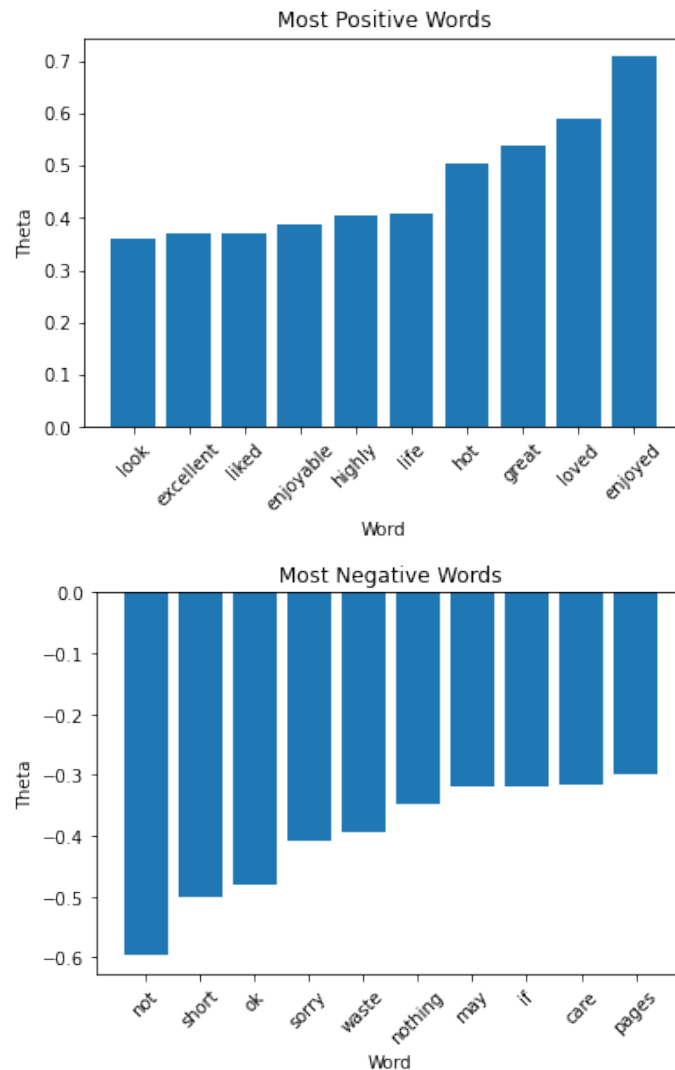
Solution:



- (h) (2 pt) Recall that each element of $\bar{\theta}$ is associated with a word. The more positive the value of the element, the more the presence of the associated word indicates that the review is positive, and similarly with negative coefficients. In this way, we can use these coefficients to find out what word-rating associations our SVM has learned.

Using $C = 0.1$ (for consistency with our results), train an SVM on `X_train`, `Y_train`. In your report, include a bar chart (coefficient vs each word) for both the ten most positive and ten most negative coefficients from the trained SVM. The words on the bar chart should be sorted by coefficient value in ascending order. As before, you may choose to use the helper function `select_classifier`.

Solution:



- (i) (1 pt) It is noteworthy that the word-rating association learned can be misleading. To illustrate this, come up with a review that is negative-sounding yet contains three of the ten words with the most positive coefficients (from your answer to the previous part).

Solution: Open-ended.

3.2 Hyperparameter Selection for a Quadratic-Kernel SVM [9 pts]

Similar to the hyperparameter selection for a linear-kernel SVM, you will perform hyperparameter selection for a quadratic-kernel SVM. Here we are assuming a kernel of the form $(\tilde{x} \cdot \tilde{x}' + r)^2$, where r is a

hyperparameter.

- (a) (5 pt) Implement the `select_param_quadratic(X, y, k=5, metric='accuracy', param_range=[])` function to choose a setting for C and r as in the previous part. Your function should call your CV function (implemented earlier) passing in instances of `SVC(kernel='poly', degree=2, C=c, coef0=r, class_weight='balanced')` with the same range of C that we use in 3.1(c). You should also use the same range for r . Make sure to test your implementation with the `test_cases.test_select_param_quadratic()` method.

