

## Homework 3: SVMs & Kernel Methods

**Due:** Wednesday, March 2, 2022 at 11:59PM EST

**Instructions:** Your answers to the questions below, including plots and mathematical work, should be submitted as a single PDF file. It's preferred that you write your answers using software that typesets mathematics (e.g. LaTeX, LyX, or MathJax via iPython), though if you need to you may scan handwritten work. You may find the minted package convenient for including source code in your LaTeX document. If you are using LyX, then the listings package tends to work better.

In this problem set we will get up to speed with SVMs and Kernels. Long at first glance, the problem set includes a lot of helpers. You will find a review of kernalization. One section will include a revision of ridge regression which you should start to be familiar with. For the second and third problem some codes are provided to save you some time. Finally, some reminders on positive (semi)definite matrices are included in the Appendix.

### 1 Support Vector Machines: SVMs with Pegasos

In this first problem we will use Support Vector Machines to predict whether the sentiment of a movie review was *positive* or *negative*. We will represent each review by a vector  $\mathbf{x} \in \mathbb{R}^d$  where  $d$  is the size of the word dictionary and  $x_i$  is equal to the number of occurrence of the  $i$ -th word in the review  $\mathbf{x}$ . The corresponding label is either  $y = 1$  for a positive review or  $y = -1$  for a negative review. In class we have seen how to transform the SVM training objective into a quadratic program using the dual formulation. Here we will use a gradient descent algorithm instead.

#### Subgradients

Recall that a vector  $g \in \mathbb{R}^d$  is a *subgradient* of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  at  $\mathbf{x}$  if for all  $\mathbf{z}$ ,

$$f(\mathbf{z}) \geq f(\mathbf{x}) + g^T(\mathbf{z} - \mathbf{x}).$$

There may be 0, 1, or infinitely many subgradients at any point. The *subdifferential* of  $f$  at a point  $\mathbf{x}$ , denoted  $\partial f(\mathbf{x})$ , is the set of all subgradients of  $f$  at  $\mathbf{x}$ . A good reference for subgradients are the course notes on Subgradients by Boyd et al. Below we derive a property that will make our life easier for finding a subgradient of the hinge loss.

1. Suppose  $f_1, \dots, f_m : \mathbb{R}^d \rightarrow \mathbb{R}$  are convex functions, and  $f(\mathbf{x}) = \max_{i=1, \dots, m} f_i(\mathbf{x})$ . Let  $k$  be any index for which  $f_k(\mathbf{x}) = f(\mathbf{x})$ , and choose  $g \in \partial f_k(\mathbf{x})$  (a convex function on  $\mathbb{R}^d$  has a non-empty subdifferential at all points). Show that  $g \in \partial f(\mathbf{x})$ .

$$\begin{aligned} f(\mathbf{z}) &\geq f(\mathbf{x}) + g^T(\mathbf{z} - \mathbf{x}) \\ f_k(\mathbf{z}) &\geq f_k(\mathbf{x}) + g_k^T(\mathbf{z} - \mathbf{x}) \\ f(\mathbf{z}) &\geq f_k(\mathbf{z}) \geq f_k(\mathbf{x}) + g_k^T(\mathbf{z} - \mathbf{x}) \\ f(\mathbf{z}) &\geq f_k(\mathbf{x}) + g_k^T(\mathbf{z} - \mathbf{x}) \end{aligned} \tag{1}$$

2. Give a subgradient of the hinge loss objective The subgradient is defined when  $y\mathbf{w}^T \mathbf{x} > 1$  or  $y\mathbf{w}^T \mathbf{x} < 1$ . We can summarize the subgradients by taking the derivative in each of those cases.  $J(\mathbf{w}) = \max\{0, 1 - y\mathbf{w}^T \mathbf{x}\}$ .

$$\text{Subgradient of } J(w) = \begin{cases} 0 & \text{if } yw^T x \geq 1 \\ -yx & \text{if } yw^T x < 1 \end{cases}$$

3. Suppose we have function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  which is sub-differentiable everywhere, i.e.  $\partial f \neq \emptyset$  for all  $x \in \mathbb{R}^n$ . Show that  $f$  is convex. Note, in the general case, a function is convex if for all  $x, y$  in the domain of  $f$  and for all  $\theta \in (0, 1)$ ,

$$\theta f(a) + (1 - \theta)f(b) \geq f(\theta a + (1 - \theta)(b))$$

Hint: Suppose  $f$  is not convex, then by definition, there exists a point in some interval:  $x_0 \in (a, b)$ , such that  $f(x_0)$  lies above the line connection  $(a, f(a)), (b, f(b))$ . Is this possible if the function is sub-differentiable everywhere?

