

## Problem Set 4

Instructor: Kamalika Chaudhuri

Due on:

**Problem 1: 8 points**

Alice, Bob and Carol have all been asked to implement the perceptron algorithm. They all have the same training and test data, and they make a single pass over the training and test data with all the algorithms (Alice's variant, Bob's variant, and Carol's variant). Carol implements the version of perceptron that we discussed in lecture.

1. Suppose Alice implements the following variant of the perceptron algorithm.

- (a) Initially:  $w_1 = 0$ .
- (b) For  $t = 1, 2, 3, 4, \dots, T$ 
  - i. If  $y_t \langle w_t, x_t \rangle < 0$  then  $w_{t+1} = w_t + y_t x_t$ .
  - ii. Otherwise:  $w_{t+1} = w_t$ .
- (c) Output  $w_{Alice} = w_{T+1}$ .

Is the test error of the classifier output by Alice's algorithm the same as the test error of Carol's algorithm, no matter what the test data is? If your answer is yes, justify your answer. If your answer is no, provide a counterexample or a brief justification.

2. Bob implements a second variant of the perceptron algorithm, as follows.

- (a) Initially:  $w_1 = 0$ .
- (b) For  $t = 1, 2, 3, 4, \dots, T$ 
  - i. If  $y_t \langle w_t, x_t \rangle \leq 0$  then  $w_{t+1} = w_t + y_t x_t / \|x_t\|$ .
  - ii. Otherwise:  $w_{t+1} = w_t$ .
- (c) Output  $w_{Bob} = w_{T+1}$ .

Is the test error of the classifier output by Bob's algorithm the same as the test error of Carol's algorithm, no matter what the test dataset is? If your answer is yes, provide a justification for your answer; if your answer is no, provide a counterexample or a brief justification.

**Solution**

1. No. Since Alice's algorithm only update  $w$  when  $y_t \langle w_t, x_t \rangle < 0$  and  $w$  is initialized as 0, it will keep the same through all iterations and the final  $w_{T+1} = 0$ . So the test error may not be the same as that of Carol's classifier.
2. No. Suppose the data lies in a 2-dimensional Euclidean space, and the underlying classifier is  $w^* = (1, 2)$ , i.e., for all  $x \in \mathbb{R}^2$ , its corresponding label is  $y = \text{sign}(\langle w^*, x \rangle)$ . Suppose our training set is  $\{((0, 2), 1), ((1, 0), 1), ((0, -2), -1), ((-1, 0), -1)\}$  and test set is  $\{((-3, 2), 1)\}$ . Then Carol's algorithm will do the following:

- See  $(x_1, y_1) = ((0, 2), 1)$  and update  $w_2 = (0, 2)$
- See  $(x_2, y_2) = ((1, 0), 1)$  and update  $w_3 = (1, 2)$
- See  $(x_3, y_3) = ((0, -2), -1)$  and keep  $w_4 = (1, 2)$

- See  $(x_4, y_4) = ((-1, 0), -1)$  and keep  $w_{Carol} = w_5 = (1, 2)$

And Bob's algorithm will do the following:

- See  $(x_1, y_1) = ((0, 2), 1)$  and update  $w_2 = (0, 1)$
- See  $(x_2, y_2) = ((1, 0), 1)$  and update  $w_3 = (1, 1)$
- See  $(x_3, y_3) = ((0, -2), -1)$  and keep  $w_4 = (1, 1)$
- See  $(x_4, y_4) = ((-1, 0), -1)$  and keep  $w_{Bob} = w_5 = (1, 1)$

$w_{Carol}$  and  $w_{Bob}$  are different. They classify  $(-3, 2)$  as 1 and  $-1$  respectively, and the test errors are not the same.

## Problem 2: 12 points

Suppose we are given a training data set  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$  where each feature vector  $x_i$  lies in  $d$ -dimensional space. For each  $x_i$ , suppose we transform it to  $z_i$  by rescaling each axis of the data by a fixed factor; that is, for every  $i = 1, \dots, n$  and every coordinate  $j = 1, \dots, d$ , we write:

$$z_i^j = \alpha^j x_i^j$$

Here  $\alpha^j$ s are real, non-zero and positive constants. Thus, our original training set  $S$  is transformed after rescaling to a new data set  $S' = \{(z_1, y_1), \dots, (z_n, y_n)\}$ .

