

27/12/2025

①

Numerical Optimization for ML

→ Learning outcomes

- understanding Convexity
- " optimization
- " Gradient Descent algorithm
- #
 - vanilla GD
 - Variants GD
- momentum-based Algorithms
 - ADAM Adaptive Momentum
 - Adam
- Newton's method
- be able to apply these algorithm.

Session 1

②

Machine Learning

e.g. Linear Regression

model ~~params~~
Parameters.
↑
k

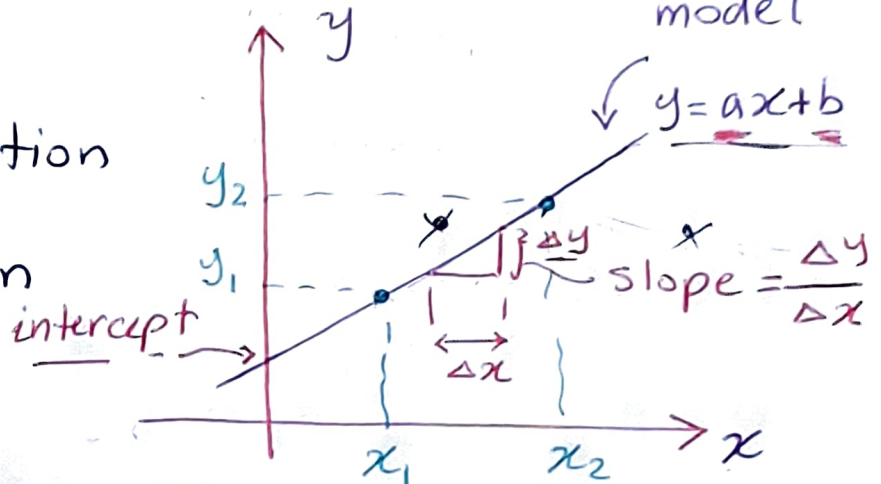
$$y = \theta_0 + \theta_1 x$$

↑
model

Case #1 exact solution

→ Analytic solution

↳ Cannot be generalized



$$y_1 = ax_1 + b$$

$$y_2 = ax_2 + b \rightarrow \begin{matrix} \text{Model} & \text{Parameters} \end{matrix} \quad \begin{matrix} 2 \text{ equations,} \\ 2 \text{ unknowns} \end{matrix}$$

⇒ unique, exact solution

"Analytic"

features

data

matrix

i	x		y
	x ₁	x ₂	y ₁
1	x ₁	y ₁	
2	x ₂	y ₂	

i	x ₁	x ₂	...	x _n	y
1					
2					
3					
⋮					
m					

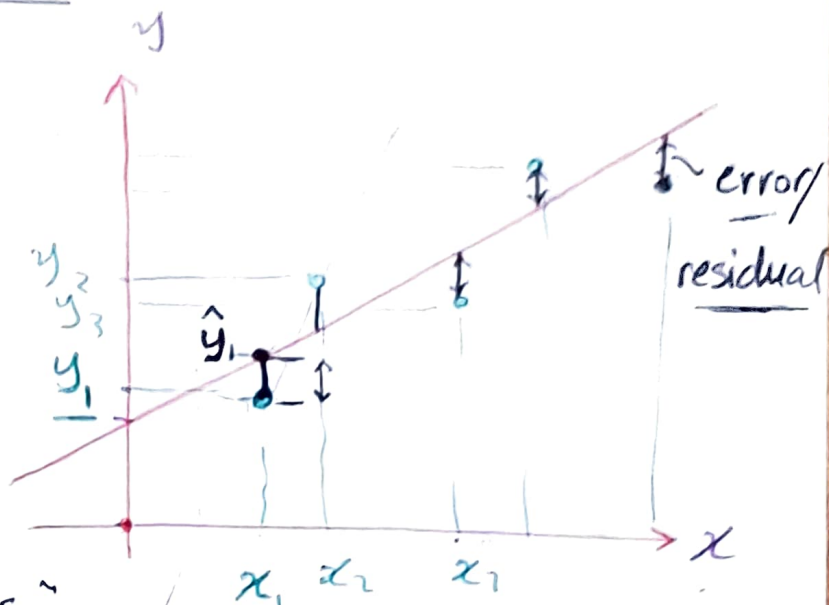
(3)

approximate solution

- Linear model

$$y = \theta_0 + \theta_1 x$$

$$y_i = \theta_0 \times 1 + \theta_1 x_i$$



e.g., "least squares"

example, using pseudo inverse.

