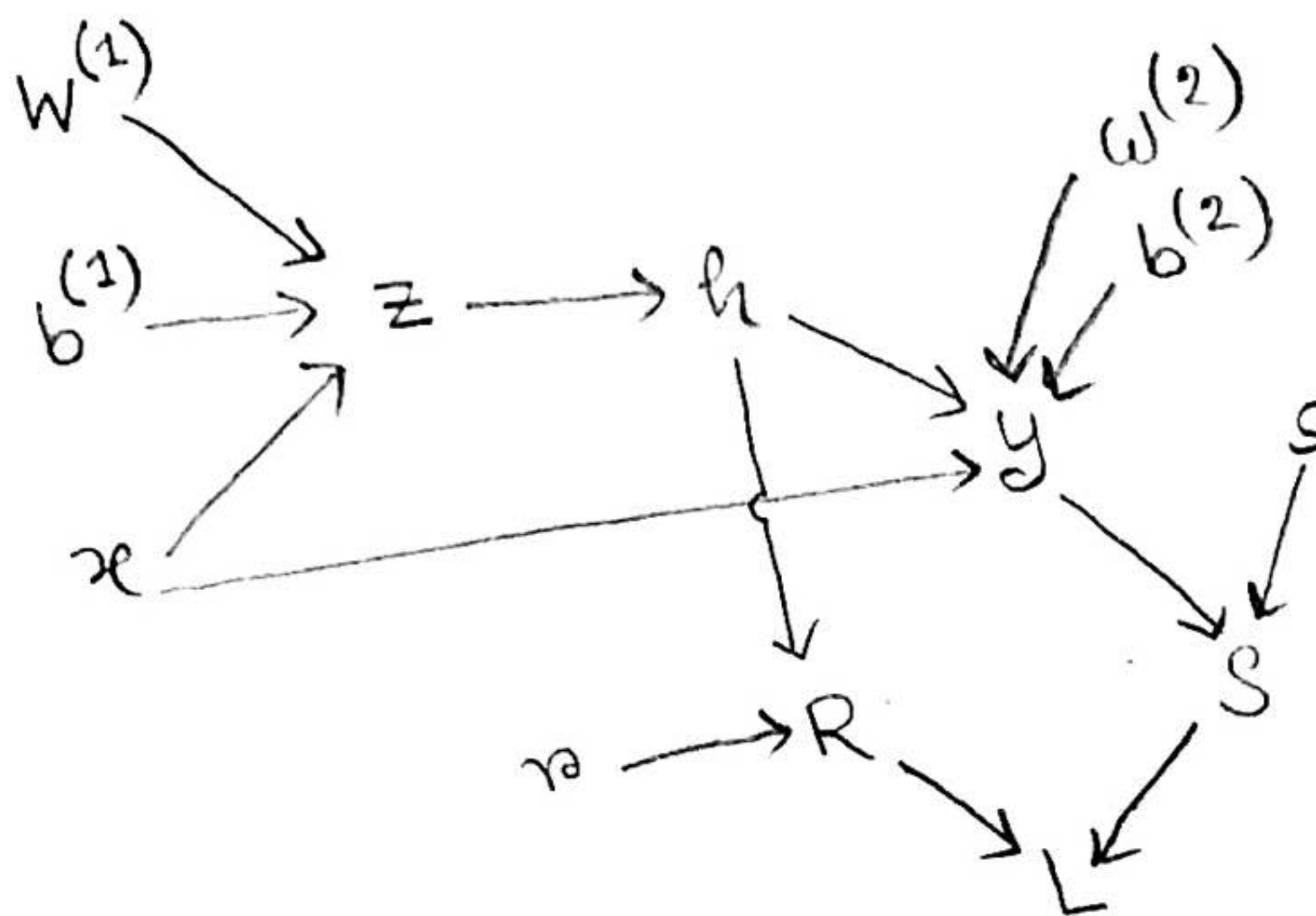Asadullah Hill Galib

CSE-891 : Homework 2

October, 11

① ·①



① ② $\dfrac{dL}{dx} = x' = ?$

$$L' = 1$$

$$S' = L'$$

$$R' = L'$$

$$y' = S'(y-s)$$

$$h' = y' \cdot \omega^{(2)} + R' \cdot v^T$$

$$z' = h' \cdot \sigma'(z)$$

$$x' = z' \omega^{(1)} + y' \cdot 1$$

(Ans)

