

## Homework 3: Max-Margin and SVM

### Introduction

This homework assignment will have you work with max-margin methods and SVM classification. The aim of the assignment is (1) to further develop your geometrical intuition behind margin-based classification and decision boundaries, (2) to explore the properties of kernels and how they provide a different form of feature development from basis functions, and finally (3) to implement a basic Kernel based classifier.

There is a mathematical component and a programming component to this homework. Please submit your PDF and Python files to Canvas, and push all of your work to your GitHub repository. If a question requires you to make any plots, like Problem 3, please include those in the writeup.

**Problem 1** (Fitting an SVM by hand, 7pts)

For this problem you will solve an SVM without the help of a computer, relying instead on principled rules and properties of these classifiers.

Consider a dataset with the following 7 data points each with  $x \in \mathbb{R}$  :

$$\{(x_i, y_i)\}_i = \{(-3, +1), (-2, +1), (-1, -1), (0, -1), (1, -1), (2, +1), (3, +1)\}$$

Consider mapping these points to 2 dimensions using the feature vector  $\phi(x) = (x, x^2)$ . The hard margin classifier training problem is:

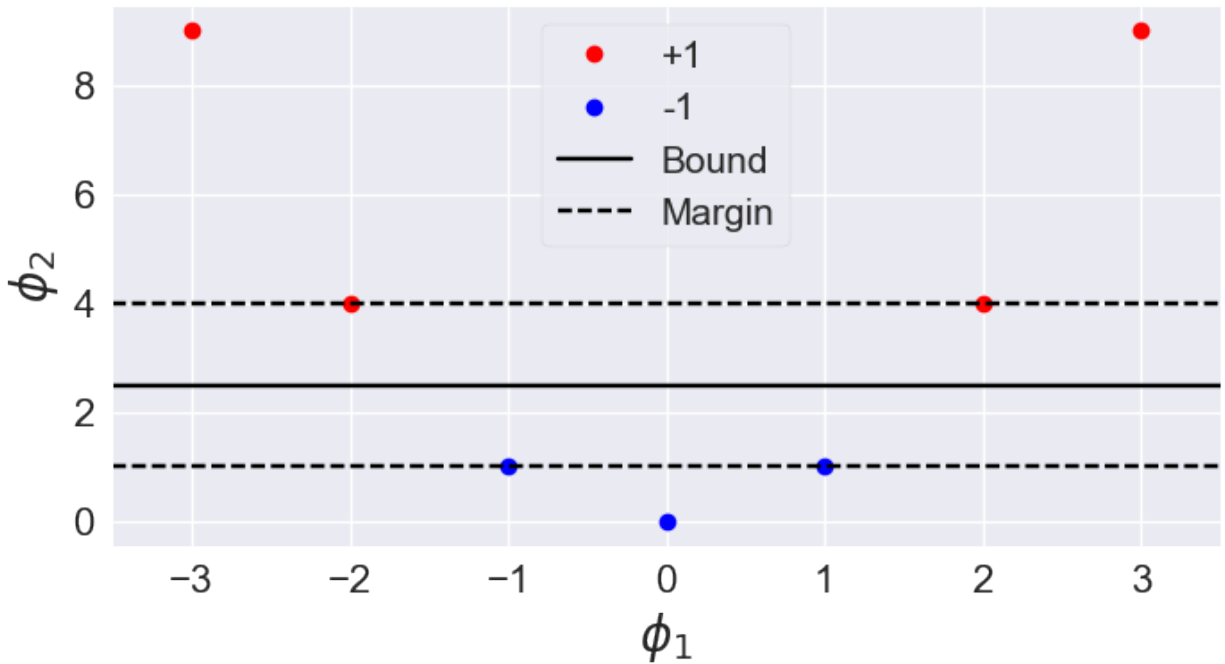
$$\begin{aligned} \min_{\mathbf{w}, w_0} \quad & \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^\top \phi(x_i) + w_0) \geq 1, \quad \forall i \in \{1, \dots, n\} \end{aligned} \tag{1}$$

The exercise has been broken down into a series of questions, each providing a part of the solution. Make sure to follow the logical structure of the exercise when composing your answer and to justify each step.

1. Plot the training data in  $\mathbb{R}^2$  and draw the decision boundary of the max margin classifier.
2. What is the value of the margin achieved by the optimal decision boundary?
3. What is a vector that is orthogonal to the decision boundary?
4. Considering discriminant  $h(\phi(x); \mathbf{w}, w_0) = \mathbf{w}^\top \phi(x) + w_0$ , give an expression for *all possible*  $(\mathbf{w}, w_0)$  that define the optimal decision boundary. Justify your answer.
5. Consider now the training problem (1). Using your answers so far, what particular solution to  $\mathbf{w}$  will be optimal for this optimization problem?
6. Now solve for the corresponding value of  $w_0$ , using your general expression from part (4.) for the optimal decision boundary. Write down the discriminant function  $h(\phi(x); \mathbf{w}, w_0)$ .
7. What are the support vectors of the classifier? Confirm that the solution in part (6.) makes the constraints in (1) binding for support vectors.

## Solution

### Problem 1



**Figure 1:** Training data in  $\mathbb{R}^2$  with decision boundary and margin.

### Problem 2

The value of the margin achieved by the optimal decision boundary is

$$\frac{y_i(\mathbf{w}^\top \mathbf{x}_i + w_0)}{\|\mathbf{w}\|} = 1.5$$

### Problem 3

The vector orthogonal to the decision boundary is given by  $\mathbf{w}$  and in this case is of arbitrary magnitude such that:

$$\mathbf{w} = \lambda[0, 1], \forall \lambda \in \mathbb{R}$$

**Problem 4**

**Problem 5**

**Problem 6**

**Problem 7**

**Problem 2** (Composing Kernel Functions, 10pts)

A key benefit of SVM training is the ability to use kernel functions  $K(\mathbf{x}, \mathbf{x}')$  as opposed to explicit basis functions  $\phi(\mathbf{x})$ . Kernels make it possible to implicitly express large or even infinite dimensional basis features. We do this by computing  $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$  directly, without ever computing  $\phi(\mathbf{x})$ .

