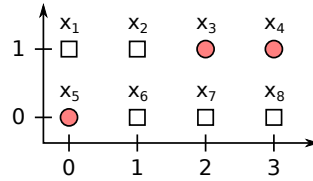


Exercise Sheet 11

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

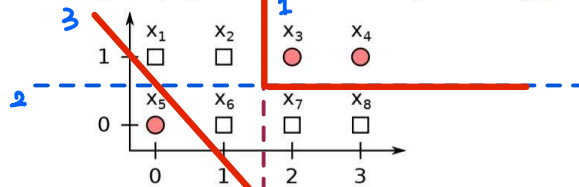
$$\mathcal{G}(\alpha) = \sum_{i=1}^N \exp(-y_i f_{\alpha,t}(x_i))$$

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 $\alpha_t = \arg \min_{\alpha_t} \mathcal{G}(\alpha)$.

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.

- (a) 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 $\text{sign}(0) = 0$.

Solution:

See the diagram above.

The dash line is the boundary of the weak classifiers.

The red line is the boundary of the final boosted classifier. ✓

- (b) 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.)

Solution:

The weighting terms represent the importance of the data points. i.e. The data points with high weight is crucial to the boundary

	1	2	3	4	5	6	7	8	
$p^{(1)} =$	0	$\frac{1}{2}$	$\frac{1}{2}$	0	0	0	0	0	$\alpha_0 = 1$
$p^{(2)} =$	0	0	$\frac{1}{4}$	$\frac{1}{4}$	0	0	$\frac{1}{4}$	$\frac{1}{4}$	$\alpha_1 = 1$ ✓
$p^{(3)} =$	$\frac{1}{3}$	0	0	0	$\frac{1}{3}$	$\frac{1}{3}$	0	0	$\alpha_2 = 1$
									$\alpha_3 = 2$

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.

(a) Show that the objective can be rewritten as $\mathcal{G}(\alpha) = N \cdot \left(\prod_{\tau=1}^{t-1} Z_{\tau} \right) \cdot \sum_{i=1}^N p_i^{(t)} \exp(-y_i \alpha_t h_t(x_i))$.

Solution:

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

$$p_i^{(2)} = \frac{1}{N} \cdot \frac{1}{Z_1} \exp(-\alpha_1 y_i h_1(x_i))$$

$$p_i^{(3)} = \frac{1}{N} \cdot \frac{1}{Z_2} \frac{1}{Z_1} \exp(-\alpha_1 y_i h_1(x_i) - \alpha_2 y_i h_2(x_i))$$

$$= \frac{1}{N} \cdot \frac{1}{Z_1} \frac{1}{Z_2} \exp\left(-\sum_{\tau=1}^2 \alpha_{\tau} y_i h_{\tau}(x_i)\right)$$

\Downarrow

$$p_i^{(t)} = \frac{1}{N} \cdot \left(\prod_{\tau=1}^{t-1} \frac{1}{Z_{\tau}} \right) \cdot \exp\left(-\sum_{\tau=1}^{t-1} \alpha_{\tau} y_i h_{\tau}(x_i)\right)$$

\Downarrow

$$N p_i^{(t)} \cdot \left(\prod_{\tau=1}^{t-1} Z_{\tau} \right) = \exp\left(-\sum_{\tau=1}^{t-1} \alpha_{\tau} y_i h_{\tau}(x_i)\right)$$

\Downarrow

$$\begin{aligned} \mathcal{G}(\alpha) &= \sum_{i=1}^N \exp(-y_i f_{\alpha, t}(x_i)) \\ &= \sum_{i=1}^N \exp\left(-y_i \sum_{\tau=1}^t \alpha_{\tau} h_{\tau}(x_i)\right) \end{aligned}$$

$$= \sum_{i=1}^N \underbrace{\exp(-y_i \sum_{\tau=1}^{t-1} \alpha_{\tau} h_{\tau}(x_i)) \cdot \exp(-y_i \alpha_t h_t(x_i))}_{\downarrow}$$

$$(N \cdot p_i^{(t)} \cdot \prod_{\tau=1}^{t-1} Z_{\tau})$$

$$= \sum_{i=1}^N (N \cdot p_i^{(t)} \cdot \prod_{\tau=1}^{t-1} Z_{\tau}) \cdot \exp(-y_i \alpha_t h_t(x_i))$$

$$= N \cdot (\prod_{\tau=1}^{t-1} Z_{\tau}) \cdot \sum_{i=1}^N p_i^{(t)} \cdot \exp(-y_i \alpha_t h_t(x_i))$$

proofed.

(b) Show that Step 3 of the AdaBoost procedure above is equivalent to computing $\alpha_t = \arg \min_{\alpha_t} \mathcal{G}(\alpha)$.

Solution:

$$\frac{\partial \mathcal{G}}{\partial \alpha_t} = \frac{\partial}{\partial \alpha_t} \left(N \cdot (\prod_{\tau=1}^{t-1} Z_{\tau}) \cdot \sum_{i=1}^N p_i^{(t)} \cdot \exp(-y_i \alpha_t h_t(x_i)) \right)$$

$$= \underbrace{N \cdot (\prod_{\tau=1}^{t-1} Z_{\tau}) \cdot \sum_{i=1}^N p_i^{(t)} \cdot \exp(-y_i \alpha_t h_t(x_i))}_{\neq 0} \cdot (-y_i h_t(x_i)) = 0$$

Since if $y_t = h_t(x_i) \Rightarrow y_t h_t(x_i) = +1$
 else $y_t \neq h_t(x_i) \Rightarrow y_t h_t(x_i) = -1$

\Downarrow

$$\underbrace{\sum_{i: y_i = h_t(x_i)} p_i^{(t)} \cdot \exp(-\alpha_t) \cdot (-1)}_{(1 - \epsilon_t)} + \underbrace{\sum_{i: y_i \neq h_t(x_i)} p_i^{(t)} \exp(\alpha_t) \cdot (+1)}_{\epsilon_t} = 0$$

\therefore we have.

$$-(1 - \epsilon_t) \cdot \exp(-\alpha_t) + \epsilon_t \cdot \exp(\alpha_t) = 0$$

爽筒H7

天

$$\frac{(1-\varepsilon_t)}{\varepsilon_t} = \frac{\exp(-\alpha_t)}{\exp(\alpha_t)} = \exp(2\alpha_t)$$

$$\dots \alpha_t^* = \frac{1}{2} \log\left(\frac{1-\varepsilon_t}{\varepsilon_t}\right)$$