

**1. Feature Selection.**

- (a) i. Let  $k = 1$ . Then since  $x_1$  and  $x_{100}$  are uncorrelated, we simply pick our feature as the feature which contributes more to the signal, namely  $x_{100}$ . Our predictor is then

$$h_w(x) = ax_{100}$$

And we pick  $a = \frac{3}{\sqrt{10}}$ . Our error is then:

$$L_D(h_w) = E[x_1^2/10] = \frac{1}{10} \text{Var}(x_1) = \frac{1}{10}$$

For  $k = 2$  and above, we simply pick  $x_1$  and  $x_{100}$  as our features,  $w = \frac{3}{\sqrt{10}}e_{100} + \frac{1}{\sqrt{10}}e_1$ . We get zero loss. Therefore, we need  $k = 2$  to get loss less than 0.01.

- ii. For  $k = 1$  we get the same predictor as the optimal feature selection.

For  $k = 2$  we would select  $x_{100}$  first and then select  $x_1$  as the second feature.

For  $k > 2$ , we would only select  $x_{100}$  and  $x_1$  as our features, since adding any feature would increase our loss.

Therefore, the optimal feature selection is the same as the feature selection with  $k = 2$ , and we require  $k = 2$  to get loss less than 0.01.

- iii. We solve the following optimization problem for a fixed  $B$ .

$$\arg \min_w E \left( \frac{3}{\sqrt{10}}x_{100} + \frac{1}{\sqrt{10}}x_1 - \langle w, x \rangle \right)^2 \quad \text{s.t.} \quad \|w\|_1 \leq B$$

Where the expectation is taken over the distribution of  $x$ . It is clear that any non-zero coefficient on  $w_i$ , for  $i \neq 1, 100$  is inefficient, since it pointlessly increases our loss and increases our  $\ell^1$  norm. Therefore, the support of  $w$  is at most  $\{1, 100\}$ , for all values of  $B$ , and we can solve the optimization problem by solving the following optimization problem:

$$\arg \min_{w_1, w_{100}} E \left( \left( \frac{3}{\sqrt{10}} - w_{100} \right) x_{100} + \left( \frac{1}{\sqrt{10}} - w_1 \right) x_1 \right)^2 \quad \text{s.t.} \quad |w_1| + |w_{100}| \leq B$$

We can expand the objective function as follows:

$$\begin{aligned} E(\ell) &= E \left( \left( \frac{3}{\sqrt{10}} - w_{100} \right) x_{100} + \left( \frac{1}{\sqrt{10}} - w_1 \right) x_1 \right)^2 \\ &= E \left( \left( \frac{3}{\sqrt{10}} - w_{100} \right)^2 x_{100}^2 + \left( \frac{1}{\sqrt{10}} - w_1 \right)^2 x_1^2 + 2 \left( \frac{3}{\sqrt{10}} - w_{100} \right) \left( \frac{1}{\sqrt{10}} - w_1 \right) x_1 x_{100} \right) \\ &= \left( \frac{3}{\sqrt{10}} - w_{100} \right)^2 E[x_{100}^2] + \left( \frac{1}{\sqrt{10}} - w_1 \right)^2 E[x_1^2] + 2 \left( \frac{3}{\sqrt{10}} - w_{100} \right) \left( \frac{1}{\sqrt{10}} - w_1 \right) E[x_1 x_{100}] \\ &= \left( \frac{3}{\sqrt{10}} - w_{100} \right)^2 + \left( \frac{1}{\sqrt{10}} - w_1 \right)^2 + 2 \left( \frac{3}{\sqrt{10}} - w_{100} \right) \left( \frac{1}{\sqrt{10}} - w_1 \right) \cdot 0 \\ &= \left( \frac{3}{\sqrt{10}} - w_{100} \right)^2 + \left( \frac{1}{\sqrt{10}} - w_1 \right)^2 \end{aligned}$$

From this it is clear that  $w_1$  and  $w_{100}$  are always positive. Additionally, for  $B \leq \frac{4}{\sqrt{10}}$ , the constraint is active, so we can write our constraints as:

$$\begin{aligned} w_1 &\geq 0 \\ w_{100} &\geq 0 \\ w_1 + w_{100} &= B \end{aligned}$$