A classifier  $C$  in the original space is said to be equal to a classifier  $C'$  in the rescaled space if for every  $x \in \mathbb{R}^d$ ,  $C(x) = C'(z)$ , where  $z$  is obtained by transforming  $x$ . For example, if  $d = 1$ , and if  $\alpha^1 = 2$ , the classifier  $C$  in the original space:

$C$ : Predict 0 if  $|x| \leq 1$  and 1 otherwise

is equal to the classifier  $C'$  in the rescaled space:

$C'$ : Predict 0 if  $|z| \leq 2$  and 1 otherwise

1. First, suppose that all the  $\alpha^j$  values are equal; that is,  $\alpha^1 = \dots = \alpha^d$ . Suppose we train a  $k$ -NN classifier  $C$  on  $S$  and a  $k$ -NN classifier  $C'$  on  $S'$ . Are these two classifiers equal? What if we trained  $C$  and  $C'$  on  $S$  and  $S'$  respectively using the ID3 Decision Tree algorithm? What if we trained  $C$  and  $C'$  on  $S$  and  $S'$  respectively using the Perceptron algorithm? If the classifiers are equal, provide a *brief* argument to justify why; if they are not equal, provide a counterexample.
2. Repeat your answers to the questions in part (1) when the  $\alpha_i$ s are different. Provide a *brief* justification for each answer if the classifiers are equal, and a counterexample if they are not.
3. From the results of parts(1)and (2), what can you conclude about how  $k$ -NN, decision trees and perceptrons behave under scaling transformations?

## Solution

For this question, we assume that ties are always broken in a consistent manner for both the  $k$ -NN and ID3 decision tree algorithms.

**$k$ -NN.** We will obtain equal  $k$ -NN classifiers before and after a space transform for an arbitrary data set, if and only if, the following *relative distance* condition holds: for any three points  $x$ ,  $x_p$  and  $x_q$  in the original space,  $d(x, x_p) \geq d(x, x_q)$  implies  $d(z, z_p) \geq d(z, z_q)$ , where  $z$ ,  $z_p$  and  $z_q$  are the points after rescaling. In other words, we need to ensure that in all cases, the nearest neighbors of a point in the original space are still the nearest neighbors in the rescaled space.

In the case of a uniform scaling factor (all  $\alpha^j = \alpha$ ), the distance between any two points  $z_1$  and  $z_2$  in the rescaled space is,

$$d(z_1, z_2) = \sqrt{\sum_j (z_1^j - z_2^j)^2} = \sqrt{\sum_j (\alpha x_1^j - \alpha x_2^j)^2} = \alpha \sqrt{\sum_j (x_1^j - x_2^j)^2} = \alpha d(x_1, x_2),$$

This is simply the distance in the original space scaled by a constant  $\alpha$ . Clearly the relative distance condition holds. In particular, this means that the training points that are the nearest neighbors of  $x$  in the original space remain the nearest neighbors of  $z$  in the rescaled space, therefore prediction for  $x$  remains the same as the prediction for  $z$ .

For nonuniform scaling factors, the relative distance condition does not necessarily hold. One extreme example is if  $\alpha^1 = 1, \alpha^2 = 0.0001$  (a very small quantity). The transform now essentially projects each point to the  $x$ -axis (although the point will not be exactly on the axis). Consider the training points  $(1, 0)$  with label 0 and  $(0, 1)$  with label 1, and a test example  $(0.1, 0)$ . In the original space,  $(0.1, 0)$  is closer to  $(1, 0)$  than  $(0, 1)$  and will be assigned label 0; in the rescaled space however, it will be rescaled to be  $(0.1, 0)$ , will be closer to  $(0, 1)$  (now rescaled as  $(0, 0.001)$ ), and thus will be assigned label 1 by the 1-NN classifier. Therefore in this case, we are not guaranteed to get the same  $k$ -NN classifier.

**ID3 Decision Tree.** The decision trees produced by the ID3 algorithm will be equal in both cases, assuming that ties are broken in a consistent manner. We can show this by induction. In what follows, we will say that a splitting rule  $(j, t)$  in the original space *is equal to* a splitting rule  $(j, \alpha^j t)$  in the rescaled space.

