

## Homework Set #3 : Kernel and SVM

Due 23rd May 2018, Wednesday, before 11:59 pm

---

### Submission

- Submit your solutions electronically on the course Gradescope site as PDF files.
- If you plan to typeset your solutions, please use the LaTeX solution template. If you must submit scanned handwritten solutions, please use a black pen on blank white paper and a high-quality scanner app.

## Problem 1 (KERNEL PROPERTIES [10 PTS])

We know that the following is a valid Kernel:  $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^d$  where  $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ ,  $d \in \mathbb{Z}^+$ . Prove that the following are valid kernels

- (a) **(5 pts)**  $K_{new}(\mathbf{x}, \mathbf{z}) = \left(\sum_{i=1}^n \sqrt{x_i} \sqrt{z_i}\right)^d$
- (b) **(5 pts)**  $K_{new}(\mathbf{x}, \mathbf{z}) = \alpha K_1(\mathbf{x}, \mathbf{z}) + \beta K_2(\mathbf{x}, \mathbf{z})$  where  $K_1$  and  $K_2$  are both valid kernels and  $\alpha, \beta \in \mathbb{R}^+$

*Hint:* If required you can prove and use the more general claim.

If  $K(\mathbf{x}, \mathbf{z})$  is a valid kernel and  $g(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  then  $K_{new}(\mathbf{x}, \mathbf{z}) = K(g(\mathbf{x}), g(\mathbf{z}))$  is also a valid kernel.

## Problem 2 (KERNELIZED SUPPORT VECTOR MACHINES [20 PTS])

Suppose you are given 6 one-dimensional points: 3 with negative labels  $x_1 = -1, x_2 = 0, x_3 = 1$  and 3 with positive labels  $x_4 = -3, x_5 = -2, x_6 = 3$ . You have seen in HW2 that if a feature map such as  $\phi(u) = (u, u^2)$  is used to transform the points then a hyperplane in the new  $\mathbb{R}^2$  feature space induced by  $\phi$  can perfectly separate the points. The kernel corresponding to the feature map is  $K(x, z) = xz(1 + xz)$

- (a) **(4 pts)** Construct a maximum margin separating hyperplane in  $\mathbb{R}^2$  which can be parameterized by its normal equation, i.e.  $w_1 y_1 + w_2 y_2 + w_0 = 0$  for appropriate choices of  $\mathbf{w} \in \mathbb{R}^3$ . Here  $(y_1, y_2) = \phi(x)$  is the result of applying the feature map  $\phi(x)$  to the original feature  $x$ . Also explicitly compute the margin for the hyperplane. *Hint:* Look at the points in the induced feature space.
- (b) **(4 pts)** Apply  $\phi(x)$  to the data and plot the points in the new  $\mathbb{R}^2$  feature space. On the plot of the transformed points, plot the separating hyperplane you found in a) and the margin and circle the support vectors (SV).
- (c) **(3 pts)** Draw the decision boundary of the separating hyperplane you found in a) in the original  $\mathbb{R}^1$  feature space.
- (d) **(6 pts)** Find the  $\alpha_i$  and  $b$  in

$$h(\mathbf{x}) = \text{sign} \left( \sum_{n=1}^{\#SV} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

for the kernel  $K$  and the support vectors chosen earlier. Do this by solving the dual form of the quadratic program.

- (e) **(3 pts)** If we add another negatively labelled training point at  $\mathbf{x} = 0.5$ , will the hyperplane or the margin change? Explain.

### Problem 3 (SUPPORT VECTOR MACHINES [18 PTS])

Suppose we are looking for a maximum-margin linear classifier *through the origin*, i.e.  $b = 0$  (also hard margin, i.e., no slack variables). In other words, we minimize  $\frac{1}{2}||\mathbf{w}||^2$  subject to  $y_n \mathbf{w}^T \mathbf{x}_n \geq 1, n = 1, \dots, N$ .

- (a) **(6 pts)** Given a single training vector  $\mathbf{x} = (a, e)^T$  with label  $y = -1$ , what is the  $\mathbf{w}^*$  that satisfies the above constrained minimization?
- (b) **(6 pts)** Suppose we have two training examples,  $\mathbf{x}_1 = (1, 1)^T$  and  $\mathbf{x}_2 = (1, 0)^T$  with labels  $y_1 = 1$  and  $y_2 = -1$ . What is  $\mathbf{w}^*$  in this case?
- (c) **(6 pts)** Suppose we now allow the offset parameter  $b$  to be non-zero. How would the classifier and the margin change in the previous question? What are  $(\mathbf{w}^*, b^*)$ ? Compare your solutions with and without offset.

