

Module 8.4 : l_2 regularization

Different forms of regularization

- l_2 regularization
- Dataset augmentation
- Parameter Sharing and tying
- Adding Noise to the inputs
- Adding Noise to the outputs
- Early stopping
- Ensemble methods
- Dropout

- For l_2 regularization we have,

$$\widetilde{\mathcal{L}}(w) = \mathcal{L}(w) + \frac{\alpha}{2} \|w\|^2$$

- For l_2 regularization we have,

$$\widetilde{\mathcal{L}}(w) = \mathcal{L}(w) + \frac{\alpha}{2} \|w\|^2$$

- For SGD (or its variants), we are interested in

$$\nabla \widetilde{\mathcal{L}}(w) = \nabla \mathcal{L}(w) + \alpha w$$

- For l_2 regularization we have,

$$\widetilde{\mathcal{L}}(w) = \mathcal{L}(w) + \frac{\alpha}{2} \|w\|^2$$

- For SGD (or its variants), we are interested in

$$\nabla \widetilde{\mathcal{L}}(w) = \nabla \mathcal{L}(w) + \alpha w$$

- Update rule:

$$w_{t+1} = w_t - \eta \nabla \mathcal{L}(w_t) - \eta \alpha w_t$$

- For l_2 regularization we have,

$$\widetilde{\mathcal{L}}(w) = \mathcal{L}(w) + \frac{\alpha}{2} \|w\|^2$$

- For SGD (or its variants), we are interested in

$$\nabla \widetilde{\mathcal{L}}(w) = \nabla \mathcal{L}(w) + \alpha w$$

- Update rule:

$$w_{t+1} = w_t - \eta \nabla \mathcal{L}(w_t) - \eta \alpha w_t$$

- Requires a very small modification to the code

- For l_2 regularization we have,

$$\widetilde{\mathcal{L}}(w) = \mathcal{L}(w) + \frac{\alpha}{2} \|w\|^2$$

- For SGD (or its variants), we are interested in

$$\nabla \widetilde{\mathcal{L}}(w) = \nabla \mathcal{L}(w) + \alpha w$$

- Update rule:

$$w_{t+1} = w_t - \eta \nabla \mathcal{L}(w_t) - \eta \alpha w_t$$

- Requires a very small modification to the code
- Let us see the geometric interpretation of this

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

$$\mathcal{L}(w^* + u) = \mathcal{L}(w^*) + u^T \nabla \mathcal{L}(w^*) + \frac{1}{2} u^T H u$$

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

$$\mathcal{L}(w^* + u) = \mathcal{L}(w^*) + u^T \nabla \mathcal{L}(w^*) + \frac{1}{2} u^T H u$$

$$\mathcal{L}(w) = \mathcal{L}(w^*) + (w - w^*)^T \nabla \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

$$\mathcal{L}(w^* + u) = \mathcal{L}(w^*) + u^T \nabla \mathcal{L}(w^*) + \frac{1}{2} u^T H u$$

$$\begin{aligned}\mathcal{L}(w) &= \mathcal{L}(w^*) + (w - w^*)^T \nabla \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*) \\ &= \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*) \quad (\because \nabla \mathcal{L}(w^*) = 0)\end{aligned}$$

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

$$\mathcal{L}(w^* + u) = \mathcal{L}(w^*) + u^T \nabla \mathcal{L}(w^*) + \frac{1}{2} u^T H u$$

$$\mathcal{L}(w) = \mathcal{L}(w^*) + (w - w^*)^T \nabla \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

$$= \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*) \quad (\because \nabla \mathcal{L}(w^*) = 0)$$

$$\nabla \mathcal{L}(w) = \nabla \mathcal{L}(w^*) + H(w - w^*)$$

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

$$\mathcal{L}(w^* + u) = \mathcal{L}(w^*) + u^T \nabla \mathcal{L}(w^*) + \frac{1}{2} u^T H u$$

$$\mathcal{L}(w) = \mathcal{L}(w^*) + (w - w^*)^T \nabla \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

$$= \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*) \quad (\because \nabla L(w^*) = 0)$$

$$\begin{aligned} \nabla \mathcal{L}(w) &= \nabla \mathcal{L}(w^*) + H(w - w^*) \\ &= H(w - w^*) \end{aligned}$$

