

Lecture 6. GD.

1. Linear Regression

Four component: ① $\{(x_i, y_i)\}_{i=1}^n$

generalized
feature vector

$$x_i = \begin{pmatrix} x_{i0} = 1 \\ x_{i1} \\ \vdots \\ x_{im} \end{pmatrix}$$

target value.
 $y_i \in \mathbb{R}$

②

model: $f(x_i) = w^T \cdot x_i$

model parameters
 $w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{pmatrix}_{(m+1) \times 1}$
 $= w_0 \cdot x_{i0} + w_1 \cdot x_{i1} + \dots + w_m \cdot x_{im}$

③

loss:

$$L(w) = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - \underline{y}_i)^2$$

prediction $f(x)$

$$= \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2$$

④ Optimization: Find w^* (best model parameters)

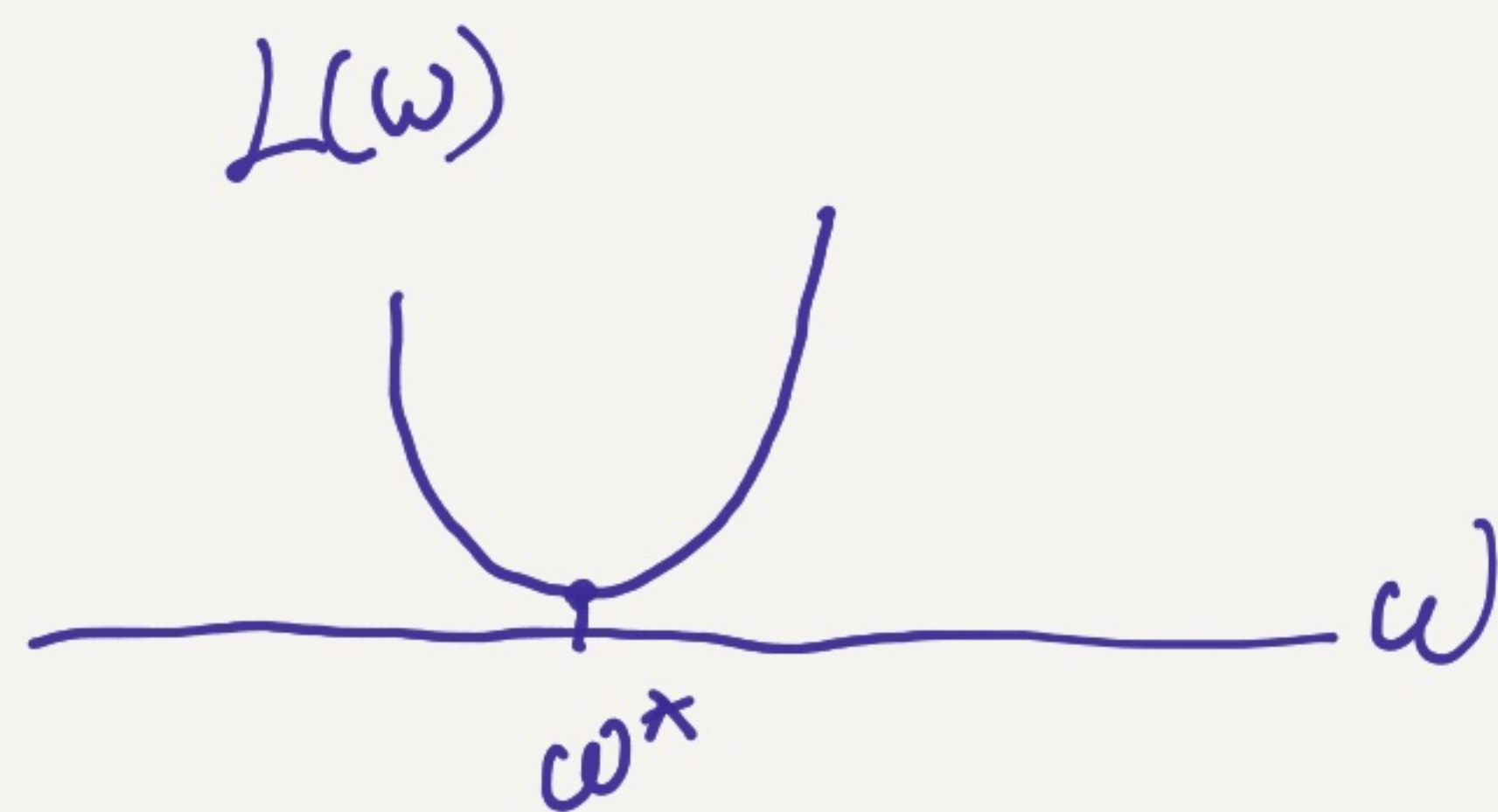
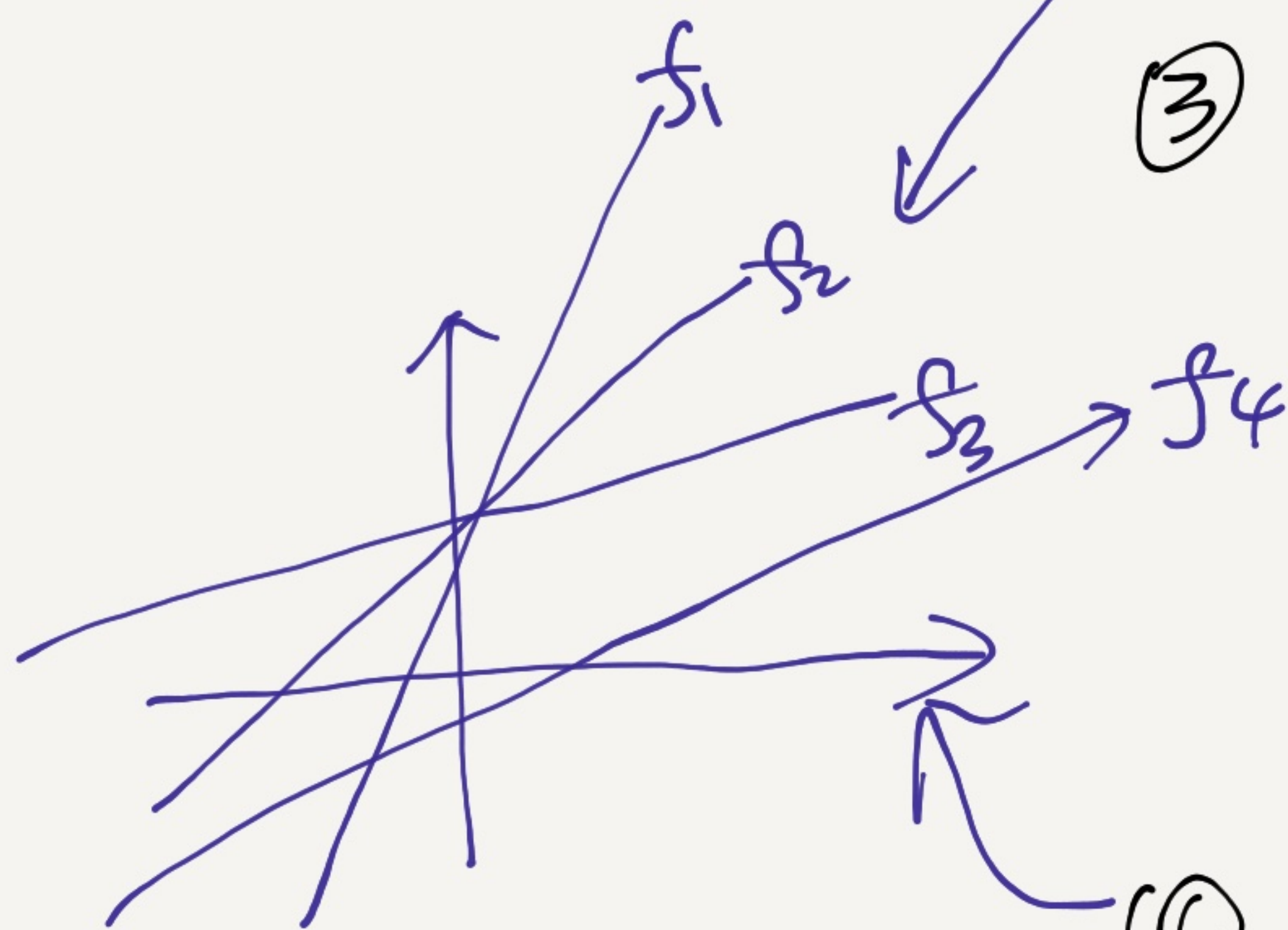
$$w^* = \arg \min_{w} L(w) = \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\sum_{i=1}^n z_i^2 = z^T \cdot z \quad z = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix}$$

