

COS 511: Theoretical Machine Learning

Homework #6
Kernels and online learning

Due:
April 10, 2014

Problem 1

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 2

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:

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

- a. [15] Suppose we run an on-line learning algorithm that uses a function F satisfying (1) 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 (1). (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. [10] (**Optional – for extra credit**) 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 $(\lg N)/2$.

Problem 3

[15] For this problem, let us suppose that labels, outcomes, expert/hypothesis predictions, etc. are all defined over the set $\{-1, +1\}$ rather than $\{0, 1\}$. Since this does not change what it means for the learner or an expert to make a mistake, this has no effect on any of the results we have discussed regarding online mistake bounds.

Let \mathcal{H} be a finite space of hypotheses $h : \mathcal{X} \rightarrow \{-1, +1\}$, and let $S = \langle x_1, \dots, x_m \rangle$ be any sequence of m distinct points in \mathcal{X} . Prove that the empirical Rademacher complexity of \mathcal{H} satisfies

$$\hat{\mathcal{R}}_S(\mathcal{H}) \leq O\left(\sqrt{\frac{\ln |\mathcal{H}|}{m}}\right)$$

by applying our analysis of online algorithms for learning with expert advice to an appropriately constructed sequence of expert prediction ξ_i and outcomes y . Give a bound with explicit constants.

Note that this bound was earlier stated without proof in class (see page 5 in the scribe notes for lecture #10), and is also a special case of Theorem 3.3 in the Mohri et al. book, although with possibly weaker constants. **Extra credit** [5] will be given for obtaining a bound of $\sqrt{(2 \ln |\mathcal{H}|)/m}$, that is, with the constant that was stated in class.