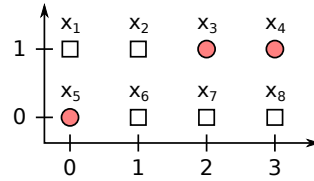


Exercise Sheet 10

Exercise 1: Boosted Classifiers (25 + 25 P)

We consider a two-dimensional dataset $x_1, \dots, x_8 \in \mathbb{R}^2$ with binary labels $y_1, \dots, y_8 \in \{-1, 1\}$.



Red circles denote the first class ($y_i = +1$) and white squares denote the second class ($y_i = -1$). We decide to classify this data with a boosted classifier and use the nearest mean classifier as a weak classifier. The boosted classifier is given by

$$f(x) = \text{sign}\left(\alpha_0 + \sum_{t=1}^T \alpha_t h_t(x)\right)$$

where $\alpha_0, \dots, \alpha_T \in \mathbb{R}$ are the boosting coefficients. The t th nearest mean classifier is given by

$$h_t(x) = \begin{cases} +1 & \|x - \mu_t^+\| < \|x - \mu_t^-\| \\ -1 & \text{else} \end{cases} \quad \text{with} \quad \mu_t^+ = \frac{\sum_{i:y_i=+1} p_i^{(t)} x_i}{\sum_{i:y_i=+1} p_i^{(t)}} \quad \text{and} \quad \mu_t^- = \frac{\sum_{i:y_i=-1} p_i^{(t)} x_i}{\sum_{i:y_i=-1} p_i^{(t)}}.$$

where $p_1^{(t)}, \dots, p_N^{(t)}$ are the data weighting terms for this classifier.

- Draw at hand a possible boosted classifier that classifies the dataset above, i.e. draw the decision boundary of the weak classifiers $h_t(x)$ and of the final boosted classifier $f(x)$. We use the convention $\text{sign}(0) = 0$.
- Write the weighting terms $p_i^{(t)}$ and the coefficients $\alpha_0, \dots, \alpha_T$ associated to the classifiers you have drawn.

(Note: In this exercise, the boosted classifier does not need to derive from a particular algorithm. Instead, the number of weak classifiers, the coefficients and the weighting terms can be picked at hand with the sole constraint that the final classifier implements the desired decision boundary.)

Exercise 2: AdaBoost as an Optimization Problem (25 + 25 P)

Consider AdaBoost for binary classification applied to some dataset $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$. The algorithm starts with uniform weighting ($\forall_{i=1}^N : p_i^{(1)} = 1/N$) and performs the following iteration:

for $t = 1 \dots T$:

- | | | |
|---------|-----------------------------------------------------------------------------------|-----------------------------------------------------------|
| Step 1: | $\mathcal{D}, p^{(t)} \mapsto h_t$ | (learn t th weak classifier using weighting $p^{(t)}$) |
| Step 2: | $\epsilon_t = \mathbb{E}_{p^{(t)}}[1_{(h_t(x) \neq y)}]$ | (compute the weighted error of the classifier) |
| Step 3: | $\alpha_t = \frac{1}{2} \log\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)$ | (set its contribution to the boosted classifier) |
| Step 4: | $\forall_{i=1}^N : p_i^{(t+1)} = Z_t^{-1} p_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))$ | (set a new weighting for the data) |

The term $\mathbb{E}_{p^{(t)}}[\cdot]$ denotes the expectation under the data weighting $p^{(t)}$, and Z_t is a normalization term. An interesting property of AdaBoost is that it can be shown to minimize some objective function

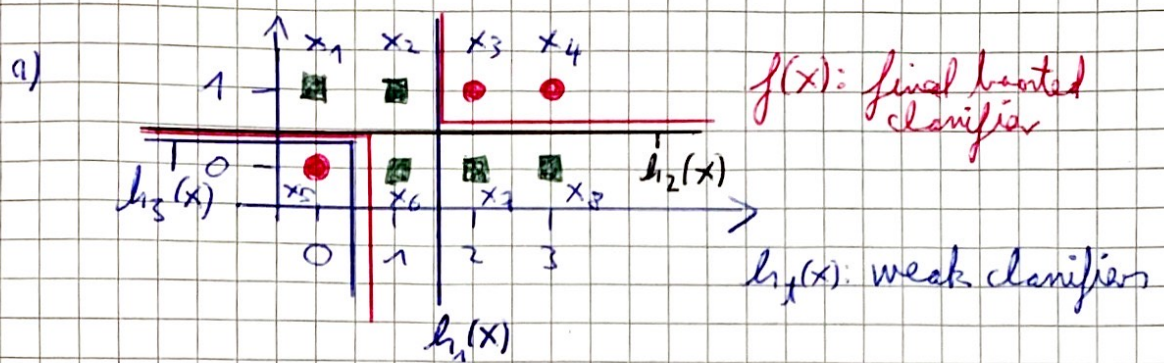
$$\mathcal{G}(\alpha) = \sum_{i=1}^N \exp(-y_i f_{\alpha,t}(x_i))$$

where $f_{\alpha,t}(x) = \sum_{\tau=1}^t \alpha_\tau h_\tau(x)$ is the output score of the boosted classifier after t iterations.

