

Machine Learning Fundamentals - Homework 3

Milena Markovic (milena.markovic@grenoble-inp.org)

Adaboost algorithm

Question 1: Explain the Adaboost algorithm seen in the course.

The Adaboost algorithm utilizes a set of weak learners, which when combined, produce an efficient final classifier. In the algorithm, we start by assigning all training points equal weights. We go through T iterations where we first train a weak classifier, pick a step size for the current iteration based on the current distribution of misclassified points, and then based on those observations, we update the weights of the training points - these weights will influence the weak learner in the next iteration to classify misclassified points correctly.

Question 2: What is the role of the distribution D_t ?

The distribution assigns weights to training points. In each iteration, if the point was misclassified, its weight increases. If it was classified correctly, its weight in the distribution decreases. When a point's weight is increased, the classifier in the next iteration is more probable to classify it correctly.

Question 3-4: After T rounds, the algorithm will learn T weak-classifiers (f_t) $1 \leq t \leq T$, with their associated weights (α_t) $1 \leq t \leq T$ where the output of each weak classifier is binary in the set $\{-1, +1\}$. How is the final classifier F obtained?

The final classifier F is obtained as $F(x) = \text{sign}(\sum_{t=1}^T \alpha_t f_t(x))$ where α_t is the step value of iteration t and f_t is the classification of the weak classifier t.

Question 5: Explain why the empirical error of the final classifier F on a training set of size m; $S = \{(x_i, y_i) | i \in \{1, \dots, m\}\}$ is bounded by the following surrogate loss:

$$\mathcal{L}(F, S) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i F(\mathbf{x}_i) \leq 0} \leq \frac{1}{m} \sum_{i=1}^m e^{-y_i \sum_{t=1}^T \alpha_t f_t(\mathbf{x}_i)}$$

where, 1_π if the predicate π is true; and 0 otherwise.

$$\mathcal{L}(F, S) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i F(\mathbf{x}_i) \leq 0} \leq \frac{1}{m} \sum_{i=1}^m e^{-y_i \sum_{t=1}^T \alpha_t f_t(\mathbf{x}_i)}$$

As per the definition of the final classifier F, the following statement can be transformed as:

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i F(x_i) \leq 0} \leq \frac{1}{m} \sum_{i=1}^m e^{-y_i F(x_i)}$$

From the constraint of the left hand side of the equation, we can determine the codomain of the expression in the sum on the right hand side:

$$y_i F(x_i) \leq 0 \Rightarrow e^{-y_i F(x_i)} \geq 1, \forall i \in [1, m]$$

As the expression from the sum on the left hand side can only be 1 or 0, and the full expressions on both sides are average values on the domain $\forall i \in [1, m]$, we can conclude that the statement that the empirical error is bounded by loss is true.

$$\mathbb{1}_{y_i F(x_i) \leq 0} \in [0, 1], \forall i \in [1, m]$$

Question 6: Show that

$$\frac{1}{m} \sum_{i=1}^m e^{-y_i F(\mathbf{x}_i)} = \sum_{i=1}^m Z_1 D_2(i) \prod_{t>1} e^{-y_i \alpha_t f_t(\mathbf{x}_i)} \quad (1)$$

$$\text{where, } \forall t, Z_t = \sum_{i=1}^m D_t(i) e^{-\alpha_t y_i f_t(\mathbf{x}_i)} \quad (2)$$

We will start off by transforming the product on the right hand side of the equation:

$$\begin{aligned} \prod_{t>1} e^{-y_i \alpha_t f_t(\mathbf{x}_i)} &= e^{-y_i \alpha_2 f_2(\mathbf{x}_i)} \cdot \dots \cdot e^{-y_i \alpha_T f_T(\mathbf{x}_i)} \\ &= e^{-y_i (\alpha_2 f_2(\mathbf{x}_i) + \dots + \alpha_T f_T(\mathbf{x}_i))} \\ &= e^{-y_i \sum_{t=2}^T \alpha_t f_t(\mathbf{x}_i)} \end{aligned}$$