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

$$\mathcal{L}(w^* + u) = \mathcal{L}(w^*) + u^T \nabla \mathcal{L}(w^*) + \frac{1}{2} u^T H u$$

$$\mathcal{L}(w) = \mathcal{L}(w^*) + (w - w^*)^T \nabla \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

$$= \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*) \quad (\because \nabla \mathcal{L}(w^*) = 0)$$

$$\begin{aligned} \nabla \mathcal{L}(w) &= \nabla \mathcal{L}(w^*) + H(w - w^*) \\ &= H(w - w^*) \end{aligned}$$

- Now,

$$\nabla \widetilde{\mathcal{L}}(w) = \nabla \mathcal{L}(w) + \alpha w$$

- Assume w^* is the optimal solution for $\mathcal{L}(w)$ [not $\widetilde{\mathcal{L}}(w)$] i.e. the solution in the absence of regularization (w^* optimal $\rightarrow \nabla \mathcal{L}(w^*) = 0$)
- Consider $u = w - w^*$. Using Taylor series approximation (upto 2^{nd} order)

$$\mathcal{L}(w^* + u) = \mathcal{L}(w^*) + u^T \nabla \mathcal{L}(w^*) + \frac{1}{2} u^T H u$$

$$\mathcal{L}(w) = \mathcal{L}(w^*) + (w - w^*)^T \nabla \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

$$= \mathcal{L}(w^*) + \frac{1}{2} (w - w^*)^T H (w - w^*) \quad (\because \nabla L(w^*) = 0)$$

$$\begin{aligned} \nabla \mathcal{L}(w) &= \nabla \mathcal{L}(w^*) + H(w - w^*) \\ &= H(w - w^*) \end{aligned}$$

- Now,

$$\begin{aligned} \nabla \widetilde{\mathcal{L}}(w) &= \nabla \mathcal{L}(w) + \alpha w \\ &= H(w - w^*) + \alpha w \end{aligned}$$

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

$$\because \nabla \tilde{L}(\tilde{w}) = 0$$

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

$$\because \nabla \tilde{L}(\tilde{w}) = 0$$

$$H(\tilde{w} - w^*) + \alpha \tilde{w} = 0$$

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

$$\therefore \nabla \tilde{L}(\tilde{w}) = 0$$

$$H(\tilde{w} - w^*) + \alpha \tilde{w} = 0$$

$$\therefore (H + \alpha \mathbb{I}) \tilde{w} = H w^*$$

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

$$\because \nabla \tilde{L}(\tilde{w}) = 0$$

$$H(\tilde{w} - w^*) + \alpha \tilde{w} = 0$$

$$\therefore (H + \alpha \mathbb{I}) \tilde{w} = H w^*$$

$$\therefore \tilde{w} = (H + \alpha \mathbb{I})^{-1} H w^*$$

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

$$\because \nabla \tilde{L}(\tilde{w}) = 0$$

$$H(\tilde{w} - w^*) + \alpha \tilde{w} = 0$$

$$\therefore (H + \alpha \mathbb{I}) \tilde{w} = H w^*$$

$$\therefore \tilde{w} = (H + \alpha \mathbb{I})^{-1} H w^*$$

- Notice that if $\alpha \rightarrow 0$ then $\tilde{w} \rightarrow w^*$ [no regularization]

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

$$\therefore \nabla \tilde{L}(\tilde{w}) = 0$$

$$H(\tilde{w} - w^*) + \alpha \tilde{w} = 0$$

$$\therefore (H + \alpha \mathbb{I}) \tilde{w} = H w^*$$

$$\therefore \tilde{w} = (H + \alpha \mathbb{I})^{-1} H w^*$$

- Notice that if $\alpha \rightarrow 0$ then $\tilde{w} \rightarrow w^*$ [no regularization]
- But we are interested in the case when $\alpha \neq 0$