$$m > n$$

"overdetermined system"Algebraic solution can be used

to find approximate solution,

$$\begin{bmatrix} y \\ y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$\vec{y} = X \vec{\theta}$$

$m \times 2$

$$X^T \vec{y} = X^T X \vec{\theta}$$

$2 \times m \quad m \times 2$

$$(X^T X)^{-1} X^T \vec{y} = (X^T X)^{-1} (X^T X) \vec{\theta}$$

$$\vec{\theta} = (X^T X)^{-1} X^T \vec{y}$$

$$\begin{array}{c} \overbrace{\quad\quad\quad}^n \\ i \\ \hline \begin{array}{cc} 1 & \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} \\ 2 & \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} \\ \vdots & \vdots \\ m & \begin{bmatrix} x_m \\ y_m \end{bmatrix} \end{array} \end{array}$$

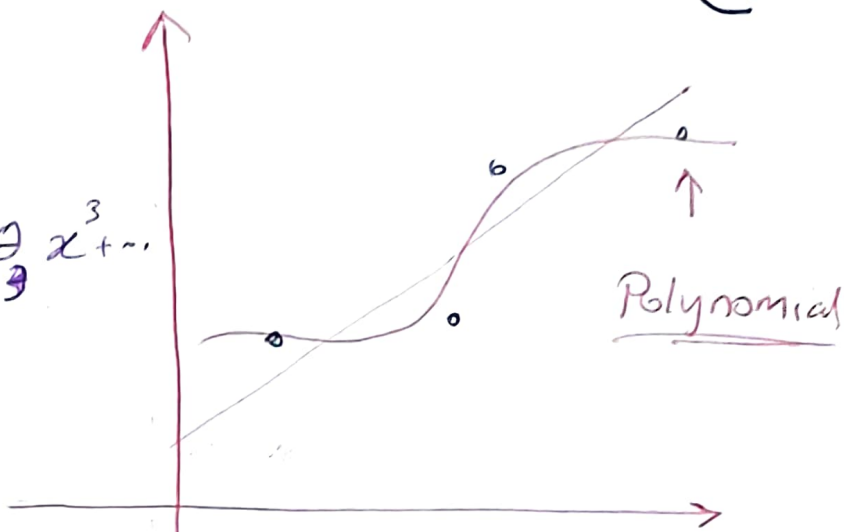
$X_{m \times 2} \quad \vec{y}_m$

polynomial model

4

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots$$

$$\vec{y} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \theta_0 \\ 0 \end{bmatrix}$$



Least norm method;

"Algebraic solution"

$$n > m$$

underdetermined system,

$$\vec{y} = \cancel{X} \vec{\theta}$$

→ infinite number of solutions

$$\boxed{\theta = X^T (X X^T)^{-1} y}$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{1n} \\ 1 & x_{21} & x_{2n} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$X_{m \times n}$$

