

$$\begin{aligned}
\text{(1a)} \quad \hat{L}(\theta) &= \frac{1}{N} \| y^{(i)} - (x^{(i)} + \delta^{(i)})^T \theta \|_2^2 \\
&= \frac{1}{N} (Y - (X + \delta)^T \theta)^T (Y - (X + \delta)^T \theta) \\
&= \frac{1}{N} [Y^T Y - Y^T (X + \delta)^T \theta - \theta^T (X + \delta)] + \\
&\quad \theta^T (X + \delta) (X + \delta)^T \theta \\
&= \frac{1}{N} [Y^T Y - 2Y^T (X + \delta)^T \theta + \theta^T (X + \delta)(X + \delta)^T \theta] \\
&= \frac{1}{N} [Y^T Y - 2Y^T (X^T + \delta^T) \theta + \theta^T (XX^T + X\delta^T + \delta\delta^T) \theta].
\end{aligned}$$

since $L(\theta) = \frac{1}{N} [Y^T Y - 2Y^T X^T \theta + \theta^T XX^T \theta]$
so $\hat{L}(\theta) = L(\theta) - \frac{1}{N} [2Y^T \delta^T \theta - \theta^T (2X\delta^T + \delta\delta^T) \theta]$

$$\begin{aligned}
&E\left[\frac{1}{N}(2Y^T \delta^T \theta - \theta^T (2X\delta^T + \delta\delta^T) \theta\right] \\
&= -E\left[\frac{1}{N}(\theta^T \delta\delta^T \theta)\right] = -\frac{1}{N}\theta^T \sigma^2 I \theta = -\frac{1}{N}\theta^T \sigma^2 \theta
\end{aligned}$$

$$\begin{aligned}
\text{so } E_{\delta \sim N}[\hat{L}(\theta)] &= L(\theta) + \frac{1}{N}\theta^T \sigma^2 \theta \\
&= L(\theta) + \frac{1}{N}\sigma^2 \|\theta\|_2^2
\end{aligned}$$

1(c) Under expectation, the loss function will be increased by $\frac{1}{N}\sigma^2\|\theta\|_2^2$ because of the noise.

1(c) $\sigma \rightarrow 0$, $E[\tilde{L}(\theta)] = L(\theta)$, the model will overfit the data.

1(d), $\sigma \rightarrow \infty$, $\tilde{L}(\theta) \approx \sigma^2\|\theta\|_2^2$, the model will underfit the data.

3. Assuming the samples $(X^{(1)}, y^{(1)}), \dots, (X^{(m)}, y^{(m)})$ are iid,

$$\begin{aligned} P(X^{(1)}, \dots, X^{(m)}, y^{(1)}, \dots, y^{(m)} | \theta) &= \prod_{i=1}^m P(X^{(i)}, y^{(i)} | \theta) \\ &= \prod_{i=1}^m (P(X^{(i)} | \theta)) P(y^{(i)} | X^{(i)}, \theta) \end{aligned}$$

$$\underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P(X^{(i)} | \theta) P(y^{(i)} | X^{(i)}, \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P(y^{(i)} | X^{(i)}, \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \text{softmax}_{y^{(i)}}(X^{(i)})$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \left[\frac{e^{a_{y^{(i)}}(X^{(i)})}}{\sum_{j=1}^C e^{a_j(X^{(i)})}} \right]$$

$$= \underset{\theta}{\operatorname{argmax}} \frac{1}{m} \sum_{i=1}^m [a_{y^{(i)}}(X^{(i)}) - \log \sum_{j=1}^C e^{a_j(X^{(i)})}]$$

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \left[\log \sum_{j=1}^C e^{a_j(X^{(i)})} - a_{y^{(i)}}(X^{(i)}) \right]$$

$$\text{so } L(\theta) = \frac{1}{m} \sum_{i=1}^m \left[\log \sum_{j=1}^C e^{a_j(X^{(i)})} - a_{y^{(i)}}(X^{(i)}) \right]$$

Expand $a_j(x) = w_j^T x + b_j$ as below :

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m \left[\log \sum_{j=1}^c e^{w_j^T x^{(i)} + b_j} - (w_j^T x^{(i)} + b_j) \right]$$

Since $L(\theta)$ is linear, we can calculate two terms separately,

For the first term, $D_{w_i} L_i = \text{softmax}_j(x^{(i)}) \cdot x^{(i)}$

For the second term, $D_{b_i} L_i = \begin{cases} -x^{(i)} & \text{if } y^{(i)} = i \\ 0 & \text{o.w.} \end{cases}$

$$\text{so } D_{w_i} L(\theta) = \frac{1}{m} \sum_{i=1}^m \left(\text{softmax}(x^{(i)}) \cdot x^{(i)} - x^{(i)} \cdot \mathbf{1}_{\{y^{(i)}=i\}} \right)$$

Similarly, $D_{b_i} L_i = \text{softmax}_j(x^{(i)})$

$$D_{b_i} L_i = \begin{cases} -1 & \text{if } y^{(i)} = i \\ 0 & \text{o.w.} \end{cases}$$

$$\text{so } D_{b_i} L(\theta) = \frac{1}{m} \sum_{i=1}^m \left(\text{softmax}(x^{(i)}) - \mathbf{1}_{\{y^{(i)}=i\}} \right)$$

Since we are doing argmin to replace argmax,

$$\therefore D_{w_i} L(\theta) = \frac{1}{m} \sum_{i=1}^m x^{(i)} \left(\mathbf{1}_{\{y^{(i)}=i\}} - \text{softmax}(x^{(i)}) \right)$$

$$D_{b_i} L(\theta) = \frac{1}{m} \sum_{i=1}^m \left(\mathbf{1}_{\{y^{(i)}=i\}} - \text{softmax}(x^{(i)}) \right).$$

$$4. \text{ hinge}_{y^{(i)}}(x^{(i)}) = \begin{cases} 0 & y^{(i)}(w^T x^{(i)} + b) > 1 \\ 1 - y^{(i)}(w^T x^{(i)} + b) & \text{o.w.} \end{cases}$$

$$\text{If } y^{(i)}(w^T x^{(i)} + b) > 1, \text{ then } \text{hinge}_{y^{(i)}}(x^{(i)}) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^d$$

$$\text{else, then } \text{hinge}_{y^{(i)}}(x^{(i)}) = -y^{(i)}x^{(i)}$$

$$\text{hinge}_{y^{(i)}}(x^{(i)}) = \mathbb{I}_{\{y^{(i)}(w^T x^{(i)} + b) < 1\}} \odot (-y^{(i)}x^{(i)})$$

$$\text{so, } \nabla_L = \frac{1}{K} \sum_{i=1}^K \mathbb{I}_{\{y^{(i)}(w^T x^{(i)} + b) < 1\}} \odot (-y^{(i)}x^{(i)})$$

$$\text{If } y^{(i)}(w^T x^{(i)} + b) > 1, \quad \nabla_b \text{hinge}_{y^{(i)}}(x^{(i)}) = 0$$

$$\text{else, } \nabla_b \text{hinge}_{y^{(i)}}(x^{(i)}) = -y^{(i)}$$

$$\text{so, } \nabla_b L = \frac{1}{K} \sum_{i=1}^K \mathbb{I}_{\{y^{(i)}(w^T x^{(i)} + b) < 1\}} \odot (-y^{(i)})$$