$$= \frac{1}{2} \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_n - y_n \end{pmatrix}^T (w^T x_1 - y_1, \dots, w^T x_n - y_n)$$

$$b = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1} \quad A = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}_{n \times m}$$

$$= \frac{1}{2} (Aw - b)^T (Aw - b)$$



$$\nabla L(w)|_{w=w^*} = 0 \quad (m+1, \text{equations})$$

$$\nabla_w L(w) = \underbrace{\left(\frac{1}{2} \cdot 2\right)}_1 A^T (Aw - b) = 0 \quad L(w) = \frac{1}{2} (Aw - b)^T (Aw - b)$$

$$A^T A w = A^T b \quad \swarrow \text{inverse}$$

$$w^* = \underbrace{(A^T A)^{-1} A^T b}_{\text{Analytic solution}}$$

challenges:

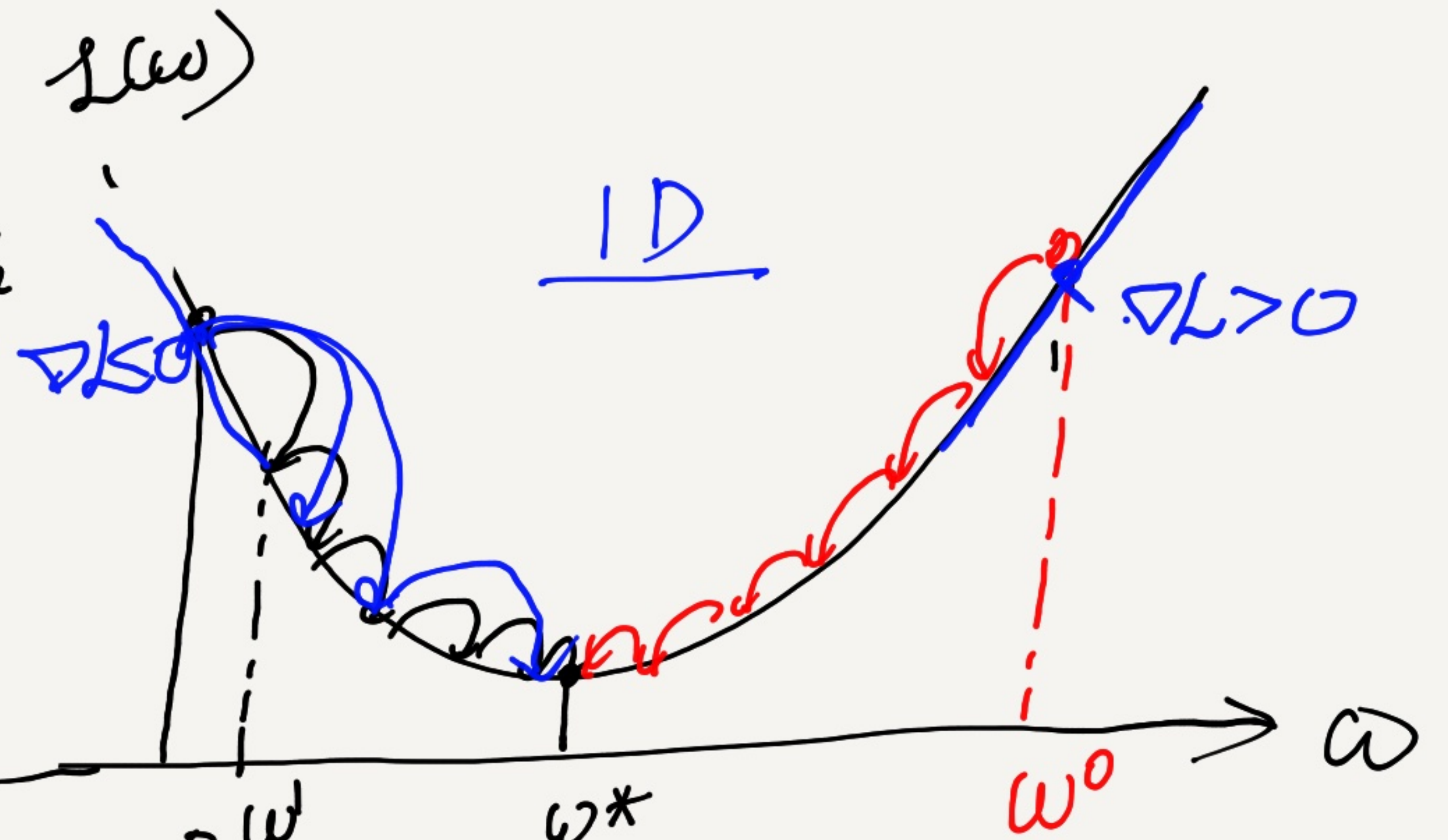
① A could be very large. we may do not have enough resources to load it and perform the calculation.