- Show that the objective can be rewritten as $\mathcal{G}(\alpha) = N \cdot \left(\prod_{\tau=1}^{t-1} Z_\tau\right) \cdot \sum_{i=1}^N p_i^{(t)} \exp(-y_i \alpha_t h_t(x_i))$.
- Show that Step 3 of the AdaBoost procedure above is equivalent to computing $\alpha_t = \arg \min_{\alpha_t} \mathcal{G}(\alpha)$.

ML10

1



b)

For the final classifier $f(x) = \text{sign}(\alpha_0 + \sum_{t=1}^T \alpha_t h_t(x))$ we need to choose proper weights α_i for each weak classifier $h_i(x)$. The first two classifiers $h_1(x)$ and $h_2(x)$ can be weighted equally:

$$\alpha_1 = \alpha_2 = 1$$

The third classifier needs a larger weight because it classifies x_5 as a single classifier (x_3 and x_4 are two classifiers). We would like values of zero to become negative such that they are correctly classified as $y_i = -1$. Therefore, we choose the bias weight negative:

$$\alpha_0 = -1$$

The data weighting term $p_i^{(t)}$ needed to select the correct linear decision boundary using the mean of the ~~corresponding~~ corresponding data points. For example, for $h_1(x)$ the correct mean is between x_2 and x_3 . So our 3 weighting vectors p_i become:

$$p_1 = [0, 1, 1, 0, 0, 0, 0, 0]$$

$$p_2 = [0, 0, 1, 1, 0, 0, 1, 1]$$

$$p_3 = [1, 0, 0, 0, 1, 1, 0, 0]$$

Our data values are

$$X = \left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 3 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 0 \\ 3 \end{pmatrix} \right\}$$

and the labels are

$$Y = \{-1, -1, +1, +1, +1, -1, -1, -1\}$$

The first classifier $h_1(x)$ is:

$$h_1(x) = \begin{cases} +1 & \|x - \mu_1^+\| < \|x - \mu_1^-\| \\ -1 & \text{else} \end{cases}$$

The mean vectors are

$$\mu_1^+ = \frac{\sum_{i: y_i = +1} p_i^{(1)} x_i}{\sum_{i: y_i = +1} p_i^{(1)}} = \frac{p_2 x_2 + p_3 x_3 + p_4 x_4}{p_2 + p_3 + p_4} = \frac{1 \cdot \begin{pmatrix} 1 \\ 2 \end{pmatrix} + 1 \cdot \begin{pmatrix} 1 \\ 3 \end{pmatrix} + 0 \cdot \begin{pmatrix} 1 \\ 3 \end{pmatrix}}{1 + 1 + 0} = \frac{1}{2} \begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

and

$$\mu_1^- = \frac{\sum_{i: y_i = -1} p_i^{(1)} x_i}{\sum_{i: y_i = -1} p_i^{(1)}} = \frac{p_1 x_1 + p_5 x_5 + p_6 x_6 + p_7 x_7 + p_8 x_8}{p_1 + p_5 + p_6 + p_7 + p_8} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\text{because } p_1 = p_5 = p_6 = p_7 = p_8 = 0$$

The first weak classifier $h_1(x)$ for the first data point x_1 becomes

$$h_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{cases} +1 & \text{if } \left\| \begin{pmatrix} 1 \\ 0 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 2 \\ 3 \end{pmatrix} \right\| < \left\| \begin{pmatrix} 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right\| \Leftrightarrow \left\| \begin{pmatrix} 0 \\ -3/2 \end{pmatrix} \right\| < \left\| \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\| \\ -1 & \text{else} \end{cases}$$

$$\Leftrightarrow h_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{cases} +1 & \text{if } \frac{3}{2} < 1 \\ -1 & \text{else} \end{cases}$$

Therefore, $h_1(x_1) = -1$.

Analogously for x_2, \dots, x_8 and $h_2(x)$ and $h_3(x)$.

