

COS 511: Theoretical Machine Learning

Homework #5
Boosting & SVM's

Due:
March 28, 2013

Problem 1

[10] Suppose, in the usual boosting set-up, that the weak learning condition is guaranteed to hold so that $\epsilon_t \leq \frac{1}{2} - \gamma$ for some $\gamma > 0$ which is not known before boosting begins. And suppose AdaBoost is run in the usual fashion, except that the algorithm is modified to halt and output the combined classifier H immediately following the first round on which it is consistent with all of the training examples (so that its training error is zero). Assume that the weak hypotheses are selected from a class of VC-dimension d . Prove that, with probability at least $1 - \delta$, the generalization error of the output combined classifier H is at most

$$\tilde{O}\left(\frac{(d/\gamma^2) + \ln(1/\delta)}{m}\right).$$

Give a bound in which all constants and log terms have been filled in explicitly.

Problem 2

Suppose AdaBoost is run for an unterminating number of rounds. In addition to our usual notation, we define for each $T \geq 1$:

$$F_T(x) = \sum_{t=1}^T \alpha_t h_t(x) \quad \text{and} \quad s_T = \sum_{t=1}^T \alpha_t.$$

Recall that each $\alpha_t \geq 0$ (since $\epsilon_t \leq \frac{1}{2}$). The *minimum margin* on round t , denoted θ_t , is the smallest margin of any training example; thus,

$$\theta_t = \min_i \frac{y_i F_t(x_i)}{s_t}.$$

Finally, we define the *smooth margin* on round t to be

$$g_t = \frac{-\ln\left(\frac{1}{m} \sum_{i=1}^m e^{-y_i F_t(x_i)}\right)}{s_t}.$$

- a. [10] Prove that

$$\theta_t \leq g_t \leq \theta_t + \frac{\ln m}{s_t}.$$

Thus, if s_t gets large, then g_t gets very close to θ_t .

- b. [10] Let us define the continuous function

$$\Upsilon(\gamma) = \frac{-\ln(1 - 4\gamma^2)}{\ln\left(\frac{1+2\gamma}{1-2\gamma}\right)}.$$

It is a fact (which you do not need to prove) that $\gamma \leq \Upsilon(\gamma) \leq 2\gamma$ for $0 \leq \gamma \leq \frac{1}{2}$.

Prove that g_T is a weighted average of the values $\Upsilon(\gamma_t)$, specifically,

$$g_T = \frac{\sum_{t=1}^T \alpha_t \Upsilon(\gamma_t)}{s_T}.$$

- c. [10] Prove that if the edges γ_t converge (as $t \rightarrow \infty$) to some value γ , where $0 < \gamma < \frac{1}{2}$, then the minimum margins θ_t converge to $\Upsilon(\gamma)$.

Problem 3

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 $\boldsymbol{\psi}(x)$ in some high dimensional space (that need not be specified) so that $K(x, u) = \boldsymbol{\psi}(x) \cdot \boldsymbol{\psi}(u)$.

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

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

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 $\mathbf{w} \cdot \boldsymbol{\psi}(x)$. 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 4 – Optional (Extra Credit)

[15] In class (as well as on Problem 1 of this homework), we showed how a weak learning algorithm that uses hypotheses from a space \mathcal{H} of bounded VC-dimension can be converted into a strong learning algorithm. However, strictly speaking, the definition of weak learnability does *not* include such a restriction on the weak hypothesis space. The purpose of this problem is to show that weak and strong learnability are equivalent, even without these restrictions.

Let \mathcal{C} be a concept class on domain X . Let A_0 be a weak learning algorithm and let $\gamma > 0$ be a (known) constant such that for every concept $c \in \mathcal{C}$ and for every distribution D on X , when given m_0 random examples x_i from D , each with its label $c(x_i)$, A_0 outputs a hypothesis h such that, with probability at least $1/2$,

$$\Pr_{x \in D} [h(x) \neq c(x)] \leq \frac{1}{2} - \gamma.$$

Here, for simplicity, we have “hard-wired” the usual parameter δ to the constant $1/2$ so that A_0 takes a fixed number of examples and only needs to succeed with fixed probability $1/2$. Note that no restrictions are made on the form of hypothesis h used by A_0 , nor on the cardinality or VC-dimension of the space from which it is chosen. For this problem, assume that A_0 is a deterministic algorithm.

Show that A_0 can be converted into a strong learning algorithm using boosting. That is, construct an algorithm A such that, for $\epsilon > 0$, $\delta > 0$, for every concept $c \in \mathcal{C}$ and for every distribution D on X , when given $m = \text{poly}(m_0, 1/\epsilon, 1/\delta, 1/\gamma)$ random examples x_i from D , each with its label $c(x_i)$, A outputs a hypothesis H such that, with probability at least $1 - \delta$,

$$\Pr_{x \in D} [H(x) \neq c(x)] \leq \epsilon.$$

Be sure to show that the number of examples needed by this algorithm is polynomial in m_0 , $1/\epsilon$, $1/\delta$ and $1/\gamma$.