### SVM with the Pegasos algorithm

You will train a Support Vector Machine using the Pegasos algorithm<sup>1</sup>. Recall the SVM objective using a linear predictor  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  and the hinge loss:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\},$$

where  $n$  is the number of training examples and  $d$  the size of the dictionary. Note that, for simplicity, we are leaving off the bias term  $b$ . Note also that we are using  $\ell_2$  regularization with a parameter  $\lambda$ . Pegasos is stochastic subgradient descent using a step size rule  $\eta_t = 1/(\lambda t)$  for iteration number  $t$ . The pseudocode is given below:

---

Input:  $\lambda > 0$ . Choose  $w_1 = 0, t = 0$   
 While termination condition not met  
   For  $j = 1, \dots, n$  (assumes data is randomly permuted)  
      $t = t + 1$   
      $\eta_t = 1/(t\lambda)$ ;  
     If  $y_j w_t^T x_j < 1$   
        $w_{t+1} = (1 - \eta_t \lambda) w_t + \eta_t y_j x_j$   
     Else  
        $w_{t+1} = (1 - \eta_t \lambda) w_t$

---

4. Consider the SVM objective function for a single training point<sup>2</sup>:  $J_i(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \max \{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\}$ . The function  $J_i(\mathbf{w})$  is not differentiable everywhere. Specify where the gradient of  $J_i(w)$  is not defined. Give an expression for the gradient where it is defined.

The gradient is defined and not defined by the following piece-wise expression:

$$\begin{cases} \text{Defined} & \text{when } y_i \mathbf{w}^T \mathbf{x}_i \neq 1 \\ \text{Not - Defined} & \text{when } y_i \mathbf{w}^T \mathbf{x}_i = 1 \end{cases}$$

<sup>1</sup>Shalev-Shwartz et al. Pegasos: Primal Estimated sub-GrAdient SOLver for SVM.

<sup>2</sup>Recall that if  $i$  is selected uniformly from the set  $\{1, \dots, n\}$ , then this objective function has the same expected value as the full SVM objective function.

5. Show that a subgradient of  $J_i(w)$  is given by

$$g w = \begin{cases} \lambda w - y_i x_i & \text{for } y_i w^T x_i < 1 \\ \lambda w & \text{for } y_i w^T x_i \geq 1. \end{cases}$$

You may use the following facts without proof: 1) If  $f_1, \dots, f_n : \mathbb{R}^d \rightarrow \mathbb{R}$  are convex functions and  $f = f_1 + \dots + f_n$ , then  $\partial f(x) = \partial f_1(x) + \dots + \partial f_n(x)$ . 2) For  $\alpha \geq 0$ ,  $\partial(\alpha f)(x) = \alpha \partial f(x)$ . (Hint: Use the first part of this problem.)

Convince yourself that if your step size rule is  $\eta_t = 1/(\lambda t)$ , then doing SGD with the subgradient direction from the previous question is the same as given in the pseudocode.

Let  $f = J_i(w) = g_i(w) + k_i(w)$ , where  $g_i(w) = \frac{1}{2} \|w\|^2$  and  $k_i(w) = \max(0, 1 - y_i w^T x_i)$ . Since the derivative of  $f$  is defined when  $y_i w^T x_i \neq 1$  we can calculate the subgradient for two cases, when  $y_i w^T x_i < 1$  and when  $y_i w^T x_i \leq 1$ .

Using the property that the property of derivatives, that the derivative of a sum is the sum of a derivatives, therefore  $\partial f_i(w) = \partial g_i(w) + \partial k_i(w)$ .

In all cases, the subgradient for  $g_i(w)$  is as follows:  $\partial g_i(w) = \lambda w$ .

Where as for  $k_i(w)$ :

$$\begin{cases} \partial k_i(w) = -y_i x_i & \text{when } y_i w^T x_i < 1 \\ \partial k_i(w) = \text{undefined} & \text{when } y_i w^T x_i = 1 \\ \partial k_i(w) = 0 & \text{otherwise} \end{cases}$$

Therefore, the sub-gradients for  $J_i(w)$  are as follows:

$$g w = \begin{cases} \lambda w - y_i x_i & \text{for } y_i w^T x_i < 1 \\ \lambda w & \text{for } y_i w^T x_i \geq 1. \end{cases}$$

## Dataset and sparse representation

We will be using the Polarity Dataset v2.0, constructed by Pang and Lee, provided in the `data_reviews` folder. It has the full text from 2000 movies reviews: 1000 reviews are classified as *positive* and 1000 as *negative*. Our goal is to predict whether a review has positive or negative sentiment from the text of the review. Each review is stored in a separate file: the positive reviews are in a folder called “pos”, and the negative reviews are in “neg”. We have provided some code in `utils_svm_reviews.py` to assist with reading these files. The code removes some special symbols from the reviews and shuffles the data. Load all the data to have an idea of what it looks like.

A usual method to represent text documents in machine learning is with *bag-of-words*. As hinted above, here every possible word in the dictionary is a feature, and the value of a word feature for a given text is the number of times that word appears in the text. As most words will not appear in any particular document, many of these counts will be zero. Rather than storing many zeros, we use a *sparse representation*, in which only the nonzero counts are tracked. The counts are stored in a key/value data structure, such as a dictionary in Python. For example, “Harry Potter and Harry Potter II” would be represented as the following Python dict: `x={'Harry':2, 'Potter':2, 'and':1, 'II':1}`.