② ①

$f(x) = vv^T x$

$J = vv^T$

$n = 3, \quad v^T = [1, 2, 3]$

$$J = vv^T = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} [1 \ 2 \ 3] = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}$$

② ②

The time and memory cost of evaluating the Jacobian is $O(n^2)$. [ There are $n \times n$ calculation & $n \times n$ grids in ʘ or in memory are needed]

② ③

$z = J^T y$

$\quad = vv^T y$     [ transpose of $J = J$]

Here, instead of evaluating $(v \cdot v^T) \cdot y$, one

should do this: $v \cdot (v^T \cdot y)$

$\underbrace{\qquad}_{first} \rightarrow$ linear in $n$ [ $1 \times 1$]

$\underbrace{\qquad}_{second} \rightarrow$ linear in $n$ [ $3 \times 1$]

So, $z = J^T y$

So, $= v \cdot (v^T y)$

$$= \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} [1 \cdot 2 \ 3] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \cdot 6 = \begin{bmatrix} 6 \\ 12 \\ 18 \end{bmatrix} \text{(Ans)}$$

③① 

Loss, $\mathcal{L} = \frac{1}{n}\|X\hat{\omega} - t\|_2^2$

$$\frac{d\mathcal{L}}{d\hat{\omega}} = \frac{2}{n}|(X\omega \frac{2}{n} X^T(X\hat{\omega} - t)$$

$$= \frac{2}{n}\left(X^T X \hat{\omega} - X^T t\right)$$

③② Under parameterized :

ⓐ

$$\frac{d\mathcal{L}}{d\hat{\omega}} = 0$$

$$\Rightarrow \frac{2}{n}\left(X^T X \hat{\omega} - X^T t\right) = 0$$

$$\Rightarrow X^T X \hat{\omega} = X^T t$$

$$\therefore \hat{\omega} = (X^T X)^{-1} X^T t \qquad [\text{as } X^T X \text{ is invertible for } n > d]$$

③②ⓑ $t_i = \omega^{*T} x_i$ .

~~$t = Xw^*$~~ $t = X\omega^*$

So, $\hat{\omega} = (X^T X)^{-1} X^T \cdot X \omega^* = \omega^*$

Therefore, $\forall x \in \mathbb{R}^d$, $(\omega^{*T} x - \hat{\omega}^T x)^2 = 0$ $[\text{when } d < n]$

and $\hat{\omega}$ achieves perfect generalization.

③ ③ Overparameterized Model: 2D

ⓐ

$n = 1, d = 2.$ $x_1 = [2 \ 1]$ $l_1 = 2$

$\hat{\omega}^T x_1 = y_1$ . Let, $\hat{\omega} = \begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix}$

So, $[\omega_1 \ \omega_2] \begin{bmatrix} 2 \\ 1 \end{bmatrix} = l_1 = 2$

$\Rightarrow 2\omega_1 + \omega_2 = 2$

$\Rightarrow \omega_2 = -2\omega_1 + 2 \longrightarrow$ equation of line

So, there exists infinitely many $\hat{\omega}$, satisfying $\hat{\omega}^T x_1 = y_1$

as ~~there's no~~ the solution can be anywhere on the line.

③ ③ ⑥

$\hat{\omega}(0) = 0$

$$\frac{d\mathcal{L}}{d\hat{\omega}} = \frac{2}{n} X^T(X\hat{\omega} - t)$$

when, $\hat{\omega}(0) = 0$ and ~~$x_0 = x_1$~~, $X = X_1$, $t = t_1$, $n = 1$, $d = 2$:

$$\frac{d\mathcal{L}}{d\hat{\omega}} = \frac{2}{n}(-X^T t) = \frac{1}{n}(-2X_1^T t_1)$$

As, $t_1 = 2$ a constant, the direction of the gradient is along $X_1$. And, it doesn't change along the trajectory as there is no $\hat{\omega}$ term in the derivative.

