

CIS 520, Machine Learning, Fall 2014: Assignment 6A

Due: Friday, November 7th, 11:59pm

[54 points]

Collaborators:

Type Collaborator Name Here

Type Collaborator Name Here

Type Collaborator Name Here

Instructions. Please write up your responses to the following problems clearly and concisely. We encourage you to write up your responses using \LaTeX ; we have provided a \LaTeX template, available on Canvas, to make this easier. **Submit your answers in PDF form to Canvas. We will not accept handwritten or paper copies of the homework**

Collaboration. You are allowed and encouraged to work together. You may discuss the homework to understand the problem and reach a solution in groups up to size **four students**. However, *each student must write down the solution independently, and without referring to written notes from the joint session. In addition, each student must write on the problem set the names of the people with whom you collaborated.* You must understand the solution well enough in order to reconstruct it by yourself. (This is for your own benefit: you have to take the exams alone.)

1 Mind Refreshing SVM [40 points]

1. **Finite Features:** You are given the data set D shown in in Figure 1 with data from a single feature X_1 in \mathbb{R}^1 and corresponding label $Y \in \{+, -\}$. The data set contains three positive examples at $X_1 = \{-3, -2, 3\}$ and three negative examples at $X_1 = \{-1, 0, 1\}$.

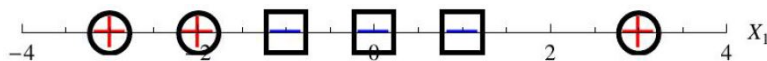


Figure 1: Dataset for SVM feature map task.

This data set, *in its current feature space*, cannot be perfectly separated using a linear separator. The data changes class twice along only one dimension; a linear classifier in one dimension can only represent a single split. In this problem, we'll investigate mapping the data to another feature space.

- (a) **[2 points]** Lets define the simple feature map $\phi(u) = (u, u^2)$ which transforms points in \mathbb{R}^1 to points in \mathbb{R}^2 . Apply ϕ to the data and plot the points in the new \mathbb{R}^2 feature space.
- (b) **[2 points]** Can a linear separator perfectly separate the points in the new \mathbb{R}^2 features space induced by ϕ ? If so, draw a linear separator that works. If not, explain why not.
- (c) **[4 points]** Give the analytic form of the kernel that corresponds to the feature map ϕ in terms of only X_1 and X'_1 . Specifically define $k(X_1, X'_1)$.

- (d) **[7 points]** Construct a maximum-margin separating hyperplane. This hyperplane will be a line in \mathbb{R}^2 , which can be parameterized by its normal equation, i.e. $w_1 Y_1 + w_2 Y_2 + c = 0$ for appropriate choices of w_1, w_2, c . Here, $(Y_1, Y_2) = \phi(X_1)$ is the result of applying the feature map ϕ to the original feature X_1 . Give the values for w_1, w_2, c . Also, explicitly compute the margin for your hyperplane. You do not need to solve a quadratic program to find the maximum margin hyperplane. Instead, let your geometric intuition guide you.
- (e) **[4 points]** On the plot of the transformed points (from part 3), plot the separating hyperplane and the margin, and circle the support vectors.
- (f) **[2 points]** Draw the decision boundary separating of the separating hyperplane, in the original \mathbb{R}^1 feature space.
- (g) **[4 points]** Compute the coefficients α and the constant b in Equation 1 for the kernel k and the support vectors $SV = \{u_1, u_2\}$ you chose in part 6. Be sure to explain how you obtained these coefficients.

$$y(x) = \text{sign} \left(\sum_{i=1}^{|SV|} \alpha_i y_i k(x, u_i) + b \right) \quad (1)$$

Think about the dual form of the quadratic program and the constraints placed on the α values.

- (h) **[2 points]** If we add another positive ($Y = +$) point to the training set at $X_1 = 5$ would the hyperplane or margin change? Why or why not?

2. **Infinite Features Spaces and Kernel Magic:** Lets define a new (infinitely) more complicated feature transformation $\phi_n : \mathbb{R}^1 \rightarrow \mathbb{R}^n$ given in Equation 2.

$$\phi_n(x) = \left\{ e^{-x^2/2}, e^{-x^2/2}x, \frac{e^{-x^2/2}x^2}{\sqrt{2}}, \dots, \frac{e^{-x^2/2}x^i}{\sqrt{i!}}, \dots, \frac{e^{-x^2/2}x^n}{\sqrt{n!}} \right\} \quad (2)$$

Suppose we let $n \rightarrow \infty$ and define new feature transformation in Equation 3. You can think of this feature transformation as taking some finite feature vector and producing an infinite dimensional feature vector rather than the simple two dimensional feature vector used in the earlier part of this problem.

$$\phi_\infty(x) = \left\{ e^{-x^2/2}, e^{-x^2/2}x, \frac{e^{-x^2/2}x^2}{\sqrt{2}}, \dots, \frac{e^{-x^2/2}x^i}{\sqrt{i!}}, \dots \right\} \quad (3)$$

- (a) **[3 points]** Assuming a *consistent* dataset, is there a finite set of points that cannot be linearly separated in this feature space? If so, give an example dataset. If not, explain why not. Note: A consistent dataset is one where there aren't any points $\mathbf{x}_i = \mathbf{x}_j$ where $y_i \neq y_j$. That is, no point can be labeled multiple different things.
- (b) **[4 points]** We know that we can express a linear classifier using only inner products of support vectors in the transformed feature space as seen in Equation 1. It would be great if we could somehow use the feature space obtained by the feature transformation ϕ_∞ . However, to do this we must be able to compute the inner product of examples in this infinite vector space. Lets define the inner product between two infinite vectors $a = \{a_1, \dots, a_i, \dots\}$ and $b = \{b_1, \dots, b_i, \dots\}$ as the infinite sum given in Equation 4.

$$k(a, b) = a \cdot b = \sum_{i=1}^{\infty} a_i b_i \quad (4)$$

We cannot explicitly compute $k(a, b)$ in this form since it contains an infinite sum. However, we can re-write $k(a, b)$ in a form that is efficiently computable. Derive the efficiently computable form. *Hint:* You may want to use the Taylor series expansion of e^x which is given in Equation 5.

$$e^x = \lim_{n \rightarrow \infty} \sum_{i=0}^n \frac{x^i}{i!} \quad (5)$$

- (c) **[6 points]** Prove or disprove each claim below.
- Suppose we translate the inputs $x'_i = x_i + x_0$ for some arbitrary x_0 before using the kernel above in an SVM. Will my kernel function change? (i.e., does $k(x_i, x_j) = k(x'_i, x'_j)$)?
 - What if we negated all the inputs $x'_i = -x_i$?
 - What about rescaling $x'_i = ax_i$ for some positive scalar a ?

2 Eigenvectors

- [7 points]** Show that if X has rank p (all its columns are linearly independent) and $n > p$ then using the p -dimensional pseudo-inverse X^+ in $\hat{w} = X^+y$ solves the least squares problem $\hat{w} = \operatorname{argmin} \sum_i (y - Xw)^2$.
- [7 points]** We want to efficiently find the largest eigenvectors of the matrix $X^T X$ where X is $n \times p$ with $p \gg n$. Show how to do this using the largest eigenvectors of XX^T