data matrix

m rows: # of data examples

n: columns: # of features

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

$$X_{m \times n}$$

e.g. 2

$$y = \theta_0 \times 1 + \theta_1 x_1 + \dots$$

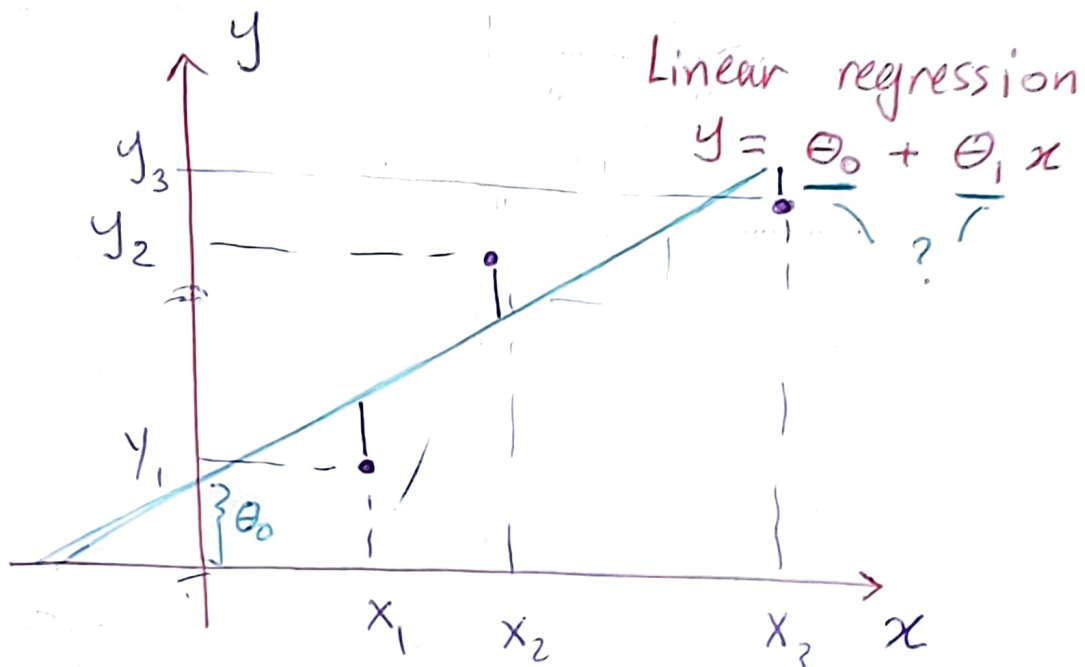
numerical solutions

(5)

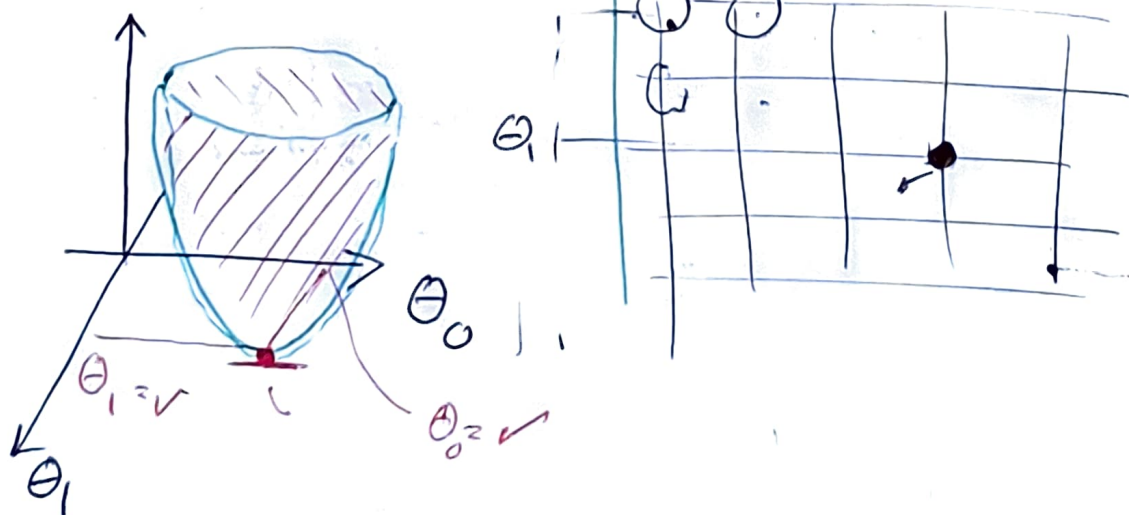
'Optimization'