Again, you may choose to use the helper function `select_classifier`. The function argument `param_range` should be a matrix with two columns with first column as values of C and second column as values of r . You will need to try out the range of parameters via two methods:

- i) Grid Search: In this case, we look at all the specified values of C in a given set and choose the best setting after tuning. For this question, the values should be considered in powers of 10 for both C (between 10^{-3} and 10^3) and r (between 10^{-3} and 10^3) [A total of 49 pairs]. This code will take a substantial time to run (our project solution runs in about 35 minutes on our test computer).
- ii) Random Search: Instead of assigning fixed values for the parameter C and r , we can also sample random values from a particular distribution. For this case, we will be sampling from a log-uniform distribution, i.e., the log of random variables follows a uniform distribution:

$$P[a \leq \log C \leq b] = k(b - a)$$

for some constant k so the distribution is valid. In other words, we sample a uniform distribution with the same range as above to yield exponents x_i , and corresponding values of $C = 10^{x_i}$.

In your case, the values should range from the powers of 10 which you used in part (i). Choose 25 pairs of such sampled pairs of (C, r) . Again, this code may take time to run (our solution runs in about 20 minutes on our test computer).

Solution: Implementations vary.

- (b) (4 pt) Find the best C and r value for AUROC using both tuning schemes mentioned above and the training data from question 2. Report your findings in tabular format. The table should have four columns: Tuning Scheme, C , r and AUROC. Again, in the case of ties, report the lower C and the lower r values that perform the best (prioritizing a lower C). Your table should look be similar to the following:

Tuning Scheme	C	r	AUROC
Grid Search			
Random Search			

How does the 5-fold CV performance vary with C and r ? Also, explain the pros and cons of grid search and random search.

Solution:

Reported CV Performance

Tuning Scheme	C	r	AUROC
Grid Search	10.0	10.0	0.9131
Random Search	6.0303	23.8505	0.9134

- How does the 5-fold CV performance vary with C and r ?
 - **Random Search** - the performances fluctuates around 0.80 to 0.92 for different pairs of C and r . Even though the result for random search varies and is not deterministic, we do expect your performances to be around 0.80 to 0.95.
 - **Grid Search** - if we fix C , the performance increases as r increases in general. If we fix r , the performance generally increases also as C increases.
- Is the use of random search better than grid search?
 - By looking at the performance, random search performed as well as grid search. As random search is not deterministic, we may not get the same set of hyperparameters if we run random search multiple times. Consequently, random search may in some cases perform better or worse than grid search. The best performance of random search is very close to that of grid search most of the time.
Given that we search through fewer pairs of parameters when using random search compared to grid search, random search tends to be more computationally efficient than grid search. However, as mentioned, the deterministic nature of grid search can be useful for comparing algorithms and models.

3.3 Learning Non-linear Classifiers with a Linear-Kernel SVM [4 pts]

Here, we will explore the use of an explicit feature mapping in place of a kernel. (Note: you do not need to write any code for question 3.3)

- (a) (2 pt) Describe a feature mapping, $\phi(\bar{x})$, that maps your data to the same feature space as the one implied by the quadratic kernel from the question above.

Solution:

$$\phi(x) = \langle \underbrace{r, \sqrt{2r}x_1, \dots, \sqrt{2r}x_n}_{\text{linear terms}}, \underbrace{\sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \dots, \sqrt{2}x_1x_n, \sqrt{2}x_2x_3, \dots, \sqrt{2}x_{n-1}x_n}_{\text{cross terms } x_i x_j, i \neq j}, \underbrace{x_1^2, \dots, x_n^2}_{\text{square terms}} \rangle$$

- (b) (2 pt) Instead of using a quadratic-kernel SVM, we could simply map the data to this higher dimensional

space via this mapping and then learn a linear classifier in this higher-dimensional space. What are the tradeoffs (pros and cons) of using an explicit feature mapping over a kernel? Explain.

Solution:

Pros: We retain interpretability and can access feature weights for each feature if we use explicit feature mapping.