### Problem 4 (TWITTER ANALYSIS USING SVM [52 PTS])

In this project, you will be working with Twitter data. Specifically, we have supplied you with a number of tweets that are reviews/reactions to movies<sup>1</sup>,

e.g., “@nickjfrost just saw *The Boat That Rocked/Pirate Radio* and I thought it was brilliant! You and the rest of the cast were fantastic! < 3”.

You will learn to automatically classify such tweets as either positive or negative reviews. To do this, you will employ Support Vector Machines (SVMs), a popular choice for a large number of classification problems.

Download the code and data sets from the course website. It contains the following data files:

- `tweets.txt` contains 630 tweets about movies. Each line in the file contains exactly one tweet, so there are 630 lines in total.
- `labels.txt` contains the corresponding labels. If a tweet praises or recommends a movie, it is classified as a positive review and labeled +1; otherwise it is classified as a negative review and labeled -1. These labels are ordered, i.e. the label for the  $i^{\text{th}}$  tweet in `tweets.txt` corresponds to the  $i^{\text{th}}$  number in `labels.txt`.

Skim through the tweets to get a sense of the data.

The python file `twitter.py` contains skeleton code for the project. Skim through the code to understand its structure.

---

<sup>1</sup>Please note that these data were selected at random and thus the content of these tweets do not reflect the views of the course staff. :-)

## 0.1 Feature Extraction [4 pts]

We will use a bag-of-words model to convert each tweet into a feature vector. A bag-of-words model treats a text file as a collection of words, disregarding word order. The first step in building a bag-of-words model involves building a “dictionary”. A dictionary contains all of the unique words in the text file. For this project, we will be including punctuations in the dictionary too. For example, a text file containing “*John likes movies. Mary likes movies2!!*” will have a dictionary `{'John':0, 'Mary':1, 'likes':2, 'movies':3, 'movies2':4, '.':5, '!':6}`. Note that the (key,value) pairs are (word, index), where the index keeps track of the number of unique words (size of the dictionary).

Given a dictionary containing  $d$  unique words, we can transform the  $n$  variable-length tweets into  $n$  feature vectors of length  $d$  by setting the  $i^{\text{th}}$  element of the  $j^{\text{th}}$  feature vector to 1 if the  $i^{\text{th}}$  dictionary word is in the  $j^{\text{th}}$  tweet, and 0 otherwise.

- (a) (2 pts) We have implemented `extract_words(...)` that processes an input string to return a list of unique words. This method takes a simplistic approach to the problem, treating any string of characters (that does not include a space) as a “word” and also extracting and including all unique punctuations.

Implement `extract_dictionary(...)` that uses `extract_words(...)` to read all unique words contained in a file into a dictionary (as in the example above). Process the tweets in the order they appear in the file to create this dictionary of  $d$  unique words/punctuations.

- (b) (2 pts) Next, implement `extract_feature_vectors(...)` that produces the bag-of-words representation of a file based on the extracted dictionary. That is, for each tweet  $i$ , construct a feature vector of length  $d$ , where the  $j^{\text{th}}$  entry in the feature vector is 1 if the  $j^{\text{th}}$  word in the dictionary is present in tweet  $i$ , or 0 otherwise. For  $n$  tweets, save the feature vectors in a feature matrix, where the rows correspond to tweets (examples) and the columns correspond to words (features). Maintain the order of the tweets as they appear in the file.

- (c) (0 pts) In `main(...)`, we have provided code to read the tweets and labels into a  $(630, d)$  feature matrix and  $(630,)$  label array. Split the feature matrix and corresponding labels into your training and test sets. **The first 560 tweets will be used for training and the last 70 tweets will be used for testing.** `**All subsequent operations will be performed on these data.**`

## 0.2 Hyperparameter Selection for a Linear-Kernel SVM [20 pts]

Next, we will learn a classifier to separate the training data into positive and negative tweets. For the classifier, we will use SVMs with two different kernels: linear and radial basis function (RBF). We will use the `sklearn.svm.SVC` class<sup>2</sup> and explicitly set only three of the initialization parameters: `kernel`, `gamma`, and `C`. As usual, we will use `SVC.fit(X,y)` to train our SVM, but in

<sup>2</sup>Note that when using SVMs with the linear kernel, it is recommended to use `sklearn.svm.LinearSVC` instead of `sklearn.svm.SVC` because the backbone of `sklearn.svm.LinearSVC` is the LIBLINEAR library, which is specifically designed for the linear kernel. For the sake of the simplicity, in this problem set we use `sklearn.svm.SVC`.

lieu of using `SVC.predict(X)` to make predictions, we will use `SVC.decision_function(X)`, which returns the (signed) distance of the samples to the separating hyperplane.