For the case where the optimal solution is on the interior of the constraint, we solve the following system of equations:

$$w_1 + w_{100} = B$$

$$w_1 = w_{100} - \frac{2}{\sqrt{10}}$$

This gives us the solution:

$$\begin{aligned} w_1 &= \frac{B}{2} - \frac{1}{\sqrt{10}} \\ w_{100} &= \frac{B}{2} + \frac{1}{\sqrt{10}} \end{aligned}$$

When the solution is an extreme point, we have the solution:

$$\begin{aligned} w_1 &= 0 \\ w_{100} &= B \end{aligned}$$

And this occurs when

$$B \leq \frac{2}{\sqrt{10}}$$

Therefore, when  $B \leq \frac{2}{\sqrt{10}}$ , we have  $k = 1$ , and when  $B > \frac{2}{\sqrt{10}}$ , we have  $k = 2$ . For  $k = 1$ , we have the same as the optimal solution, namely  $w = \frac{3}{\sqrt{10}}e_{100}$ , and the loss is  $\frac{1}{10}$ . For  $\frac{2}{\sqrt{10}} \leq B \leq \frac{4}{\sqrt{10}}$ , we have:

$$w = \left(\frac{B}{2} + \frac{1}{\sqrt{10}}\right)e_{100} + \left(\frac{B}{2} - \frac{1}{\sqrt{10}}\right)e_1$$

And the loss is:

$$\begin{aligned} L_D(h_w) &= \left(\frac{3}{\sqrt{10}} - \frac{B}{2} - \frac{1}{\sqrt{10}}\right)^2 + \left(\frac{1}{\sqrt{10}} - \frac{B}{2} + \frac{1}{\sqrt{10}}\right)^2 \\ &= 2\left(\frac{B}{2} - \frac{2}{\sqrt{10}}\right)^2 \end{aligned}$$

When  $B > \frac{4}{\sqrt{10}}$ , we have  $k = 2$ , and the loss is zero.

iv. We calculate the correlation coefficient for each  $x_i$  and  $y$ . For  $i = 1$ , we have:

$$\begin{aligned} \rho_1 &= \frac{E[x_1 y]}{\sqrt{E[x_1^2]E[(y - E(y))^2]}} \\ &= \frac{E\left[\frac{1}{\sqrt{10}}x_1^2 + \frac{3}{\sqrt{10}}x_{100}x_1\right]}{\sqrt{\text{Var}(x_1)\text{Var}(y)}} \\ &= \frac{\frac{1}{\sqrt{10}}E[x_1^2] + \frac{3}{\sqrt{10}}E[x_{100}]E[x_1]}{\sqrt{\text{Var}\left(\frac{1}{\sqrt{10}}x_1 + \frac{3}{\sqrt{10}}x_{100}\right)}} \\ &= \frac{\frac{1}{\sqrt{10}}}{\sqrt{\frac{1}{10} + \frac{9}{10}}} \\ &= \frac{1}{\sqrt{10}} \end{aligned}$$

Note that for each  $i$ , the denominator is the same, 1.

For  $i = 100$ , we have:

$$\begin{aligned} \rho_{100} &= \frac{E[x_{100} y]}{\sqrt{E[x_{100}^2]E[(y - E(y))^2]}} \\ &= E\left[\frac{1}{\sqrt{10}}x_1x_{100} + \frac{3}{\sqrt{10}}x_{100}^2\right] \\ &= \frac{3}{\sqrt{10}} \end{aligned}$$

for  $i \neq 1, 100$ , we have:

$$\begin{aligned}\rho_i &= \frac{E[x_i y]}{\sqrt{E[x_i^2]E[(y - E(y))^2]}} \\ &= E\left[\frac{1}{\sqrt{10}}x_i x_1 + \frac{3}{\sqrt{10}}x_i x_{100}\right] \\ &= \frac{3}{\sqrt{10}}E[x_i x_{100}] \\ &= \frac{3}{\sqrt{10}} \cdot \frac{9}{10} = \frac{2.7}{\sqrt{10}}\end{aligned}$$

Therefore, if  $k = 1$ , we select  $x_{100}$  as our feature, which is the ideal feature. If  $k = 2$ , we select  $x_{100}$  and any other feature different from  $x_1$ . For  $k = 3, \dots, 99$ , we select any features different from  $x_1$ . Only for  $k = 100$ , do we finally select  $x_1$ .