Cons: For d features in original data, we have $(d^2 + d + 1)$ features in the transformed data, which is very computationally expensive

3.4 Linear-Kernel SVM with L1 Penalty and Squared Hinge Loss [4 pts]

In this part of the project, you will explore the use of a different penalty (i.e., regularization term) and a different loss function. In particular, we will use the L1 penalty and squared hinge loss which corresponds to the following optimization problem.

$$\underset{\bar{\theta}, b}{\text{minimize}} \|\bar{\theta}\|_1 + C \sum_{i=1}^n \text{loss}(y^{(i)}(\bar{\theta} \cdot \bar{x}^{(i)} + b))$$

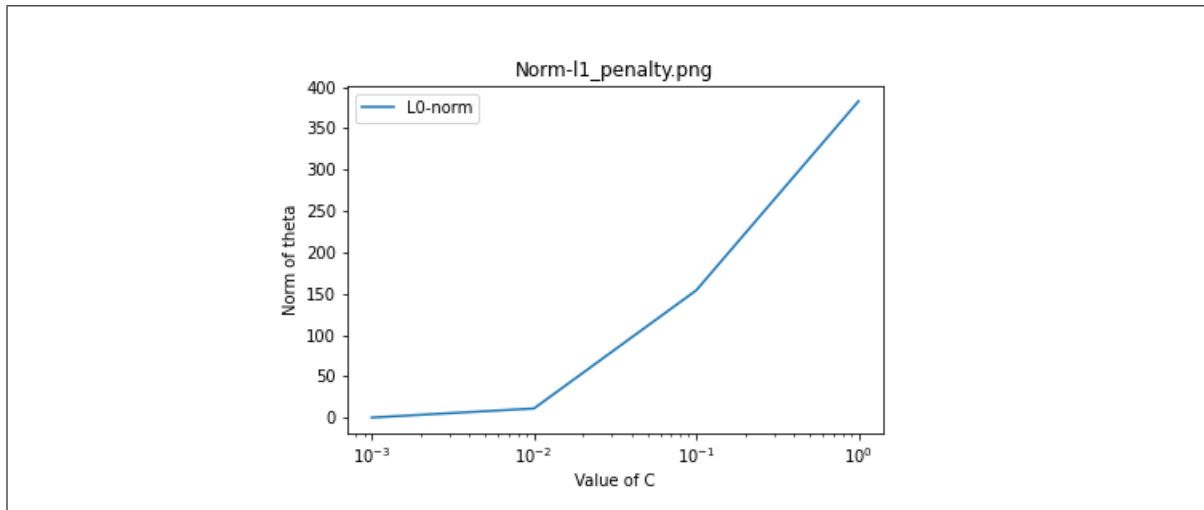
where $\text{loss}(z) = \max\{0, (1 - z)\}^2$. We will make use of the `LinearSVC()` class, which uses the squared hinge loss and allows us to specify the penalty. We will consider only a linear-kernel SVM and the original (untransformed) features. When calling `LinearSVC` please use the following settings: `LinearSVC(penalty='l1', dual=False, C=c, class_weight='balanced')`. As always, you may implement and use the helper function `select_classifier` to instantiate your SVM classifier.

- (a) (1 pt) Using the training data from question 2 and 5-fold CV, find the best setting for C that maximizes the AUROC given that $C \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0\}$. When we say "grid search" here, we mean searching a one-dimensional grid where the only hyperparameter that is changing is C . In the case of ties, report the lower C value. Report your findings.

Solution: The value of C that maximizes CV performance is $C = 1.0$ with AUROC performance using cross-validation of 0.8898 and AUROC performance on the entire test set of 0.8870.

- (b) (1 pt) Similar to 3.1(f), plot the L0-norm of the learned parameter $\bar{\theta}$ against C using complete training data and no cross-validation. You should be able to re-use the function `plot_weight` with different input parameters without writing additional code. Include the plot in your write-up.

Solution:



- (c) (1 pt) Compare and contrast the graphs that you generated for both the L1 and L2 penalty. Why does this difference occur? (Hint: Think about the gradients)

Solution:

- First, Looking at scale, we can see that the L1 penalty produces much sparser theta vectors for all settings of C. Since the gradient of the L1 norm is constant with respect to each element, there is just as much value in reducing a small value as there is a large value. This ends up pushing the coefficients of the unimportant features all the way to zero thus resulting in a sparser graph. Likewise, for the L2 norm, the gradient is linear with respect to each element meaning that it more value is gained from reducing large elements and less is gained from reducing small elements. This incentivizes balancing the weights of all features in order to not have expensive large values.

