

## PROBLEM SET 1

Due: Tuesday, February 7, 3 p.m. by email

Please title your email "CSCI599\_PS1".

## 1. Relationships between concept classes

Assume each example  $x$  is given by  $n$  boolean variables. A decision list is a function of the form: "if  $\ell_1$  then  $b_1$ , else if  $\ell_2$  then  $b_2$ , else if  $\ell_3$  then  $b_3$ , ..., else  $b_m$ ," where each  $\ell_i$  is a literal (either a variable or its negation) and each  $b_i \in \{0, 1\}$ .

- (a) Show that conjunctions (like  $x_1 \wedge \bar{x}_2 \wedge x_3$ ) and disjunctions (like  $x_1 \vee \bar{x}_2 \vee x_3$ ) are special cases of decision lists.
- (b) Show that decision lists are a special case of linear threshold functions. That is, any function that can be expressed as a decision list can also be expressed as a linear threshold function " $f(x) = +$  iff  $w_1x_1 + \dots + w_nx_n \geq w_0$ ", for some values  $w_0, w_1, \dots, w_n$ .

## 2. VC-dimension of simple concept classes

- (a) *Parity Functions:* For a set  $I \subseteq \{1, 2, \dots, n\}$ , we define a parity function  $h_I : \{0, 1\}^n \rightarrow \{0, 1\}$  as follows: On a binary vector  $x = (x_1, \dots, x_n) \in \{0, 1\}^n$ , we have:

$$h_I(x) = \left( \sum_{i \in I} x_i \right) \bmod 2.$$

(That is,  $h_I$  computes the parity of bits in  $I$ .) What is the VC-dimension of the class of all such parity functions,  $\mathcal{H}_{\text{parity}}^n = \{h_I : I \subseteq \{1, 2, \dots, n\}\}$ ?

- (b) *Axis-aligned rectangles:* For two vectors  $a, b \in \mathbb{R}^n$  with  $a \leq b$  (coordinate-wise), we define an axis-aligned rectangle  $h_{\text{rec}}^{a,b} : \mathbb{R}^n \rightarrow \{0, 1\}$  as  $h_{\text{rec}}^{a,b}(x) = 1$  if  $a \leq x \leq b$  and  $h_{\text{rec}}^{a,b} = 0$  otherwise. Let  $\mathcal{H}_{\text{rec}}^n$  be the class of all axis-aligned rectangles in  $\mathbb{R}^n$ . What is the VC-dimension of the class  $\mathcal{H}_{\text{rec}}^n$ ?
- (c) *Boolean conjunctions:* Let  $\mathcal{H}_{\text{con}}^n$  be the class of Boolean conjunctions over the variables  $x_1, \dots, x_n$ . What is the VC-dimension of the class  $\mathcal{H}_{\text{con}}^n$ ?

3. An online mistake-bound algorithm is said to be conservative if it changes its hypothesis only when a mistake is made.

All of the online mistake-bound algorithms we have seen in class (the elimination algorithm, the decision list algorithm, Winnow, etc.) are conservative. In this problem, you will show that this is not a coincidence:

Let  $C$  be any concept class, and let  $A$  be any online learning algorithm (not necessarily conservative) which has a finite mistake bound  $M$  for  $C$ . Prove that there exists a conservative online learning algorithm  $A'$  for  $C$  which also has mistake bound  $M$ .

4. The Perceptron Convergence Theorem shows that the Perceptron algorithm will not make too many mistakes as long as every example is “far” from the separating hyperplane of the target halfspace. In this problem you will explore a variant of the Perceptron algorithm and show that it performs well (given a little help in the form of a good initial hypothesis) as long as every example is “far” (in terms of angle) from the separating hyperplane of the *current hypothesis*.

Consider the following variant of Perceptron:

- Start with an initial hypothesis vector  $w = w^{\text{init}}$ .
- Given example  $x \in \mathbb{R}^n$ , predict according to the linear threshold function  $w \cdot x \geq 0$ .
- Given the true label of  $x$ , update hypothesis vector  $w$  as follows:
  - If the prediction is correct, leave  $w$  unchanged.
  - If the prediction is incorrect, set  $w \leftarrow w - (w \cdot x)x$ .

So the update step differs from that of Perceptron shown in class in that  $(w \cdot x)x$  (rather than  $x$ ) is added or subtracted to  $w$ . (Note that if  $\|x\|_2 = 1$ , then this update causes vector  $w$  to become orthogonal to  $x$ , i.e., we add or subtract the multiple of  $x$  that shrinks  $w$  as much as possible.)

Suppose that we run this algorithm on a sequence of examples that are labeled according to some linear threshold function  $v \cdot x \geq 0$  for which  $\|v\|_2 = 1$ . Suppose moreover that

- each example vector  $x$  has  $\|x\|_2 = 1$ ;
- the initial hypothesis vector  $w^{\text{init}}$  satisfies  $\|w^{\text{init}}\|_2 = 1$  and  $w^{\text{init}} \cdot v \geq \gamma$  for some fixed  $\gamma > 0$ ;
- each example vector  $x$  satisfies  $\frac{|w \cdot x|}{\|w\|_2} \geq \delta$ , where  $w$  is the current hypothesis vector when  $x$  is received. (Note that for a unit vector  $x$ , this quantity  $\frac{|w \cdot x|}{\|w\|_2}$  is the cosine of the angle between vectors  $w$  and  $x$ .)

Show that under these assumptions the algorithm described above will make at most  $\frac{2}{\delta^2} \ln(1/\gamma)$  many mistakes.

5. Describe a set  $S$  of  $O(n)$  examples over  $\{0, 1\}^n$  that are linearly separable by a hyperplane through the origin, but where the Perceptron algorithm takes exponential time for learning (i.e., time  $2^{\Omega(n)}$ ). Specifically, imagine we repeatedly cycle the online perceptron algorithm through the examples until we have  $w \cdot x \geq 0$  for every positive example  $x \in S$  and we have  $w \cdot x < 0$  for every negative  $x \in S$ . For simplicity, use the version of the Perceptron algorithm (shown in class) that does not normalize examples (i.e., when a mistake is made on a positive example  $x$ , it sets  $w = w + x$  and similarly for mistakes on negatives). Show that your set of examples has the desired property.