When training SVMs, we begin by computing the kernel matrix  $\mathbf{K}$ , over our training data  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . The kernel matrix, defined as  $K_{i,i'} = K(\mathbf{x}_i, \mathbf{x}_{i'})$ , expresses the kernel function applied between all pairs of training points.

In class, we saw Mercer's theorem, which tells us that any function  $K$  that yields a positive semi-definite kernel matrix forms a valid kernel, i.e. corresponds to a matrix of dot-products under *some* basis  $\phi$ . Therefore instead of using an explicit basis, we can build kernel functions directly that fulfill this property.

A particularly nice benefit of this theorem is that it allows us to build more expressive kernels by composition. In this problem, you are tasked with using Mercer's theorem and the definition of a kernel matrix to prove that the following compositions are valid kernels, assuming  $K^{(1)}$  and  $K^{(2)}$  are valid kernels. Recall that a positive semi-definite matrix  $\mathbf{K}$  requires  $\mathbf{z}^\top \mathbf{K} \mathbf{z} \geq 0$ ,  $\forall \mathbf{z} \in \mathbb{R}^n$ .

1.  $K(\mathbf{x}, \mathbf{x}') = c K^{(1)}(\mathbf{x}, \mathbf{x}')$  for  $c > 0$
2.  $K(\mathbf{x}, \mathbf{x}') = K^{(1)}(\mathbf{x}, \mathbf{x}') + K^{(2)}(\mathbf{x}, \mathbf{x}')$
3.  $K(\mathbf{x}, \mathbf{x}') = f(\mathbf{x}) K^{(1)}(\mathbf{x}, \mathbf{x}') f(\mathbf{x}')$  where  $f$  is any function from  $\mathbb{R}^m$  to  $\mathbb{R}$
4.  $K(\mathbf{x}, \mathbf{x}') = K^{(1)}(\mathbf{x}, \mathbf{x}') K^{(2)}(\mathbf{x}, \mathbf{x}')$

[Hint: Use the property that for any  $\phi(\mathbf{x})$ ,  $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$  forms a positive semi-definite kernel matrix. ]

5. (a) The exp function can be written as,

$$\exp(x) = \lim_{i \rightarrow \infty} \left( 1 + x + \dots + \frac{x^i}{i!} \right).$$

Use this to show that  $\exp(xx')$  (here  $x, x' \in \mathbb{R}$ ) can be written as  $\phi(x)^\top \phi(x')$  for some basis function  $\phi(x)$ . Derive this basis function, and explain why this would be hard to use as a basis in standard logistic regression.

- (b) Using the previous identities, show that  $K(\mathbf{x}, \mathbf{x}') = \exp(K^{(1)}(\mathbf{x}, \mathbf{x}'))$  is a valid kernel.

6. Finally use this analysis and previous identities to prove the validity of the Gaussian kernel:

$$K(\mathbf{x}, \mathbf{x}') = \exp \left( \frac{-\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2} \right)$$

**Solution**

**Problem 3** (Scaling up your SVM solver, 10pts (+opportunity for extra credit))

For this problem you will build a simple SVM classifier for a binary classification problem. We have provided you two files for experimentation: training *data.csv* and validation *val.csv*.

- First read the paper at <http://www.jmlr.org/papers/volume6/bordes05a/bordes05a.pdf> and implement the Kernel Perceptron algorithm and the Budget Kernel Perceptron algorithm. Aim to make the optimization as fast as possible. Implement this algorithm in *problem3.py*.

[Hint: For this problem, efficiency will be an issue. Instead of directly implementing this algorithm using numpy matrices, you should utilize Python dictionaries to represent sparse matrices. This will be necessary to have the algorithm run in a reasonable amount of time. ]

- Next experiment with the hyperparameters for each of these models. Try seeing if you can identify some patterns by changing  $\beta$ ,  $N$  (the maximum number of support vectors), or the number of random training samples taken during the Randomized Search procedure (Section 4.3). Note the training time, training and validation accuracy, and number of support vectors for various setups.
- Lastly, compare the classification to the naive SVM imported from scikit-learn by reporting accuracy on the provided validation data. *For extra credit, implement the SMO algorithm and implement the LASVM process and do the same as above.*<sup>a</sup>

We are intentionally leaving this problem open-ended to allow for experimentation, and so we will be looking for your thought process and not a particular graph. Visualizations should be generated using the provided code. You can use the trivial  $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$  kernel for this problem, though you are welcome to experiment with more interesting kernels too.

In addition, provide answers the following reading questions **in one or two sentences for each**.

1. In one short sentence, state the main purpose of the paper.
2. Describe each of the parameters in Eq. 1 in the paper
3. State, informally, one guarantee about the Kernel perceptron algorithm described in the paper.
4. What is the main way the budget kernel perceptron algorithm tries to improve on the perceptron algorithm?
5. (*if you did the extra credit*) In simple words, what is the theoretical guarantee of LASVM algorithm? How does it compare to its practical performance?

---

<sup>a</sup>Extra credit only makes a difference to your grade at the end of the semester if you are on a grade boundary.

## Solution

### **Calibration [1pt]**

Approximately how long did this homework take you to complete?