## Question 3.1

Given: $\tilde{\omega}$ for L2 reg. case is:

$$\hat{\omega} = Q(\Lambda + \alpha I)^{-1} \Lambda Q^T \omega^*$$

Second order Taylor expansion to approximate reg. cost $\hat{J}(\theta)$.

$$\hat{J}(\omega) = J(\omega^*) + \frac{1}{2}(\omega - \omega^*)^T H (\omega - \omega^*) \qquad \text{where } \omega^* \text{ is optimal sol. on } J.$$

$\downarrow$ find gradient $\nabla_\omega \hat{J}(\omega)$

$\nabla_\omega J(\omega^*) = 0 \qquad \omega^*$ is a minimum on $J$, so gradient is zero.

$\nabla_\omega \hat{J}(\theta) = H(\omega - \omega^*)$

$\downarrow$ find weight update at iteration $t$

$$\omega^t = \omega^{t-1} - \varepsilon \underline{H(\omega^{t-1} - \omega^*)}$$
$$\text{gradient evaluated at } \omega^{t-1}$$

$$\omega^t - \omega^* = (I - \varepsilon H)(\omega^{t-1} - \omega^*)$$

$\downarrow$ decompose $H \rightarrow Q \Lambda Q^T$

$$\omega^t - \omega^* = (I - \varepsilon Q \Lambda Q^T)(\omega^{t-1} - \omega^*)$$
$$= Q(I - \varepsilon \Lambda) Q^T (\omega^{t-1} - \omega^*)$$

$\downarrow$
- $\omega^{(0)} = [0 ... 0]$, initialization form origin.
- learning rate is small

$$\omega^t = Q[I - (I - \varepsilon \Lambda)^t] Q^T \omega^*$$

$$\hat{\omega} = Q(\Lambda + \alpha I)^{-1} \Lambda Q^T \omega^* \qquad \text{(L2 reg.)}$$
$$= Q[I - (\Lambda + \alpha I)^{-1} \alpha] Q^T \omega^*$$

condition for $\omega^t = \tilde{\omega}$:

$$(I - \varepsilon \Lambda)^t = (\Lambda + \alpha I)^{-1} \alpha$$

$\downarrow$ $\varepsilon \lambda_i$ is small

$$t \approx \frac{1}{\varepsilon \alpha} \qquad \therefore t \text{ and } \alpha \text{ are inversely related.}$$

early stopping is equivalent to using a large regularization constant.

## Question 3.2

1) Verify $\nabla_{\omega^{(k)}} \mathcal{J} = g \, h^{(k-1)T} + \lambda \nabla_{\omega^{(k)}} \Omega(\theta)$

$\underbrace{\quad}_{m \times n} \quad \underbrace{\quad}_{m \times 1} \underbrace{\quad}_{1 \times n} \quad \underbrace{\quad}_{m \times n}$

### non-regularized

$\nabla_{\omega^{(k)}} \mathcal{J} = \underbrace{\nabla_{a^{(k)}} \mathcal{J}}_{m \times 1} \cdot \underbrace{\frac{\partial a^{(k)}}{\partial \omega^{(k)}}}_{1 \times n}$  (chain rule)

$\underbrace{\quad}_{m \times n}$

↘ this is $g$        ↓ ⓐ : $m \times 1$ ⇒ ⓦ : $m \times n$

this is vector to matrix diff,
but rows in ⓦ are independent.
Let's break it down to scalar to
vector operation as follows:

Take $a_1^{(k)}$, the first element in ⓐ ; and $\omega_{1,:}^{(k)}$, the first (scalar)
row of matrix $\omega^{(k)}$  $(1 \times n)$

$\frac{\partial a_1^{(k)}}{\partial \omega_{1,:}^{(k)}} = \begin{bmatrix} \frac{\partial f}{\partial \omega_{1,1}} \\ \vdots \\ \frac{\partial f}{\partial \omega_{1,n}} \end{bmatrix}$, where $f = \sum_{i=1}^{n} \omega_{1,i}^{(k)} \cdot h_i^{(k-1)}$

$\frac{\partial a_1^{(k)}}{\partial \omega_{1,:}^{(k)}} = \begin{bmatrix} h_1^{(k-1)} \\ h_2^{(k-1)} \\ \vdots \\ h_n^{(k-1)} \end{bmatrix}$

$\frac{\partial a_i^{(k)}}{\partial \omega_{i,:}^{(k)}}$ for all $i \in 1 \dots m$ generates the same result.

Hence, $\frac{\partial a_i^{(k)}}{\partial \omega_{i,:}^{(k)}} = \begin{bmatrix} h_1^{(k-1)} \\ h_2^{(k-1)} \\ \vdots \\ h_n^{(k-1)} \end{bmatrix}$

∴ this does not depend on $i$
∴ it can be applied to each and every element in ⓐ
(outer product)

∴ $\nabla_{\omega^{(k)}} \mathcal{J} = \nabla_{a^{(k)}} \mathcal{J} \cdot h^{(k-1)T}$

### regularized

independent

$\mathcal{J}\_reg = \lambda \Omega(\omega_1, \omega_2 \cdots \omega_i)$   ∴ no chain rule

$\frac{\partial \mathcal{J}\_reg}{\omega_k} = \lambda \nabla_{\omega^{(k)}}(\theta)$

### overall  $\nabla_{\omega^{(k)}} \mathcal{J} = g \, h^{(k-1)T} + \lambda \nabla_{\omega^{(k)}} \Omega(\theta)$

---

2) Verify $\nabla_{h^{(k-1)}} \mathcal{J} = \omega^{(k)T} \cdot g$

$\underbrace{\quad}_{n \times 1} \quad \underbrace{\quad}_{n \times m} \underbrace{\quad}_{m \times 1}$

$\nabla_{h^{(k-1)}} \mathcal{J} = \left( \frac{\partial a^{(k)}}{\partial h^{(k-1)}} \right)^T \cdot \nabla_{a^{(k)}} \mathcal{J}$

↘ this is $g$

$\frac{\partial a^{(k)}}{\partial h^{(k-1)}} = \begin{bmatrix} \frac{\partial a_1}{\partial h_1} & \frac{\partial a_1}{\partial h_2} & \cdots & \frac{\partial a_1}{\partial h_n} \\ \frac{\partial a_2}{\partial h_1} & & \ddots & \\ \vdots & & & \\ \frac{\partial a_m}{\partial h_1} & & & \end{bmatrix} \Big\} m$

$\underbrace{\qquad\qquad\qquad}_{n}$

During forward propagation, $\underline{a^{(k)} = \omega^{(k)} h^{(k-1)}}$, so $\frac{\partial a^{(k)}}{\partial h^{(k-1)}}$ is exactly $\omega^{(k)}$.

∴ $\nabla_{h^{(k-1)}} \mathcal{J} = \omega^{(k)T} \cdot g$