

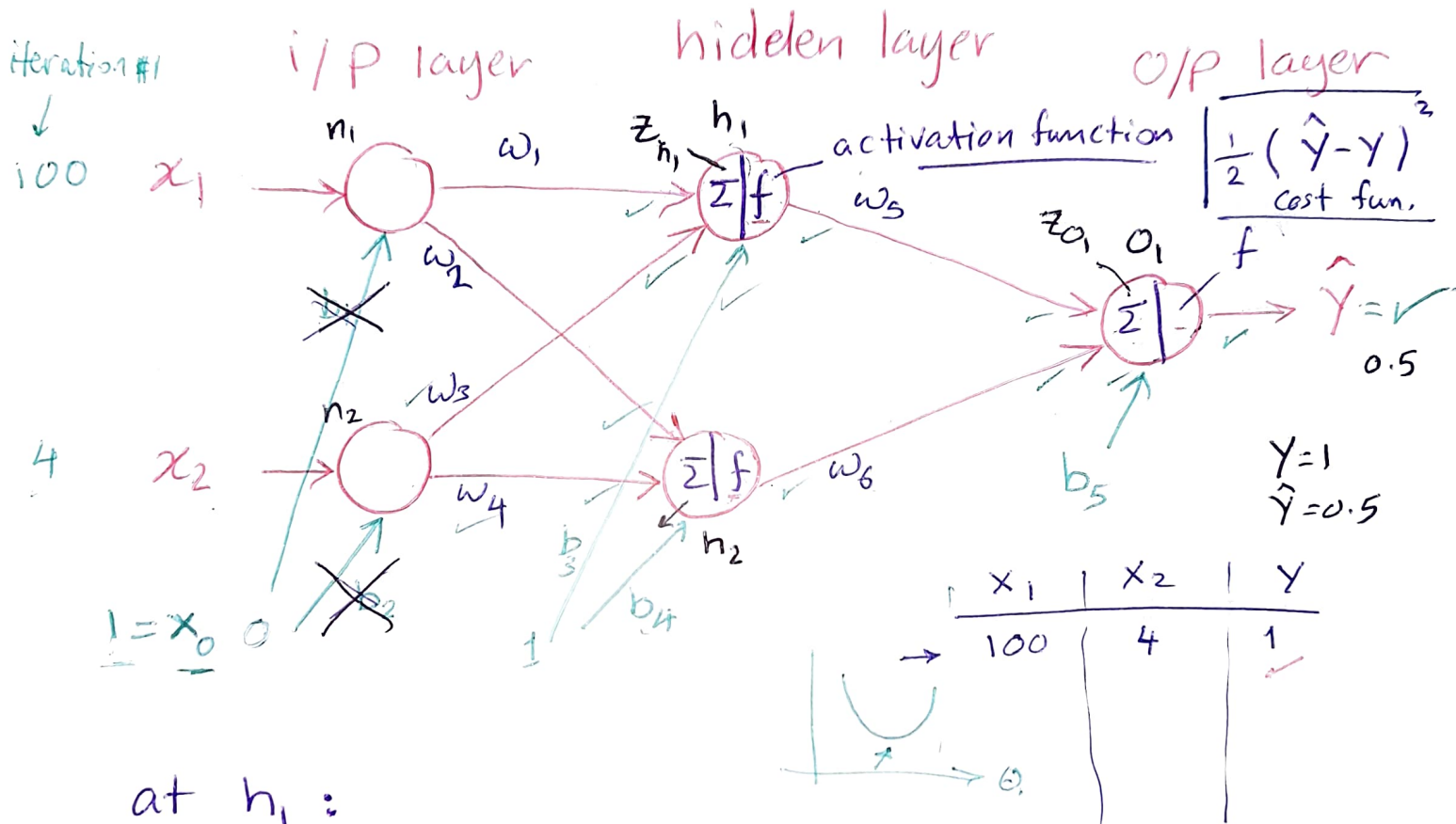
①

Numerical Optimization, AI45 Mans. session 5

11/2/2025

→ Neural Network

$$\hat{x}_0 \theta_0 + \theta_1 x_1 + \theta_2 x_2$$



$$z_{h_1} = x_1 w_1 + x_2 w_3 + b_3$$

$$\text{Out}_{h_1} = f(z_{h_1}) = f(x_1 w_1 + x_2 w_3 + b_3)$$

$$\text{Out}_{h_1} = \frac{1}{1 + e^{-z_{h_1}}} = \frac{1}{1 + e^{-(x_1 w_1 + x_2 w_3 + b_3)}}$$

$$\text{Out}_{h_2} = f(x_1 w_2 + x_2 w_4 + b_4)$$

(2)

$$z_{0,1} = \text{Out}_{h_1} w_5 + \text{Out}_{h_2} w_6 + b_5$$

$$\hat{y} = f(z_{0,1}) = \frac{1}{1 + e^{-(\text{Out}_{h_1} w_5 + \text{Out}_{h_2} w_6 + b_5)}}$$

$$= \frac{1}{1 + e^{-z_{0,1}}}$$

ex.