minimization

maximization



grid search



Linear regression problem

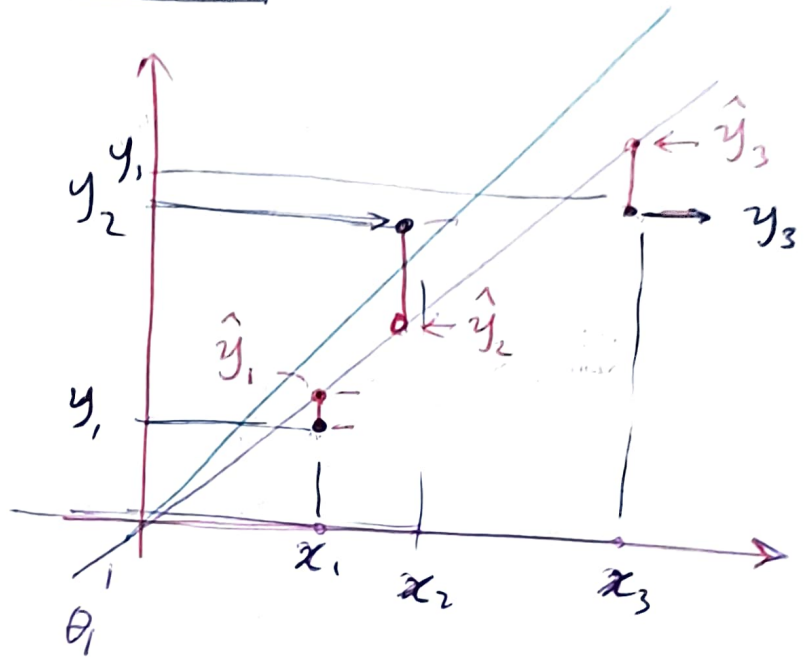
(6)

1- parameter

intercept = 0

$$\theta_1 = ?$$

$$\hat{y} = \theta_1 x$$



Error ; residuals

$$y_i - \hat{y}_i = y_i - \theta_1 x_i$$

Sum of Absolute values of Residuals

Sum of Squared values of Residuals.

Sum of Squared error
↓

Sum of Squared Residuals = ~~MA~~ SE

RSS

Residual squares sum.

$$\sum_{i=1}^m (y_i - \hat{y}_i)^2 = \sum_{i=1}^m (y_i - \underbrace{\theta_1 x_i}_{\text{loss}})^2$$

cost function

Optimization problem

$$\underset{\theta_1}{\text{minimize}} (RSS) = \underset{\theta_1}{\text{min}} \underbrace{\sum_{i=1}^m (y_i - \theta_1 x_i)^2}_{\text{objective fun}}$$

(7)

- For linear regression model

$$\hat{y}_i = \theta_0 + \theta_1 x_i$$

find optimum values of parameters θ

the minimize prediction error $|\vec{y} - \vec{\hat{y}}|$

→ Steps involve:

~~Assume~~ ^{select} certain model (linear regression)

→ select certain error function (SSR, MAE, MSE, SAT)

1) Parameter initialization (assume $\theta = 0$)

or any other initial value.

2) Predict output \hat{y} using model

3) Evaluate prediction error $|y_i - \hat{y}_i|$

4) Find direction & magnitude of changing
-ve, or +ve
of model parameters (θ)

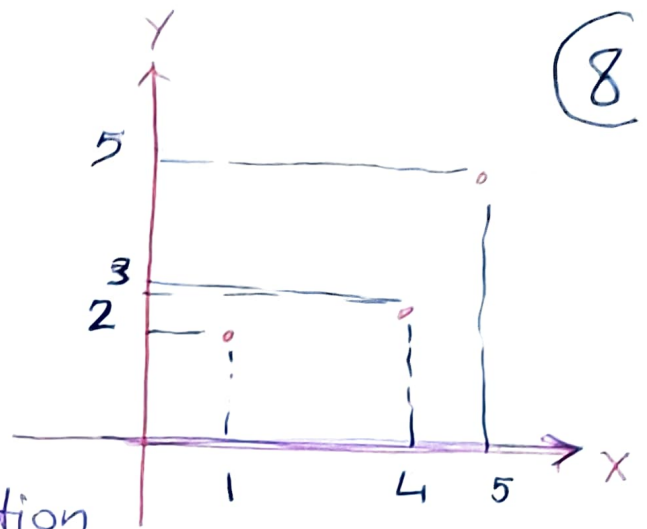
5) update model parameters

6) go to step 2 (repeat untill convergence)