For  $k < 100$ , we have the same loss as if we only selected  $x_{100}$ , since the coefficient of the features we select, other than  $x_{100}$ , is zero. We thus get loss  $\frac{1}{10}$ , and only for  $k = 100$  do we get loss 0.

- (b) i. Let  $k = 1$ . Then we compute the loss for  $x_1, x_2$ , and  $x_i$ , for  $i \neq 1, 2$  as the chosen features.

$$\begin{aligned}E[(y - az_1)^2] &= E[(z_2 - az_1)^2] = E[(z_1 - az_2)^2] - E[z_1 - az_2]^2 = \text{Var}[z_1 - az_2] \\ &= \text{Var}[z_1] + a^2 \text{Var}[z_2] = 1 + a^2 \\ \implies a &= 0, \quad L_D(h_{w_1}) = 1\end{aligned}$$

For  $k = 2$ , we replace the coefficient of  $z_2$  with  $b = 0.0001$ , so we can reuse the same calculations for  $z_i$ .

$$\begin{aligned}E[(y - az_1 - abz_2)^2] &= E[(z_2 - az_1 - abz_2)^2] = \text{Var}(z_2 - az_1 - abz_2) \\ &= (1 - ab)^2 + a^2 \\ \frac{\partial L}{\partial a}(1 - ab)^2 + a^2 &= 0 \\ 0 &= 2(1 - ab)(-b) + 2a \\ a &= b - ab^2 \\ b &= a(1 + b^2) \\ a &= \frac{1}{b^{-1} + b} \\ \implies L &= \left(1 - \frac{b}{b^{-1} + b}\right)^2 + \left(\frac{b}{b^{-1} + b}\right)^2 \\ &= \frac{1}{b^{-2} + 1}\end{aligned}$$

Therefore, selecting  $x_2$  our feature we have that our loss is:

$$L_D(h_{w_2}) = \frac{1}{10^{-8} + 1}$$

Since the  $z_i$  have no covariance between each other, we have the same logic for  $x_i$ . Selecting  $x_i$  for  $i \neq 1, 2$ , we have:

$$L_D(h_{w_i}) = \frac{1}{10^{-6} + 1}$$

Therefore, for  $k = 1$ , we would select any  $x_i$  for  $i \neq 1, 2$ .

For  $k = 2$ , we would select  $x_1$  and  $x_2$  as our features, since we get zero loss by selecting  $w = -10^4 e_1 + 10^4 e_2$ . This gives us:

$$\begin{aligned}h_w(x) &= -10^4 x_1 + 10^4 x_2 = -10^4 z_1 + 10^4 (z_1 + 10^{-4} z_2) \\ &= z_2\end{aligned}$$

Therefore, we have:

$$L_D(h_w) = E[(y - h_w(x))^2] = E[(z_2 - z_2)^2] = 0$$

Thus, for  $k > 2$ , we similarly select  $x_1$  and  $x_2$  as our features.

- ii. For greedy feature selection we select  $x_i$ ,  $i > 2$  as our first feature, since they have the lowest loss. We now calculate the loss for  $k = 2$ , given that we've selected some  $x_i$  as our first feature.

We begin with  $x_1$  as our first second feature.

$$\begin{aligned}\ell &= E[(z_2 - az_1 - bz_i - b10^{-3}z_2)^2] \\ &= \text{Var}((1 - b10^{-3})z_2 - az_1 - bz_i) \\ &= (1 - b10^{-3})^2 + a^2 + b^2\end{aligned}$$

Since  $a^2$  is non-negative, we set  $a = 0$  to minimize the loss, and we choose  $b$  as in the optimal case for  $k = 1$ . We check  $x_2$  as our second feature.

$$\begin{aligned}\ell &= E[(z_2 - az_1 - a10^{-4}z_2 - bz_i - b10^{-3}z_2)^2] \\ &= \text{Var}((1 - a10^{-4} - b10^{-3})z_2 - az_1 - bz_i) \\ &= (1 - a10^{-4} - b10^{-3})^2 + a^2 + b^2\end{aligned}$$