- (d) (1 pt) Note that using the Squared Hinge Loss (as opposed to the Hinge Loss) changes the objective function as shown above. What effect do you expect this will have on the optimal solution?

Solution:

- Squared hinge loss is lenient on the cases where a data point is classified correctly but it is within the margin compared to hinge loss. On the other hand, squared hinge loss punishes misclassifications more harshly. So, we see a comparatively wider margin and more support vectors when using squared hinge loss.

3.5 Perceptron Classifier [5 Points]

Here we will compare the SVM classifier with one of the earliest classification algorithms– the Perceptron.

- (a) (2 pt) We've provided you the `train_perceptron()` function, which returns $\bar{\theta}$ and b values corresponding to the decision boundary found by the Perceptron algorithm. Write the code necessary to evaluate this decision boundary's accuracy. (Note: ensure that your data is from `get_split_binary(class_size=500)` when calling `train_perceptron()`).

Solution:

- Answers may vary

- (b) (3 pt) Compare your Perceptron's accuracy on the book review data set to your SVM from part 3.1.e. Make sure that you are using the same dataset on both models. Report the accuracy of both models. Why do you think Perceptron performs better/worse than the SVM?

Solution:

- Perceptron accuracy of 82.2%
- Linear SVC accuracy of 83.6%, 82.6%, 72.4%
- Since dataset is linearly separable, linear SVC will converge to the maximum margin separator while the Perceptron will only find a valid decision boundary. This results in the SVC performing better at test time.

4 Asymmetric Cost Functions and Class Imbalance [20 pts]

The training data we have provided you with so far is *balanced*: the data contain an equal number of positive and negative ratings. But this is not always the case. In many situations, you are given imbalanced data, where one class may be represented more than the others.

In this section, you will investigate the objective function of the SVM, and how we can modify it to fit situations with class imbalance. Recall that the objective function for an SVM in scikit-learn is as follows:

$$\begin{aligned} & \underset{\bar{\theta}, b, \xi_i}{\text{minimize}} \quad \frac{\|\bar{\theta}\|^2}{2} + C \sum_{i=1}^n \xi_i \\ & \text{subject to} \quad y^{(i)} (\bar{\theta} \cdot \phi(\bar{x}^{(i)}) + b) \geq 1 - \xi_i \\ & \quad \quad \quad \xi_i \geq 0, \forall i = 1, 2, 3, \dots, n \end{aligned}$$

We can modify it in the following way:

$$\begin{aligned} & \underset{\bar{\theta}, b, \xi_i}{\text{minimize}} \quad \frac{\|\bar{\theta}\|^2}{2} + W_p * C \sum_{i|y^{(i)}=1} \xi_i + W_n * C \sum_{i|y^{(i)}=-1} \xi_i \\ & \text{subject to} \quad y^{(i)} (\bar{\theta} \cdot \phi(\bar{x}^{(i)}) + b) \geq 1 - \xi_i \\ & \quad \quad \quad \xi_i \geq 0, \forall i = 1, 2, 3, \dots, n \end{aligned}$$

where $\sum_{i|y^{(i)}=1}$ is a sum over all indices i where the corresponding point is positive $y^{(i)} = 1$. Similarly, $\sum_{i|y^{(i)}=-1}$ is a sum over all indices i where the corresponding point is negative $y^{(i)} = -1$.

4.1 Arbitrary class weights [6 pts]

W_p and W_n are called “class weights” and are built-in parameters in scikit-learn.

- (a) (1 pt) Describe how this modification will change the solution. If W_n is much greater than W_p , what does this mean in terms of classifying positive and negative points? Refer to the weighted SVM formulation for a brief justification of your reasoning.

Solution: This modification changes the objective function by effectively changing the importance of different points in the minimization process. That is, points in a class with a higher weight (in this case, points with negative labels) are more likely to correspond to a lower slack term ξ_i as the slack term ξ_i is more highly weighted in the minimization. Therefore, negative-labeled points are more likely to be correctly classified.

