
Theoretical Foundations of Machine Learning - Homework #4

Jacob Abernethy and Chansoo Lee

Due: 11/20/2015

Homework Policy: Working in groups is fine, but *every student* must submit their own writeup. Please write the members of your group on your solutions. There is no strict limit to the size of the group but we may find it a bit suspicious if there are more than 4 to a team. Questions labelled with **(Challenge)** are not strictly required, but you'll get some participation credit if you have something interesting to add, even if it's only a partial answer.

1) **Solving Feasibility Problems with Perceptron.** We can define a special class of Linear Programming Feasibility (LFP) problems as follows. Given a set of m vectors $\mathbf{a}^1, \dots, \mathbf{a}^m \in \mathbb{R}^d$, we want to find a *feasible solution*, which is a vector $\mathbf{y} \in \mathbb{R}^d$ such that $\mathbf{a}^i \cdot \mathbf{y} > 0$ for all $i = 1, \dots, m$. For any vector \mathbf{x} , we can define $\nu(\mathbf{x}) := \min_i \frac{\mathbf{a}^i \cdot \mathbf{x}}{\|\mathbf{a}^i\|_2 \|\mathbf{x}\|_2}$, which is positive if \mathbf{x} is a feasible solution. Define the *wiggle room* of the feasibility problem as $\nu^* := \sup_{\mathbf{x} \in \mathbb{R}^d} \nu(\mathbf{x})$.

Assume we are given an LFP defined by $\mathbf{a}^1, \dots, \mathbf{a}^m \in \mathbb{R}^d$ with a positive wiggle room $\nu^* > 0$. Using Perceptron as a subroutine, give an algorithm that finds a feasible solution requiring no more than $(\frac{1}{\nu^*})^2$ updates to Perceptron. Prove that your solution works.

2) **Tuning Parameters.** We are going to imagine we have some algorithm \mathcal{A} with a performance bound that depends on some input values (which can not be adjusted) and some tuning parameters (which can be optimized). We will use greek letters (α, η, ζ , etc.) for the tuning parameters and capital letters (T, D, N , etc.) for inputs. We would like the bound to be the tightest possible, up to multiplicative constants. For each of the following, tune the parameters to obtain the optimal bound. Using *big-Oh* notation is fine to hide constants, but you *must not* ignore the dependence on the input parameters. For example, assuming $M, T > 0$, imagine we have a performance guarantee of the form:

$$\text{Performance}(\mathcal{A}; M, T, \epsilon) \leq M\epsilon + \frac{T}{\epsilon} \tag{1}$$

and we know $\epsilon > 0$. Then by optimizing the above expression with respect to the free parameter we can set $\epsilon = \sqrt{\frac{T}{M}}$. With this value we obtain $\text{Performance}(\mathcal{A}; M, T, \epsilon) = O(\sqrt{MT})$

NOTE: I didn't have to make up this problem, I actually pulled all the bounds below from different papers!

(a) $\text{Performance}(\mathcal{A}; T, N, \eta) \leq \frac{\log N + \eta T}{1 - \exp(-\eta)}$

(b) $\text{Performance}(\mathcal{A}; T, \eta) \leq \max(T\eta, \eta^{-2})$

(c) $\text{Performance}(\mathcal{A}; T, \eta) \leq \frac{T}{\eta} + \exp(\eta)$. (Note: you needn't obtain the optimal choice of η here or the tightest possible bound, but try to tune in order to get a bound that is $o(T)$ – i.e. the bound should grow strictly slower than linear in T .)

(d) $\text{Performance}(\mathcal{A}; N, T, \eta, \epsilon) \leq \frac{T\epsilon}{\eta} + \frac{N}{\epsilon} + T\eta$

3) **The Doubling Trick.** Let's imagine that an online learning algorithm that runs in T rounds has the bound in Eq. (1). By optimally tuning ϵ we obtain a bound of the form $O(\sqrt{TM})$. The problem with this approach is that it requires us to know T in advance. Is there a way around this issue?

Imagine constructing a modified algorithm \mathcal{A}' that does the following iterative procedure. \mathcal{A}' starts with an initial parameter ϵ_1 , implements algorithm \mathcal{A} using this parameter for T_1 rounds, then adjusts the parameter to ϵ_2 , and runs \mathcal{A} for T_2 rounds, and so forth. Let's say ϵ gets updated k times, where $T_1 + T_2 + \dots + T_k = T$. You can also assume that $\text{Performance}(\mathcal{A}'; T, M) = \sum_{i=1}^k \text{Performance}(\mathcal{A}; T_i, M, \epsilon_i)$.