We check  $x_j$  where  $j \neq i$  and  $j > 2$  as our second feature. With the same logic, we have:

$$\ell = (1 - (a + b)10^{-3})^2 + a^2 + b^2$$

Clearly, for every choice of  $a$  and  $b$ , we have that the loss is minimized when we select  $x_j$  as our next feature. Continuing this logic, for  $k = 3, \dots, 98$ , we select  $x_i$  where  $i > 2$  as our features. Our loss is then:

$$\ell = \left(1 - \left(\sum_{i=1}^k a_{i+2}\right)10^{-3}\right)^2 + \sum_{i=1}^k a_{i+2}^2$$

Taking the critical points, we have:

$$\begin{aligned}[a_j] : & -2 \left(1 - 10^{-3} \sum_{i=1}^k a_{i+2}\right) 10^{-3} + 2a_j = 0 \\ \implies a_j &= 10^{-3} \left(1 - 10^{-3} \sum_{i=1}^k a_{i+2}\right) \\ a_j &= 10^{-3} - 10^{-6}ka_j \\ a_j &= \frac{10^{-3}}{1 + 10^{-6}k}\end{aligned}$$

Where we have all the  $a_j$ s are equal by the second line. Therefore, the loss is:

$$\begin{aligned}L_D(h_w) &= \left(1 - 10^{-3} \cdot k \cdot \frac{10^{-3}}{1 + 10^{-6}k}\right)^2 + k \cdot \left(\frac{10^{-3}}{1 + 10^{-6}k}\right)^2 \\ &= \left(1 - \frac{10^{-6}k}{1 + 10^{-6}k}\right)^2 + \frac{10^{-6}k}{(1 + 10^{-6}k)^2} \\ &= \frac{1}{(1 + 10^{-6}k)^2} + \frac{10^{-6}k}{(1 + 10^{-6}k)^2} \\ &= \frac{1 + 10^{-6}k}{(1 + 10^{-6}k)^2} \\ &= \frac{1}{1 + 10^{-6}k}\end{aligned}$$

For  $k = 99$ , we select  $x_2$  as our feature, since adding  $x_1$  wouldn't decrease our loss. We get the following loss:

$$\begin{aligned}
 L_D(h_w) &= E \left[ \left( z_2 - 10^{-3} \sum_{i=3}^{100} a_i z_2 - 10^{-4} a_2 z_2 + \sum_{i=3}^{100} a_i z_i + a_2 z_1 \right)^2 \right] \\
 &= \text{Var} \left( \left( 1 - 10^{-3} \sum_{i=3}^{100} a_i - 10^{-4} a_2 \right) z_2 + \sum_{i=3}^{100} a_i z_i + a_2 z_1 \right) \\
 &= \left( 1 - 10^{-3} \sum_{i=3}^{100} a_i - 10^{-4} a_2 \right)^2 + \sum_{i=2}^{100} a_i^2
 \end{aligned}$$

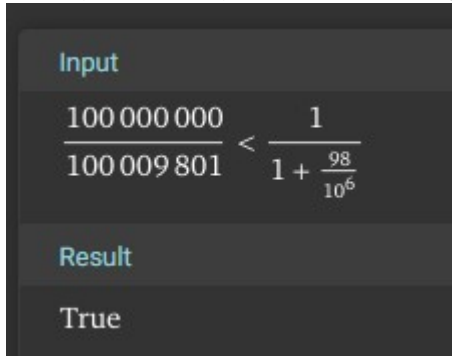
Taking critical points with  $a_2$  and  $a_i$ ,  $i > 2$ , we have:

$$\begin{aligned}
 [a_i] : & -2 \left( 1 - 10^{-3} \sum_{i=3}^{100} a_i - 10^{-4} a_2 \right) 10^{-3} + 2a_i = 0 \\
 \Rightarrow & a_i = 10^{-3} \left( 1 - 10^{-3} \sum_{i=3}^{100} a_i - 10^{-4} a_2 \right) \\
 [a_2] : & -2 \left( 1 - 10^{-3} \sum_{i=3}^{100} a_i - 10^{-4} a_2 \right) 10^{-4} + 2a_2 = 0 \\
 \Rightarrow & a_2 = 10^{-4} \left( 1 - 10^{-3} \sum_{i=3}^{100} a_i - 10^{-4} a_2 \right) \\
 \Rightarrow & a_2 = 10^{-4} (1 - 10^{-3} 98a_i - 10^{-4} a_2) \\
 \Rightarrow & a_i = 10^{-3} (1 - 10^{-3} 98a_i - 10^{-4} a_2) \\
 \Rightarrow & a_2 = \frac{10^5}{10^9 + 98 \cdot 10^2 + 1} \\
 \Rightarrow & a_i = \frac{10^6}{10^9 + 98 \cdot 10^2 + 1}
 \end{aligned}$$