6. Write a function that converts an example (a list of words) into a sparse bag-of-words representation. You may find Python's Counter<sup>3</sup> class to be useful here. Note that a Counter is itself a dictionary.
7. Load all the data and split it into 1500 training examples and 500 validation examples. Format the training data as a list `X_train` of dictionaries and `y_train` as the list of corresponding 1 or -1 labels. Format the test set similarly.

We will be using linear classifiers of the form  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ , and we can store the  $\mathbf{w}$  vector in a sparse format as well, such as `w={'minimal':1.3, 'Harry':-1.1, 'viable':-4.2, 'and':2.2, 'product':9.1}`. The inner product between  $\mathbf{w}$  and  $\mathbf{x}$  would only involve the features that appear in both  $\mathbf{x}$  and  $\mathbf{w}$ , since whatever doesn't appear is assumed to be zero. For this example, the inner product would be `x(Harry) * w(Harry) + x(and) * w(and) = 2*(-1.1) + 1*(2.2)`. To help you along, `utils.svm_reviews.py` includes two functions for working with sparse vectors: 1) a dot product between two vectors represented as dictionaries and 2) a function that increments one sparse vector by a scaled multiple of another vector, which is a very common operation. It is worth reading the code, even if you intend to implement it yourself. You may get some ideas on how to make things faster.

8. Implement the Pegasos algorithm to run on a sparse data representation. The output should be a sparse weight vector  $\mathbf{w}$  represented as a dictionary. Note that our Pegasos algorithm starts at  $w = 0$ , which corresponds to an empty dictionary. **Note:** With this problem, you will need to take some care to code things efficiently. In particular, be aware that making copies of the weight dictionary can slow down your code significantly. If you want to make a copy of your weights (e.g. for checking for convergence), make sure you don't do this more than once per epoch. **Also:** If you normalize your data in some way, be sure not to destroy the sparsity of your data. Anything that starts as 0 should stay at 0.

Note that in every step of the Pegasos algorithm, we rescale every entry of  $w_t$  by the factor  $(1 - \eta_t \lambda)$ . Implementing this directly with dictionaries is very slow. We can make things significantly faster by representing  $w$  as  $w = sW$ , where  $s \in \mathbb{R}$  and  $W \in \mathbb{R}^d$ . You can start with  $s = 1$  and  $W$  all zeros (i.e. an empty dictionary). Note that both updates (i.e. whether or not we have a margin error) start with rescaling  $w_t$ , which we can do simply by setting  $s_{t+1} = (1 - \eta_t \lambda) s_t$ .

9. If the update is  $w_{t+1} = (1 - \eta_t \lambda)w_t + \eta_t y_j x_j$ , then verify that the Pegasos update step is equivalent to:

$$\begin{aligned} s_{t+1} &= (1 - \eta_t \lambda) s_t \\ W_{t+1} &= W_t + \frac{1}{s_{t+1}} \eta_t y_j x_j. \end{aligned}$$

Implement the Pegasos algorithm with the  $(s, W)$  representation described above.

Using the definitions given, we can substitute in the necessary values and manipulate

---

<sup>3</sup><https://docs.python.org/2/library/collections.html>

the expression to reach our equivalency:

$$\begin{aligned} w_{t+1} &= s_{t+1} W_{t+1} \\ w_{t+1} &= ((1 - \eta_t \lambda) s_t) (W_t + \frac{1}{((1 - \eta_t \lambda) s_t)} \eta_t y_j x_j) \\ w_{t+1} &= ((1 - \eta_t \lambda) s_t) W_t + \eta_t y_j x_j \end{aligned} \tag{2}$$

4

10. Run both implementations of Pegasos on the training data for a couple epochs. Make sure your implementations are correct by verifying that the two approaches give essentially the same result. Report on the time taken to run each approach.
11. Write a function `classification_error` that takes a sparse weight vector  $\mathbf{w}$ , a list of sparse vectors  $\mathbf{X}$  and the corresponding list of labels  $\mathbf{y}$ , and returns the fraction of errors when predicting  $y_i$  using  $\text{sign}(\mathbf{w}^T \mathbf{x}_i)$ . In other words, the function reports the 0-1 loss of the linear predictor  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ .
12. Search for the regularization parameter that gives the minimal percent error on your test set. You should now use your faster Pegasos implementation, and run it to convergence. A good search strategy is to start with a set of regularization parameters spanning a broad range of orders of magnitude. Then, continue to zoom in until you're convinced that additional search will not significantly improve your test performance. Plot the test errors you obtained as a function of the parameters  $\lambda$  you tested. (Hint: the error you get with the best regularization should be closer to 15% than 20%. If not, maybe you did not train to convergence.)

## Error Analysis

Recall that the *score* is the value of the prediction  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ . We like to think that the magnitude of the score represents the confidence of the prediction. This is something we can directly verify or refute.

13. Break the predictions on the test set into groups based on the score (you can play with the size of the groups to get a result you think is informative). For each group, examine the percentage error. You can make a table or graph. Summarize the results. Is there a correlation between higher magnitude scores and accuracy?

In natural language processing one can often interpret why a model has performed well or poorly on a specific example. The first step in this process is to look closely at the errors that the model makes.