The corresponding unit norm vector of $X_1$: $\frac{1}{\sqrt{5}}\begin{bmatrix} 2 \\ 1 \end{bmatrix}$

Using the $d$ above, we get

$$\hat{\omega} = 2 \cdot \frac{1}{(\sqrt{5})^2}\begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 0.4 \end{bmatrix}$$ (using squared-norm)

Here,

the gradient descent finds the
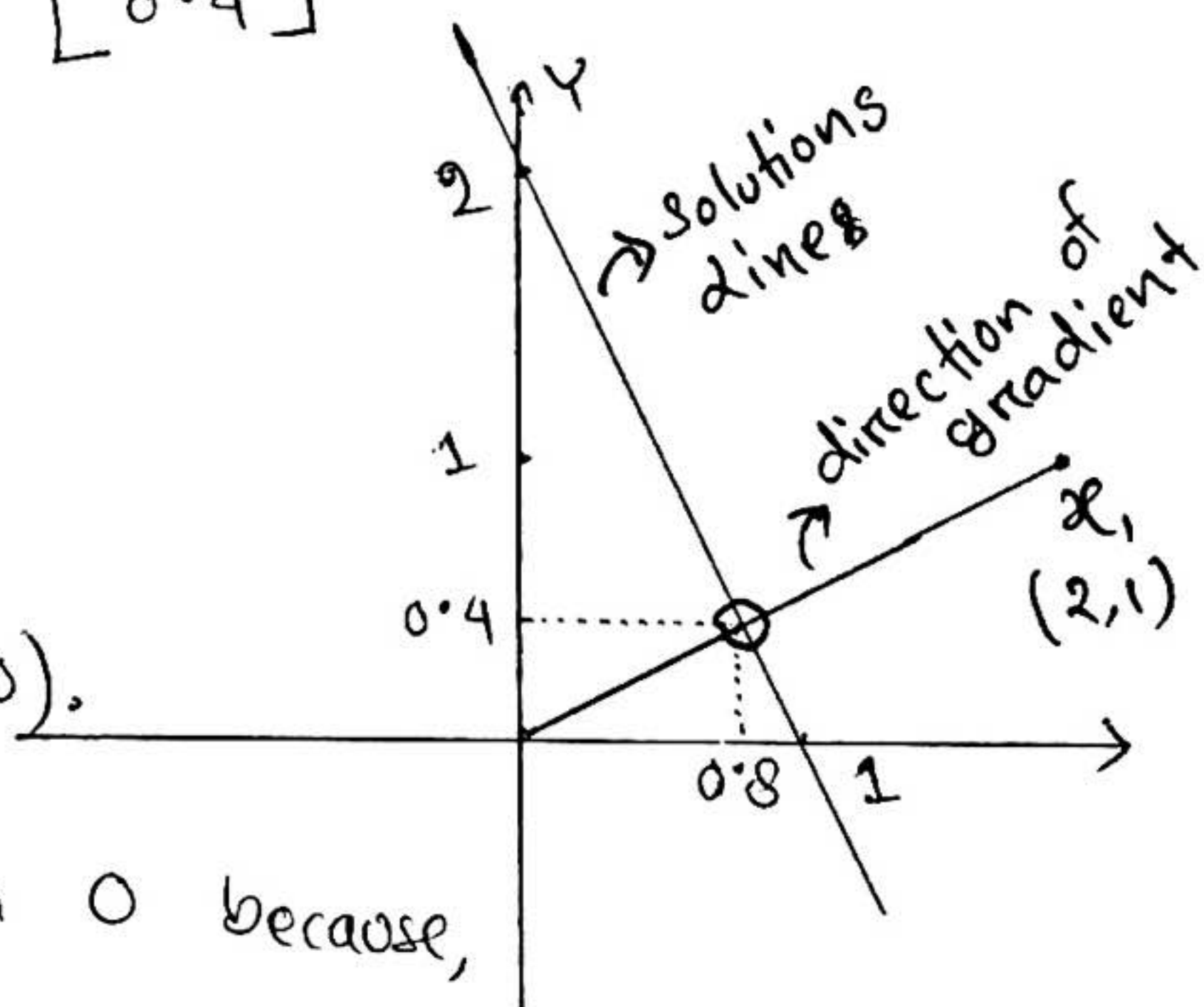
closest solution from $W_0$ (0,0).