$$A^T \cdot A \quad . \quad (A^T \cdot A)^T$$

② $(A^T A)^{-1}$ may not exist.

2. Gradient-based approaches.

1). main idea: start from a random location, e.g., ω^0 , update ω^0 to approach ω^* . Generate a sequence of model parameters.



$$\omega^0 \rightarrow \omega^1 \rightarrow \omega^2 \rightarrow \dots \rightarrow \omega^k = \omega^*$$

$$L(\omega^0) > L(\omega^1) > L(\omega^2) > \dots > L(\omega^k)$$

To design this optimization, we need to determine

① moving direction.

② step size

③ stopping condition. (k)

2) Moving direction (use the sign of the gradient to determine the trend direction of the $L(\omega)$)

∇L is positive > 0 , $-\nabla L(-)$
 ∇L is < 0 , $-\nabla L(+)$

N.D:

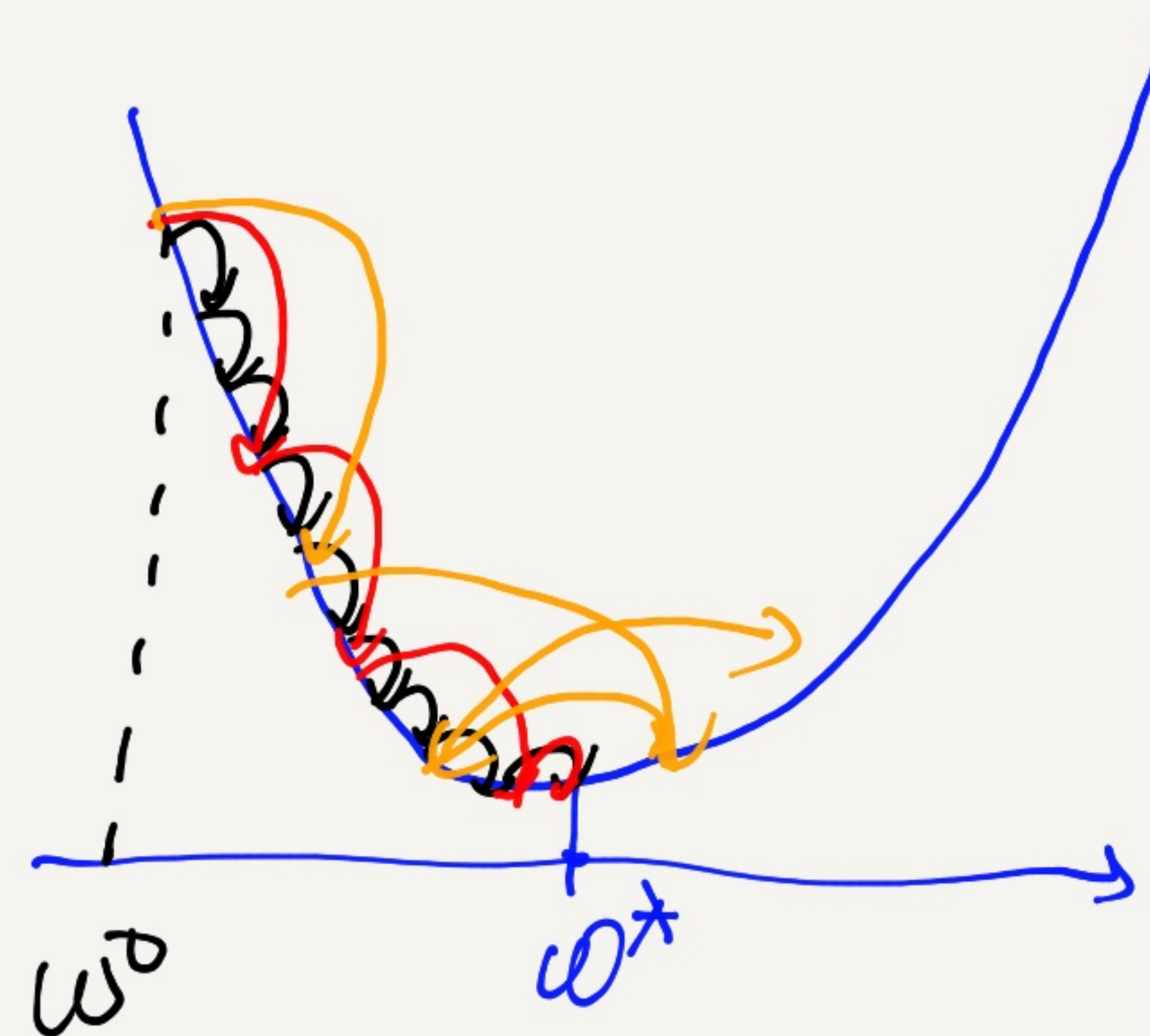
$$\omega = \begin{pmatrix} \omega_0 \\ \omega_1 \\ \vdots \\ \omega_m \end{pmatrix} \quad \nabla L = \begin{pmatrix} \frac{\partial L}{\partial \omega_0} \\ \frac{\partial L}{\partial \omega_1} \\ \vdots \\ \frac{\partial L}{\partial \omega_m} \end{pmatrix} \longrightarrow$$