- linear regression

$$\rightarrow \hat{y} = \theta_1 x$$

\rightarrow using SSR



1) parameter initialization

let $\theta_1 = \text{zero}$ (or any initial value(s))

$$\hat{y} = \text{zero} \times x = 0$$

\rightarrow 2) Predict output

3) error = $|y_i - \hat{y}_i|$

$$= |\hat{y}_i - y_i| \quad \sim \text{SSR}$$

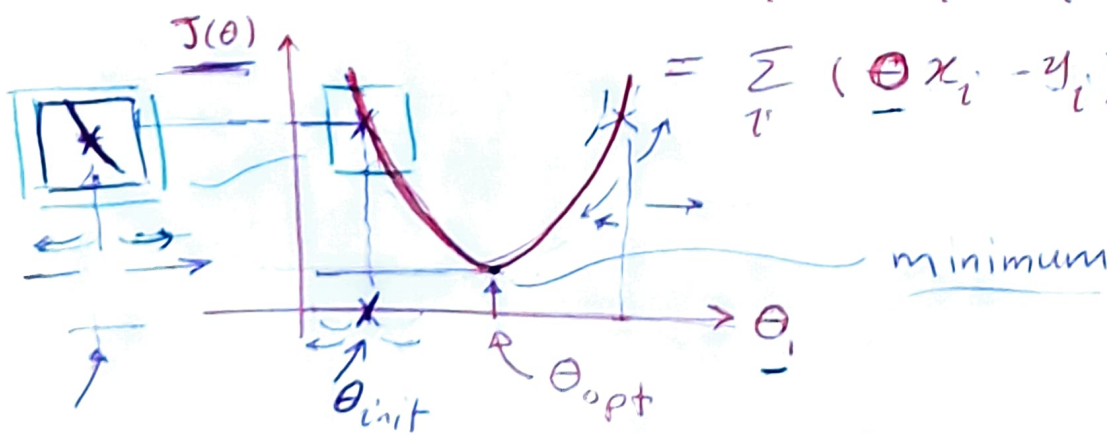
$\theta_1 = 0$				
i	x	y	$\hat{y} = \theta_1 x$	error
1	1	2	0	2
2	4	3	0	3
3	5	5	0	5

4) $J(\theta) = \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \underline{2^2 + 3^2 + 5^2}$

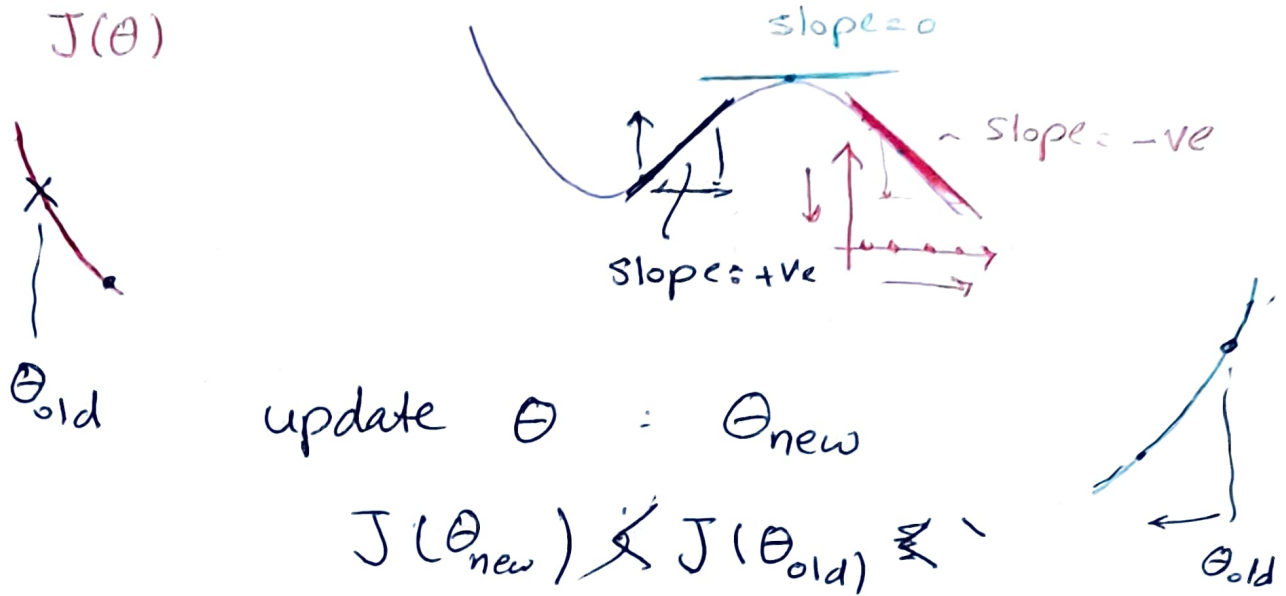
5) update θ ; $\theta_{\text{new}} = \theta_{\text{old}} - \underbrace{\text{const.}}_{\text{learning rate}} \cdot \text{gradient}$

