

Optimization algorithms (part 1)

$$w_{t+1} = w_t - \alpha \nabla l(w_t) \quad \textcircled{\eta} \Rightarrow \underline{\underline{lr}}$$

$$\underline{\underline{w}} l(w) = 4w^2 - 12w + 9$$

$$w = 1.5$$

$$\text{we have } l'(w) = 8w - 12 \quad \checkmark$$

$$l''(w) = \underline{\underline{8}}$$

$$\nabla l(w_t) = 8w_t - 12$$

we start with $\underline{\underline{w_0 = 3}}$ & follow the update rule

$$w_{k+1} = w_k - \eta l'(w_k)$$

$$\Rightarrow w_{k+1} = w_k - \eta (8w_k - 12)$$

$$\Rightarrow \boxed{w_{k+1} = w_k - 0.2 (8w_k - 12)} \quad \text{we assume } \eta = \underline{\underline{0.2}}$$

$k \Rightarrow \underline{\underline{\text{No. of steps}}}$

The Taylor series expansion of a function $f(x)$ that is differentiable at a is given by

$$f(x) = f(a) + \frac{f'(a)}{1!} (x-a) + \frac{f''(a)}{2!} (x-a)^2 + \dots$$

We consider a loss function $l(w)$ & replace $f(\cdot)$ by $l(\cdot)$, x by w_{k+1} & a by w_k . We get -

$$l(w_{k+1}) = l(w_k) + \frac{l'(w_k)}{1!} [w_{k+1} - w_k] + \frac{l''(w_k)}{2!} (w_{k+1} - w_k)^2 + \dots$$

Assumption :- Loss function in any bounded region can be approximated with a quadratic function.

$$\Rightarrow l(w_{k+1}) = l(w_k) + (w_{k+1} - w_k) l'(w_k) + \frac{(w_{k+1} - w_k)^2}{2} l''(w_k) \quad \checkmark$$

We need to find w_{k+1} & we partially differentiate w.r.t w_{k+1}

$$\frac{d}{dw_{k+1}} l(w_{k+1}) = \frac{d}{dw_{k+1}} l(w_k) + \frac{d}{dw_{k+1}} [(w_{k+1} - w_k) l'(w_k)] +$$

$$= \frac{d}{dw_{k+1}} \left[\frac{(w_{k+1} - w_k)^2}{2} l''(w_k) \right]$$

$$= 0 + l'(w_k) \frac{d}{dw_{k+1}} (w_{k+1} - w_k) + \frac{l''(w_k)}{2} \frac{d}{dw_{k+1}} (w_{k+1} - w_k)^2$$

$\hookrightarrow 2(w_{k+1} - w_k)$

$$= l'(w_k) + \frac{l''(w_k)}{2} [w_{k+1} - w_k]$$

$$\boxed{\frac{d}{dw_{k+1}} l(w_{k+1}) = l'(w_k) + \frac{l''(w_k)}{2} [w_{k+1} - w_k]} \quad \checkmark$$

$$= l'(\omega_k) + \frac{l''(\omega_k)}{2} [\omega_{k+1} - \omega_k]$$

$$\boxed{\frac{d l(\omega_{k+1})}{d \omega_{k+1}} = l'(\omega_k) + l''(\omega_k) [\omega_{k+1} - \omega_k]} \quad \checkmark$$

For extremum points; $\frac{d l(\omega_{k+1})}{d \omega_{k+1}} = 0$

$$\Rightarrow l'(\omega_k) + l''(\omega_k) [\omega_{k+1} - \omega_k] = 0$$

$$\Rightarrow \underline{l''(\omega_k)} [\omega_{k+1} - \omega_k] = -l'(\omega_k)$$

$$\Rightarrow \omega_{k+1} - \omega_k = -\frac{1}{l''(\omega_k)} l'(\omega_k)$$

$$\Rightarrow \omega_{k+1} = \omega_k - \frac{1}{l''(\omega_k)} l'(\omega_k) \quad \checkmark$$

$$\Rightarrow \boxed{\omega_{k+1} = \omega_k - \underline{\eta} l'(\omega_k)} \quad \text{where } \eta = \frac{1}{l''(\omega_k)}$$

