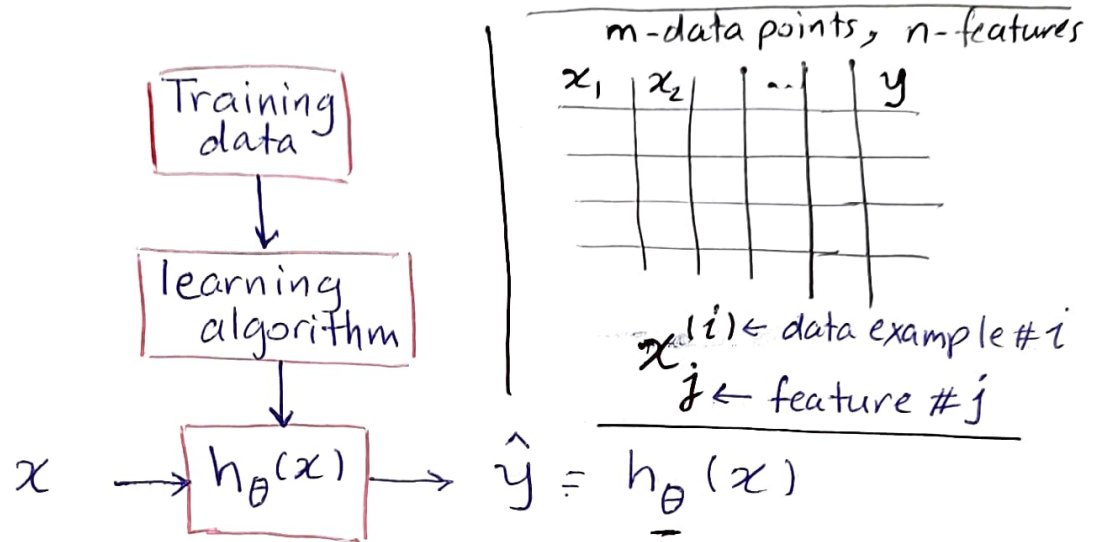


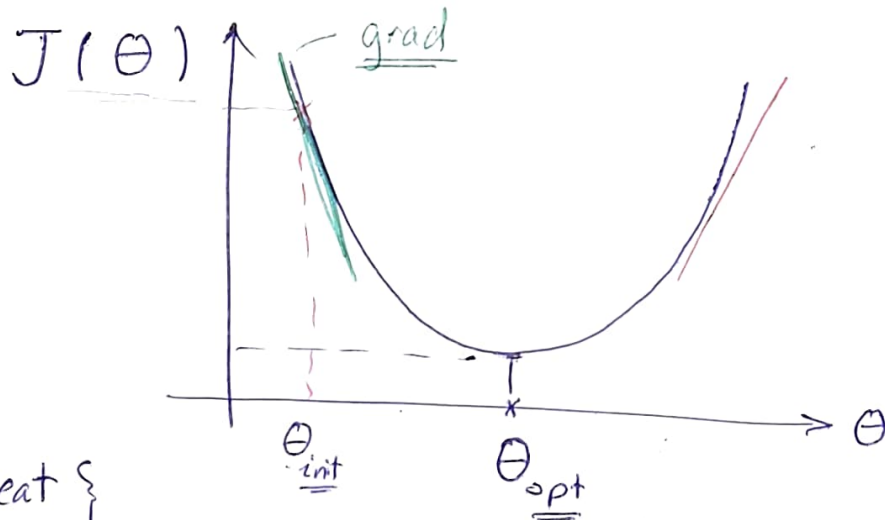
①

→ Numerical Optimization for ML, session 4, Mans.
AI45, 10/2/2025

Review



→ Cost function of $(\hat{y} - y)$ "error"



Repeat {

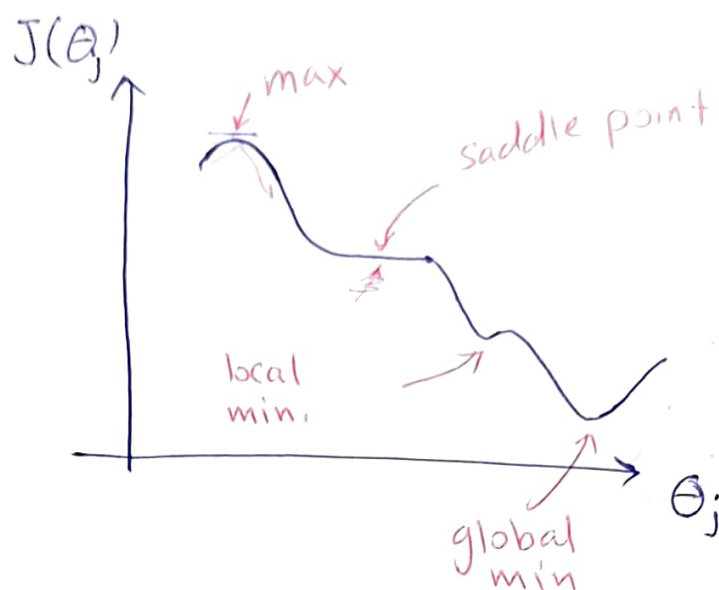
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n)$$

assign new value to the parameter θ_j

} until convergence

(2)