$$\hat{y} - y = 0.5 = \text{error} \Rightarrow J(\quad)$$

$$\frac{\partial J(\quad)}{\partial w_5}$$

$$\text{update } w_{5_{\text{new}}} = w_{5_{\text{old}}} - \alpha \frac{\partial J(\quad)}{\partial w_5}$$

$$\text{update } w_{i_{\text{new}}} = w_{i_{\text{old}}} - \alpha \frac{\partial J(\quad)}{\partial w_i}$$

for all parameters (w 's and b 's)

e.g., $J(w_i) = \frac{1}{2} (\hat{y} - y)^2$

$$\hat{y} = h_{w,b}(\vec{x}) = \boxed{\quad}$$

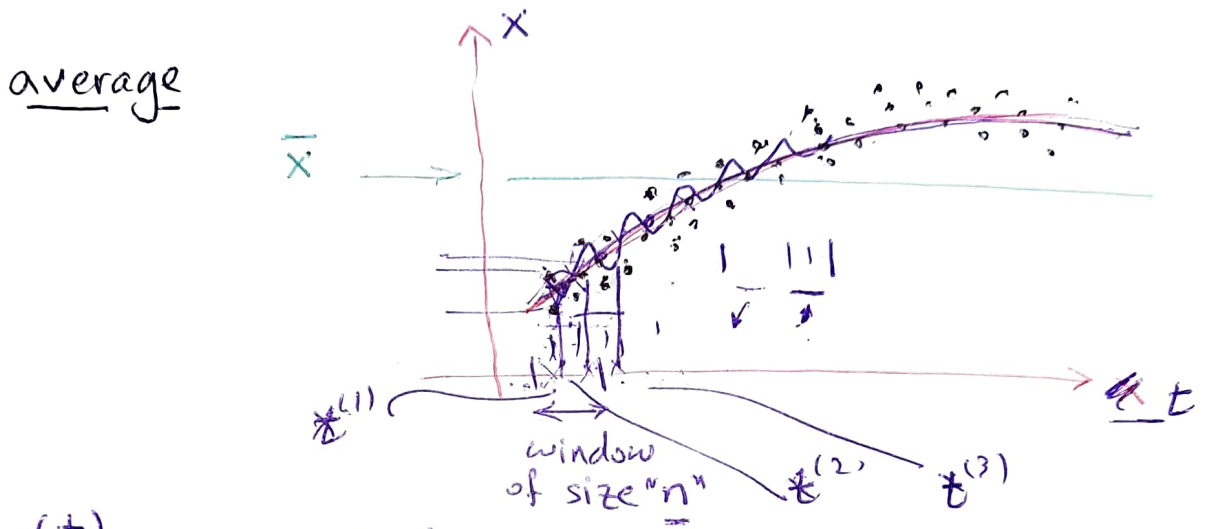
$$\frac{\partial J(\quad)}{\partial w_5} = \frac{\partial J(\quad)}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_{0,1}} \times \frac{\partial z_{0,1}}{\partial w_5} = \underline{0.1}$$

$$\Rightarrow \text{update } w_{5_{\text{new}}} = w_{5_{\text{old}}} - \alpha 0.1 = w_{5_{\text{new}}} \checkmark$$

$$\frac{\partial J(\quad)}{\partial w_1} = \frac{\partial J}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_{0,1}} \times \frac{\partial z_{0,1}}{\partial \text{Out}_{h_1}} \times \frac{\partial \text{Out}_{h_1}}{\partial z_{h_1}} \times \frac{\partial z_{h_1}}{\partial w_1}$$

