Massachusetts Institute of Technology

Department of Electrical Engineering and Computer Science

6.S064 Introduction to Machine Learning

**Midterm exam (March 21, 2013)**

**Your name & ID:** ⸻⸻⸻⸻

- This is a closed book exam

- You do not need nor are permitted to use calculators

- The value of each question – number of points awarded for full credit – is shown in parenthesis

- The problems are not necessarily in any order of difficulty. We recommend that you read through all the problems first, then do the problems in whatever order suits you best.

- Record all your answers in the places provided

| Problem 1.1 | Problem 2.1 | Problem 2.2 | Problem 3.1 | Problem 4.1 | Total |
|:-----------:|:-----------:|:-----------:|:-----------:|:-----------:|:-----:|
| 12 | 12 | 10 | 14 | 9 | 57 |
| | | | | | |

**(1.1)**   The perceptron algorithm remains surprisingly popular for a 50+ year old (algorithm). Hard to find a simpler algorithm for training linear or non-linear classifiers. Worth knowing this algorithm by heart, yes?

(a) **(6 points)** Let $\theta = 0$ (vector) initially. We run the perceptron algorithm to train $y = \text{sign}(\theta \cdot x)$, where $x \in \mathcal{R}^2$, using the three labeled 2-dimensional points in the figure below. In each figure, approximately draw both $\theta$ and the decision boundary *after* updating the parameters (if needed) based on the corresponding point.



(b) **(6 points)** We were wondering what would happen if we normalized all the input examples. In other words, instead of running the algorithm using $x$, we would run it with feature vectors $\phi(x) = x/\|x\|$. Let's explore this with linear classifiers that include offset, i.e., we use $y = \text{sign}(\theta \cdot x + \theta_0)$ or $y = \text{sign}(\theta \cdot \phi(x) + \theta_0)$ after feature mapping. In the figure below, label *all the points* such that 1) the perceptron algorithm wouldn't converge if they are given as original points, 2) the algorithm would converge if run with $\phi(x) = x/\|x\|$.



2

**(2.1)** Support vector machines (SVMs) are so popular that we decided to try them out a bit in a simple setting where $x \in \mathcal{R}$. We used both linear and non-linear SVMs:

$$\min w^2/2 \quad \text{subject to} \quad y^{(t)}(wx^{(t)} + w_0) \geq 1, \quad t = 1, \ldots, n \tag{1}$$
$$\min \|\theta\|^2/2 \quad \text{subject to} \quad y^{(t)}(\theta \cdot \phi(x^{(t)}) + \theta_0) \geq 1, \quad t = 1, \ldots, n \tag{2}$$

where $\phi(x)$ is a feature vector constructed from the real valued input $x$. We wish to compare the resulting classifiers when $\phi(x) = [x, x^2]^T$.

(a) **(3 points)** Provide three input points $x^{(1)}$, $x^{(2)}$, and $x^{(3)}$, where $x^{(i)} \in [0, 4]$, and their associated $\pm 1$ labels such that 1) they cannot be separated with the linear classifier, but 2) are separable by the non-linear classifier with $\phi(x) = [x, x^2]^T$. You may find Figure 2.1 helpful in answering this question.
$(x^{(1)}, y^{(1)}) = ( \quad , \quad )$
$(x^{(2)}, y^{(2)}) = ( \quad , \quad )$
$(x^{(3)}, y^{(3)}) = ( \quad , \quad )$

(b) **(4 points)** Map your three labeled points as labeled feature vectors in Figure 2.1. Approximately draw the resulting *decision boundary* in the feature space of the non-linear SVM classifier with $\phi(x) = [x, x^2]^T$.
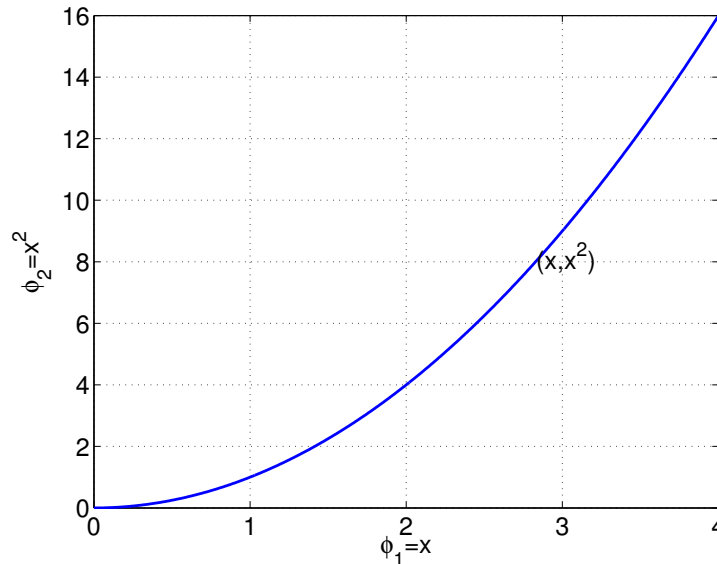


Figure 2.1. Feature space.

