

## Homework 3

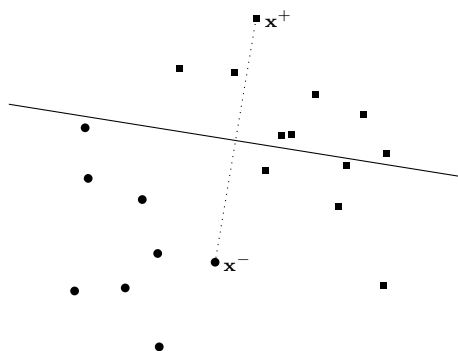
Due: 5 April 2018

**Homework submission:** We will collect your homework **at the beginning of class** on the due date. If you cannot attend class that day, you can leave your solution in Phyllis Wan's postbox in the Department of Statistics, 10th floor SSW, at any time before then.

We do not accept homework submitted late. There will be no exceptions.

### Problem 1 (A primitive ensemble)

This problem considers a very simple ensemble classifier for a two-class problem, based on the idea of "randomly throwing out hyperplanes" discussed in class.



The basic idea is to define a "weak classifier" as follows:

- Choose two training data points  $x^-$  and  $x^+$ , one in each class.
- Place an affine plane "in the middle" between the two, i.e. perpendicular to the connecting line. The connecting line is drawn in the figure as a dotted line, the affine plane as a solid line. In formulas, the plane is given by a normal vector  $w$  and an offset  $c_w$  computed as

$$w := \frac{x^+ - x^-}{\|x^+ - x^-\|} \quad \text{and} \quad c_w := \langle w, x^- + \frac{1}{2}(x^+ - x^-) \rangle$$

- Try out each of the two possible orientations of the plane; select the one with the smaller training error. That is, the weak classifier is defined as

$$f(\cdot) = \text{sgn}(\langle \cdot, v \rangle - c)$$

where

$$\text{either } v := w \text{ or } v := -w \quad \text{and} \quad c := \langle v, x^- + \frac{1}{2}(x^+ - x^-) \rangle,$$

depending on which choice achieves smaller error.

We combine many of these weak classifiers into a single ensemble classifier.

More precisely:

- Split the available data into two equally sized parts (training and test).
- Select  $m$  pairs of points  $(\mathbf{x}_1^-, \mathbf{x}_1^+), \dots, (\mathbf{x}_m^-, \mathbf{x}_m^+)$ , where each  $\mathbf{x}_i^-$  is drawn uniformly at random from class  $-1$ , and  $\mathbf{x}_i^+$  uniformly from class  $+1$ . Sample with replacement, i.e. training data points can appear more than once.
- For each such pair  $(\mathbf{x}_i^-, \mathbf{x}_i^+)$ , compute the classifier  $f_i$  given by  $(\mathbf{v}_i, c_i)$  as described above.
- The overall classifier  $g_m$  is defined as the majority vote

$$g_m(\mathbf{x}) = \text{sgn}\left(\sum_{j=1}^m f_j(\mathbf{x})\right) = \text{sgn}\left(\sum_{j=1}^m \text{sgn}(\langle \mathbf{v}_j, \mathbf{x} \rangle - c_j)\right).$$

We will try to understand how well this classifier performs depending on how many weak classifiers we include in the computation.

**Questions:**

1. Implement the classifier  $g_m$  above in R: Write a training function `[V, C] = train.g(m, trainingData)` that returns a  $(m \times d)$ -matrix  $\mathbf{V}$  and a length  $m$ -vector  $\mathbf{C}$  ( $d$  is the length of  $\mathbf{x}$ ), where

$$\mathbf{V} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \dots \\ \mathbf{v}_m \end{bmatrix} \quad \text{and} \quad \mathbf{C} = \begin{pmatrix} c_1 \\ c_2 \\ \dots \\ c_m \end{pmatrix}.$$

2. Write a function `classify(x, V, C)` that applies the classifier to a new data point  $\mathbf{x}$ . This function should return  $\pm 1$ .

We perform experiments on two data sets: The USPS digit data we have already used before, and another data set often used as an example, the Wisconsin Breast Cancer Data (available on the course page). For each data set:

3. Randomly and equally separate the data into a training and a test data set. (Subsetting the data does not perform well in the WBCD dataset. I tried using loops. The time is not that bad. )
4. For each  $m = 1, 3, \dots, 199$ :
  - (a) Train a classifier  $g_m$ .
  - (b) Classify each point in the test set.
  - (c) Compute the error rate on the test set.
5. Record all errors and plot the error rate as a function of  $m$ .

Optional:

6. You can also further increase the range of  $m$  to see what happens.