

Homework No. 5

ALEXANDER SHENDER 328626114, SNIR HORDAN 205689581

Question 1.

1. Prove that when running AdaBoost, the distribution is updated such that the error of the chosen weak classifier h_t , w.r.t the updated distribution $D_i^{(t+1)}$, is exactly $\frac{1}{2}$.

That is, prove that $\sum_i D_i^{(t+1)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i} = \frac{1}{2}$.

Hint: You can fill the missing steps in the following derivation:

$$\sum_i D_i^{(t+1)} \cdot \mathbf{1}_{h_t(x_i) \neq y_i} = \dots = \frac{\epsilon_t}{\epsilon_t + (1 - \epsilon_t) \exp\{-2w_t\}} = \dots = \frac{1}{2}.$$

We start by indeed from writing the expression for the updated error value with the updated data weights (distribution) for the last chosen weak classifier:

$$E_{t+1} = \sum_i D_i^{t+1} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

By using:

$$D_i^{t+1} = D_i^t \cdot \frac{\exp(-w_t y_i h_t(x_i))}{\sum_j D_j^t \exp(-w_t y_j h_t(x_j))} = D_i^t \cdot \frac{\exp(-w_t y_i h_t(x_i))}{Z_t}$$

Putting back:

$$E_{t+1} = \sum_i D_i^{t+1} \cdot \mathbf{1}_{h_t(x_i) \neq y_i} = \sum_i \frac{D_i^t \exp(-w_t y_i h_t(x_i))}{Z_t} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

Z_t is a normalization factor, so we can put it outside of the sum:

$$E_{t+1} = \frac{\sum_i D_i^t \exp(-w_t y_i h_t(x_i))}{Z_t} \cdot \mathbf{1}_{h_t(x_i) \neq y_i}$$

We divide into 2 cases:

$$\begin{cases} E : h_t(x_i) = y_i \\ C : h_t(x_i) \neq y_i \end{cases}$$

For each case,

$$\begin{cases} E : \mathbf{1}_{h_t(x_i) \neq y_i} = 0 ; y_i h_t(x_i) = 1 \\ C : \mathbf{1}_{h_t(x_i) \neq y_i} = 1 ; y_i h_t(x_i) = -1 \end{cases}$$

So we get:

$$E_{t+1} = \frac{\sum_{i \in E} D_i^t \exp(-w_t y_i h_t(x_i))}{Z_t} \cdot 0 + \frac{\sum_{i \in C} D_i^t \exp(-w_t y_i h_t(x_i))}{Z_t} \cdot 1 = \frac{\sum_{i \in C} D_i^t \exp(w_t)}{Z_t}$$

The numerator contains the expression for the error value, since we have isolated for case $\{E : h_t(x_i) = y_i\}$:

$$E_t = \sum_{i \in C} D_i^t$$

$$E_{t+1} = \frac{E_t \exp(w_t)}{Z_t}$$

Opening the denominator using same 2 cases:

$$Z_t = \sum_{i \in E} D_i^t \exp(-w_t) + \sum_{i \in C} D_i^t \exp(w_t)$$

Putting back:

$$E_{t+1} = \frac{E_t \exp(w_t)}{\sum_{i \in E} D_i^t \exp(-w_t) + \sum_{i \in C} D_i^t \exp(w_t)} \quad ; \quad / \exp(w_t)$$

$$E_{t+1} = \frac{E_t}{\sum_{i \in E} D_i^t \exp(-2w_t) + \sum_{i \in C} D_i^t} = \frac{E_t}{\sum_{i \in E} D_i^t \exp(-2w_t) + E_t}$$

Taking $\exp(-2w_t)$ out of the sum:

$$E_{t+1} = \frac{E_t}{\exp(-2w_t) \sum_{i \in E} D_i^t + E_t}$$

The sum over weights of the correct predictions E, is 1-sum over incorrect predictions, since they sum to 1:

$$\sum_{i \in E} D_i^t = 1 - \sum_{i \in C} D_i^t = 1 - E_t$$

So we get:

$$E_{t+1} = \frac{E_t}{\exp(-2w_t) (1 - E_t) + E_t}$$

Using the expression for the weight of the weak classifier:

$$w_t = \frac{1}{2} \log \left(\frac{1}{E_t} - 1 \right)$$

$$\begin{aligned} \exp(-2w_t) &= \exp \left(-\log \left(\frac{1}{E_t} - 1 \right) \right) = \exp \left(\log \left(\left(\frac{1}{E_t} - 1 \right)^{-1} \right) \right) = \left(\frac{1}{E_t} - 1 \right)^{-1} = \left(\frac{1 - E_t}{E_t} \right)^{-1} \\ &= \frac{E_t}{1 - E_t} \end{aligned}$$

Putting everything back:

$$E_{t+1} = \frac{E_t}{\frac{E_t}{1 - E_t} (1 - E_t) + E_t} = \frac{E_t}{E_t + E_t} = \frac{1}{2}$$

Question 2

2.a.