Plugging these values into our loss function, we get:

$$\begin{aligned}
 L_D(h_w) &= (1 - 10^{-3} 98 \cdot a_i - 10^{-4} a_2)^2 + 98 \cdot a_i^2 + a_2^2 \\
 &= 0.999901...
 \end{aligned}$$

Which, to be clear, is better than the previous results (verified by wolfram alpha).



The screenshot shows a Wolfram Alpha interface. Under the 'Input' section, the expression  $\frac{100\,000\,000}{100\,009\,801} < \frac{1}{1 + \frac{98}{10^6}}$  is entered. Under the 'Result' section, the output is 'True'.

This is still (obviously) worse than the target of 0.01. Therefore, the only  $k$  for which we get loss lower than 0.01 is  $k = 100$ , in which we finally are able to select  $x_1$  and  $x_2$ , and we let  $w = -10^4 e_1 + 10^4 e_2$ . As shown above, this achieves zero loss.

iii. We can write the loss function as follows:

$$\begin{aligned}
 L_D(h_w) &= E \left[ \left( z_2 - w_1 z_1 - w_2 z_1 - w_2 10^{-4} z_2 - 10^{-3} \sum_{k=3}^{100} w_k z_k \right)^2 \right] \\
 &= \text{Var} \left[ \left( 1 - 10^{-4} w_2 - 10^{-3} \sum_{k=3}^{100} w_k \right) z_2 - (w_1 + w_2) z_1 - \sum_{k=3}^{100} w_k z_k \right] \\
 &= \left( 1 - 10^{-4} w_2 - 10^{-3} \sum_{k=3}^{100} w_k \right)^2 + (w_1 + w_2)^2 + \sum_{k=3}^{100} w_k^2
 \end{aligned}$$

Notice that for  $k > 2$ , each of the  $w_k$ 's contribute the same to decreasing the loss, and for increasing the  $\ell^1$  norm of  $w$ . Therefore, without loss of generality, we can assume that they are all equal. Let  $w_1 = x$ ,  $w_2 = y$ , and  $w_3 = z$ . Then we have the following

$$L_D(h_w) = (1 - 10^{-4}y - 10^{-3}98z)^2 + (x + y)^2 + 98z^2$$

We set up the following Lagrangian:

$$\mathcal{L}(x, y, z, \lambda) = (1 - 10^{-4}y - 10^{-3}98z)^2 + (x + y)^2 + 98z^2 - \lambda(|x| + |y| + 98|z| - B)$$

In order to find the optimal point, we have to investigate all 8 cases for the signs of  $x, y, z$ . However, we can do some simplification. Notice that the loss is minimized when  $x$  is the opposite sign as  $y$ , since the term

$$(x + y)^2$$

is minimized when  $x$  and  $y$  are opposite signs. Additionally, the only time when  $z$  decreases the loss is when  $z$  is positive. Similarly, the only time that  $y$  minimizes the loss is when  $y$  is positive. Therefore, we really only have one case to consider: that is when  $x$  is negative,  $y$  is positive, and  $z$  is positive. Thus, we have the following Lagrangian:

$$\mathcal{L}(x, y, z, \lambda) = (1 - 10^{-4}y - 10^{-3}98z)^2 + (x + y)^2 + 98z^2 - \lambda(-x + y + 98z - B)$$

We can solve this system of equations to get the solutions, but another way we can think about this is that when  $B > 2 \cdot 10^4$ , we have zero loss, since we can select  $w_1 = -10^4$  and  $w_2 = 10^4$ . When  $B < 2 \cdot 10^4$ , we still want to select  $w_1 = -w_2$ , since this minimizes the loss, because we get an unbiased estimate of  $z_2$ . When we select  $w_1 = -w_2$  as our only features, we have the following loss:

$$L_D(h_w) = (1 - 10^{-4}y)^2$$

And since

$$2|y| \leq B$$

We can manage zero loss as long as  $1 - 10^{-4}y = 0$  is feasible, and this occurs when  $B > 2 \cdot 10^4$ . When  $B < 2 \cdot 10^4$  and we keep the same features, we have the following loss:

$$\begin{aligned}
 L_D(h_w) &= (1 - 10^{-4}y)^2 \\
 &= \left( 1 - 10^{-4} \frac{B}{2} \right)^2
 \end{aligned}$$

