

Solution for Homework Assignment 4

1. Let \mathbf{x} denote the vector of points and \mathbf{y} the vector of labels. Note the $\mathbf{x}^2 = \mathbf{x}^\top \mathbf{x} = \sum_i x_i^2 = \|\mathbf{x}\|_2^2$.

$$\begin{aligned}
 \mathbb{E}[L(w_i)] &= \sum_i \frac{x_i^2}{\mathbf{x}^2} \sum_j \underbrace{\left(\frac{y_i}{x_i} x_j - y_j\right)}_{w_i^*}^2 \\
 &= \frac{1}{\mathbf{x}^2} \sum_i \sum_j (y_i x_j - y_j x_i)^2 \\
 &= \frac{1}{\mathbf{x}^2} \sum_i \sum_j (y_i^2 x_j^2 - 2y_i x_j y_j x_i + y_j^2 x_i^2) \\
 &= \frac{1}{\mathbf{x}^2} \sum_i \sum_j (2y_i^2 x_j^2 - 2y_i x_j y_j x_i) \\
 &= 2\mathbf{y}^2 - \frac{2}{\mathbf{x}^2} \sum_i \sum_j y_i y_j x_i x_j \\
 L(w^*) &= \left(\mathbf{x} \underbrace{\frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^2}}_{w^*} - \mathbf{y}\right)^2 \\
 &= \left(\frac{\mathbf{x}\mathbf{x}^\top}{\mathbf{x}^2} \mathbf{y} - \mathbf{y}\right)^\top \left(\frac{\mathbf{x}\mathbf{x}^\top}{\mathbf{x}^2} \mathbf{y} - \mathbf{y}\right) \\
 &= \mathbf{y}^\top \frac{\mathbf{x}\mathbf{x}^\top}{\mathbf{x}^2} \frac{\mathbf{x}\mathbf{x}^\top}{\mathbf{x}^2} \mathbf{y} - 2\mathbf{y}^\top \frac{\mathbf{x}\mathbf{x}^\top}{\mathbf{x}^2} \mathbf{y} + \mathbf{y}^2 \\
 &= \mathbf{y}^2 - \mathbf{y}^\top \frac{\mathbf{x}\mathbf{x}^\top}{\mathbf{x}^2} \mathbf{y} \\
 &= \mathbf{y}^2 - \frac{1}{\mathbf{x}^2} \sum_i \sum_j y_i y_j x_i x_j
 \end{aligned}$$

2. Let us define the notations:

- t_f : True label.
- y_f : Prediction result.
- E_n : Loss function of the n th data point.
- a_f : Linear summation of the hidden layer output plus the bias.
- w_{fj} : Weights between the hidden layer and output layer.
- b_f : Bias item of the hidden layer output.
- z_f : One final output of output layer after non-linear transformation.
- a_j : Linear summation of the input layer output plus the bias.
- w_{ji} : Weights between the input layer and hidden layer.
- b_j : Bias item of the input layer output.
- z_j : The j th node's output of hidden layer after non-linear transformation.

Based on the chain rule of derivative and the notations mentioned above, we have:

$$a_f = \sum_j w_{fj} z_j + b_f \quad z_f = \Phi(a_f)$$

$$a_j = \sum_i w_{ji} x_i + b_j \quad z_j = \Phi(a_j)$$

So we can calculate the derivatives of weights and bias:

$$(1) : \quad \frac{\partial E_n}{\partial w_{fj}} = \frac{\partial E_n}{\partial a_f} \cdot \frac{\partial a_f}{\partial w_{fj}}$$

$$\delta_f = \frac{\partial E_n}{\partial a_f} = \frac{\frac{1}{2}(z_f - t_f)^2}{\frac{\partial a_f}{\partial z_f}} = \frac{\frac{1}{2}(z_f - t_f)^2}{\frac{\partial z_f}{\partial a_f}} \cdot \frac{\partial z_f}{\partial a_f} = (z_f - t_f) \cdot \Phi'(a_f) = (y_f - t_f) \cdot \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{a_f^2}{2}}$$

$$(\Phi'(a) = (\int_{-\infty}^a \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{z^2}{2}} dz)' = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{a^2}{2}} \quad \text{and} \quad z_f = y_f)$$

$$\frac{\partial a_f}{\partial w_{fj}} = z_j, \quad \text{So} \quad \frac{\partial E_n}{\partial w_{fj}} = (y_f - t_f) \cdot \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{a_f^2}{2}} \cdot z_j$$

$$(2) : \quad \frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \cdot \frac{\partial a_j}{\partial w_{ji}}$$

$$\delta_j = \frac{\partial E_n}{\partial a_j} = \frac{\partial E_n}{\partial a_f} \cdot \frac{\partial a_f}{\partial a_j} = \delta_f \cdot \frac{\partial a_f}{\partial a_j} = \delta_f \cdot \frac{\partial a_f}{\partial z_j} \cdot \frac{\partial z_j}{\partial a_j} = \delta_f \cdot w_{fj} \cdot \sqrt{2\pi} \cdot e^{-\frac{a_f^2}{2}}$$

$$(\Phi'(a) = (\int_{-\infty}^a \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{z^2}{2}} dz)' = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{a^2}{2}})$$

$$\frac{\partial a_j}{\partial w_{ji}} = x_i, \quad \text{So} \quad \frac{\partial E_n}{\partial w_{ji}} = \delta_f \cdot w_{fj} \cdot \sqrt{2\pi} \cdot e^{-\frac{a_f^2}{2}} \cdot x_i$$

$$(3) : \quad \frac{\partial E_n}{\partial b_f} = \frac{\partial E_n}{\partial a_f} \cdot \frac{\partial a_f}{\partial b_f} = \delta_f = (y_f - t_f) \cdot \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{a_f^2}{2}}$$

$$(4) : \quad \frac{\partial E_n}{\partial b_j} = \frac{\partial E_n}{\partial a_j} \cdot \frac{\partial a_j}{\partial b_j} = \delta_j = \delta_f \cdot w_{fj} \cdot \sqrt{2\pi} \cdot e^{-\frac{a_f^2}{2}}$$

3. The definition of max function is as following:

$$h(a) = \max(0, a) = \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{if } a < 0 \end{cases}$$

And the integral function of max function is as following:

$$H(a) = \int \max(0, a) = \begin{cases} \frac{1}{2}a^2 & \text{if } a \geq 0 \\ 0 & \text{if } a < 0 \end{cases}$$

$$= \frac{1}{2} \max(0, a) a$$

Based on the 16/23 of "Bregman divergences for constructing regularizers and losses" in the lecture 5, the formula of matching loss is as following:

$$\triangle_H(w \cdot x, h^{-1}) = H(w \cdot x) - H(h^{-1}(y)) - (w \cdot x - h^{-1}(y))y$$

Since we have $h^{-1}(y) = y$, so we have :

$$\triangle_H(w \cdot x, y) = H(w \cdot x) - H(y) - (w \cdot x - y)y$$

$$\text{So we have : Matching loss} = \frac{1}{2} \max(0, w \cdot x) w \cdot x - \frac{1}{2} \max(0, y) y - (w \cdot x - y)y$$