- (b) (2 pt) Create a linear-kernel SVM with hinge loss and L2-penalty with $C = 0.1$. This time, when calling SVC, set `class_weight = {-1: 1, 1: 10}`, or implement and use your `select_classifier helper function`. This corresponds to setting $W_n = 1$ and $W_p = 10$. Train this SVM on the training data `X_train, y_train`.

- (c) (2 pt) Report the performance of the modified SVM on the test data X_{test} , y_{test} for each metric below.

Note: You should be using SVC, not LinearSVC.

Performance Measures	Performance
Accuracy	
F1-Score	
AUROC	
Precision	
Sensitivity	
Specificity	

Solution:

Performance Measures	Performance
Accuracy	0.8240
F1-Score	0.8281
AUROC	0.9042
Precision	0.8092
Sensitivity	0.8480
Specificity	0.8000

- (d) (1 pt) Also, answer the following: Compared to your work in question 3.1(e), which performance measures were affected the most by the new class weights? Why do you suspect this is the case?

Note: We set $C = 0.1$ to ensure that interesting trends can be found regardless of your work in question 3. This may mean that your value for C differs in 4 and 3.1(e). In a real machine learning setting, you'd have to be more careful about how you compare models.

Solution: Answers will depend based on what metric was optimized for in Question 3.1(e). In general, we would expect sensitivity to increase and specificity to decrease due to the classifier penalizing misclassified positive points more heavily than negative ones.

Note: This was not required, but to see this trend more clearly, consider setting $C = 0.01$. Think about why a lower C value might accentuate this trend. One possible explanation is that by prioritizing a simpler solution (which this setting of C encourages), the optimal solution now sacrifices the lower weighted negative points in favor of a less complex decision boundary.

4.2 Imbalanced data [5 pts]

You just saw the effect of arbitrarily setting the class weights when our training set is already balanced. Let's return to the class weights you are used to: $W_n = W_p$. We turn our attention to class imbalance. Using the

functions you wrote in part 2, we have provided you with a second feature matrix and vector of binary labels `IMB_features`, `IMB_labels`. This class-imbalanced data set has 800 positive points and 200 negative points. It also comes with a corresponding test feature matrix and label vector pair `IMB_test_features`, `IMB_test_labels`, which have the same class imbalances.

- (a) (2 pt) Create a linear-kernel SVM with hinge loss, L2-penalty and as before, $C = 0.1$. Set `class_weight={-1: 1, 1: 1}`, which returns the SVM to the formulation you have seen in class. Now train this SVM on the class-imbalanced data `IMB_features`, `IMB_labels` provided.
- (b) (2 pt) Use this classifier to predict the provided test data `IMB_test_features`, `IMB_test_labels` and report the accuracy, specificity, sensitivity, precision, AUROC, and F1-Score of your predictions:

Class Weights	Performance Measures	Performance
$W_n = 1, W_p = 1$	Accuracy	
$W_n = 1, W_p = 1$	F1-Score	
$W_n = 1, W_p = 1$	Auroc	
$W_n = 1, W_p = 1$	Precision	
$W_n = 1, W_p = 1$	Sensitivity	
$W_n = 1, W_p = 1$	Specificity	

Solution:

Class Weights	Performance Measures	Performance
$W_n = 1, W_p = 1$	Accuracy	0.8400
$W_n = 1, W_p = 1$	F1-Score	0.9029
$W_n = 1, W_p = 1$	Auroc	0.8802
$W_n = 1, W_p = 1$	Precision	0.8774
$W_n = 1, W_p = 1$	Sensitivity	0.9300
$W_n = 1, W_p = 1$	Specificity	0.4800

- (c) (1 pt) How has training on an imbalanced data set affected performance?

Solution: Training on an imbalanced data set has affected sensitivity, which increased, along with specificity, which decreased significantly. Changing the data set to consist of far more positive samples than negative samples results in the classifier trained on this data tending to classify many more positive samples correctly, as shown by the sensitivity and specificity metrics. The sensitivity (true positive rate) is 0.9300 which indicates almost all positive samples were classified correctly, while the specificity (true negative rate) is 0.4800 which indicates about half of the negative samples were correctly classified as negative.