SSR : $J(\theta) = \sum_i (\hat{y}_i - y_i)^2$

$$= \sum_i (\theta x_i - y_i)^2$$



slope \equiv derivative of curve \equiv grad. (9)

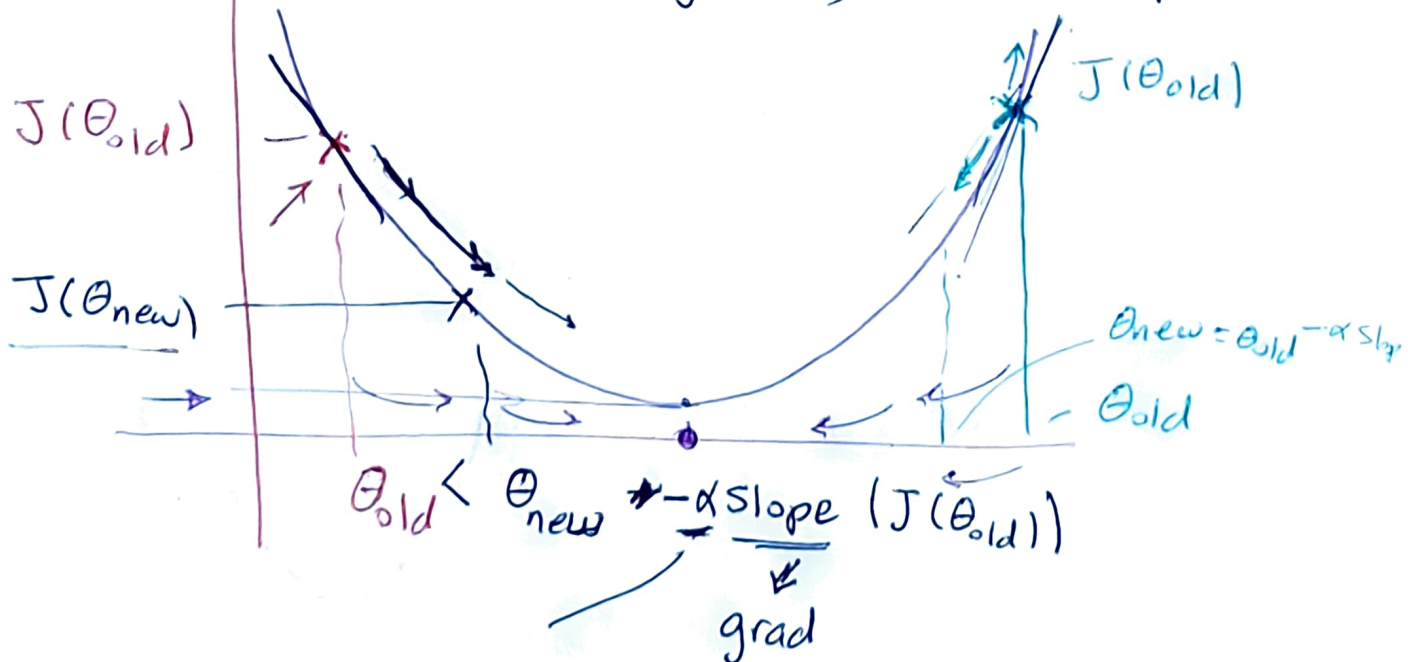


update θ :

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \text{grad}(J(\theta_{\text{old}}))$$

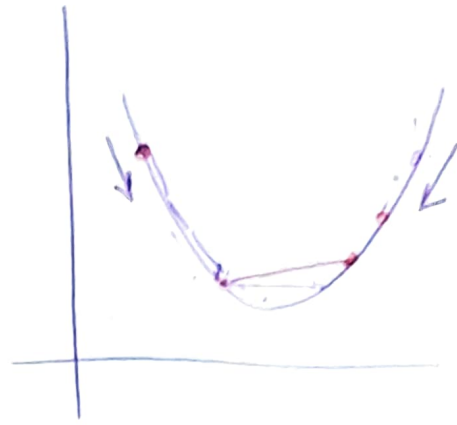
Const.

