

CS 6375 – ASSIGNMENT 3

Names of students in your group:

Adithya Sundararajan Iyer (asi2000000)

Siddhant Suresh Medar (ssm2000002)

Number of free late days used: 2

Note: You are allowed a **total** of 4 free late days for the **entire semester**. You can use at most 2 for each assignment. After that, there will be a penalty of 10% for each late day.

Please list clearly all the sources/references that you have used in this assignment.

PART I

(10 points)

1. Consider a regression problem of trying to estimate the function $f: X \rightarrow y$ where X is a vector of feature attributes and y is a continuous real valued output variable. You would like to use a bagging model where you first create M bootstrap samples and then use them to create M different models – h_1, h_2, \dots, h_M . You can assume that all models are of the same type.

The error for each of the models would be described as:

$$\epsilon_i(x) = f(x) - h_i(x)$$

where x is the input data and h_i is the model created using i^{th} bootstrap sample

The expected value of the squared error for any of the models will be defined as:

$$E(\epsilon_i(x)^2) = E[(f(x) - h_i(x))^2]$$

The average value of the expected squared error for each of the models acting individually is defined as:

$$E_{avg} = \frac{1}{M} \sum_{i=1}^M E(\epsilon_i(x)^2)$$

Now, you decide to aggregate the models using a committee approach as follows:

$$h_{agg}(x) = \frac{1}{M} \sum_{i=1}^M h_i(x)$$

The error using the aggregated model is defined as:

$$E_{agg}(x) = E\left[\left\{\frac{1}{M} \sum_{i=1}^M h_i(x) - f(x)\right\}^2\right]$$

which can be simplified as:

$$E_{agg}(x) = E\left[\left\{\frac{1}{M} \sum_{i=1}^M \epsilon_i(x)\right\}^2\right]$$

where we used the value of ϵ_i is defined above.

Prove that

$$E_{agg} = \frac{1}{M} E_{avg}$$

provided you make the following assumptions:

1. Each of the errors have a 0 mean

$$E(\epsilon_i(x)) = 0 \text{ for all } i$$

2. Errors are uncorrelated

$$E(\epsilon_i(x)\epsilon_j(x)) = 0 \text{ for all } i \neq j$$

To prove:

$$E_{agg} = \frac{1}{M} E_{avg}$$

We have,

$$E_{agg}(x) = E \left[\left\{ \frac{1}{M} \sum_{i=1}^M E_i(x) \right\}^2 \right] = \frac{1}{M^2} E \left[\left\{ \sum_{i=1}^M E_i(x) \right\}^2 \right]$$

$$E_{agg}(x) = \frac{1}{M^2} E [E_1^2(x) + E_2^2(x) + \dots + E_M^2(x) + 2E_1(x)E_2(x) + \dots]$$

Since errors are uncorrelated:

$$E[E_1(x)E_2(x)] = 0$$

$$E_{agg}(x) = \frac{1}{M^2} E [E_1^2(x) + E_2^2(x) + \dots + E_M^2(x)] = \frac{1}{M^2} E \left[\sum_{i=1}^M E_i(x)^2 \right]$$

But we have:

$$E_{avg}(x) = \frac{1}{M} E \left[\sum_{i=1}^M E_i(x)^2 \right]$$

We can rewrite:

$$E_{agg}(x) = \frac{1}{M} \left(\frac{1}{M} E \left[\sum_{i=1}^M E_i(x)^2 \right] \right)$$

$$\therefore E_{agg}(x) = \frac{1}{M} E_{avg}(x)$$

(10 points)

2. Jensen's inequality states that for any *convex* function f :

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i)$$

In question 1, we had assumed that each of the errors are uncorrelated i.e.

$$E(\epsilon_i(x)\epsilon_j(x)) = 0 \text{ for all } i \neq j$$

This is not really true, as the models are created using bootstrap samples and have correlation with each other. Now, let's remove that assumption. Show that using Jensen's inequality, it is still possible to prove that:

$$E_{agg} \leq E_{avg}$$

To prove:

$$E_{agg} \leq E_{avg}$$

We have,

$$E_{agg}(x) = E\left[\frac{1}{M} \sum_{i=1}^M E_i(x)\right]^2$$

$$E_{avg}(x) = \frac{1}{M} E\left[\sum_{i=1}^M E_i(x)^2\right]$$

According to Jensen's inequality:

$$\lambda_i = \frac{1}{M} F(x_i) = E(E_i(x)^2)$$

$$E\left[\sum_{i=1}^M \frac{1}{M} E_i(x)^2\right] \leq \frac{1}{M} E\left[\sum_{i=1}^M E_i(x)^2\right]$$

Using Cauchy-Schwarz inequality:

$$\left|\sum_{i=1}^M u_i v_i\right|^2 \leq \sum_{j=1}^n |v_j|^2 \sum_{k=1}^n |v_k|^2$$

$$\frac{1}{M} \left(E_1(x) + E_2(x) + \dots + E_M(x) \right)^2 \leq E_1^2(x) + E_2^2(x) + \dots + E_M^2(x)$$

$$M \cdot E \left[\frac{1}{M} \sum_{i=1}^M E_i(x) \right]^2 \leq E \left(E_1^2(x) + E_2^2(x) + \dots + E_M^2(x) \right)$$

