

CSE 512 - Machine Learning HW 3
Mihir Chakradeo 111462188

I. Boosting

ML-HW3

[I] Boosting -

(1) Show That:
$$E_{\text{train}} = \frac{1}{N} \sum_{j=1}^N \delta(H(x^j) \neq y^j) \leq \frac{1}{N} \sum_{j=1}^N \exp(-f(x^j) y^j)$$

LHS =
$$E_{\text{train}} = \frac{1}{N} \sum_{j=1}^N \delta(H(x^j) \neq y^j)$$

It is given that,
$$\delta(H(x^j) \neq y^j) = \begin{cases} 1, & \text{if } H(x^j) \neq y^j \\ 0, & \text{otherwise} \end{cases}$$

$$\therefore E_{\text{train}} = \frac{1}{N} \sum_{j=1}^N \begin{cases} 1, & \text{if } H(x^j) \neq y^j \\ 0 & \text{otherwise} \end{cases}$$

We also know that $H(x^j) = \text{sign}\{f(x^j)\}$ & $y^j \in \{-1, 1\}$

• If we have true $y^j = -1$ & we predict $f(x^j) = +1$.
 then $y^j f(x^j) = -1$.

• when true $y^j = +1$ & $f(x^j) = -1$, $y^j f(x^j) = -1$.

• when true $y^j = -1$ & $f(x^j) = +1$, $y^j f(x^j) = 1$

• when $y^j = +1$ & $f(x^j) = +1$, $y^j f(x^j) = 1$.

We can see, that $y^j f(x^j) \geq 0$ means no error and
 $y^j f(x^j) < 0$ means there is some error.

$$\therefore E_{\text{train}} = \frac{1}{N} \sum_{j=1}^N \begin{cases} 1, & \text{if } y^j f(x^j) \leq 0 \\ 0, & \text{otherwise} \end{cases}$$

Now, for some z , $\exp(-z) \geq 1$ if $z \leq 0$

~~and $\exp(-z) < 1$ if $z > 0$~~

$$\therefore \exp(-y^j f(x^j)) \geq 1 \text{ when } y^j f(x^j) \leq 0.$$

$$\therefore E_{\text{train}} \leq \frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j))$$

(2) Show that: $\frac{1}{N} \sum_{j=1}^N \exp(-f(x_j) y_j) = \prod_{t=1}^T Z_t$.

→ It is given that:

$$w_j^{(t+1)} = \frac{w_j^{(t)} \exp(-\alpha_t y_j h_t(x_j))}{Z_t}$$

$$\text{where } Z_t = \sum_{j=1}^N w_j^{(t)} \exp(-\alpha_t y_j h_t(x_j))$$

Let us consider 1 sample.

As y_j and $h_t(x_j) \in \{-1, 1\}$.

$$w_j^{(t+1)} = w_j^{(1)} \cdot \frac{\exp(-\alpha_1 y_j h_1(x_j))}{Z_1} \dots \frac{\exp(-\alpha_T y_j h_T(x_j))}{Z_T}$$

But, we initialize $w_j^{(1)} = \frac{1}{N}$.

$$\begin{aligned} \therefore w_j^{(t+1)} &= \frac{1}{N} \cdot \frac{\exp(-\alpha_1 y_j h_1(x_j))}{Z_1} \dots \frac{\exp(-\alpha_T y_j h_T(x_j))}{Z_T} \\ &= \frac{1}{N} \frac{\exp\left(-y_j \sum_{t=1}^T \alpha_t h_t(x_j)\right)}{\prod_{t=1}^T Z_t} \end{aligned}$$

Given that, $f(x_j) = \sum_{t=1}^T \alpha_t h_t(x_j)$.

$$\therefore w_j^{(t+1)} = \frac{1}{N} \frac{\exp(-y_j f(x_j))}{\prod_{t=1}^T Z_t}$$

Now, let us extend for N data points:

$$\sum_{j=1}^N w_j^{(t+1)} = \left[\frac{1}{N} \right] \left[\frac{\sum_{j=1}^N \exp(-y_j f(x_j))}{\prod_{t=1}^T Z_t} \right]$$

Now, $w_j^{(t+1)}$ is a normalizing factor. i.e. $\sum_{j=1}^N w_j^{(t+1)} = 1$.

$$\therefore 1 = \frac{\frac{1}{N} \sum_{j=1}^N \exp(-y^j f(x^j))}{\prod_{j=1}^T Z_t}$$

$$\text{ i.e. } \frac{1}{N} \sum_{j=1}^N \exp(-f(x^j) y^j) = \prod_{j=1}^T Z_t$$

(3)

(a) Find the value of α_t that minimizes Z_t , and then

$$\text{ST: } Z_t^{\text{opt}} = 2\sqrt{\epsilon_t(1-\epsilon_t)}.$$

→ We know: ~~Z_t~~ $Z_t = (1-\epsilon_t) \exp(-\alpha_t) + \epsilon_t \exp(\alpha_t)$

(i) Finding optimal α_t .