SVMs have hyperparameters that must be set by the user. For both linear and RBF-kernel SVMs, we will select the hyperparameters using 5-fold cross-validation (CV). Using 5-fold CV, we will select the hyperparameters that lead to the ‘best’ mean performance across all 5 folds.

- (a) **(6 pts)** The result of a 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**.<sup>3</sup>

Implement `performance(...)`. All measures, except sensitivity and specificity, are implemented in `sklearn.metrics` library. You can use `sklearn.metrics.confusion_matrix(...)` to calculate the other two.

- (b) **(4 pts)** Next, implement `cv_performance(...)` to return the mean  $k$ -fold CV performance for the performance metric passed into the function. Here, you will make use of `SVC.fit(X,y)` and `SVC.decision_function(X)`, as well as your `performance(...)` function.

You may have noticed that the proportion of the two classes (positive and negative) are not 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 your write-up, briefly describe why it might be beneficial to maintain class proportions across folds. Then, in `main(...)`, use `sklearn.cross_validation.StratifiedKFold(...)` to split the data for 5-fold CV, making sure to stratify using only the training labels.

- (c) **(4 pts)** Now, implement `select_param_linear(...)` to choose a setting for  $C$  for a linear SVM based on the training data and the specified metric. Your function should call `cv_performance(...)`, passing in instances of `SVC(kernel='linear', C=c)` with different values for  $C$ , e.g.,  $C = 10^{-3}, 10^{-2}, \dots, 10^2$ .
- (d) **(6 pts)** Finally, using the training data from Section 0.1 and the functions implemented here, find the best setting for  $C$  for each performance measure mentioned above. Report your findings in tabular format (up to the fourth decimal place):

$C$	accuracy	F1-score	AUROC	precision	sensitivity	specificity
$10^{-3}$						
$10^{-2}$						
$10^{-1}$						
$10^0$						
$10^1$						
$10^2$						
best $C$						

Your `select_param_linear(...)` function returns the ‘best’  $C$  given a range of values. How does the 5-fold CV performance vary with  $C$  and the performance metric?

<sup>3</sup>Read menu [link](#) to understand the meaning of these evaluation metrics.

### 0.3 Hyperparameter Selection for an RBF-kernel SVM [16 pts]

Similar to the hyperparameter selection for a linear-kernel SVM, you will perform hyperparameter selection for an RBF-kernel SVM.

- (a) **(4 pts)** Describe the role of the additional hyperparameter  $\gamma$  for an RBF-kernel SVM. How does  $\gamma$  affect generalization error?
- (b) **(6 pts)** Implement `select_param_rbf(...)` to choose a setting for  $C$  and  $\gamma$  via a grid search. Your function should call `cv_performance(...)`, passing in instances of `SVC(kernel='rbf', C=c, gamma=gamma)` with different values for  $C$  and  $\gamma$ . Explain what kind of grid you used and why.
- (c) **(6 pts)** Finally, using the training data from Section 0.1 and the function implemented here, find the best setting for  $C$  and  $\gamma$  for each performance measure mentioned above. Report your findings in tabular format. This time, because we have a two-dimensional grid search, report only the best score for each metric, along with the accompanying  $C$  and  $\gamma$  setting.

metric	score	$C$	$\gamma$
accuracy			
F1-score			
AUROC			
precision			
sensitivity			
specificity			

How does the CV performance vary with the hyperparameters of the RBF-kernel SVM?

### 0.4 Test Set Performance [12 pts]

In this section, you will apply the two classifiers learned in the previous sections to the test data from Section 0.1. Once you have predicted labels for the test data, you will measure performance.

- (a) **(4 pts)** Based on the results you obtained in Section 0.2 and Section 0.3, choose a hyperparameter setting for the linear-kernel SVM and a hyperparameter setting for the RBF-kernel SVM. Explain your choice.  
Then, in `main(...)`, using the training data extracted in Section 0.1 and `SVC.fit(...)`, train a linear- and an RBF-kernel SVM with your chosen settings.
- (b) **(2 pts)** Implement `performance_test(...)` which returns the value of a performance measure, given the test data and a trained classifier.
- (c) **(6 pts)** For each performance metric, use `performance_test(...)` and the two trained linear- and RBF-kernel SVM classifiers to measure performance on the test data. Report the results. Be sure to include the name of the performance metric employed, and the performance on the test data. How do the test performance of your two classifiers compare?