→ Chain Rule ←

→ Exponentially Weighted Moving Average (EWMA)



$$\underline{v(t)} = \beta \underline{v(t-1)} + (1-\beta)x$$

$$v^{(0)} = 0$$

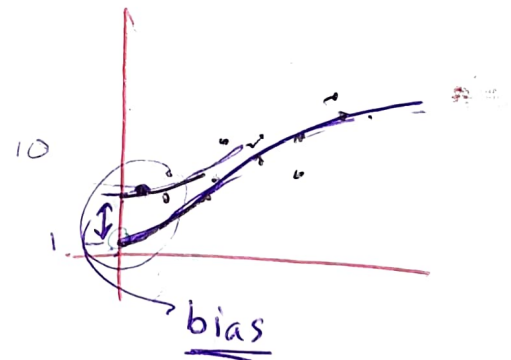
$$\underline{v^{(1)}} = \beta \underline{v^{(0)}} + (1-\beta)x^{(1)}$$

$$\begin{aligned} v^{(2)} &= \beta v^{(1)} + (1-\beta)x^{(2)} \\ &= \beta(1-\beta)x^{(1)} + (1-\beta)x^{(2)} \end{aligned}$$

$$v^{(3)} = \beta v^{(2)} + (1-\beta)x^{(3)}$$

$$v^{(3)} = \beta^2(1-\beta)x^{(1)} + \beta(1-\beta)x^{(2)} + (1-\beta)x^{(3)}$$

$$\underline{v^{(4)}} = \beta^3(1-\beta)x^{(1)} + \beta^2(1-\beta)x^{(2)} + \beta(1-\beta)x^{(3)} + (1-\beta)x^{(4)}$$



for bias
Correction

$$\hat{v}(t) = \frac{v(t)}{1 - \beta^t}$$

(4)

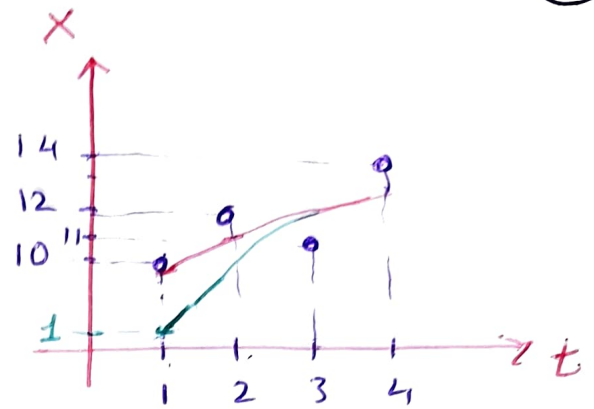
Bias correction

$$\text{let } \underline{\beta = 0.9}$$

$$v^{(0)} = 0 \quad \leftarrow$$

$$\begin{aligned} v^{(1)} &= \beta v^{(0)} + (1-\beta) x^{(1)} \\ &= 0.9 \times 0 + 0.1 \times 10 \\ &= 1 \end{aligned}$$

$$\begin{aligned} v^{(2)} &= \beta v^{(1)} + 0.1 \times 12 \\ &= 1 + 1.2 \end{aligned}$$

Bias corrected \leftarrow

$$\begin{aligned} \hat{v}^{(1)} &= \frac{v^{(1)}}{1-\beta^1} \\ &= \frac{1}{1-0.9^1} = \frac{1}{0.1} = 10 \end{aligned}$$

⋮

$$\hat{v}^{(5)} = \frac{v^{(5)}}{1-\beta^5}$$

=

for large time index.

$$1-\beta^t \approx 1$$

$$\hat{v}^{(t)} \approx v^{(t)}$$

(5)

Review Momentum

Standard form

$$v^{(t)} = \beta v^{(t-1)} + \alpha \nabla (J(\theta^{(t)}))$$

$$\theta^{(t+1)} = \theta^{(t)} - v^{(t)}$$

$$\theta^{(t+1)} = \theta^{(t)} - \beta v^{(t-1)} - \alpha \cdot \underline{\text{grad.}}$$

EWMA form

velocity, $v^{(t)}$

$$v^{(t)} = \beta v^{(t+1)} + (1-\beta) \nabla (J(\theta^{(t)}))$$

momentum, $m^{(t)}$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha v^{(t)}$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \beta v^{(t+1)} - \alpha (1-\beta) \cdot \underline{\text{grad.}}$$

scaled
learning
rate

$$\tilde{\alpha} = \frac{\alpha}{1-\beta}$$

(6)

Review RMSprop.

$$v^{(+)} = \beta v^{(+ - 1)} + (1 - \beta) (\text{grad } L)^2$$

$$\theta^{(++)} = \theta^{(+)} - \frac{\alpha}{\sqrt{v^{(+)} + \epsilon}} \nabla (J(\theta^{(+)}))$$

Adam (Adaptive Momentum)

EWMA momentum

$$m^{(+)} \equiv v_1^{(+)} = \beta_1 v_1^{(+ - 1)} + (1 - \beta_1) \underbrace{\nabla (J(\theta^{(+)}))}_{\text{grad.}}$$

$$v_2^{(+)} = \beta_2 v_2^{(+ - 1)} + (1 - \beta_2) \left(\nabla J(\theta^{(+)} \right)^2$$