$$\frac{\partial Z_t}{\partial \alpha} = 0 \Rightarrow 0 = -(1-\epsilon_t) \exp(-\alpha_t) + \epsilon_t \exp(\alpha_t)$$

$$(1-\epsilon_t) \exp(-\alpha_t) = \epsilon_t \exp(\alpha_t)$$

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

$$\therefore \alpha_t = \ln \sqrt{\frac{1-\epsilon_t}{\epsilon_t}}$$

(ii) Subs. optimal α_t in Z_t to get Z_t^{opt} .

$$\therefore Z_t^{\text{opt}} = (1-\epsilon_t) \exp\left(-\ln \sqrt{\frac{1-\epsilon_t}{\epsilon_t}}\right) + \epsilon_t \exp\left(\ln \sqrt{\frac{1-\epsilon_t}{\epsilon_t}}\right)$$

$$= \frac{(1-\epsilon_t)}{\sqrt{1-\epsilon_t}} \sqrt{\epsilon_t} + \epsilon_t \frac{\sqrt{1-\epsilon_t}}{\sqrt{\epsilon_t}}$$

$$= \frac{(1-\epsilon_t)\epsilon_t + \epsilon_t(1-\epsilon_t)}{(\sqrt{1-\epsilon_t})(\sqrt{\epsilon_t})}$$

$$Z_t^{\text{opt}} = \frac{2\epsilon_t(1-\epsilon_t)}{\sqrt{\epsilon_t(1-\epsilon_t)}} = \underline{2\sqrt{\epsilon_t(1-\epsilon_t)}}$$

(b) From part (a)

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

But, $\epsilon_t = \frac{1}{2} - r_t$.

$$\begin{aligned}\therefore Z_t &= 2\sqrt{\cancel{\epsilon_t(1-\epsilon_t)}} \sqrt{\left(\frac{1}{2}-r_t\right)\left(1-\frac{1}{2}+r_t\right)} \\ &= 2\sqrt{\left(\frac{1}{2}-r_t\right)\left(\frac{1}{2}+r_t\right)}\end{aligned}$$

Now $(a+b)(a-b) = a^2 - b^2$.

$$\begin{aligned}\therefore Z_t &= 2\sqrt{\left(\frac{1}{2}\right)^2 - r_t^2} \\ &= 2\sqrt{\left(\frac{1}{4}\right) - r_t^2} \\ &= \sqrt{\frac{4}{4} - 4r_t^2} \\ &= \sqrt{1 - 4r_t^2}.\end{aligned}$$

But, $\log(1-x) \leq -x$ for $0 \leq x < 1$

$$\therefore 1-x \leq e^{-x}.$$

$$\therefore \sqrt{1-x} \leq \sqrt{e^{-x}}$$

$$\therefore \sqrt{1-x} \leq e^{-x/2}$$

$$\begin{aligned}\therefore \sqrt{1-4r_t^2} &\leq e^{-\frac{4r_t^2}{2}} \\ &\leq \exp(-2r_t^2)\end{aligned}$$

$$\therefore \underline{Z_t \leq \exp(-2r_t^2)}.$$

(c) We have:

$$\epsilon_{\text{train}} \leq \exp \left(-2 \sum_{t=1}^T r_t^2 \right)$$

we want to show that if each classifier is better than random, then: $\epsilon_{\text{train}} \leq \exp(-2Tr^2)$ $\begin{cases} r_t \geq r, \forall t \\ r > 0 \end{cases}$

→ We know that $r_t > 0$ implies better than random.

i.e. $r_t \geq r$, where $r > 0, \forall t$.

i.e. $r_t^2 \geq r^2$.

Now, we want an upper bound on ϵ_{train} .

∴ we can replace r_t^2 by r^2 .

$$\text{i.e. } \epsilon_{\text{train}} \leq \exp \left(-2 \sum_{t=1}^T r^2 \right)$$

$$\underline{\epsilon_{\text{train}} \leq \exp(-2Tr^2)}$$

II. Action Recognition with CNN

Submitted ipython notebook in zip

1. **First TODO: define model** - in notebook
2. **Second TODO: define loss function and optimizer** - in notebook
3. **Third TODO: Train model and show accuracy**- in notebook
4. **Fourth TODO: Tell us what you did:**
I tried the VGG-19 from ResNet paper: <https://arxiv.org/pdf/1512.03385.pdf>. The convolutional layers have a kernel_size of 3 throughout. After successive conv, activations, there are maxpool layers, with a kernel size of 2, which reduce the size of output by 2. After every Max pooling layer, the number of filters is doubled. At the end I have used three fully connected layers
5. **Fifth TODO: Submit score to Kaggle**
Kaggle Score: 61.03975
6. **Sixth TODO: Flatten 3d** - In notebook
7. **Seventh TODO: Design a network using 3D convolution on videos for video classification** - In notebook
8. **Eighth TODO: Tell us what you did:**
I have used a two layer network with a fixed filter size of 3. The number of filters for 1st conv layer is 8 with padding 2, and then it is doubled (16) for the next conv layer, again with padding 2. After each convolutional layer, I have used BatchNorm. After RELU activations, I have used Max Pooling with a kernel size of 2 and stride 2.

Kaggle Score: 58.10397