4.3 Choosing appropriate class weights [6 pts]

- (a) (4 pt) Now we will return to setting the class weights given the situation we explored in Part 4.2. Using what you have done in the preceding parts, **find an appropriate setting for the class weights** that mitigates the situation in part 4.2 and improves the classifier trained on the imbalanced data set. That is, find class weights that give a good mean cross-validation performance (Think: which performance metric(s) are informative in this situation, and which metric(s) are less meaningful? Make sure the metric you use for cross-validation is a good choice given the imbalanced class weights). **Report here your strategy for choosing an appropriate performance metric and weight parameters.** This question requires you to choose hyperparameters based on cross-validation; you should not be using the test data to choose hyperparameters.

Solution: **F1-score** is a good metric to use when we want to balance precision and recall in the presence of class imbalance. This means that it considers both the number of positive labels that we classified correctly, and the negative labels that we classified incorrectly. Note that optimizing for this metric might heighten the higher positive bias because it is not adversely affected by false positives (due to their relatively low number). Weights found using Cross Validation with respect to F1-score across a reasonably large range of weights ($W_n=2$ to 7, $W_p=2$ to 7) are: **$W_n=2$, $W_p=5$** (CV Score: 0.9047815785726068).

Alternative solution: We could also use specificity in order to make sure that we counteract the bias towards positive labels. Weights found using Cross Validation with respect to **specificity** across a range of weights $W_n=1$ to 7, $W_p=2$ to 8 are: **$W_n=7$, $W_p=2$** (CV Score: 0.555)

Alternative solution: We could also use AUROC since it measures the trade off between the True Positive and False Positive Rate which is what we're looking for here. Weights found using Cross Validation with respect to **AUROC** across a range of weights $W_n=1$ to 7, $W_p=2$ to 8 are: **$W_n=2$, $W_p=5$** (CV Score: 0.85596875)

- (b) (2 pt) Use your custom classifier to predict the provided test data `IMB_test_features`, `IMB_test_labels` again, and **report the accuracy, specificity, sensitivity, precision, AUROC, and F1-Score of your predictions:**

Class Weights	Performance Measures	Performance
$W_n = ?, W_p = ?$	Accuracy	
$W_n = ?, W_p = ?$	F1-Score	
$W_n = ?, W_p = ?$	Auroc	
$W_n = ?, W_p = ?$	Precision	
$W_n = ?, W_p = ?$	Sensitivity	
$W_n = ?, W_p = ?$	Specificity	

Solution:

Class Weights	Performance Measures	Performance
$W_n = 3, W_p = 5$	Accuracy	0.8480
$W_n = 3, W_p = 5$	F1-Score	0.9064
$W_n = 3, W_p = 5$	Auroc	0.8639
$W_n = 3, W_p = 5$	Precision	0.8932
$W_n = 3, W_p = 3$	Sensitivity	0.9200
$W_n = 3, W_p = 5$	Specificity	0.5600