- Let \tilde{w} be the optimal solution for $\tilde{L}(w)$ [i.e regularized loss]

$$\therefore \nabla \tilde{L}(\tilde{w}) = 0$$

$$H(\tilde{w} - w^*) + \alpha \tilde{w} = 0$$

$$\therefore (H + \alpha \mathbb{I}) \tilde{w} = H w^*$$

$$\therefore \tilde{w} = (H + \alpha \mathbb{I})^{-1} H w^*$$

- Notice that if $\alpha \rightarrow 0$ then $\tilde{w} \rightarrow w^*$ [no regularization]
- But we are interested in the case when $\alpha \neq 0$
- Let us analyse the case when $\alpha \neq 0$

- If H is symmetric Positive Semi Definite

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\tilde{w} = (H + \alpha\mathbb{I})^{-1} H w^*$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\begin{aligned}\tilde{w} &= (H + \alpha\mathbb{I})^{-1} H w^* \\ &= (Q\Lambda Q^T + \alpha\mathbb{I})^{-1} Q\Lambda Q^T w^*\end{aligned}$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\begin{aligned}\tilde{w} &= (H + \alpha\mathbb{I})^{-1} H w^* \\ &= (Q\Lambda Q^T + \alpha\mathbb{I})^{-1} Q\Lambda Q^T w^* \\ &= (Q\Lambda Q^T + \alpha Q\mathbb{I}Q^T)^{-1} Q\Lambda Q^T w^*\end{aligned}$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\begin{aligned}\tilde{w} &= (H + \alpha\mathbb{I})^{-1} H w^* \\ &= (Q\Lambda Q^T + \alpha\mathbb{I})^{-1} Q\Lambda Q^T w^* \\ &= (Q\Lambda Q^T + \alpha Q\mathbb{I}Q^T)^{-1} Q\Lambda Q^T w^* \\ &= [Q(\Lambda + \alpha\mathbb{I})Q^T]^{-1} Q\Lambda Q^T w^*\end{aligned}$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\begin{aligned}\tilde{w} &= (H + \alpha\mathbb{I})^{-1} H w^* \\ &= (Q\Lambda Q^T + \alpha\mathbb{I})^{-1} Q\Lambda Q^T w^* \\ &= (Q\Lambda Q^T + \alpha Q\mathbb{I}Q^T)^{-1} Q\Lambda Q^T w^* \\ &= [Q(\Lambda + \alpha\mathbb{I})Q^T]^{-1} Q\Lambda Q^T w^* \\ &= Q^{T^{-1}} (\Lambda + \alpha\mathbb{I})^{-1} Q^{-1} Q\Lambda Q^T w^*\end{aligned}$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\begin{aligned}\tilde{w} &= (H + \alpha\mathbb{I})^{-1} H w^* \\ &= (Q\Lambda Q^T + \alpha\mathbb{I})^{-1} Q\Lambda Q^T w^* \\ &= (Q\Lambda Q^T + \alpha Q\mathbb{I}Q^T)^{-1} Q\Lambda Q^T w^* \\ &= [Q(\Lambda + \alpha\mathbb{I})Q^T]^{-1} Q\Lambda Q^T w^* \\ &= Q^{T^{-1}} (\Lambda + \alpha\mathbb{I})^{-1} Q^{-1} Q\Lambda Q^T w^* \\ &= Q(\Lambda + \alpha\mathbb{I})^{-1} \Lambda Q^T w^* \quad (\because Q^{T^{-1}} = Q)\end{aligned}$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\begin{aligned} \tilde{w} &= (H + \alpha\mathbb{I})^{-1} H w^* \\ &= (Q\Lambda Q^T + \alpha\mathbb{I})^{-1} Q\Lambda Q^T w^* \\ &= (Q\Lambda Q^T + \alpha Q\mathbb{I}Q^T)^{-1} Q\Lambda Q^T w^* \\ &= [Q(\Lambda + \alpha\mathbb{I})Q^T]^{-1} Q\Lambda Q^T w^* \\ &= Q^{T^{-1}} (\Lambda + \alpha\mathbb{I})^{-1} Q^{-1} Q\Lambda Q^T w^* \\ &= Q(\Lambda + \alpha\mathbb{I})^{-1} \Lambda Q^T w^* \quad (\because Q^{T^{-1}} = Q) \\ \tilde{w} &= QDQ^T w^* \end{aligned}$$