Next, using the formula for the distribution

$$D_{t+1}(i) = \frac{D_t(i) e^{-\alpha_t y_i f_t(\mathbf{x}_i)}}{Z_t}$$

we will start transforming the equation by substituting the product and D_2 :

$$\begin{aligned}
\frac{1}{m} \sum_{i=1}^m e^{-y_i F(x_i)} &= \sum_{i=1}^m Z_1 D_2(i) \prod_{t>1} e^{-y_i \alpha_t f_t(x_i)} \\
&= \sum_{i=1}^m Z_1 \frac{D_1(i) e^{-y_i \alpha_1 f_1(x_i)}}{\sum_{j=1}^m D_1(j)} e^{-y_i \sum_{t=2}^T \alpha_t f_t(x_i)}
\end{aligned}$$

Because of the initial value of D_1 given by the algorithm definition, we can continue

$$\forall i \in \{1, \dots, m\}, D_1(i) = \frac{1}{m}$$

$$\begin{aligned}
&= \frac{1}{m} \sum_{i=1}^m e^{-y_i \sum_{t=1}^T \alpha_t f_t(x_i)} \\
&= \frac{1}{m} \sum_{i=1}^m e^{-y_i F(x_i)}
\end{aligned}$$

which proves the equation (1).

Question 7: By induction deduce that :

$$\frac{1}{m} \sum_{i=1}^m e^{-y_i F(\mathbf{x}_i)} = \prod_{t=1}^T Z_t$$

As the normalization terms are all positive, the minimization of the surrogate loss (Eq. 1) is then equivalent to the minimization of the normalization factors Z_t , at each iteration.

$$\begin{aligned}
 D_{T+1}(i) &= \frac{D_T(i) e^{-y_i \alpha_T f_T(x_i)}}{\sum_{j=1}^m D_T(j) e^{-y_j \alpha_T f_T(x_j)}} \\
 D_{T+1}(i) &= D_1(i) \frac{e^{-y_i \alpha_1 f_1(x_i)}}{Z_1} \cdot \dots \cdot \frac{e^{-y_i \alpha_T f_T(x_i)}}{Z_T} \\
 D_{T+1}(i) &= \frac{e^{-y_i \sum_{t=1}^T \alpha_t f_t(x_i)}}{m \prod_{t=1}^T Z_t} = \frac{e^{-y_i F(\mathbf{x}_i)}}{m \prod_{t=1}^T Z_t} \\
 \sum_{i=1}^m D_{T+1}(i) &= \sum_{i=1}^m \frac{e^{-y_i F(\mathbf{x}_i)}}{m \prod_{t=1}^T Z_t} = 1 \\
 \frac{1}{m} \sum_{i=1}^m e^{-y_i F(\mathbf{x}_i)} &= \prod_{t=1}^T Z_t
 \end{aligned}$$

Question 8: Considering the equation (Eq. 2); for which value of α_t – expressed with respect to the error $\epsilon_t = \sum_{i:y_i \neq f_t(x_i)} D_t(i)$ – the factor Z_t is minimized?

We will start by transforming the definition of Z_t through the error ϵ_t .

$$\begin{aligned}
 \epsilon_t &= \sum_{i:y_i \neq f_t(x_i)} D_t(i) \Rightarrow \sum_{i:y_i = f_t(x_i)} D_t(i) = 1 - \epsilon_t \\
 Z_t &= \sum_{i=1}^m D_t(i) e^{-\alpha_t y_i f_t(x_i)} \\
 &= \sum_{i=1}^m \left(\prod_{y_i \neq f_t(x_i)} D_t(i) e^{-\alpha_t \cdot (-1)} \right. \\
 &\quad \left. + \prod_{y_i = f_t(x_i)} D_t(i) e^{-\alpha_t \cdot (+1)} \right) \\
 Z_t &= \epsilon_t \cdot e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t}
 \end{aligned}$$

In order to minimize Z_t , we want the difference to be zero. After expressing α_t through ϵ_t , we get the same step size formula that is used in the algorithm.

$$\begin{aligned}
 \epsilon_t e^{\alpha_t} &= (1 - \epsilon_t) e^{-\alpha_t} \quad | \cdot e^{\alpha_t} \\
 \epsilon_t e^{2\alpha_t} &= 1 - \epsilon_t \quad | : \epsilon_t \\
 e^{2\alpha_t} &= \frac{1 - \epsilon_t}{\epsilon_t} \quad | \ln \\
 2\alpha_t &= \ln \frac{1 - \epsilon_t}{\epsilon_t} \\
 \alpha_t &= \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}
 \end{aligned}$$

Question 9: For this particular value of αt , what is the minimum value of Z_t ?