---

<sup>4</sup>There is one subtle issue with the approach described above: if we ever have  $1 - \eta_t \lambda = 0$ , then  $s_{t+1} = 0$ , and we'll have a divide by 0 in the calculation for  $W_{t+1}$ . This only happens when  $\eta_t = 1/\lambda$ . With our step-size rule of  $\eta_t = 1/(\lambda t)$ , it happens exactly when  $t = 1$ . So one approach is to just start at  $t = 2$ . More generically, note that if  $s_{t+1} = 0$ , then  $w_{t+1} = 0$ . Thus an equivalent representation is  $s_{t+1} = 1$  and  $W = 0$ . Thus if we ever get  $s_{t+1} = 0$ , simply set it back to 1 and reset  $W_{t+1}$  to zero, which is an empty dictionary in a sparse representation.

14. (Optional) Choose an input example  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$  that the model got wrong. We want to investigate what features contributed to this incorrect prediction. One way to rank the importance of the features to the decision is to sort them by the size of their contributions to the score. That is, for each feature we compute  $|w_i x_i|$ , where  $w_i$  is the weight of the  $i$ th feature in the prediction function, and  $x_i$  is the value of the  $i$ th feature in the input  $\mathbf{x}$ . Create a table of the most important features, sorted by  $|w_i x_i|$ , including the feature name, the feature value  $x_i$ , the feature weight  $w_i$ , and the product  $w_i x_i$ . Attempt to explain why the model was incorrect. Can you think of a new feature that might be able to fix the issue? Include a short analysis for at least 2 incorrect examples. Can you think of new features that might help fix a problem? (Think of making groups of words.)

## 2 Kernel Methods

### 2.1 Kernelization review

Consider the following optimization problem on a data set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \mathcal{Y}$ :

$$\min_{\mathbf{w} \in \mathbb{R}^d} R\left(\sqrt{\langle \mathbf{w}, \mathbf{w} \rangle}\right) + L(\langle \mathbf{w}, \mathbf{x}_1 \rangle, \dots, \langle \mathbf{w}, \mathbf{x}_n \rangle),$$

where  $\mathbf{w}, \mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ , and  $\langle \cdot, \cdot \rangle$  is the standard inner product on  $\mathbb{R}^d$ . The function  $R : [0, \infty) \rightarrow \mathbb{R}$  is nondecreasing and gives us our regularization term, while  $L : \mathbb{R}^n \rightarrow \mathbb{R}$  is arbitrary<sup>5</sup> and gives us our loss term. We noted in lecture that this general form includes soft-margin SVM and ridge regression, though not lasso regression. Using the representer theorem, we showed if the optimization problem has a solution, there is always a solution of the form  $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$ , for some  $\alpha \in \mathbb{R}^n$ . Plugging this into the our original problem, we get the following “kernelized” optimization problem:

$$\min_{\alpha \in \mathbb{R}^n} R\left(\sqrt{\alpha^T K \alpha}\right) + L(K\alpha),$$

where  $K \in \mathbb{R}^{n \times n}$  is the Gram matrix (or “kernel matrix”) defined by  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$ . Predictions are given by

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x}),$$

and we can recover the original  $\mathbf{w} \in \mathbb{R}^d$  by  $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$ .

The *kernel trick* is to swap out occurrences of the kernel  $k$  (and the corresponding Gram matrix  $K$ ) with another kernel. For example, we could replace  $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$  by  $k'(\mathbf{x}_i, \mathbf{x}_j) = \langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_j) \rangle$  for an arbitrary feature mapping  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . In this case, the recovered  $\mathbf{w} \in \mathbb{R}^d$  would be  $\mathbf{w} = \sum_{i=1}^n \alpha_i \psi(\mathbf{x}_i)$  and predictions would be  $\langle \mathbf{w}, \psi(\mathbf{x}_i) \rangle$ .

More interestingly, we can replace  $k$  by another kernel  $k''(\mathbf{x}_i, \mathbf{x}_j)$  for which we do not even know or cannot explicitly write down a corresponding feature map  $\psi$ . Our main example of this is the RBF kernel

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

<sup>5</sup>You may be wondering “Where are the  $y_i$ ’s?”. They’re built into the function  $L$ . For example, a square loss on a training set of size 3 could be represented as  $L(s_1, s_2, s_3) = \frac{1}{3} \left[ (s_1 - y_1)^2 + (s_2 - y_2)^2 + (s_3 - y_3)^2 \right]$ , where each  $s_i$  stands for the  $i$ th prediction  $\langle \mathbf{w}, \mathbf{x}_i \rangle$ .

for which the corresponding feature map  $\psi$  is infinite dimensional. In this case, we cannot recover  $w$  since it would be infinite dimensional. Predictions must be done using  $\alpha \in \mathbb{R}^n$ , with  $f(x) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$ .

Your implementation of kernelized methods below should not make any reference to  $w$  or to a feature map  $\psi$ . Your learning routine should return  $\alpha$ , rather than  $w$ , and your prediction function should also use  $\alpha$  rather than  $w$ . This will allow us to work with kernels that correspond to infinite-dimensional feature vectors.