- If H is symmetric Positive Semi Definite

$$H = Q\Lambda Q^T \quad [Q \text{ is orthogonal, } QQ^T = Q^T Q = \mathbb{I}]$$

$$\begin{aligned} \tilde{w} &= (H + \alpha\mathbb{I})^{-1} H w^* \\ &= (Q\Lambda Q^T + \alpha\mathbb{I})^{-1} Q\Lambda Q^T w^* \\ &= (Q\Lambda Q^T + \alpha Q\mathbb{I}Q^T)^{-1} Q\Lambda Q^T w^* \\ &= [Q(\Lambda + \alpha\mathbb{I})Q^T]^{-1} Q\Lambda Q^T w^* \\ &= Q^{T^{-1}} (\Lambda + \alpha\mathbb{I})^{-1} Q^{-1} Q\Lambda Q^T w^* \\ &= Q(\Lambda + \alpha\mathbb{I})^{-1} \Lambda Q^T w^* \quad (\because Q^{T^{-1}} = Q) \\ \tilde{w} &= QDQ^T w^* \end{aligned}$$

where $D = (\Lambda + \alpha\mathbb{I})^{-1} \Lambda$, is a diagonal matrix which we will see in more detail soon

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

- So what is happening here?

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha\mathbb{I})^{-1}\Lambda Q^T w^* \\ &= QDQ^T w^*\end{aligned}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} \\ \vdots \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & \\ & \frac{1}{\lambda_2 + \alpha} \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha\mathbb{I})^{-1}\Lambda Q^T w^* \\ &= QDQ^T w^*\end{aligned}$$

$$(\Lambda + \alpha\mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha\mathbb{I})^{-1}\Lambda Q^T w^* \\ &= QDQ^T w^*\end{aligned}$$

$$(\Lambda + \alpha\mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha\mathbb{I})^{-1}\Lambda$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha\mathbb{I})^{-1}\Lambda Q^T w^* \\ &= QDQ^T w^*\end{aligned}$$

$$(\Lambda + \alpha\mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha\mathbb{I})^{-1}\Lambda$$

$$(\Lambda + \alpha\mathbb{I})^{-1}\Lambda = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha \mathbb{I})^{-1} \Lambda$$

$$(\Lambda + \alpha \mathbb{I})^{-1} \Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha \mathbb{I})^{-1} \Lambda$$

$$(\Lambda + \alpha \mathbb{I})^{-1} \Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & \\ & & \ddots \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha \mathbb{I})^{-1} \Lambda$$

$$(\Lambda + \alpha \mathbb{I})^{-1} \Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha \mathbb{I})^{-1} \Lambda$$

$$(\Lambda + \alpha \mathbb{I})^{-1} \Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha\mathbb{I})^{-1}\Lambda Q^T w^* \\ &= QDQ^T w^*\end{aligned}$$

$$(\Lambda + \alpha\mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha\mathbb{I})^{-1}\Lambda$$

$$(\Lambda + \alpha\mathbb{I})^{-1}\Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

- So what is happening here?
- w^* first gets rotated by Q^T to give $Q^T w^*$
- However if $\alpha = 0$ then Q rotates $Q^T w^*$ back to give w^*
- If $\alpha \neq 0$ then let us see what D looks like
- So what is happening now?

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha\mathbb{I})^{-1}\Lambda Q^T w^* \\ &= QDQ^T w^*\end{aligned}$$

$$(\Lambda + \alpha\mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha\mathbb{I})^{-1}\Lambda$$

$$(\Lambda + \alpha\mathbb{I})^{-1}\Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

- Each element i of $Q^T w^*$ gets scaled by $\frac{\lambda_i}{\lambda_i + \alpha}$ before it is rotated back by Q

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha \mathbb{I})^{-1} \Lambda$$

$$(\Lambda + \alpha \mathbb{I})^{-1} \Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