If we select any  $x_i$ ,  $i > 2$ , we we have the following problem:

$$\begin{aligned}
 L_D(h_w) &= (1 - 10^{-3}z)^2 + z^2 \\
 B &\geq |z|
 \end{aligned}$$

This is minimized always at  $z = \frac{10^3}{10^7+1}$ , and we get the loss  $\frac{1}{1+10^{-7}}$ . However, for  $B < 2 \cdot 10^4$ , in the other case we get that the loss is always less than 1. The only time when we get loss less than 0.01 is when  $B > 9000$ , and we select  $x_1$  and  $x_2$  as our features.

iv. We compute the various correlation coefficients. First, notice that:

$$\begin{aligned}\text{Var}(y) &= \text{Var}(z_2) = 1 \\ \text{Var}(x_1) &= \text{Var}(z_1) = 1 \\ \text{Var}(x_2) &= \text{Var}(z_1 + 10^{-4}z_2) = 1 + 10^{-8} \\ \text{Var}(x_i) &= \text{Var}(z_i + 10^{-3}z_2) = 1 + 10^{-6}\end{aligned}$$

Therefore, we have:

$$\begin{aligned}\rho_1 &= \frac{E[x_1 y]}{\sqrt{\text{Var}(x_1) \text{Var}(y)}} = \text{Cov}(x_1, y) \\ &= \text{Cov}(z_1, z_2) = 0 \\ \rho_2 &= \frac{E[x_2 y]}{\sqrt{\text{Var}(x_2) \text{Var}(y)}} = \frac{\text{Cov}(x_2, y)}{\sqrt{1 + 10^{-8}}} \\ &= \frac{\text{Cov}(z_1 + 10^{-4}z_2, z_2)}{\sqrt{1 + 10^{-8}}} \\ &= \frac{10^{-4}}{\sqrt{1 + 10^{-8}}} \\ \rho_i &= \frac{E[x_i y]}{\sqrt{\text{Var}(x_i) \text{Var}(y)}} = \frac{\text{Cov}(x_i, y)}{\sqrt{1 + 10^{-6}}} \\ &= \frac{\text{Cov}(z_i + 10^{-3}z_2, z_2)}{\sqrt{1 + 10^{-6}}} \\ &= \frac{10^{-3}}{\sqrt{1 + 10^{-6}}}\end{aligned}$$

Therefore, for  $k < 99$  we will always choose  $x_3, \dots, x_{100}$  as our features, since they have the highest correlation with  $y$ . As before, we have to wait until  $k = 99$  to select  $x_2$ , and we have to wait until  $k = 100$  to select  $x_1$ .

Therefore, we have the same loss as in part (b), and so we have to wait until  $k = 100$  to get loss less than 0.01.

## 2. Boosting as Coordinate Descent.

(a) We can write the loss as follows:

$$\begin{aligned}L_{D^{(t+1)}}(h_t) &= \sum_{i=1}^m D_i^{(t+1)} [[y_i \neq h_t(x_i)]] \\ &= \frac{\sum_{i=1}^m D_i^{(t)} \exp(\alpha_t) [[y_i \neq h_t(x_i)]]}{\sum_{j=1}^m D_j^{(t)} \exp(-\alpha_t y_j h_t(x_j))} \\ &= \frac{\sum_{i=1}^m \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} [[y_i \neq h_t(x_i)]]}{\sum_{j=1}^m D_j^{(t)} \exp(-\alpha_t y_j h_t(x_j))} \\ &= \frac{\sum_{i=1}^m \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} [[y_i \neq h_t(x_i)]]}{\sum_{j=1}^m \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} [[y_j \neq h_t(x_j)]] + \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} [[y_j = h_t(x_j)]]} \\ &= \frac{1}{1 + \frac{\sum_{j=1}^m \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} [[y_j = h_t(x_j)]]}{\sum_{i=1}^m \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} [[y_i \neq h_t(x_i)]]}} \\ &= \frac{\sum_{j=1}^m \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} [[y_j = h_t(x_j)]]}{\sum_{i=1}^m \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} [[y_i \neq h_t(x_i)]]} = \frac{(1 - \epsilon_t) \cdot m \sqrt{\frac{\epsilon_t}{1-\epsilon_t}}}{\epsilon_t \cdot m \sqrt{\frac{1-\epsilon_t}{\epsilon_t}}} \\ &= \frac{1-\epsilon_t}{\epsilon_t}\end{aligned}$$