Further, $\frac{d^2 l(\omega_{k+1})}{d \omega_{k+1}^2} = l''(\omega_k) > 0$

$$\Rightarrow \omega_{k+1} \text{ is } \underline{\text{local minima.}} \quad l(\omega_{k+1})$$

$$\left[\text{Note :- learning rate } (\eta) = \frac{1}{l''(\omega_k)} \right]$$

In multivariate case;

$$\checkmark \underline{\underline{l'(\omega_k)}} = \underline{\underline{\nabla l(\omega_k)}} = g$$

$$\checkmark \underline{\underline{l''(\omega_k)}} = \underline{\underline{H}} \quad [\text{Hessian matrix of } l \text{ at } \omega_k]$$

$\omega_k \rightarrow \text{vector}$

$H \rightarrow \text{matrix}$

$$\text{Also, } (\omega_{k+1} - \omega_k)^T l''(\omega_k) = (\omega_{k+1} - \omega_k)^T H (\omega_{k+1} - \omega_k)$$

under the same set of assumptions;

$$l(\omega_{k+1}) = \underline{\underline{l(\omega_k)}} + (\omega_{k+1} - \omega_k)^T g + \frac{1}{2} (\omega_{k+1} - \omega_k)^T H (\omega_{k+1} - \omega_k) + (\omega_{k+1} - \omega_k)^T l''(\omega_k)$$

Partially differentiating w.r.t ω_{k+1} ;

$$\nabla_{\omega_{k+1}} l(\omega_{k+1}) = \nabla_{\omega_{k+1}} l(\omega_k) + g [\nabla (\omega_{k+1} - \omega_k)^T] + \nabla_{\omega_{k+1}} \left[\frac{1}{2} (\omega_{k+1} - \omega_k)^T H (\omega_{k+1} - \omega_k) \right] \rightarrow (\omega_{k+1} - \omega_k)^T l'(\omega_k)$$

$$= \underline{\underline{g}} + \underline{\underline{H (\omega_{k+1} - \omega_k)}}$$

$$\nabla_{\omega_{k+1}} l(\omega_k) = 0$$

$$\left[\begin{array}{l} \text{Reason:-} \\ 1) \nabla_{\omega_{k+1}} (\omega_{k+1} - \omega_k)^T = 1 \\ 2) \nabla_{\omega_{k+1}} [(\omega_{k+1} - \omega_k)^T H (\omega_{k+1} - \omega_k)] = \\ H (\omega_{k+1} - \omega_k) \end{array} \right]$$

$$\Rightarrow \boxed{\nabla_{\omega_{k+1}} \underline{\underline{l(\omega_{k+1})}} = g + H(\omega_{k+1} - \omega_k)}$$

For extremum points; $\nabla_{\omega_{k+1}} l(\omega_{k+1}) = 0$ ✓

$$\Rightarrow g + H(\omega_{k+1} - \omega_k) = 0 \quad \checkmark$$

$$\Rightarrow \underline{H}(\omega_{k+1} - \omega_k) = -g$$

$$\Rightarrow H^{-1}H(\omega_{k+1} - \omega_k) = -H^{-1}g \quad \checkmark$$

$$\Rightarrow \omega_{k+1} - \omega_k = -H^{-1}g$$

$$\Rightarrow \omega_{k+1} = \omega_k - H^{-1}g$$

$$\Rightarrow \boxed{\omega_{k+1} = \omega_k - H^{-1} \nabla l(\omega_k)}$$

(2)

$$\frac{1}{l''(\omega_k)} = \eta$$

$$\underline{\eta} = \underline{H^{-1}}$$

we have;

