

$$1) \quad \omega = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad v = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \alpha = 0.2 \\ \beta = 0.8$$

$$X = \begin{bmatrix} \begin{bmatrix} 2 \\ -1 \end{bmatrix} & \begin{bmatrix} -1 \\ 3 \end{bmatrix} \end{bmatrix} \quad y = \begin{bmatrix} 3 \\ 1 \end{bmatrix} \quad B = 2$$

Iteration 1:

$$\rightarrow \hat{y}_i = \omega^T X_i = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ -1 \end{bmatrix} = 2 \quad \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} -1 \\ 3 \end{bmatrix} = -1$$

$$MSE = \frac{1}{B} \sum_{i=1}^B (\hat{y}_i - y_i)^2 = \frac{1}{2} \left[ (-2-3)^2 + (-1-1)^2 \right]$$

$$\hat{y} = X \omega = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

$$\hat{y} - y = \begin{bmatrix} 2 \\ -1 \end{bmatrix} - \begin{bmatrix} 3 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -2 \end{bmatrix}$$

$$\nabla_{\omega} = \frac{2}{B} \frac{dL}{d\hat{y}} \frac{d\hat{y}}{d\omega} = \frac{1}{2} \cdot 2 \cdot X^T (\hat{y} - y)$$

$$\nabla_{\omega} = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix}^T \begin{bmatrix} -1 \\ -2 \end{bmatrix} = \begin{bmatrix} -2 \\ -5 \end{bmatrix}$$

$$v = \beta v - \alpha \nabla_{\omega} = 0.8 \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.2 \begin{bmatrix} -2 \\ -5 \end{bmatrix} = \begin{bmatrix} 0.4 \\ 1.0 \end{bmatrix}$$

$$\omega' = \omega + v = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0.4 \\ 1.0 \end{bmatrix} = \begin{bmatrix} 1.4 \\ 1.0 \end{bmatrix}$$

Iteration 2:

$$\hat{y} = Xw = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix} \begin{bmatrix} 1.4 \\ 1 \end{bmatrix} = \begin{bmatrix} 1.8 \\ 1.6 \end{bmatrix}$$

$$\hat{y} - y = \begin{bmatrix} 1.8 \\ 1.6 \end{bmatrix} - \begin{bmatrix} 3 \\ 1 \end{bmatrix} = \begin{bmatrix} -1.2 \\ 0.6 \end{bmatrix}$$

$$\Delta w = \frac{1}{B} \cdot 2 \cdot X^T \cdot (\hat{y} - y) = \frac{1}{2} \cdot 2 \cdot \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix}^T \cdot \begin{bmatrix} -1.2 \\ 0.6 \end{bmatrix}$$

$$= \begin{bmatrix} -3 \\ 2.4 \end{bmatrix}$$

$$v = 0.8 \begin{bmatrix} 0.4 \\ 1 \end{bmatrix} - 0.2 \begin{bmatrix} -3 \\ 2.4 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.32 \end{bmatrix}$$

$$w = w + v = \begin{bmatrix} 1.4 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0.32 \end{bmatrix} = \underline{\underline{\begin{bmatrix} 2.4 \\ 1.32 \end{bmatrix}}}$$

$$2) s_t = p_1 s_{t-1} + (1-p_1)g$$

$g \rightarrow$  constant gradient

$$s = 0$$

$p_1 \rightarrow$  decay rate for first moment

$$\hat{s}_t = \frac{s_t}{1 - p_1^t} \quad (\text{bias correction})$$

Prove  $s_t = g$  at every step



$t=1$

$$s_1 = p_1 \cdot s_0 + (1-p_1)g = 0 + (1-p_1)g$$

$$\hat{s}_1 = \frac{s_1}{1-p_1^1} = \frac{(1-\cancel{p_1})g}{1-\cancel{p_1}} = g$$

Inductive Step:

Assume  $\hat{s}_t = g$  for  $t=k$  and prove for  $t=k+1$

→ step  $t=k$  we have  $\hat{s}_k = g$

$$\hat{s}_k = \frac{s_k}{1-p_1^k} = g$$

$$s_{k+1} = p_1 s_k + (1-p_1)g$$

$$= p_1 (1-p_1^k)g + (1-p_1)g$$

$$= (p_1 (1-p_1^k) + (1-p_1))g$$

$$s_{k+1} = (p_1 - p_1^{k+1} + 1 - p_1)g = (1-p_1^{k+1})g$$

$$\hat{s}_{k+1} = \frac{s_{k+1}}{1-p_1^{k+1}} = \frac{(1-\cancel{p_1^{k+1}})g}{1-\cancel{p_1^{k+1}}} = g$$

Answer for the reason of zig-zag behaviour at 2nd ill-conditioned plot:

Because along the directions that have large eigenvalues, gradient takes big steps and overshoot which causes the zig-zag behaviour.

Momentum: Accumulates past gradients and smooth the zig-zag behaviour a little bit, enabling faster convergence (fewer steps taken)

Preconditioning: Reduces the difference between eigenvalues which help preventing the ill-condition and zig-zag behaviour. Steps are more consistent and directly to the minimum.