moving direction
 $-\nabla L$ (opposite direction)

3) Step size. ϵ / learning rate

if step size is small, it takes longer time to reach ω^* .

Too large step size may not converge to ω^* .



① define ϵ as a small positive value.

$$10^{-3}, 10^{-4}, 10^{-5}$$

② line search.

$$\epsilon^* = \arg \min_{\epsilon} L(\overbrace{\omega + \epsilon \cdot (-\nabla L)}^{\text{new } \omega})$$

③ Dynamic ϵ :

$$\epsilon_1 > \epsilon_2 > \epsilon_3 \cdots > \epsilon_k$$

4) stopping conditions

① $\|\nabla L(w)\|_2 < \epsilon$ (small positive value, e.g. 10^{-5}) $\|\nabla L(w)\|_2 = 0$

② define the max number of iterations, $K=1000$.

5) GD Algorithm.

Input: $\{(x_i, y_i)\}_{i=1}^n$, ϵ , $\delta = 10^{-5}$, max-iter.

Output: w^*

1. Generate w^0 randomly: $w = w^0$, $i = 0$

2. while $\|\nabla L(w)\|_2 > \delta$ and $i < \text{max-iter}$

① calculate the gradient. $\hat{g} = \nabla L(w) = A^T(Aw - b)$

② update w : $w = w + \epsilon \cdot (-\hat{g})$

3. $i = i + 1$
 $w^* = w$. return w^*