We run the ID3 algorithm on  $S$  and  $S'$  simultaneously, and maintain the following invariants at each step of the algorithm. If  $T$  and  $T'$  are the trees built based on  $S$  and  $S'$  respectively, then, (a)  $T$  and  $T'$  have the same structure, (b) for each internal node  $v$  in  $T$ , the splitting rule at  $v$  is equal to the splitting rule at the corresponding internal node  $v'$  in  $T'$  and (c) if  $D$  is the dataset associated with a leaf node in  $T$ , then the dataset associated with the corresponding leaf node in  $T'$  is the rescaled version of points in  $D$ .

The invariant holds at the beginning of the algorithm, as the only (leaf) node is the root, which is associated with  $S$  in  $T$  and  $S'$  in  $T'$ . Suppose the invariant holds at step  $t$  of the algorithm, and at step  $t + 1$  we split a node  $v$  in  $T$  such that the dataset associated with  $v$  is  $D$ . If the splitting rule used is  $(j, t)$ , then, this splitting rule has the highest information gain among all the possible splitting rules. Observe that as the corresponding node  $v'$  in  $T'$  is associated with a scaled version  $D'$  of  $D$ , for any  $j$  and  $t$ , the information gain of a splitting rule  $(j, \alpha^j t)$  at  $v'$  is equal to the information gain of the splitting rule  $(j, t)$  in  $v$ . Thus, assuming that ties are broken consistently, we will pick the splitting rule  $(j, \alpha^j t)$  to split node  $v'$ . Thus invariants (a) and (b) are maintained after step  $t + 1$ . Finally, invariant (c) is also maintained as the subset of  $D$  for which feature  $j$  is  $\leq t$  is exactly equal to the subset of  $D'$  for which feature  $j$  is  $\leq \alpha^j t$ .

Thus, at the end of the ID3 decision tree algorithm, we arrive at two trees  $T$  and  $T'$  which have exactly the same structure, where the corresponding nodes  $v$  and  $v'$  have equal splitting rules. Thus if a test example  $x$  follows a path  $P$  in  $T$  from the root to the leaf, its rescaled version  $z$  will follow exactly the same path in  $T'$  from root to leaf and will be classified the same way. Therefore the two decision trees will be equal.

**Perceptron.** For a uniform scaling factor  $\alpha$ , we claim that at any step, if the hyperplane normal in the original space is  $w$ , then the hyperplane normal in the rescaled space must be  $\alpha w$ . If this claim is true, then the classifiers in the two spaces will be equal, because as  $\alpha > 0$ ,  $\text{sign}(\langle w, x_t \rangle) = \text{sign}(\langle \alpha w, z_t \rangle) = \text{sign}(\langle \alpha w, \alpha x_t \rangle)$ .

We prove this by induction. The base case is trivial because  $w$  is initialized to 0 in both spaces. Then suppose our claim is true for step  $t - 1$ , we show that the claim still holds at step  $t$ . At step  $t$  the algorithm predicts the label for the training data  $(x_t, y_t)$  in the original space and training data

$(z_t = \alpha x_t, y_t)$  in the rescaled space. It is easy to see that the prediction result is the same for the classifiers in both spaces as  $\alpha > 0$ . If the result is correct, then no change is made to either  $w$ . If the result is wrong, the normal in the original space is updated to  $w + y_t x_t$ , while in the rescaled space, the normal is updated to  $\alpha w + y_t z_t = \alpha(w + y_t x_t)$ . Thus the claim still holds at this step. Therefore the Perceptron algorithm produces equal classifiers in both spaces.

For non-uniform  $\alpha$ 's, the two classifiers are not equal. One counter-example is given below. There is only one positive training data  $(2, -2)$ , which becomes  $(1, -2)$  in the rescaled space. Consider the test data  $(2, 1)$ , and the rescaled version  $(1, 1)$ . The resulting classifier classifies them into different labels.

**Behavior under scaling transformations.** In case of uniform scaling transformations (same  $\alpha^i$ ) across all features/dimensions, all the 3 algorithms are equally robust. However, in case of non-uniform scaling transformations (different  $\alpha^i$ ), ID3 Decision Trees are more robust to compared to  $k$ -NN and Perceptrons.