Since $M > 0$, multiplying both sides by $1/M$ won't change the inequality sign

$$E \left[\frac{1}{M} \sum_{i=1}^M E_i(x) \right]^2 \leq E \left[\sum_{i=1}^M \frac{1}{M} E_i(x)^2 \right]$$

$$E \left[\frac{1}{M} \sum_{i=1}^M E_i(x) \right]^2 \leq \frac{1}{M} E \left[\sum_{i=1}^M E_i(x)^2 \right]$$

$$\therefore \mathbf{E}_{agg} \leq \mathbf{E}_{avg}$$

(10 points)

3. Deriving the training error for AdaBoost:

In class, we discussed the steps of Adaboost algorithm. Recall that the final hypothesis for a Boolean classification problem at the end of T iterations is given by:

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

The above equation says that the final hypothesis is the weighted hypothesis generated at the end of each individual step.

Also recall that the weight for the point i at step $t+1$ is given by:

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times e^{-\alpha_t h_t(i) y(i)}$$

where:

$D_t(i)$ is the normalized weight of point i in step t

$h_t(i)$ is the hypothesis (prediction) at step t for point i

α_t is the final "voting power" of hypothesis h_t

$y(i)$ is the true label for point i

Z_t is the normalization factor at step t (it ensures that the weights sum up to 1.0)

Note that at step 1, the points have equal weight

$$D_1 = \frac{1}{N}$$

where N is the total number of data points.

At each of the steps, the total error of h_t will be defined as $\epsilon_t = \frac{1}{2} - \gamma_t$, which is a way of saying that the error will be better than 50% by a value γ_t .

Prove that at the end of T steps, the overall training error will be bounded by:

$$\exp(-2 \sum_{t=1}^T \gamma_t^2)$$

That is, the overall training error of the hypothesis H will be less than or equal to the amount indicated above.

To prove:

$$E_{Train} \leq e^{-2 \sum_{i=1}^T \gamma_t^2}$$

We have,

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times e^{-\alpha_t h_t(i) y(i)}$$

$$D_{t+1}(i) = D_1(i) \prod_{t=1}^T \frac{e^{-y_i \alpha_t h_t(x_i)}}{Z_t} = D_1(i) \frac{e^{-y_i \sum_{t=1}^T \alpha_t h_t(x_i)}}{\prod_{t=1}^T Z_t}$$

Let: $\sum_{t=1}^T \alpha_t h_t(x_i) = F(x_i)$

And we know: $D_1(i) = 1/N$

$$D_{t+1}(i) = D_1(i) \prod_{t=1}^T \frac{e^{-y_i \alpha_t h_t(x_i)}}{Z_t} = \frac{1}{N} \frac{e^{-y_i F(x_i)}}{\prod_{t=1}^T Z_t}$$

Given,

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right) = \text{sign}(F(x_i))$$

Since it is a sign function,

$$y_i \neq H(x_i) \rightarrow y_i \cdot F(x_i) \leq 0$$

This condition indicates an error. For an i that satisfies this condition,

$$E_{Train}(i) = 1$$

Overall training error:

$$E_{Train} \leq \frac{1}{N} e^{-y_i F(x_i)}$$

This can be rewritten as:

$$E_{Train} \leq \sum_i D_{t+1}(i) \prod_{t=1}^T Z_t$$

$D_{t+1}(i)$ is the weight distribution and its summation over $i = 1$.

Thus we get:

$$E_{Train} \leq \prod_{t=1}^T Z_t$$

Now,

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$$

$$Z_t = \sum_i D_t(i) \cdot \begin{cases} e^{-\alpha_t}, & \text{if } y_i \neq h_t(x_i) \\ e^{\alpha_t}, & \text{if } y_i = h_t(x_i) \end{cases}$$

$$Z_t = \sum_{i, y_i \neq h_t(x_i)} D_t(i) \cdot e^{-\alpha_t} + \sum_{i, y_i = h_t(x_i)} D_t(i) \cdot e^{\alpha_t}$$

$$Z_t = e^{-\alpha_t}(1 - \varepsilon_t) + e^{\alpha_t}(\varepsilon_t)$$

Substituting for the value of α_t we get,

$$Z_t = 2\sqrt{\varepsilon_t(1 - \varepsilon_t)}$$

But we know,

$$\varepsilon_t = 1/2 - \gamma_t$$

Solving, we get:

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

For all $x \geq 0$, we have $e^x \geq 1+x$

Therefore, for the ensemble, we have:

$$E_{Train} \leq \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2} \leq \prod_{t=1}^T (e^{-4\gamma_t^2})^{1/2}$$

$$E_{Train} \leq \prod_{t=1}^T e^{-2\gamma_t^2}$$

$$\prod_{t=1}^T e^{-2\gamma_t^2} = e^{-2 \sum_{i=1}^T \gamma_t^2}$$

$$\therefore E_{Train} \leq e^{-2 \sum_{i=1}^T \gamma_t^2}$$

Hence, we can say that the Adaboost Training Error is bounded by:

$$\exp\left(-2 \sum_{t=1}^T \gamma_t^2\right)$$