

COS 511: Theoretical Machine Learning

Homework #5
SVM's and On-line Learning

Due: April 3, 2008

Problem 1

- a. [10] In class, we argued that if a function L satisfies the “minmax property”

$$\min_{\mathbf{w}} \max_{\boldsymbol{\alpha}} L(\mathbf{w}, \boldsymbol{\alpha}) = \max_{\boldsymbol{\alpha}} \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha}), \quad (1)$$

and if $(\mathbf{w}^*, \boldsymbol{\alpha}^*)$ are the desired solutions

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \max_{\boldsymbol{\alpha}} L(\mathbf{w}, \boldsymbol{\alpha}) \quad (2)$$

$$\boldsymbol{\alpha}^* = \arg \max_{\boldsymbol{\alpha}} \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha}), \quad (3)$$

then $(\mathbf{w}^*, \boldsymbol{\alpha}^*)$ is a saddle point:

$$L(\mathbf{w}^*, \boldsymbol{\alpha}^*) = \max_{\boldsymbol{\alpha}} L(\mathbf{w}^*, \boldsymbol{\alpha}) = \min_{\mathbf{w}} L(\mathbf{w}, \boldsymbol{\alpha}^*). \quad (4)$$

(Here, it is understood that \mathbf{w} and $\boldsymbol{\alpha}$ may belong to a restricted space (e.g., $\boldsymbol{\alpha} \geq 0$) which we omit for brevity.)

Prove the converse of what was shown in class. That is, prove that if $(\mathbf{w}^*, \boldsymbol{\alpha}^*)$ satisfies Eq. (4), then Eqs. (1), (2) and (3) are also satisfied. You should not assume anything special about L (such as convexity), but you can assume all of the relevant minima and maxima exist.

- b. [10] Let a_1, \dots, a_n be nonnegative real numbers, not all equal to zero, and let b_1, \dots, b_n and c all be positive real numbers. Use the method of Lagrange multipliers to find the values of x_1, \dots, x_n which minimize

$$-\sum_{i=1}^n a_i \ln x_i$$

subject to the constraint that

$$\sum_{i=1}^n b_i x_i \leq c.$$

Show how this implies that relative entropy is nonnegative.

Problem 2

Suppose we use support-vector machines with the kernel:

$$K(x, u) = \begin{cases} 1 & \text{if } x = u \\ 0 & \text{otherwise.} \end{cases}$$

As we discussed in class, this corresponds to mapping each x to a vector $\psi(x)$ in some high dimensional space (that need not be specified) so that $K(x, u) = \psi(x) \cdot \psi(u)$.

As usual, we are given m examples $(x_1, y_1), \dots, (x_m, y_m)$ where $y_i \in \{-1, +1\}$. Assume for simplicity that all the x_i 's are distinct (i.e., $x_i \neq x_j$ for $i \neq j$).

- a. [10] Recall that the weight vector \mathbf{w} used in SVM's has the form

$$\mathbf{w} = \sum_i \alpha_i y_i \psi(x_i).$$

Compute the α_i 's explicitly that would be found using SVM's with this kernel.

- b. [6] Recall that the SVM algorithm outputs a classifier that, on input x , computes the sign of $\mathbf{w} \cdot \psi(x)$. What is the value of this inner product on training example x_i ? What is the value of this inner product on any example x not seen during training? Based on these answers, what kind of generalization error do you expect will be achieved by SVM's using this kernel?
- c. [6] Recall that the generalization error of SVM's can be bounded using the margin δ (which is equal to $1/\|\mathbf{w}\|$), or using the number of support vectors. What is δ in this case? How many support vectors are there in this case? How are these answers consistent with your answer in part (b)?

Problem 3

Consider the problem of learning with expert advice when one of the experts gives perfect predictions. On some round t , let q be the fraction of surviving experts that predict 1. (A surviving expert is one that has not made any mistakes so far.) In class, we talked about the halving algorithm which predicts with the majority vote of the expert predictions, and we talked about the randomized weighted majority algorithm (with β set to zero) which predicts with one randomly selected expert.

In general, we can predict 1 with probability $F(q)$ and 0 with probability $1 - F(q)$ for some function F . For instance, for the halving algorithm, $F(q)$ is 1 if $q > 1/2$ and 0 if $q < 1/2$ (and arbitrary if $q = 1/2$). For the randomized weighted majority algorithm (again, with $\beta = 0$), $F(q) = q$.

Consider now a function $F : [0, 1] \rightarrow [0, 1]$ satisfying the following property:

$$1 + \frac{\lg q}{2} \leq F(q) \leq -\frac{\lg(1 - q)}{2}. \quad (5)$$

- a. [15] Suppose we run an on-line learning algorithm that uses a function F satisfying (5) as described above. Show that the expected number of mistakes made by the learning algorithm is at most $(\lg N)/2$, where N is the number of experts.
- b. [10] Show that the function

$$F(q) = \frac{\lg(1 - q)}{\lg q + \lg(1 - q)}$$

has range $[0, 1]$ and satisfies (5). (At the endpoints, we define $F(0) = 0$ and $F(1) = 1$ to make F continuous, but you *don't* need to worry about these.)

- c. [15] **(Optional)** Suppose now that there are $k \geq 2$ possible outcomes rather than just 2. In other words, the outcome y_t is now in the set $\{1, \dots, k\}$ (rather than $\{0, 1\}$ as we have considered up until now), and likewise, both experts and the learning algorithm make predictions in this set. Assume one of the experts makes perfect predictions. On some round t , let q_j be the fraction of surviving experts predicting outcome $j \in \{1, \dots, k\}$. Suppose that the learning algorithm predicts each outcome j with probability

$$\frac{\lg(1 - q_j)}{\sum_{i=1}^k \lg(1 - q_i)}.$$

Show that the expected number of mistakes of this learning algorithm is at most $(\lg N)/2$.