$$= \frac{1 - \epsilon_t}{\epsilon_t} \cdot \sqrt{\frac{\epsilon_t}{1 - \epsilon_t}} = 1$$

Therefore,

$$L_{D^{(t+1)}}(h_t) = \frac{1}{2}$$

(b) Let  $h_t[h] := h \sum_{k \leq t | h_k = h} \alpha_k$

$$\begin{aligned} \frac{\partial L_S^{exp}}{\partial w[h]}(h_{w^{(t)}}) &= \frac{\partial}{\partial w[h]} \frac{1}{m} \sum_{i=1}^m \exp(-y_i h_{w^{(t)}}(x_i)) \\ &= \frac{1}{m} \sum_{i=1}^m \exp(-y_i h_{w^{(t)}}(x_i)) \left( -y_i \frac{\partial h_{w^{(t)}}(x_i)}{\partial w[h]} \right) \\ &= -\frac{1}{m} \sum_{i=1}^m \exp(-y_i h_{w^{(t)}}(x_i)) y_i h_t[h](x_i) \\ &= -\frac{1}{m} \sum_{i=1}^m \exp\left(-y_i \sum_{k=1}^{t-1} \alpha_k h_k(x_i)\right) y_i h_t[h](x_i) \\ &= -\frac{1}{m} \sum_{i=1}^m \prod_{k=1}^{t-1} \exp(-y_i \alpha_k h_k(x_i)) y_i h_t[h](x_i) \end{aligned}$$

Recall the definition of  $D_i^{(t)}$ :

$$\begin{aligned} D_i^{(t)} &= \frac{D_i^{(t-1)} \exp(-y_i \alpha_t h_t(x_i))}{\sum_{j=1}^m D_j^{(t-1)} \exp(-\alpha_t y_j h_t(x_j))} \\ &= \frac{D_i^{(t-1)} \exp(-y_i \alpha_t h_t(x_i))}{C^{(t-1)}} \\ D_i^{(t)} \prod_{k=1}^{t-1} C^{(k)} &= \prod_{k=1}^{t-1} D_i^{(1)} \exp(-y_i \alpha_k h_k(x_i)) = \frac{1}{m} \prod_{k=1}^{t-1} \exp(-y_i \alpha_k h_k(x_i)) \end{aligned}$$

Therefore, writing  $\prod_{k=1}^{t-1} C^{(k)} = C_t$ , we have:

$$\frac{\partial L_S^{exp}}{\partial w[h]}(h_{w^{(t)}}) = -C_t \sum_{i=1}^m D_i^{(t)} (y_i h_t[h](x_i))$$

We then have:

$$\begin{aligned} \sum_{i=1}^m D_i^{(t)} y_i h_t[h](x_i) &= \sum_{i=1}^m D_i^{(t)} ([y_i = h_t(x_i)] - [y_i \neq h_t(x_i)]) \\ \sum_{i=1}^m D_i^{(t)} y_i h_t[h](x_i) &= \sum_{i=1}^m D_i^{(t)} (1 - 2[y_i \neq h_t(x_i)]) \\ &= 1 - 2L_{D^{(t)}}(h_t[h]) \end{aligned}$$

Therefore, we have:

$$\begin{aligned} \frac{\partial L_S^{exp}}{\partial w[h]}(h_{w^{(t)}}) &= C_t (2L_{D^{(t)}}(h_t[h]) - 1) \\ &= 2C_t (L_{D^{(t)}}(h_t[h]) - \frac{1}{2}) \end{aligned}$$



$$\propto L_{D^{(t)}}(h_t[h]) - \frac{1}{2}$$