RMS prop.

$$\theta^{(++)} = \theta^{(+)} - \frac{\alpha}{\sqrt{\hat{v}_2^{(+)} + \epsilon}} \hat{v}_1^{(+)}$$

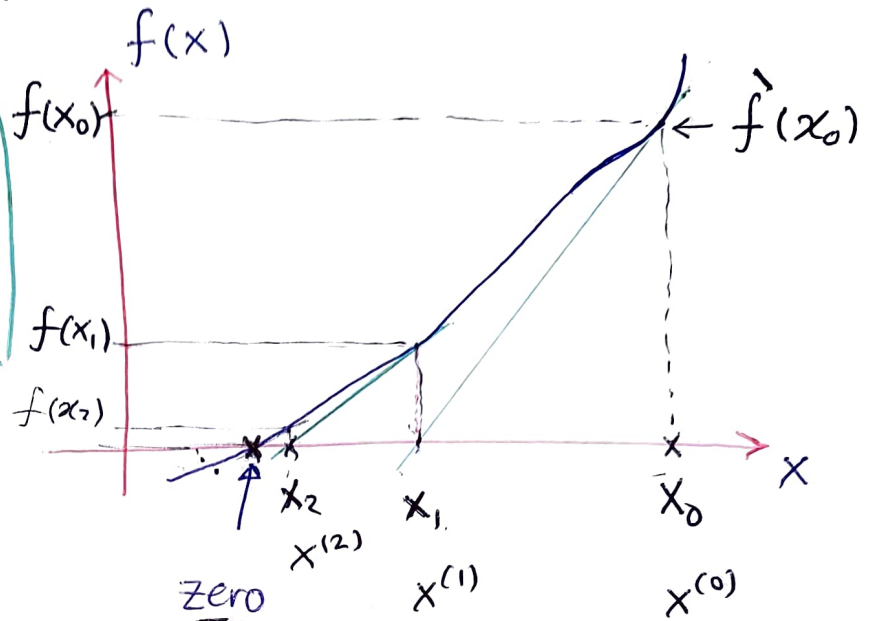
bias-corrected $v_2^{(+)}$ bias-corrected $\hat{v}_1^{(+)}$

commonly $\beta_1 = 0.9$ $\beta_2 = 0.999$ $\epsilon = 1 \times 10^{-8}$

Newton's Method

7

Zero of a function;
value of x s.t. $f(x)=0$



$$\Rightarrow f'(x_0) = \frac{\Delta y}{\Delta x} = \frac{f(x_0) - 0}{x_0 - x_1}$$

$$f'(x_0) = \frac{f(x_0)}{x_0 - x_1}$$

$$\Rightarrow x_0 - x_1 = \frac{f(x_0)}{f'(x_0)} \Rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

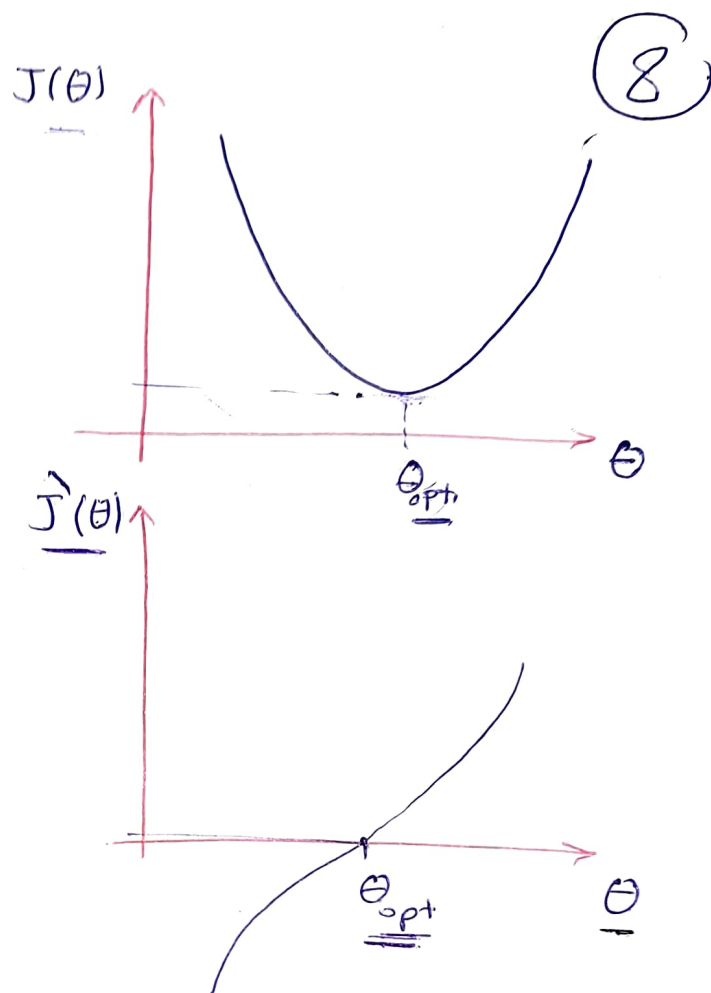
$$x^{(+)} = x^{(+1)} - \frac{f(x^{(+1)})}{f'(x^{(+1)})}$$