We can get the minimum value of Z_t by using the value of αt acquired in question 8.

$$\begin{aligned} Z_t &= \epsilon_t e^{\alpha t} + (1-\epsilon_t) e^{-\alpha t} \quad \alpha t = \frac{1}{2} \ln \frac{1-\epsilon}{\epsilon} \\ &= \epsilon_t \cdot \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} + (1-\epsilon_t) \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} \\ Z_t &= 2 \cdot \sqrt{\epsilon_t (1-\epsilon_t)} \end{aligned}$$

Question 10: Considering the following variable change $\gamma_t = \frac{1}{2} - \epsilon_t$; show that :

$$\forall t, Z_t = \sqrt{1 - 4\gamma_t^2}$$

$$\begin{aligned} Z_t &= 2 \cdot \sqrt{\epsilon_t (1-\epsilon_t)} \quad | \epsilon = \frac{1}{2} - \gamma \\ &= \sqrt{4((1/2 - \gamma)(1/2 + \gamma))} = \sqrt{4(1/4 - \gamma^2)} = \\ Z_t &= \sqrt{1 - 4\gamma^2} \end{aligned}$$

Question 11: For $\gamma_t < 1/2$, we have $\sqrt{1 - 4\gamma_t^2} \leq e^{-2\gamma_t^2}$. In this case show that :

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i \neq F(x_i)} \leq \prod_{t=1}^T Z_t \leq e^{-2 \sum_{t=1}^T \gamma_t^2}$$

We will split this expression in two parts. First off, we will prove the left part the same way as we have in question 5.

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i \neq F(x_i)} &= \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i F(x_i) \leq 0} \\ \prod_{t=1}^T Z_t &= \frac{1}{m} \sum_{i=1}^m e^{-y_i F(x_i)} \end{aligned}$$

①

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i F(x_i) \leq 0} \leq \frac{1}{m} \sum_{i=1}^m e^{-y_i F(x_i)}$$

\Downarrow

$$y_i F(x_i) \leq 0 \Rightarrow e^{-y_i F(x_i)} \geq 1 \Rightarrow \text{True}$$

Second, we will transform the right-most part of the expression into a product:

$$\begin{aligned} \textcircled{2} \quad \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2} &\leq e^{-2 \sum_{t=1}^T \gamma_t^2} \\ e^{-2 \sum_{t=1}^T \gamma_t^2} &= e^{-2\gamma_1^2} \cdot e^{-2\gamma_2^2} \cdots e^{-2\gamma_T^2} = \prod_{t=1}^T e^{-2\gamma_t^2} \\ \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2} &\leq \prod_{t=1}^T e^{-2\gamma_t^2} \quad | \quad \sqrt{1 - 4\gamma_t^2} \leq e^{-2\gamma_t^2} \\ \Rightarrow \textcircled{2} \end{aligned}$$

When we take the condition from the question, we see that the expression is valid.

Question 12: Explain why the misclassification error of the final classifier F is ensured to converge to 0 when the number of iterations T tends to infinity.

From the equation in the previous question, we can establish that the limit of the misclassification is bounded by the following:

$$\lim_{T \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{y_i \neq F(x_i)\}} \leq \lim_{T \rightarrow \infty} e^{-2 \sum_{t=1}^T \gamma_t^2}$$

To find the limit of the left side, we need to establish that the value of γ is positive. As the limit of the left size is 0, the limit of the right side has to be less or equal to 0 as well which means the error does converge to 0.

$$\begin{aligned} \lim_{T \rightarrow \infty} e^{-2 \sum_{t=1}^T \gamma_t^2} &= e^{\lim_{T \rightarrow \infty} -2 \sum_{t=1}^T \gamma_t^2} \quad |\gamma > 0 \\ &= e^{-\infty} = 0 \\ \lim_{T \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{y_i \neq F(x_i)\}} &\leq 0 \end{aligned}$$