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0$$



$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\vec{\theta})$$

$$\theta_{j\text{-new}} = \theta_{j\text{-old}} - \alpha \frac{\partial}{\partial \theta_j} J(\vec{\theta})$$

for any θ_j

(note that we update all θ_j 's simultaneously)

n -features ; x_1, x_2, \dots, x_n

$\Rightarrow (n+1)$ -parameters $\theta_0, \theta_1, \dots, \theta_n$

e.g., linear regression

$$h_{\vec{\theta}}(\vec{x}) = \theta_0 \overset{x_0=1}{x_0} + \theta_1 x_1 + \dots + \theta_n x_n$$

e.g., logistic regression

$$h_{\vec{\theta}}(\vec{x}) = \frac{1}{1 + e^{-(\theta_0 \overset{x_0=1}{x_0} + \theta_1 x_1 + \dots + \theta_n x_n)}}$$

③

$$\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\vec{\theta} := \vec{\theta} - \underline{\alpha} \nabla (J(\vec{\theta}))$$

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

$$\vec{\theta}_{k+1} = \vec{\theta}_k - \underline{\alpha} \nabla (J(\vec{\theta}_k))$$

→ Vanilla gradient descent

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

↖ for all data points m

→ Variants :

- SGD

- mini-batch GD

accelerating GD

Momentum-based methods → Momentum GD

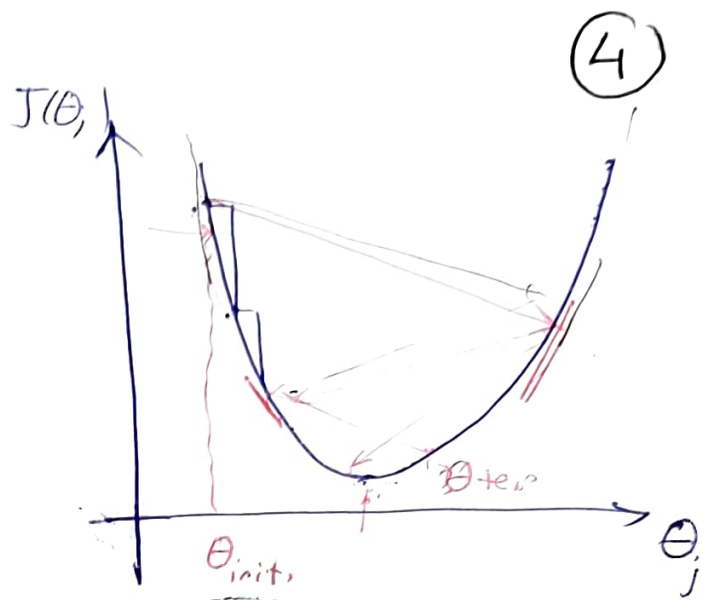
→ Nesterov method

- Momentum GD

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

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

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



$$\theta^{(1)} = \theta_{init}$$

①

$$v^{(1)} = 0 + \alpha \nabla J(\theta^{(1)})$$

$$\theta^{(2)} = \theta^{(1)} + v^{(1)}$$

$$\theta^{(2)} = \theta^{(1)} - \alpha \nabla J(\theta^{(1)})$$

②

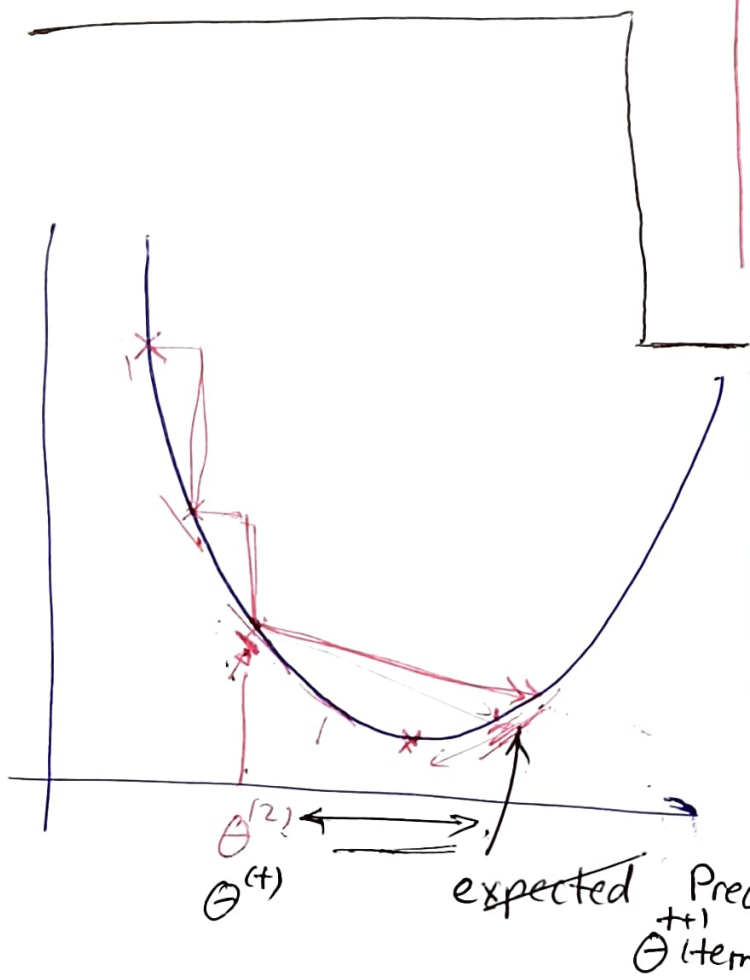
$$v^{(2)} = \beta v^{(1)} + \alpha \nabla J(\theta^{(2)})$$

$$v^{(2)} = \beta (\alpha \nabla J(\theta^{(1)})) + \alpha \nabla J(\theta^{(2)})$$

$$\theta^{(3)} = \theta^{(2)} -$$

(exponentially decaying
weighted average)

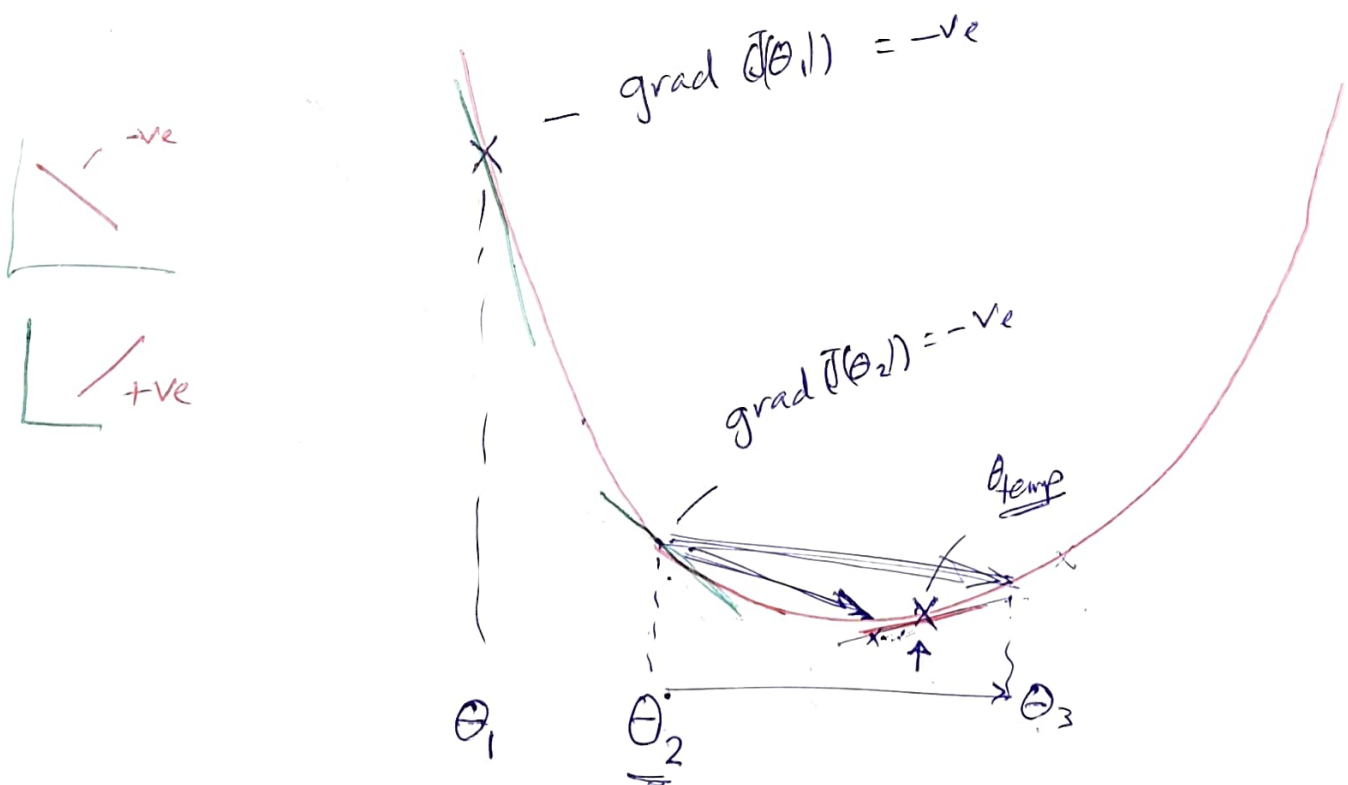
Nesterov's ?