Zero of $J'(\theta)$

\Rightarrow minimum of $J(\theta)$

(or maximum, or
saddle point, or
local minimum!)

\rightarrow use Newton's method
to find zero of $J'(\theta)$



$$\theta^{(t+1)} = \theta^{(t)} - \frac{J'(\theta^{(t)})}{J''(\theta^{(t)})}$$

for multivariate cases

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \frac{\nabla J(\theta^{(t)})}{?}$$

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - (H^{-1}) \nabla J(\theta^{(t)})$$

GD $O(n)$

complexity? $O(n^3)$

(9)

Remember

gradient $\nabla f = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ \vdots \\ \frac{\partial f}{\partial \theta_n} \end{bmatrix}$

e.g., $f(x, y) \Rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$

Hessian

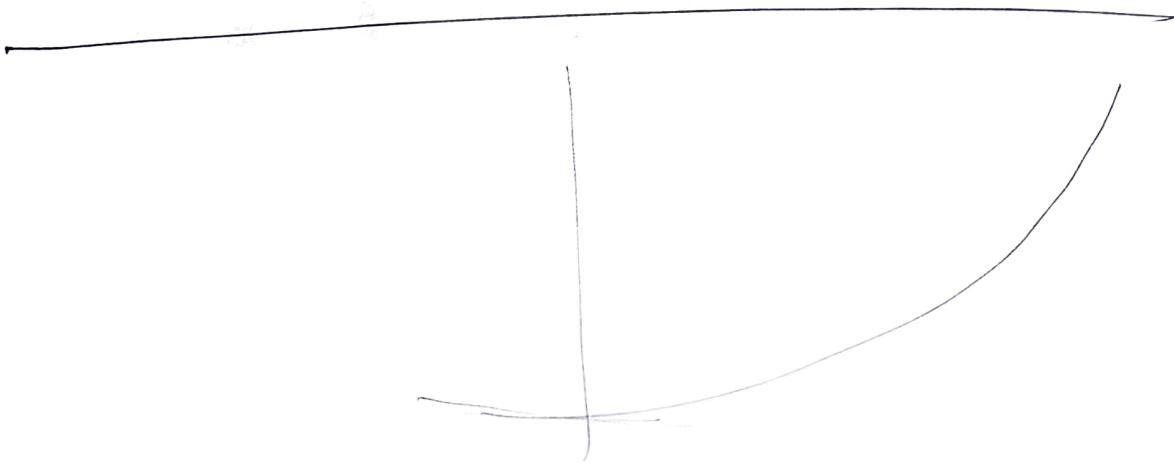
$$H = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$

$$H = \begin{bmatrix} H_{i,j} \end{bmatrix}$$

$$H_{i,j} = \frac{\partial^2 f}{\partial \theta_i \partial \theta_j}$$

$$\nabla J(\theta_1, \dots, \theta_n) = \begin{bmatrix} \partial \end{bmatrix}$$

- Quasi-Newton* methods
- Secant method
- BFGS (Broyden, Fletcher, Goldfarb, Shanno,)



note

note that the same method can be expressed in different forms.

e.g., momentum method (standard form)

$$v^{(t)} = \beta v^{(t-1)} + \alpha \nabla J(\theta^{(t)})$$

$$\theta^{(t+1)} = \theta^{(t)} - v^{(t)}$$

OR

$$v^{(t+1)} = \beta v^{(t)} + \alpha \nabla J(\theta^{(t)})$$

$$\theta^{(t+1)} = \theta^{(t)} - v^{(t+1)}$$

OR

$$v^{(t+1)} = \beta v^{(t)} - \alpha \nabla J(\theta^{(t)})$$

$$\theta^{(t+1)} = \theta^{(t)} + v^{(t+1)}$$

and so on, ...

Regardless of how it is written, the parameters are updated in the same manner.