"learning rate" ; control step size

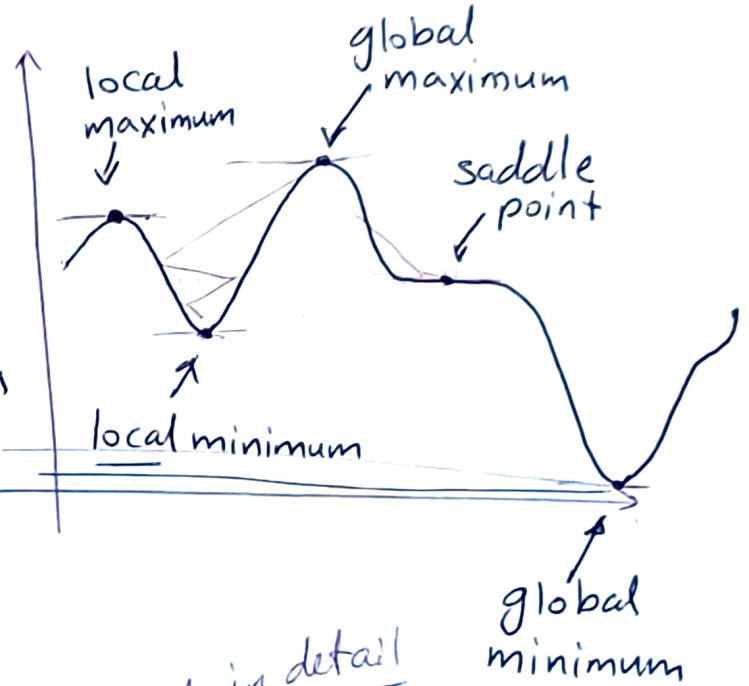


→ Convex function

- global minimum
- "local minimum"



→ non convex



- minimum ; Plural : minima
- maximum ; plural : maxima

step?

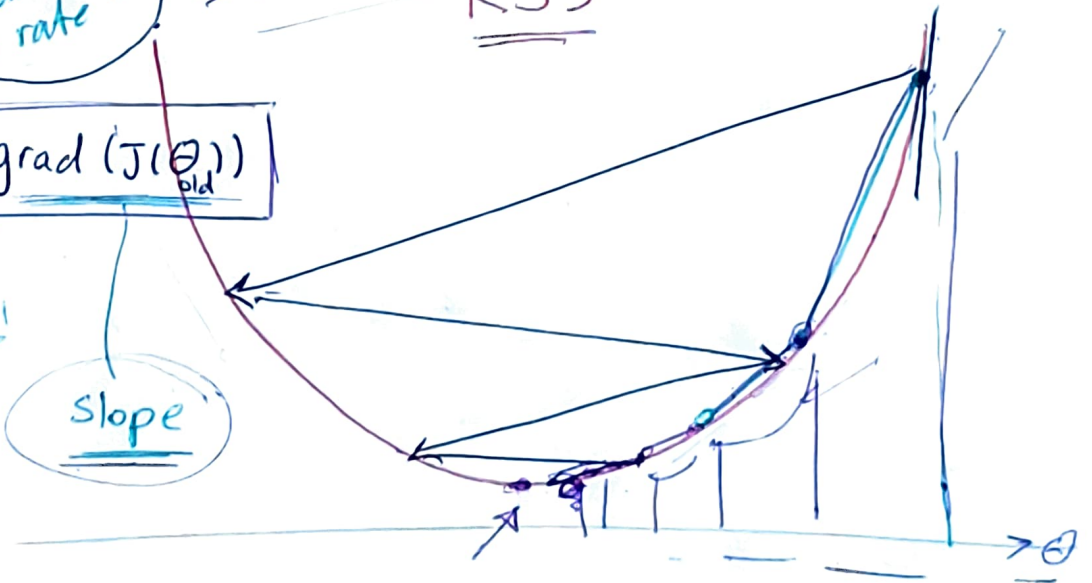
learning rate

to be discussed in detail
RSS

$$\Theta_{\text{new}} = \Theta_{\text{old}} - \alpha \text{grad}(\mathcal{J}(\Theta_{\text{old}}))$$