## 2.2 Kernel problems

### Ridge Regression: Theory

Suppose our input space is  $\mathcal{X} = \mathbb{R}^d$  and our output space is  $\mathcal{Y} = \mathbb{R}$ . Let  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  be a training set from  $\mathcal{X} \times \mathcal{Y}$ . We'll use the "design matrix"  $X \in \mathbb{R}^{n \times d}$ , which has the input vectors as rows:

$$X = \begin{pmatrix} -\mathbf{x}_1 - \\ \vdots \\ -\mathbf{x}_n - \end{pmatrix}.$$

Recall the ridge regression objective function:

$$J(w) = \|Xw - y\|^2 + \lambda \|w\|^2,$$

for  $\lambda > 0$ .

15. Show that for  $w$  to be a minimizer of  $J(w)$ , we must have  $X^T X w + \lambda I w = X^T y$ . Show that the minimizer of  $J(w)$  is  $w = (X^T X + \lambda I)^{-1} X^T y$ . Justify that the matrix  $X^T X + \lambda I$  is invertible, for  $\lambda > 0$ . (You should use properties of positive (semi)definite matrices. If you need a reminder look up the Appendix.)

We can minimize  $J(w)$  by taking its gradient and setting it to 0. The gradient is as follows:

$$\nabla J(w) = \nabla(\|Xw - y\|^2 + \lambda \|w\|^2) = 2X^T X w + 2\lambda I w - 2X^T y$$

Setting it to 0 and moving over the  $-X^T y$  term we have:

$$X^T X w + \lambda I w = X^T y$$

Lastly, we can factor out  $w$  and take the inverse of the resulting matrix to arrive at our solution:

$$X^T X w + \lambda I w = X^T y \rightarrow (X^T X + \lambda I) w = X^T y \rightarrow w = (X^T X + \lambda I)^{-1} X^T y$$

By definition, square symmetric matrices are PSD. Therefore, it will have a number positive eigenvalues equal to the matrices rank, and eigenvalues equivalent to 0 equal to the dimension of its kernel. Formalized, we can say it will have  $d$  eigenvalues greater than 0, where  $d = \text{Rank}(X^T X)$  and eigenvalues equal to 0 for  $n - d = \dim(\ker(X^T X))$  where  $n$  is the number of columns of the matrix  $X^T X$ . Also, due to rank nullity, we know that  $\text{Rank}(X) = \text{Rank}(X^T) = \text{Rank}(X^T X)$ . We also know that when we add any number,  $\alpha$  to the diagonal of a square matrix, its eigenvalues are shifted by a degree of  $\alpha$ . Therefore, there must exist some  $\alpha > 0$  such that the eigenvalues of  $X^T X$  are all positive, and therefore, the matrix will be full rank and therefore invertible.

16. Rewrite  $X^T X w + \lambda I w = X^T y$  as  $\mathbf{w} = \frac{1}{\lambda}(X^T y - X^T X w)$ . Based on this, show that we can write  $w = X^T \alpha$  for some  $\alpha$ , and give an expression for  $\alpha$ .

$$\begin{aligned}
 X^T X w + \lambda I w &= X^T y \\
 \lambda I w &= X^T y - X^T X w \\
 w &= \frac{1}{\lambda}(X^T y - X^T X w) \\
 w &= X^T \frac{1}{\lambda}(y - X w) \\
 \text{let } \alpha &= \frac{1}{\lambda}(y - X w) \\
 w &= X^T \alpha
 \end{aligned} \tag{3}$$

We have shown that  $X^T X w + \lambda I w = X^T y$  can be re expressed as  $w = X^T \alpha$  when we let  $\alpha = \frac{1}{\lambda}(y - X w)$

17. Based on the fact that  $\mathbf{w} = X^T \alpha$ , explain why we say  $\mathbf{w}$  is “in the span of the data.”

Since  $X$  is our data, and  $w$  is defined as a linear combination of our columns of  $X$ , (formalized mathematically by  $w = X^T \alpha$ ) by definition it  $w$  is in the span of our data.

18. Show that  $\alpha = (\lambda I + X X^T)^{-1} y$ . Note that  $X X^T$  is the kernel matrix for the standard vector dot product. (Hint: Replace  $\mathbf{w}$  by  $X^T \alpha$  in the expression for  $\alpha$ , and then solve for  $\alpha$ .)

$$\begin{aligned}
 \alpha &= \frac{1}{\lambda}(y - X w) \\
 \alpha &= \frac{1}{\lambda}(y - X X^T \alpha) \\
 \lambda \alpha &= y - X X^T \alpha \\
 \lambda \alpha + X X^T \alpha &= y \\
 (Id_n \lambda + X X^T) \alpha &= y \\
 \alpha &= (Id_n \lambda + X X^T)^{-1} y
 \end{aligned} \tag{4}$$

19. Give a kernelized expression for the  $X \mathbf{w}$ , the predicted values on the training points. (Hint: Replace  $\mathbf{w}$  by  $X^T \alpha$  and  $\alpha$  by its expression in terms of the kernel matrix  $X X^T$ .)