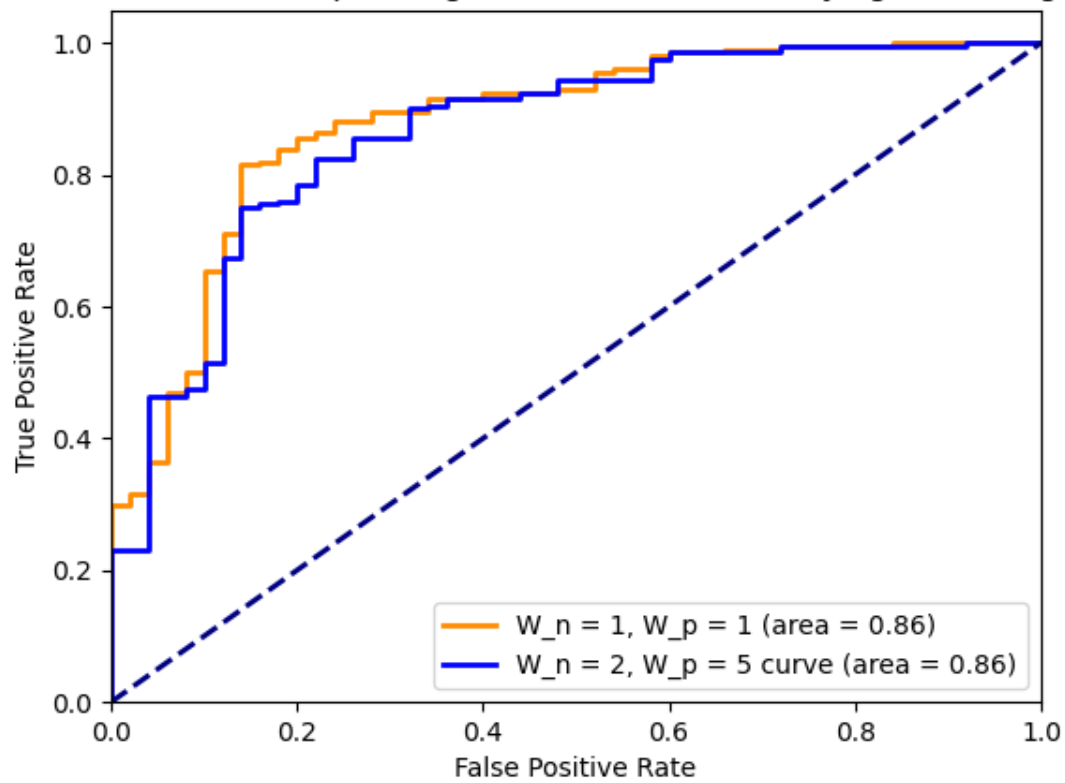
Answers may vary based on W_n and W_p weights determined by student in 4.3(a).

4.4 The ROC curve [3 pts]

Given the above results, we are interested in investigating the AUROC metric more. First, provide a plot of the ROC curve with labeled axes for both $W_n = 1, W_p = 1$ and your custom setting of W_n, W_p from above. Put both curves on the same set of axes. Make sure to label the plot in a way that indicates which curve is which.

Solution:

4.4(a) - Receiver operating characteristic with varying class weights



5 Challenge [20 pts]

Now, a challenge: in the previous problems, we had transformed the data into a binary dataset by combining multiple labels to generate two labels.

For this challenge, you will consider the original multiclass labels of the reviews. We have already prepared a held-out test set `heldout_features` for this challenge, and multiclass training data `multiclass_features`, `multiclass_labels`. This training data has 2250 reviews, 750 of each class. **You must work only with the provided data; acquiring new data to train your model is not permitted.** (Notice, if you look into the dataset, there are additional information that could be leveraged. Your goal is to train a multiclass classifier using the `SVC` or `LinearSVC` classes to predict the true ratings of the held-out test set, i.e., you will train your model on `multiclass_features` and test on `heldout_features`. If you wish to take advantage of the other features in the dataset, you will have to modify either the `get_multiclass_training_data` function in `helper.py` or the `generate_feature_matrix` function in `project1.py`.

Note that the class balance of this training set matches the class balance of the heldout set. Also note that, given the size of the data and the feature matrix, training may take several minutes.

In order to attempt this challenge, we encourage you to apply what you have learned about hyperparameter selection and consider the following extensions:

1. **Try different feature engineering methods.** The bag-of-words models we have used so far are simplistic. There are other methods to extract different features from the raw data, such as:
 - (a) Using a different method for extracting words from the ratings
 - (b) Using only a subset of the raw features
 - (c) Using the number of times a word occurs in a ratings as a feature (rather than binary 0, 1 features indicating presence)
 - (d) Include phrases from ratings in addition to words.
 - (e) Scaling or normalizing the data
 - (f) Alternative feature representations
2. **Read about one-vs-one and one-vs-all.** These are the two standard methods to implement multiclass classifier using binary classifiers. You should understand the differences between them and implement at least one.

You will have to save the output of your classifier into a `csv` file using the helper function `generate_challenge_labels(y, username)` we have provided. The base name of the output file must be your username followed by the extension `csv`. For example, the output filename for a user with username `foo` would be `foo.csv`. This file will be submitted according to the instructions at the end of the file. You may use the file `test_output.py` to ensure that your output has the correct format. To run this file, simply run `python test_output.py -i username.csv`, replacing the file `username.csv` with your generated output file.