$$\checkmark \underline{l(\omega_{k+1})} = l(\omega_k) + \underbrace{(\omega_{k+1} - \omega_k)^T}_{\underline{\eta \nabla l(\omega_k)}} \nabla l(\omega_k) + \frac{1}{2} \underbrace{(\omega_{k+1} - \omega_k)^T}_{\underline{\eta \nabla l(\omega_k)}} H (\omega_{k+1} - \omega_k)$$

we assume $\underline{\omega_{k+1}} = \omega_k - \underline{\eta \nabla l(\omega_k)}$ where η is unknown learning rate.

$$\omega_{k+1} - \omega_k = -\eta \nabla l(\omega_k) = \underline{\underline{\eta g}}$$

$$\Rightarrow l(\omega_k - \eta \nabla l(\omega_k)) = l(\omega_k) + \underbrace{(\omega_k - \eta \nabla l(\omega_k))^T}_{\underline{\underline{\eta \nabla l(\omega_k)}}} \nabla l(\omega_k) + \frac{1}{2} \underbrace{(\omega_k - \eta \nabla l(\omega_k))^T}_{\underline{\underline{\eta \nabla l(\omega_k)}}} H (\omega_k - \eta \nabla l(\omega_k) - \omega_k)$$

$$= l(\omega_k) - \underline{\underline{\eta}} \underline{\underline{g^T}} \underline{\underline{g}} + \frac{1}{2} \underline{\underline{(-\eta g^T)}} \underline{\underline{H}} \underline{\underline{(-\eta g)}}$$

$$= \underline{\underline{l(\omega_k)}} - \underline{\underline{\eta}} \underline{\underline{g^T}} \underline{\underline{g}} + \frac{\eta^2}{2} \underline{\underline{g^T H g}}$$

We partially differentiate wrt. η .

$$\frac{\partial}{\partial \eta} \ell(\omega_k - \eta \nabla \ell(\omega_k)) = \underbrace{-g^T g}_{\text{}} + \eta \underbrace{g^T H g}_{\text{}}$$

for extremum points,

$$\frac{\partial}{\partial \eta} \ell(\omega_k - \eta \nabla \ell(\omega_k)) = 0$$

$$\Rightarrow -g^T g + \eta g^T H g = 0$$

$$\Rightarrow \boxed{\eta = \frac{g^T g}{g^T H g}} \quad \checkmark$$

Ans $\underline{H} = \underline{\Phi} \underline{D} \underline{\Phi}^T$ where $\underline{D} = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \lambda_n \end{pmatrix}$ \checkmark
Eigen dec. Th.

$\underline{\Phi}$ is eigenvector matrix where each column of Φ corresponds to one eigenvector of H .

Now. $\underline{g^T H g} = \underline{g^T \Phi} \underline{D} \underline{\Phi^T g}$
 $= \underline{X^T} \underline{D} \underline{X}$ [we assume $\underline{X} = \underline{\Phi^T g}$]

$$\boxed{g^T H g = \sum_{i=1}^n \lambda_i x_i^2} \quad \checkmark$$

Let $\underline{\lambda_{\max}} = \max_i \underline{\lambda_i}$

$$\Rightarrow \lambda_i \leq \underline{\lambda_{\max}}$$

$$x_i^2 > 0$$

$$\Rightarrow \lambda_i x_i^2 \leq \underline{\lambda_{\max}} x_i^2$$

$$\Rightarrow \sum_{i=1}^n \lambda_i x_i^2 \leq \underline{\lambda_{\max}} \sum_{i=1}^n x_i^2$$

$$\Rightarrow \sum_{i=1}^n \lambda_i x_i^2 \leq \underline{\lambda_{\max}} \sum_{i=1}^n x_i^2$$