$$\begin{aligned}
 X w &= X(X^T \alpha) \\
 &= X(X^T (Id_n \lambda + X X^T)^{-1} y) \\
 &= X X^T (Id_n \lambda + X X^T)^{-1} y
 \end{aligned} \tag{5}$$



20. Give an expression for the prediction  $f(x) = x^T \mathbf{w}^*$  for a new point  $x$ , not in the training set. The expression should only involve  $x$  via inner products with other  $\mathbf{x}$ 's. (Hint: It is often convenient to define the column vector

$$\mathbf{k}_x = \begin{pmatrix} \mathbf{x}^T \mathbf{x}_1 \\ \vdots \\ \mathbf{x}^T \mathbf{x}_n \end{pmatrix}$$

to simplify the expression.)

$$\begin{aligned} f(x) &= x^T \mathbf{w} \\ &= x^T (X^T \alpha) \\ &= \sum_{i=1}^n \alpha_i \langle x^T, x_i \rangle \\ &= \sum_{i=1}^n \alpha_i k_{xi} \end{aligned} \tag{6}$$

## Kernels and Kernel Machines

There are many different families of kernels. So far we spoken about linear kernels, RBF/Gaussian kernels, and polynomial kernels. The last two kernel types have parameters. In this section, we'll implement these kernels in a way that will be convenient for implementing our kernelized ridge regression later on. For simplicity, we will assume that our input space is  $\mathcal{X} = \mathbb{R}$ . This allows us to represent a collection of  $n$  inputs in a matrix  $X \in \mathbb{R}^{n \times 1}$ . You should now refer to the jupyter notebook `skeleton_code_kernels.ipynb`.

21. Write functions that compute the RBF kernel  $k_{\text{RBF}(\sigma)}(x, x') = \exp(-\|x - x'\|^2 / (2\sigma^2))$  and the polynomial kernel  $k_{\text{poly}(a,d)}(x, x') = (a + \langle x, x' \rangle)^d$ . The linear kernel  $k_{\text{linear}}(x, x') = \langle x, x' \rangle$ , has been done for you in the support code. Your functions should take as input two matrices  $W \in \mathbb{R}^{n_1 \times d}$  and  $X \in \mathbb{R}^{n_2 \times d}$  and should return a matrix  $M \in \mathbb{R}^{n_1 \times n_2}$  where  $M_{ij} = k(W_i, X_j)$ . In words, the  $(i, j)$ 'th entry of  $M$  should be kernel evaluation between  $w_i$  (the  $i$ th row of  $W$ ) and  $x_j$  (the  $j$ th row of  $X$ ). For the RBF kernel, you may use the scipy function `cdist(X1, X2, 'sqeuclidean')` in the package `scipy.spatial.distance`.
22. Use the linear kernel function defined in the code to compute the kernel matrix on the set of points  $x_0 \in \mathcal{D}_X = \{-4, -1, 0, 2\}$ . Include both the code and the output.
23. Suppose we have the data set  $\mathcal{D}_{X,y} = \{(-4, 2), (-1, 0), (0, 3), (2, 5)\}$  (in each set of parentheses, the first number is the value of  $x_i$  and the second number the corresponding value of the target  $y_i$ ). Then by the representer theorem, the final prediction function will be in the span of the functions  $x \mapsto k(x_0, x)$  for  $x_0 \in \mathcal{D}_X = \{-4, -1, 0, 2\}$ . This set of functions will look quite different depending on the kernel function we use. The set of functions  $x \mapsto k_{\text{linear}}(x_0, x)$  for  $x_0 \in \mathcal{X}$  and for  $x \in [-6, 6]$  has been provided for the linear kernel.
- Plot the set of functions  $x \mapsto k_{\text{poly}(1,3)}(x_0, x)$  for  $x_0 \in \mathcal{D}_X$  and for  $x \in [-6, 6]$ .
  - Plot the set of functions  $x \mapsto k_{\text{RBF}(1)}(x_0, x)$  for  $x_0 \in \mathcal{X}$  and for  $x \in [-6, 6]$ .

Note that the values of the parameters of the kernels you should use are given in their definitions in (a) and (b).

24. By the representer theorem, the final prediction function will be of the form  $f(x) = \sum_{i=1}^n \alpha_i k(x_i, x)$ , where  $x_1, \dots, x_n \in \mathcal{X}$  are the inputs in the training set. We will use the class `KernelMachine` in the skeleton code to make prediction with different kernels. Complete the `predict` function of the class `KernelMachine`. Construct a `KernelMachine` object with the RBF kernel (`sigma=1`), with prototype points at  $-1, 0, 1$  and corresponding weights  $\alpha_i$   $1, -1, 1$ . Plot the resulting function.

Note: For this last problem, and for other problems below, it may be helpful to use partial application on your kernel functions. For example, if your polynomial kernel function has signature `polynomial_kernel(W, X, offset, degree)`, you can write `k = functools.partial(polynomial_kernel, offset=2, degree=2)`, and then a call to `k(W,X)` is equivalent to `polynomial_kernel(W, X, offset=2, degree=2)`, the advantage being that the extra parameter settings are built into `k(W,X)`. This can be convenient so that you can have a function that just takes a kernel function `k(W,X)` and doesn't have to worry about the parameter settings for the kernel.