2) a)

Show that the objective

$$1) \quad \mathcal{G}(\alpha) = \sum_{i=1}^N \exp\left(-y_i \sum_{c=1}^T \alpha_c h_c(x_i)\right)$$

can be rewritten as

$$2) \quad \mathcal{G}(\alpha) = N \cdot \left(\prod_{c=1}^{t-1} Z_c \right) \sum_{i=1}^N p_i^{(t)} \exp\left(-y_i \alpha_t h_t(x_i)\right)$$

we rewrite 1):

$$3) \quad \mathcal{G}(\alpha) = \sum_{i=1}^N \exp\left(-y_i \sum_{c=1}^{t-1} \alpha_c h_c(x_i)\right) \cdot \exp\left(-y_i \alpha_t h_t(x_i)\right)$$

Using the AdaBoost algorithm we can iteratively calculate p_i :

$$p_i^{(1)} = \frac{1}{N}$$

$$p_i^{(2)} = Z_1^{-1} p_i^{(1)} \exp\left(-\alpha_1 y_i h_1(x_i)\right)$$

$$p_i^{(3)} = Z_2^{-1} p_i^{(2)} \exp\left(-\alpha_2 y_i h_2(x_i)\right)$$

$$= Z_2^{-1} Z_1^{-1} \frac{1}{N} \exp\left(-\alpha_1 y_i h_1(x_i)\right) \exp\left(-\alpha_2 y_i h_2(x_i)\right)$$

$$p_i^{(t)} = \prod_{c=1}^{t-1} Z_c^{-1} \frac{1}{N} \exp\left(-y_i \sum_{c=1}^{t-1} \alpha_c h_c(x_i)\right)$$

$$4) \quad \Leftrightarrow N \left(\prod_{c=1}^{t-1} Z_c \right) p_i^{(t)} = \exp\left(-y_i \sum_{c=1}^{t-1} \alpha_c h_c(x_i)\right)$$

Inserting 4) in 3) yields

$$\mathcal{G}(\alpha) = \sum_{i=1}^N N \left(\prod_{c=1}^{t-1} Z_c \right) p_i \exp\left(-y_i \alpha_t h_t(x_i)\right)$$

$$= N \left(\prod_{c=1}^{t-1} Z_c \right) \sum_{i=1}^N p_i \exp\left(-y_i \alpha_t h_t(x_i)\right)$$

b) The value α_+ which minimizes the objective $\mathcal{G}(\alpha)$ can be obtained by setting the partial derivative to 0:

$$\alpha_+ = \underset{\alpha_+}{\operatorname{argmin}} \mathcal{G}(\alpha)$$

$$\frac{\partial \mathcal{G}(\alpha)}{\partial \alpha_+} = \frac{\partial}{\partial \alpha_+} \left(N \left(\prod_{c=1}^{t-1} z_c \right) \sum_{i=1}^N p_i^{(t)} \exp(-\gamma_i \alpha_+ h_+(x_i)) \right)$$

$$= \underbrace{N \left(\prod_{c=1}^{t-1} z_c \right)}_{\text{const.}} \sum_{i=1}^N p_i^{(t)} (-\gamma_i h_+(x_i)) \exp(-\gamma_i \alpha_+ h_+(x_i)) \stackrel{!}{=} 0$$

Since γ_i and $h_+(x_i)$ are $\in \{-1, +1\}$

Therefore we can split the previous term into

$$\frac{\partial \mathcal{G}(\alpha)}{\partial \alpha_+} = \underbrace{\sum_{i \in \mathcal{I}^{++}} p_i^{(t)} \exp(-\alpha_+) (-1)}_{\gamma_i \cdot h_+(x_i) = +1} + \underbrace{\sum_{i \in \mathcal{I}^{+-}} p_i^{(t)} \exp(\alpha_+) \cdot 1}_{\gamma_i \cdot h_+(x_i) = -1} = 0$$

$$\Leftrightarrow \underbrace{\left(\sum_{i \in \mathcal{I}^{++}} p_i^{(t)} \right)}_{= 1 - \varepsilon_+} \exp(-\alpha_+) = \underbrace{\left(\sum_{i \in \mathcal{I}^{+-}} p_i^{(t)} \right)}_{= \varepsilon_+} \exp(\alpha_+)$$

$$\Leftrightarrow \frac{1 - \varepsilon_+}{\varepsilon_+} = \frac{\exp(\alpha_+)}{\exp(-\alpha_+)}$$

$$\Leftrightarrow \ln \left(\frac{1 - \varepsilon_+}{\varepsilon_+} \right) = \alpha_+ + \alpha_+$$

$$\Leftrightarrow \alpha_+ = \frac{1}{2} \ln \left(\frac{1 - \varepsilon_+}{\varepsilon_+} \right)$$