step

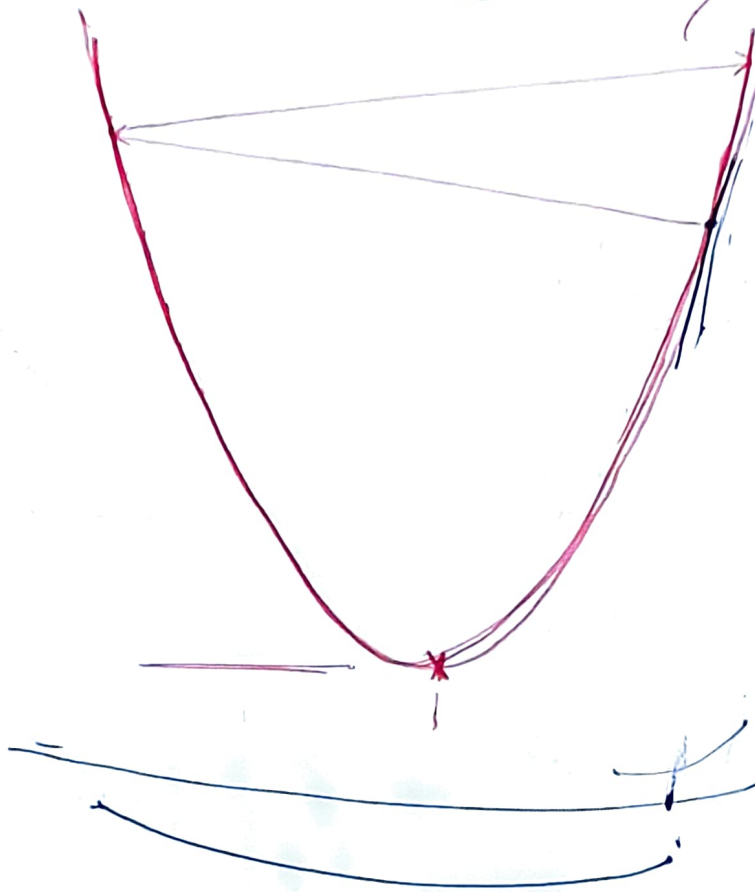
slope



too big learning rate

divergence (11)
not
convergence

α : too large



SSR

(12)

$$= \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

i	X	Y	\hat{y}
1	x_1	y_1	\hat{y}_1
2	x_2	y_2	\hat{y}_2
3	\vdots	\vdots	\vdots
m	x_m	y_m	\hat{y}_m

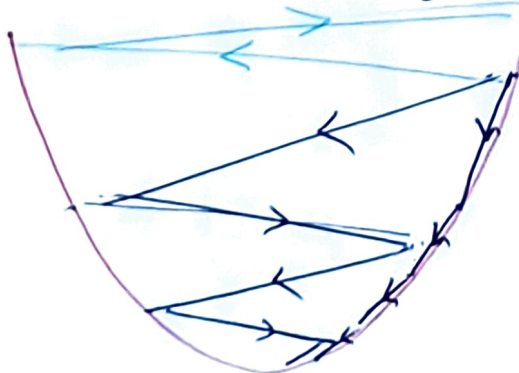
$$\|\vec{\text{error}}\| \equiv \|\vec{\text{Residual}}\| = \|\vec{\hat{y}} - \vec{y}\|$$

l_2 -norm of a vector \Rightarrow ~~l_1 -minimization~~
 l_2 -minimization

$$\|\vec{\text{error}}\| = \left\| \begin{bmatrix} \hat{y}_1 - y_1 \\ \hat{y}_2 - y_2 \\ \vdots \\ \hat{y}_m - y_m \end{bmatrix} \right\| = \sqrt{\sum_{i=1}^m (\hat{y}_i - y_i)^2}$$

Mean squared error MSE

$$MSE = \frac{SSR}{m} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$



l_1 - minimization

(13)

l_1 - norm of a vector $\vec{e} = \vec{\hat{y}} - \vec{y}$

$$= \sum_{i=1}^n |e_i| \quad \vec{v} = \begin{bmatrix} 1 \\ -1 \\ -5 \\ 2 \\ 4 \end{bmatrix}$$

ex

$$l_1\text{-norm}(\vec{v}) = |1| + |-1| + |-5| + |2| + |4| = 13$$

MAE : mean absolute error

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (|\hat{y}_i - y_i|)$$

$$= \sum_{i=1}^n |(\theta x_i - y_i)|$$