We will evaluate your performance in this challenge based on three components:

1. Write-Up and Code [8 pts]: We will evaluate how much effort you have applied to attempt this challenge based on your write-up and code. **Ensure that both are present.** Within your write-up, you must provide discussions of the choices you made when designing the following components of your classifier:

- Feature engineering
- Hyperparameter selection
- Algorithm selection (e.g., quadratic vs. linear kernel)
- Multiclass method (e.g., one-vs-rest vs. one-vs-all)
- Any techniques you used that go above and beyond current course material

Solution: Answers will vary.

2. Test Scores [8 pts]: We will evaluate your classifier based on accuracy. Consider the following confusion matrix:

	-1	0	1
-1	x_1		
0		x_2	
1	y_1		x_3

where each column corresponds to the actual class and each row corresponds to the predicted class. For instance, y_1 in the matrix above is the number of reviews with true rating -1 (poor), but are classified as a review with rating 1 (good) by your model. The accuracy for a multiclass classification problem is defined as follows:

$$\text{accuracy} = \frac{x_1 + x_2 + x_3}{n}$$

where n is the number of samples.

Solution: Answers will vary.

3. (4 pt) Now that you've implemented a variety of powerful ML and NLP techniques, it's important to remember that with great power comes great responsibility. When ML solutions are deployed in the real world, there are very real effects on people. Whether through irresponsible usage or preexisting biases in your dataset, these effects can be detrimental to people's lives. In your write-up, please answer the following questions:

- (2 pt) Imagine that you decide to now use this model to determine whether a news article was good or not by taking the average sentiment of all the comments made on the article. What is a potential societal bias that might be present in the learned model?
- (2 pt) Explain how you could potentially fix the issue that you identified.

Solution: [Answers will vary.](#)

NOTE: You may know that a model typically performs better with a larger number of data, and may have consequently concluded that a good strategy involves scraping Amazon for more data. **Please do NOT do so**, as it violates Amazon's terms of use and is therefore not permitted.

Include your code (copy and pasted) in your report. Please try to format lines of code so they are visible within the pages. **Make sure to match your code with the correct problems on Gradescope.** Failure to do so, may result in a 0 for the corresponding problems.

Upload your file `username.csv` containing the label predictions for the held-out data to the canvas assignment named Project 1 Challenge Submissions.

Appendix A: Approximate Run-times for Programming Problems

- **Problem 3.1 c:** around 30 minutes
- **Problem 3.2 a i:** around 35 minutes
- **Problem 3.2 a ii:** around 20 minutes
- **Problem 3.4 a:** around 10 seconds
- **Problem 3.4 b:** around 3 seconds
- **Problem 4.3:** around 15 minutes

N.B. these are approximate times, not exact. Different computers will result in different run-times, so do not panic if yours is a little different. Algorithmic optimization can also improve run-time noticeably in certain cases. However, if it is taking more than twice as long, something might be wrong. To help save time, we recommend testing your implementations first with the provided test cases when possible.

Appendix B: Topics and Concepts

The relevant topics for each section are as follows:

- **Problem 3.1 a, b, e, f, Problem 3.4 c, d**
 - Support Vector Machines; Primal Formulation; Geometric Margin; Loss Functions and Regularization
- **Problem 3.2 a, Problem 3.3 a, Problem 4.1 a**

- Dual Formulation; Kernels
- **Problem 3.1** c, d, **Problem 3.2** b, **Problem 3.3** b, **Problem 3.4** a, b, **Problem 4.1** b, c, **Problem 4.2** a, b, **Problem 4.3** a, b, **Problem 4.4** a, b
- Performance Measures

Appendix C: Further Reading

Below are some topics (in no particular order) you may find useful to research for the challenge portion of this project. This is not an exhaustive list, nor do we know for certain that they will improve your classifier performance, but they are avenues for you to explore.

1. Term Frequency - Inverse Document Frequency (TF-IDF)
2. Topic Modeling (Latent Dirichlet Allocation)
3. Data Augmentation¹
4. Stemming and Lemmatization
5. Part-of-speech tagging and position²
6. N-grams

¹<https://www.aclweb.org/anthology/D19-1670.pdf>

²<https://www.aclweb.org/anthology/W02-1011.pdf>