Nesterov method ;

$$v^{(+)} = \beta v^{(+-1)} + \alpha \nabla J(\underbrace{\theta^{(+)} - \beta v^{(+-1)}}_{\theta_{temp}})$$

$$\begin{aligned} \theta^{(++1)} &= \theta^{(+)} - v^{(+)} \\ &= \boxed{\theta^{(+)} - \beta v^{(+-1)}}_{\theta_{temp}} - \alpha \nabla J(\theta_{temp}) \end{aligned}$$



Momentum

$$\begin{aligned} v^{(+)} &= \beta (\underbrace{\text{grad}(J(\theta_1))}_{-ve} + \underbrace{\text{grad}(J(\theta_2))}_{-ve}) \\ &= -ve \uparrow \end{aligned}$$

$$\theta_3 \approx \theta_2 - (-ve \uparrow)$$

$$\theta_3 \gg \theta_2$$

Nesterov

$$v^{(+)} = \beta (\underbrace{\text{grad}(J(\theta_1))}_{-ve} + \underbrace{\text{grad}(\theta_{temp})}_{-ve})$$

$$\theta_{temp} = \theta_2 - \beta \text{grad}(J(\theta_1))$$

$$\theta_{temp} \gg \theta_2$$

$$\theta_3 = \theta_2 - (-ve/+ve) \downarrow$$

6

Accelerating GD

→ learning rate

$$\vec{\theta} := \vec{\theta} - \alpha \nabla J(\vec{\theta})$$

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix}$$

- Can we use different α for each parameter?

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \vec{\alpha} = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \alpha_j \frac{\partial}{\partial \theta_j} J(\vec{\theta}^{(t)})$$

for $\vec{\alpha}$

$$\vec{\theta} := \vec{\theta} - \vec{\alpha} \odot \nabla J(\vec{\theta})$$

vector
vector
vector $\in \mathbb{R}^{n+1}$

sparse

| x_5 |
|-------|
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 1 |

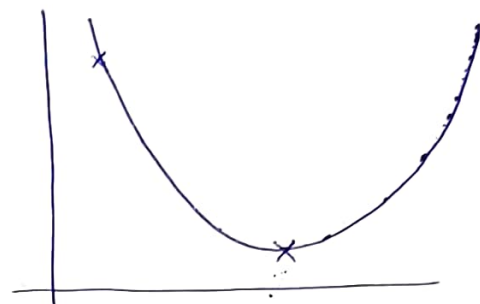
Hadamard product (element-wise)

$$\vec{X} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \vec{Y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \quad \vec{X} \odot \vec{Y} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \\ x_3 y_3 \end{bmatrix}$$

(7)

→ Decaying learning rate

Starting with large α ,
decrease α with time.



I- adaptive gradient (Adagrad)

$$\Theta_j^{(t+1)} = \Theta_j^{(t)} - \frac{\alpha_j}{\epsilon + \sqrt{\sum_{k=1}^t (\text{grad}(\Theta_j^{(k)}))^2}} \text{grad}(\quad)$$

↑
numerical stabilizer

expressed with velocity term $v^{(t+1)} = v^{(t)} + (\text{grad}(\quad))^2$

$$\Theta^{(t+1)} = \Theta^{(t)} - \frac{\alpha}{\epsilon + \sqrt{v^{(t)}}} \text{grad}(\quad)$$

(8)

II- RMS Prop (root mean square propagation)

$$\Theta^{(t+1)} = \Theta^{(t)} - \left[\frac{\alpha}{\epsilon + \sqrt{v^{(t)}}} \right] \underbrace{\text{grad}(\Theta^{(t)})}_{\text{adaptive learning rate}}$$

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

III - Adam (A d a p t i v e m o m e n t u m)