- Each element i of $Q^T w^*$ gets scaled by $\frac{\lambda_i}{\lambda_i + \alpha}$ before it is rotated back by Q
- if $\lambda_i \gg \alpha$ then $\frac{\lambda_i}{\lambda_i + \alpha} = 1$

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha\mathbb{I})^{-1}\Lambda Q^T w^* \\ &= QDQ^T w^*\end{aligned}$$

$$(\Lambda + \alpha\mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

$$D = (\Lambda + \alpha\mathbb{I})^{-1}\Lambda$$

$$(\Lambda + \alpha\mathbb{I})^{-1}\Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

- Each element i of $Q^T w^*$ gets scaled by $\frac{\lambda_i}{\lambda_i + \alpha}$ before it is rotated back by Q
- if $\lambda_i \gg \alpha$ then $\frac{\lambda_i}{\lambda_i + \alpha} = 1$
- if $\lambda_i \ll \alpha$ then $\frac{\lambda_i}{\lambda_i + \alpha} = 0$

$$\begin{aligned}\tilde{w} &= Q(\Lambda + \alpha \mathbb{I})^{-1} \Lambda Q^T w^* \\ &= Q D Q^T w^*\end{aligned}$$

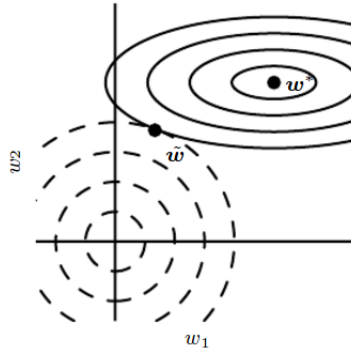
$$(\Lambda + \alpha \mathbb{I})^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$$

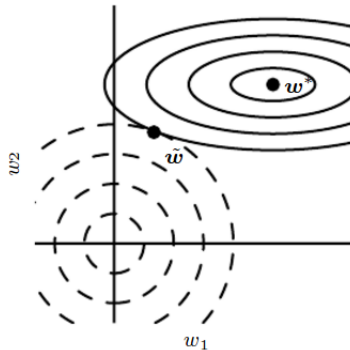
$$D = (\Lambda + \alpha \mathbb{I})^{-1} \Lambda$$

$$(\Lambda + \alpha \mathbb{I})^{-1} \Lambda = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

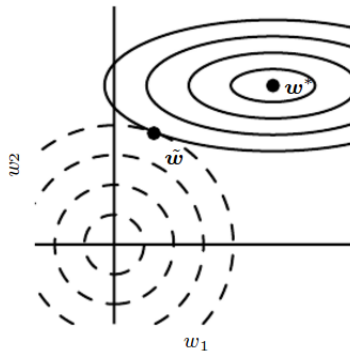
- Each element i of $Q^T w^*$ gets scaled by $\frac{\lambda_i}{\lambda_i + \alpha}$ before it is rotated back by Q
- if $\lambda_i \gg \alpha$ then $\frac{\lambda_i}{\lambda_i + \alpha} = 1$
- if $\lambda_i \ll \alpha$ then $\frac{\lambda_i}{\lambda_i + \alpha} = 0$
- Thus only significant directions (larger eigen values) will be retained.

$$\text{Effective parameters} = \sum_{i=1}^n \frac{\lambda_i}{\lambda_i + \alpha} < n$$

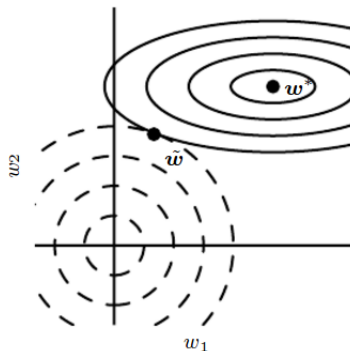




- The weight vector(w^*) is getting rotated to (\tilde{w})



- The weight vector(w^*) is getting rotated to (\tilde{w})
- All of its elements are shrinking but some are shrinking more than the others



- The weight vector(w^*) is getting rotated to (\tilde{w})
- All of its elements are shrinking but some are shrinking more than the others
- This ensures that only important features are given high weights