Let $\alpha \in \mathbb{R}_{>0}$

Notice σ is positive-homogeneous, i.e. $\sigma(\alpha x) = \max\{0, \alpha x\} = \alpha \max\{0, x\} = \alpha \sigma(x)$

Denote $\alpha\theta = (\alpha W^{(1)}, \dots, \alpha W^{(L)})$.

Let h be defined on θ as in question and \tilde{h} defined on $\alpha\theta$ as in question.

Claim. $\tilde{h}_L(x) = \alpha^L h_L(x)$

Proof. Proof by induction on L .

Base : $\tilde{h}_1(x) = \sigma(\alpha W^{(1)T} x) \underbrace{=}_{\alpha \text{ positive homogeneous}} \alpha \sigma(W^{(1)T} x) = \alpha h_1(x)$

Assume for $n < L$, $\tilde{h}_n(x) = \alpha^n h_n(x)$

Induction step:

By definition, $\tilde{h}_L(x) = \sigma(\alpha W^{(L)T} \tilde{h}_{L-1}(x)) \underbrace{=}_{\text{induction hypothesis}} \sigma(\alpha W^{(L)T} \alpha^{L-1} h_{L-1}(x)) =$
 $\underbrace{=}_{\alpha \text{ positive-homogeneous}} \alpha^L \sigma(W^{(L)T} h_{L-1}(x)) = \alpha^L h_L(x)$ □

By definition, $F_{\alpha\theta}(x) = W^{(L)T} \tilde{h}_{L-1}(x) = W^{(L)T} \alpha^{L-1} h_{L-1}(x) = \alpha^{L-1} F_\theta(x)$

Then $C = \alpha^{L-1}$.

2.b.

$$c_{y_i} = \frac{\exp((F_{\alpha\theta}(x))_i)}{\sum_{j=1}^k \exp((F_{\alpha\theta}(x))_j)} = \frac{\exp(\alpha^L (F_\theta(x))_i)}{\sum_{j=1}^k \exp(\alpha^L (F_\theta(x))_j)} \xrightarrow{\alpha \rightarrow 0} \frac{1}{k}$$

Then the distribution is uniform.

2.c.

Case 1 $(F_\theta(x))_i > (F_\theta(x))_j, \forall j \in \{m_1, \dots, m_{k-s}\}$ and $(F_\theta(x))_i = (F_\theta(x))_j, \forall j \in \{m_{k-s+1}, \dots, m_k\}, (\{m_1, \dots, m_k\} = \{1, \dots, k\})$

$$c_{y_i} = \frac{\exp((F_{\alpha\theta}(x))_i)}{\sum_{j=1}^k \exp((F_{\alpha\theta}(x))_j)} = \frac{\exp(\alpha^L (F_\theta(x))_i)}{\sum_{j=1}^k \exp(\alpha^L (F_\theta(x))_j)} = \frac{1}{\sum_{j=1}^k \exp(\alpha^L ((F_\theta(x))_j - (F_\theta(x))_i))} = \dots$$

$$\dots = \frac{1}{s + \underbrace{\sum_{j \in \{m_1, \dots, m_{k-s}\}} \exp(\alpha^L ((F_\theta(x))_j - (F_\theta(x))_i))}_{\xrightarrow{\alpha \rightarrow \infty} 0}} \xrightarrow{\alpha \rightarrow \infty} \frac{1}{s}$$

($s \geq 1$ because exists at least one maximum value then limit is well-defined)

Case 2 $\exists 1 \leq j \leq k$ s.t. $(F_\theta(x))_i < (F_\theta(x))_j$

$$c_{y_i} = \frac{1}{1 + \underbrace{\sum_{j \neq i} \exp(\alpha^L ((F_\theta(x))_j - (F_\theta(x))_i))}_{\xrightarrow{\alpha \rightarrow \infty} \infty}} \xrightarrow{\alpha \rightarrow \infty} 0$$

Then distribution is uniform on the $s \leq k$ entries with maximum value and zero otherwise.

$$\text{i.e. } \left(0, \dots, \underbrace{\frac{1}{s}}_m, \dots, 0 \right) \text{ where } m \in \operatorname{argmax}_{1 \leq i \leq k} ((F_\theta(x))_i).$$