### Kernel Ridge Regression: Practice

In the zip file for this assignment, we provide a training `krr-train.txt` and test set `krr-test.txt` for a one-dimensional regression problem, in which  $\mathcal{X} = \mathcal{Y} = \mathcal{A} = \mathbb{R}$ . Fitting this data using kernelized ridge regression, we will compare the results using several different kernel functions. Because the input space is one-dimensional, we can easily visualize the results.

25. Plot the training data. You should note that while there is a clear relationship between  $x$  and  $y$ , the relationship is not linear.
26. In a previous problem, we showed that in kernelized ridge regression, the final prediction function is  $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$ , where  $\alpha = (\lambda I + K)^{-1} y$  and  $K \in \mathbb{R}^{n \times n}$  is the kernel matrix of the training data:  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ , for  $\mathbf{x}_1, \dots, \mathbf{x}_n$ . In terms of kernel machines,  $\alpha_i$  is the weight on the kernel function evaluated at the training point  $\mathbf{x}_i$ . Complete the function `train_kernel_ridge_regression` so that it performs kernel ridge regression and returns a `KernelMachine` object that can be used for predicting on new points.
27. Use the code provided to plot your fits to the training data for the RBF kernel with a fixed regularization parameter of 0.0001 for 3 different values of `sigma`: 0.01, 0.1, and 1.0. What values of `sigma` do you think would be more likely to over fit, and which less?
28. Use the code provided to plot your fits to the training data for the RBF kernel with a fixed `sigma` of 0.02 and 4 different values of the regularization parameter  $\lambda$ : 0.0001, 0.01, 0.1, and 2.0. What happens to the prediction function as  $\lambda \rightarrow \infty$ ?
29. Find the best hyperparameter settings (including kernel parameters and the regularization parameter) for each of the kernel types. Summarize your results in a table, which gives training error and test error for each setting. Include in your table the best settings for each kernel type, as well as nearby settings that show that making small change in any one of the hyperparameters in either direction will cause the performance to get worse. You should use average square loss on the test set to rank the parameter settings. To make

things easier for you, we have provided an sklearn wrapper for the kernel ridge regression function we have created so that you can use sklearn's GridSearchCV. Note: Because of the small dataset size, these models can be fit extremely fast, so there is no excuse for not doing extensive hyperparameter tuning.

30. Plot your best fitting prediction functions using the polynomial kernel and the RBF kernel. Use the domain  $x \in (-0.5, 1.5)$ . Comment on the results.
31. The data for this problem was generated as follows: A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  was chosen. Then to generate a point  $(x, y)$ , we sampled  $x$  uniformly from  $(0, 1)$  and we sampled  $\epsilon \sim \mathcal{N}(0, 0.1^2)$  (so  $\text{var}(\epsilon) = 0.1^2$ ). The final point is  $(x, f(x) + \epsilon)$ . What is the Bayes decision function and the Bayes risk for the loss function  $\ell(\hat{y}, y) = (\hat{y} - y)^2$ .

### 3 Kernel SVMs with Kernelized Pegasos (Optional)

32. Load the SVM training `svm-train.txt` and `svm-test.txt` test data from the zip file. Plot the training data using the code supplied. Are the data linearly separable? Quadratically separable? What if we used some RBF kernel?
33. Unlike for kernel ridge regression, there is no closed-form solution for SVM classification (kernelized or not). Implement kernelized Pegasos. Because we are not using a sparse representation for this data, you will probably not see much gain by implementing the “optimized” versions described in the problems above.
34. Find the best hyperparameter settings (including kernel parameters and the regularization parameter) for each of the kernel types. Summarize your results in a table, which gives training error and test error (i.e. average 0/1 loss) for each setting. Include in your table the best settings for each kernel type, as well as nearby settings that show that making small change in any one of the hyperparameters in either direction will cause the performance to get worse. You should use the 0/1 loss on the test set to rank the parameter settings.
35. Plot your best fitting prediction functions using the linear, polynomial, and the RBF kernel. The code provided may help.

## Appendix

Here we are recalling important properties of positive (semi)definite matrices. The exercises below are for revisions for student who may not feel comfortable with these notions. None of the appendix is for credit.

### A Positive Semidefinite Matrices (not for credit)

In statistics and machine learning, we use positive semidefinite matrices a lot. Let's recall some definitions from linear algebra that will be useful here:

**Definition.** A set of vectors  $\{x_1, \dots, x_n\}$  is **orthonormal** if  $\langle x_i, x_i \rangle = 1$  for any  $i \in \{1, \dots, n\}$  (i.e.  $x_i$  has unit norm), and for any  $i, j \in \{1, \dots, n\}$  with  $i \neq j$  we have  $\langle x_i, x_j \rangle = 0$  (i.e.  $x_i$  and  $x_j$  are orthogonal).

Note that if the vectors are column vectors in a Euclidean space, we can write this as  $x_i^T x_j = \mathbb{1}_{i \neq j}$  for all  $i, j \in \{1, \dots, n\}$ .