$$\checkmark \quad \checkmark \quad \underline{X^T X} \quad \checkmark$$

$$X = \underline{\Phi^T g}$$

[Reason :- H is a symmetric matrix
 $\Rightarrow \underline{\Phi}$ is orthogonal. [eigenvectors of symmetric matrix]
 $\Rightarrow \underline{\Phi^T \Phi} = \underline{I} \Rightarrow \underline{g^T \Phi \Phi^T g} = \underline{g^T g} = \underline{g^T \Phi \Phi^T g}$]

$$\underline{\Phi^T \Phi} = \underline{\Phi \Phi^T} = \underline{I}$$

$$1. \Rightarrow \Phi \cdot \Phi = I \Rightarrow \Phi^T \cdot \Phi = I$$

$$\Rightarrow g^T H g \leq \lambda_{\max} g^T \underbrace{\Phi \Phi^T}_{I} g$$

$$\Rightarrow g^T H g \leq \lambda_{\max} \underline{\underline{g^T g}}$$

$$\Rightarrow \frac{g^T H g}{g^T H g} \leq \lambda_{\max} \frac{g^T g}{g^T H g}$$

$$\Rightarrow \frac{g^T g}{g^T H g} \geq \frac{1}{\lambda_{\max}}$$

$$\Rightarrow \boxed{\eta \geq \frac{1}{\lambda_{\max}}}$$

$$\Phi^T \Phi = \Phi \Phi^T = I$$

$$g^T \Phi^T \Phi = g^T \Phi \Phi^T = g^T$$

$$g^T \Phi^T \Phi g = g^T \Phi \Phi^T g = g^T g$$

$$\underline{\underline{\frac{1}{\lambda_{\max}} \text{ is the v. of } H.}}$$

$$l(\omega) = 7 + (\omega_1 - 1)^2 + (\omega_2 - 1)^6$$

$$\text{Now } \nabla_{\omega} l(\omega) = \begin{pmatrix} \frac{\partial l(\omega)}{\partial \omega_1} \\ \frac{\partial l(\omega)}{\partial \omega_2} \end{pmatrix} = \begin{pmatrix} 2(\omega_1 - 1) \\ 6(\omega_2 - 1)^5 \end{pmatrix}$$

$$H(l(\omega)) = \begin{pmatrix} \frac{\partial^2 l(\omega)}{\partial \omega_1^2} & \frac{\partial^2 l(\omega)}{\partial \omega_2 \partial \omega_1} \\ \frac{\partial^2 l(\omega)}{\partial \omega_1 \partial \omega_2} & \frac{\partial^2 l(\omega)}{\partial \omega_2^2} \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 6 \end{pmatrix}$$

We have $H = \begin{pmatrix} 2 & 0 \\ 0 & 6 \end{pmatrix}$ The eigenvalues of H are 2, 6.

$$\text{As } \eta \geq \frac{1}{\lambda_{\max}} \Rightarrow \boxed{\eta \geq \frac{1}{6}}$$

We consider optimal learning rate $(\eta) = 1/6$.

Gradient update rule

$$\omega_{k+1} = \omega_k - \frac{1}{6} \nabla_{\omega_k} l(\omega_k)$$

$$= \begin{pmatrix} \omega_1^k \\ \omega_2^k \end{pmatrix} - \frac{1}{6} \begin{pmatrix} 2\omega_1^k - 2 \\ 6\omega_2^k - 6 \end{pmatrix}$$

$$= \begin{pmatrix} \omega_1^k - \frac{1}{3}\omega_1^k + \frac{1}{3} \\ \omega_2^k - \omega_2^k + \frac{6}{6} \end{pmatrix} = \begin{pmatrix} \frac{2}{3}\omega_1^k + \frac{1}{3} \\ 1 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \omega_1^{k+1} \\ \omega_2^{k+1} \end{pmatrix} = \begin{pmatrix} \frac{2}{3}\omega_1^k + \frac{1}{3} \\ 1 \end{pmatrix}$$