3

(c) **(3 points)** Consider two labeled points $(x = 1, y = 1)$ and $(x = 3, y = -1)$. Is the resulting geometric margin we attain in the feature space using feature vectors $\phi(x) = [x, x^2]^T$

( ) greater, ( ) equal, or ( ) smaller

than the geometric margin resulting from using the input $x$ directly?

(d) **(2 points)** In general, is the geometric margin we would attain using scaled feature vectors $\phi(x) = 2[x, x^2]^T$

( ) greater, ( ) equal, ( ) smaller, or ( ) could be any of these

in comparison to the geometric margin resulting from using $\phi(x) = [x, x^2]^T$?

**(2.2)** **(10 points)** We were given a dataset with $n = 10$ labeled two-dimensional examples. Having learned about the support vector machine, we wanted to try it out on this data. However, we were unsure which kernel function to use. So we asked students in the class for possible kernel functions, and tried each one of them:

$$(1)\ K(x, x') = x \cdot x'$$
$$(2)\ K(x, x') = (x \cdot x') + 2(x \cdot x')^2$$
$$(3)\ K(x, x') = 1 - 2(x \cdot x')^2$$
$$(4)\ K(x, x') = (1 + x \cdot x')^5$$

(a) For any kernel $K(x, x') = \phi(x) \cdot \phi(x')$. What can you say about $K(x, x)$ in general?
( ) $K(x, x)$ is never 0, ( ) $K(x, x) < 0$ never, ( ) $K(x, x) > 1$ always

(b) The optimization routine for the dual SVM problem complained about one of the kernels, and it couldn't be used further. Which one (1-4)? ( )

(c) Only one of the kernels resulted in zero training error. Which one (1-4)? ( )

(d) As far as the training error is concerned, one of the kernels was by far the worst. Which one (1-4)? ( )

(e) After training the three classifiers (one we couldn't use), we obtained additional data to assess their test errors. Which one would you expect to have the lowest test error (1-4)? ( )

**(3.1)** We are faced with a content filtering problem where the idea is to rank new songs by trying to predict how they might be rated by a particular user. Each song $x$ is represented by a feature vector $\phi(x)$ whose coordinates capture specific acoustical properties. The ratings are binary valued $y \in \{0, 1\}$ ("need earplugs" or "more like this"). Given $n$ already rated songs, we decided to use regularized linear regression to predict the binary ratings. The training criterion is

$$J(\theta) = \frac{\lambda}{2}\|\theta\|^2 + \frac{1}{n}\sum_{t=1}^{n}(y^{(t)} - \theta \cdot \phi(x^{(t)}))^2/2 \tag{3}$$

(a) **(8 points)** Let $\hat{\theta}$ be the optimal setting of the parameters with respect to the above criterion, which of the following conditions must be true (check all that apply)

( ) $\lambda\hat{\theta} - \frac{1}{n}\sum_{t=1}^{n}(y^{(t)} - \hat{\theta} \cdot \phi(x^{(t)}))\phi(x^{(t)}) = 0$

( ) $J(\hat{\theta}) \geq J(\theta)$, for all $\theta \in \mathcal{R}^d$

( ) If we increase $\lambda$, the resulting $\|\hat{\theta}\|$ will decrease

( ) If we add features to $\phi(x)$ (whatever they may be), the resulting squared training error will NOT increase

(b) **(3 points)** Once we have the estimated parameters $\hat{\theta}$, we must decide how to predict ratings for new songs. Note that the possible rating values are 0 or 1. When do we choose rating $y = 1$ for a new song $x$? Please write the corresponding expression.

(c) **(3 points)** If we change $\lambda$, we obtain different $\hat{\theta}$, and therefore different rating predictions according to your rule above. What will happen to your predicted ratings when we increase the regularization parameter $\lambda$?

**(4.1)** Consider the simple K-means algorithm for clustering. Each iteration of the algorithm consists of two steps, assigning points to the centroids, and updating the centroids based on the points assigned to them. We will assume that $k = 2$.

(a) **(2 points)** If we initialize the centroids to be the means of the two well-separated clusters, will the centroids change after the first iteration? (Y/N) ( )

(b) **(3 points)** If we initialize the centroids by drawing a random point from each of the two well-separated clusters, how many iterations does it take for the k-means to converge? ( )

(c) **(4 points)** Consider two spherical (circle-like) clusters of radius $\delta$ as shown below. The clusters are centered at locations $-x$ and $x$. For which values of $x$ would the k-means algorithm fail to find the centers of the two clusters regardless of the initialization?