**Definition.** A matrix is **orthogonal** if it is a square matrix with orthonormal columns.

It follows from the definition that if a matrix  $M \in \mathbb{R}^{n \times n}$  is orthogonal, then  $M^T M = I$ , where  $I$  is the  $n \times n$  identity matrix. Thus  $M^T = M^{-1}$ , and so  $MM^T = I$  as well.

**Definition.** A matrix  $M$  is **symmetric** if  $M = M^T$ .

**Definition.** For a square matrix  $M$ , if  $Mv = \lambda v$  for some column vector  $v$  and scalar  $\lambda$ , then  $v$  is called an **eigenvector** of  $M$  and  $\lambda$  is the corresponding **eigenvalue**.

**Theorem.** [Spectral Theorem] A real, symmetric matrix  $M \in \mathbb{R}^{n \times n}$  can be diagonalized as  $M = Q\Sigma Q^T$ , where  $Q \in \mathbb{R}^{n \times n}$  is an orthogonal matrix whose columns are a set of orthonormal eigenvectors of  $M$ , and  $\Sigma$  is a diagonal matrix of the corresponding eigenvalues.

**Definition.** A real, symmetric matrix  $M \in \mathbb{R}^{n \times n}$  is **positive semidefinite (psd)** if for any  $x \in \mathbb{R}^n$ ,

$$x^T M x \geq 0.$$

Note that unless otherwise specified, when a matrix is described as positive semidefinite, we are implicitly assuming it is real and symmetric (or complex and Hermitian in certain contexts, though not here).

As an exercise in matrix multiplication, note that for any matrix  $A$  with columns  $a_1, \dots, a_d$ , that is

$$A = \begin{pmatrix} | & & | \\ a_1 & \cdots & a_d \\ | & & | \end{pmatrix} \in \mathbb{R}^{n \times d},$$

we have

$$A^T M A = \begin{pmatrix} a_1^T M a_1 & a_1^T M a_2 & \cdots & a_1^T M a_d \\ a_2^T M a_1 & a_2^T M a_2 & \cdots & a_2^T M a_d \\ \vdots & \vdots & \cdots & \vdots \\ a_d^T M a_1 & a_d^T M a_2 & \cdots & a_d^T M a_d \end{pmatrix}.$$

So  $M$  is psd if and only if for any  $A \in \mathbb{R}^{n \times d}$ , we have  $\text{diag}(A^T M A) = (a_1^T M a_1, \dots, a_d^T M a_d)^T \succeq 0$ , where  $\succeq$  is elementwise inequality, and  $0$  is a  $d \times 1$  column vector of 0's.

1. Use the definition of a psd matrix and the spectral theorem to show that all eigenvalues of a positive semidefinite matrix  $M$  are non-negative. [Hint: By Spectral theorem,  $\Sigma = Q^T M Q$  for some  $Q$ . What if you take  $A = Q$  in the “exercise in matrix multiplication” described above?]
2. In this problem, we show that a psd matrix is a matrix version of a non-negative scalar, in that they both have a “square root”. Show that a symmetric matrix  $M$  can be expressed as  $M = BB^T$  for some matrix  $B$ , if and only if  $M$  is psd. [Hint: To show  $M = BB^T$  implies  $M$  is psd, use the fact that for any vector  $v$ ,  $v^T v \geq 0$ . To show that  $M$  psd implies  $M = BB^T$  for some  $B$ , use the Spectral Theorem.]

## B Positive Definite Matrices (not for credit)

**Definition.** A real, symmetric matrix  $M \in \mathbb{R}^{n \times n}$  is **positive definite (spd)** if for any  $x \in \mathbb{R}^n$

with  $x \neq 0$ ,

$$x^T M x > 0.$$

1. Show that all eigenvalues of a symmetric positive definite matrix are positive. [Hint: You can use the same method as you used for psd matrices above.]
2. Let  $M$  be a symmetric positive definite matrix. By the spectral theorem,  $M = Q\Sigma Q^T$ , where  $\Sigma$  is a diagonal matrix of the eigenvalues of  $M$ . By the previous problem, all diagonal entries of  $\Sigma$  are positive. If  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ , then  $\Sigma^{-1} = \text{diag}(\sigma_1^{-1}, \dots, \sigma_n^{-1})$ . Show that the matrix  $Q\Sigma^{-1}Q^T$  is the inverse of  $M$ .
3. Since positive semidefinite matrices may have eigenvalues that are zero, we see by the previous problem that not all psd matrices are invertible. Show that if  $M$  is a psd matrix and  $I$  is the identity matrix, then  $M + \lambda I$  is symmetric positive definite for any  $\lambda > 0$ , and give an expression for the inverse of  $M + \lambda I$ .
4. Let  $M$  and  $N$  be symmetric matrices, with  $M$  positive semidefinite and  $N$  positive definite. Use the definitions of psd and spd to show that  $M + N$  is symmetric positive definite. Thus  $M + N$  is invertible. (Hint: For any  $x \neq 0$ , show that  $x^T(M + N)x > 0$ . Also note that  $x^T(M + N)x = x^T M x + x^T N x$ .)