Can you construct a schedule for the sequence of (ϵ_i, T_i) that achieves the same bound (up to a multiplicative constant factor) as the optimally tuned bound (namely, $O(\sqrt{MT})$), even though T is unknown in advance? In other words, you want to choose the sequence of sequence of T_1, T_2, \dots , with the associated $\epsilon_1, \epsilon_2, \dots$ so that whenever the game truly ends, at a previously unknown T , the bound $\text{Performance}(\mathcal{A}'; T, M) = O(\sqrt{MT})$ will always hold. (Note: this is referred to as a “doubling trick”.)

4) **Exponential Weights Algorithm with a Prior.** The Exponential Weights Algorithm had the initial weights w_1^1, \dots, w_n^1 all set to 1. What if instead we imagine an algorithm \mathcal{A}' where we set these weights according to $w_i^1 = p_i$ where \vec{p} is some distribution (i.e. $p_i \geq 0$ for all i and $\sum_i p_i = 1$). We will do the same multiplicative update rule as before.

Prove that with this modified algorithm we achieve the following bound: for any expert i we have that

$$\text{Loss}_T(\mathcal{A}) \leq \frac{\log p_i^{-1} + \eta \text{Loss}_T(\text{expert } i)}{1 - e^{-\eta}}$$

5) **Subsets as Experts.** We saw that when we wanted to do “predictive sorting” it's not easy to apply the halving algorithm to the class of all permutations (rankings) as experts. But this isn't the case for all classes of “complex experts”.

Imagine a setting where we have N experts but our goal is not to choose one but $k < N$ of them on each round! We can imagine having a “hyperexpert” for each subset $S \subset [N]$, with $|S| = k$, of which there are clearly $\binom{N}{k}$. Let \mathcal{S}_k^N be the set of all k -sized subsets of $[N]$. On round t , each “base” expert i suffers loss ℓ_i^t which implies that the hyperexpert S suffers loss

$$\text{loss}_t(S) := \sum_{i \in S} \ell_i^t,$$

that is, the hyperexpert loss is additive across the base experts in the subset. Our aim is to run the Exponential Weights Algorithm on these subset hyperexperts. With this well-defined loss on each hyperexpert and a given parameter η , we can define the weight update $w_S^{t+1} = w_S^t \exp(-\eta \text{loss}_t(S))$.

In many scenarios in which we are dealing with hyperexperts, it's suitable to compute the “projected” weights for each base expert i . That is, assume we can run our algorithm by simply knowing $u_i^t := \sum_{S \in \mathcal{S}_k^N: i \in S} w_S^t$. In other words, if our weights define a distribution over \mathcal{S}_k^N , then the value u_i^t corresponds to the (unnormalized) marginal probability of i being in a randomly drawn subset. Can we obtain these values efficiently? If so, how?

Hint: Given a vector of positive values $v_{1:n} := \langle v_1, \dots, v_n \rangle$, we can define

$$\text{SumProd}(v_{1:n}, k) := \sum_{S \in \mathcal{S}_k^n} \prod_{i \in S} v_i$$

Naively this requires $O(n^k)$ computation, but perhaps there is a faster way to compute SumProd?

6) **Reduction to the Halving Algorithm. (Challenge)** In the 0-noise expert setting, where there is one expert that makes no mistakes, we showed that a simple strategy, the *Halving algorithm*, is guaranteed to make no more than $\log N$ mistakes. But what if we know that the best expert will make some errors, but definitely no more than k mistakes? Try to find a master algorithm, via a reduction to the 0-noise case, such that the number of mistakes of the algorithm is never more than

$$\max \left\{ m : m \leq \log_2 N + \log_2 \binom{m}{\leq k} \right\}$$

where $\binom{m}{\leq k}$ is the sum of binomial coefficients $\binom{m}{0} + \dots + \binom{m}{k}$. (*Hint:* You need to find a large set of hyperexperts on which to perform the halving algorithm. You may need to argue that, without loss of generality, the majority of these hyperexperts are incorrect on every round.)