

$$\begin{aligned}
 1) a. \quad \bar{L}(\theta) &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)} + \delta^{(i)})^T \theta)^2 \\
 &= \frac{1}{N} \sum_{i=1}^N (y^{(i)} - x^{(i)T} \theta + \delta^{(i)T} \theta)^2 \\
 &= \frac{1}{N} \sum_{i=1}^N [(y^{(i)} - x^{(i)T} \theta) + \delta^{(i)T} \theta]^2 \\
 &= \frac{1}{N} \sum_{i=1}^N [(y^{(i)} - x^{(i)T} \theta)^2 + 2 \delta^{(i)T} \theta (y^{(i)} - x^{(i)T} \theta) \\
 &\quad + \theta^T \delta^{(i)} \delta^{(i)T} \theta]
 \end{aligned}$$

Then, by the linearity of the expectation function,

$$\begin{aligned}
 E[\bar{L}(\theta)] &= \frac{1}{N} \sum_{i=1}^N (E[(y^{(i)} - x^{(i)T} \theta)^2] + E[2 \delta^{(i)T} \theta (y^{(i)} - x^{(i)T} \theta)] \\
 &\quad + E[\theta^T \delta^{(i)} \delta^{(i)T} \theta])
 \end{aligned}$$

Since $\theta, y^{(i)}, x^{(i)}$ are independent of $\delta^{(i)}$,

$$\begin{aligned}
 &= \frac{1}{N} \sum_{i=1}^N [(y^{(i)} - x^{(i)T} \theta)^2 + 2 E[\delta^{(i)T}] \theta (y^{(i)} - x^{(i)T} \theta) \\
 &\quad + \theta^T E[\delta^{(i)} \delta^{(i)T}] \theta)
 \end{aligned}$$

For a zero centered Gaussian distribution, $E[\delta^{(i)}] = 0$, and we know $E[\delta^{(i)} \delta^{(i)T}] = \sigma^2 I$, so

$$\begin{aligned}
 &= \bar{L}(\theta) + \frac{1}{N} \sum_{i=1}^N (2 \cdot 0 \cdot \theta (y^{(i)} - x^{(i)T} \theta) + \theta^T (\sigma^2 I) \theta) \\
 &= \bar{L}(\theta) + \sigma^2 \theta^T \theta \\
 &= \bar{L}(\theta) + \sigma^2 \|\theta\|^2
 \end{aligned}$$

b. The addition of noise would regularize by attempting to also minimize the magnitude of θ .

c. If $\sigma \rightarrow 0$, $\bar{L}(\theta) = \bar{L}(\theta) + \sigma^2 \|\theta\|^2 \rightarrow \bar{L}(\theta)$, so there would be no regularization effect.

d. if $\sigma \rightarrow \infty$, $\bar{L}(\theta) = \bar{L}(\theta) + \sigma^2 \|\theta\|^2 \rightarrow \infty$, so it would be impossible to minimize the cost function.

3) We have that

$$\text{softmax}_i(x) = \frac{e^{\tilde{w}_i^T \tilde{x}}}{\sum_{j=1}^c e^{\tilde{w}_j^T \tilde{x}}}$$

For $\tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix}$, $\tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix}$.

Then, let \mathcal{L}_i be the gradient of the softmax for the i -th observation in the data set. Since \mathcal{L} is the negative log likelihood, it is additive. Then

$$\mathcal{L}_i(\tilde{x}^{(i)}) = -\log \text{softmax}_{y^{(i)}}(\tilde{x}^{(i)})$$

$$= \log \sum_{j=1}^c e^{\tilde{w}_j^T \tilde{x}^{(i)}} - \log e^{\tilde{w}_{y^{(i)}}^T \tilde{x}^{(i)}}$$

$$= \log \sum_{j=1}^c e^{\tilde{w}_j^T \tilde{x}^{(i)}} - \tilde{w}_{y^{(i)}}^T \tilde{x}^{(i)}$$

Then $\frac{\partial \mathcal{L}_i}{\partial w_k} = \frac{1}{\sum_{j=1}^c e^{\tilde{w}_j^T \tilde{x}^{(i)}}} e^{\tilde{w}_k^T \tilde{x}^{(i)}} \tilde{x}^{(i)} - \frac{\partial}{\partial w_k} (\tilde{w}_{y^{(i)}}^T \tilde{x}^{(i)})$

$$= \text{softmax}_k(\tilde{x}^{(i)}) \tilde{x}^{(i)} - \mathbb{I}_{\{y^{(i)}=k\}} \tilde{x}^{(i)}$$

$$= (\text{softmax}_k(\tilde{x}^{(i)}) - \mathbb{I}_{\{y^{(i)}=k\}}) \tilde{x}^{(i)}$$

Ans

$$\frac{\partial \mathcal{L}}{\partial w_k} = \sum_{i=1}^n \frac{\partial \mathcal{L}_i}{\partial w_k}$$

$$4) \mathcal{L}(\omega, b) = \frac{1}{k} \sum_{i=1}^k \text{hinge}_{y^{(i)}}(x^{(i)})$$

$$= \frac{1}{k} \sum_{i=1}^k \max(0, 1 - y^{(i)}(\omega^T x^{(i)} + b))$$

$$\frac{\partial \mathcal{L}}{\partial \omega} = \frac{1}{k} \sum_{i=1}^k \mathbb{I}_{\{y^{(i)}(\omega^T x^{(i)} + b) < 1\}} \left(\frac{\partial}{\partial \omega} (1 - y^{(i)}(\omega^T x^{(i)} + b)) \right)$$

$$= \frac{1}{k} \sum_{i=1}^k \mathbb{I}_{\{y^{(i)}(\omega^T x^{(i)} + b) < 1\}} (-y^{(i)} x^{(i)})$$

$$= \frac{1}{k} X^T \begin{bmatrix} \mathbb{I}_{\{y^{(1)}(\omega^T x^{(1)} + b) < 1\}} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \mathbb{I}_{\{y^{(k)}(\omega^T x^{(k)} + b) < 1\}} \end{bmatrix} Y$$

Where $X = \begin{bmatrix} - & x^{(1)} & - \\ & \vdots & \\ - & x^{(k)} & - \end{bmatrix}$ and $Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(k)} \end{bmatrix}$

Also,

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{1}{k} \sum_{i=1}^k \mathbb{I}_{\{y^{(i)}(\omega^T x^{(i)} + b) < 1\}} \left(\frac{\partial}{\partial b} (1 - y^{(i)}(\omega^T x^{(i)} + b)) \right)$$

$$= \frac{1}{k} \sum_{i=1}^k \mathbb{I}_{\{y^{(i)}(\omega^T x^{(i)} + b) < 1\}} (-y^{(i)})$$

$$= -\frac{1}{k} \sum_{i=1}^k \mathbb{I}_{\{y^{(i)}(\omega^T x^{(i)} + b) < 1\}} y^{(i)}$$