It is the closest distance from O because,

the direction of gradient and solutions line are orthogonal

to each other (i.e. slope of the solutions line, $m_1 = -2$

"        "      " gradient des. ", $m_2 = \frac{1}{2}$

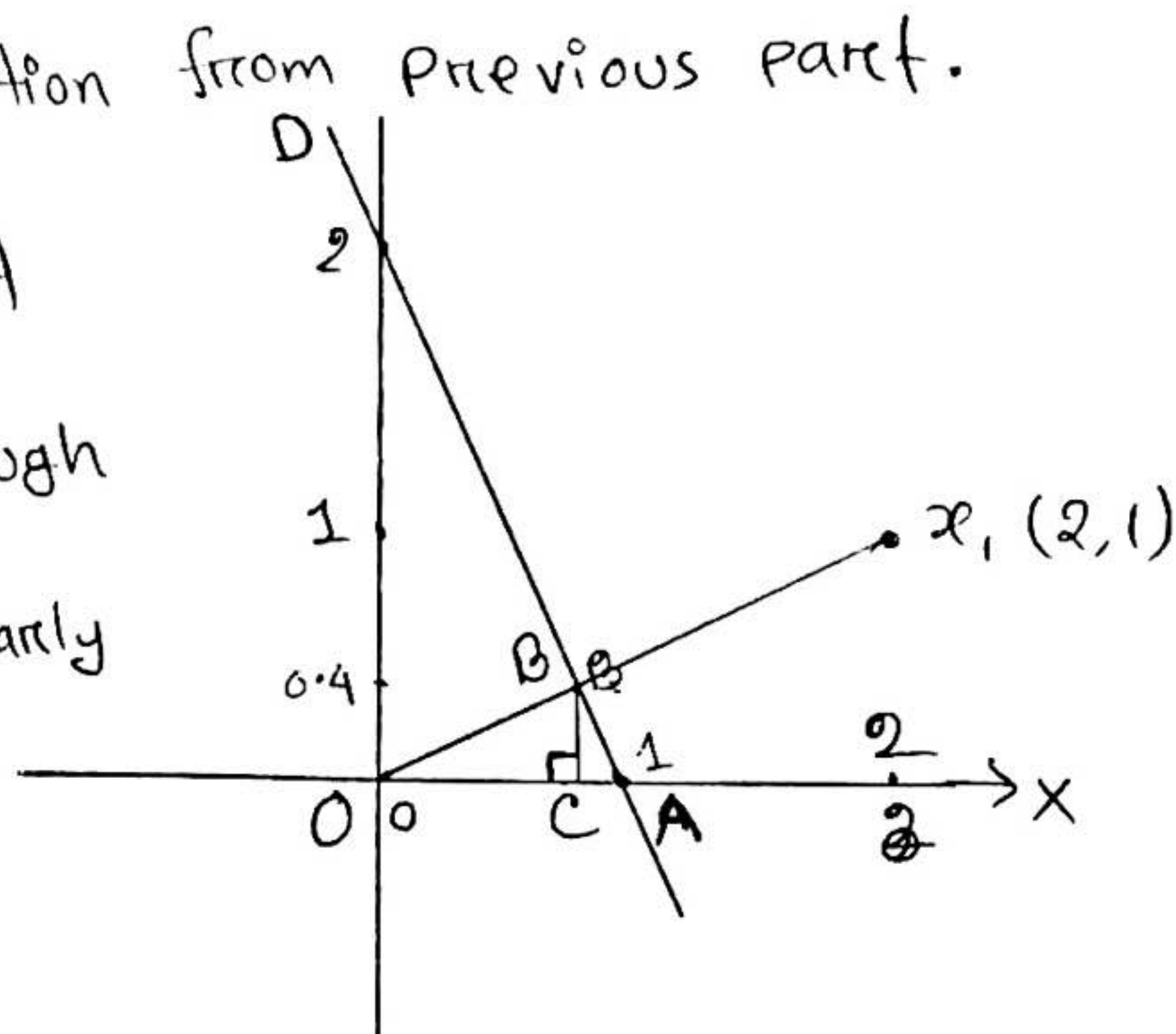$m_1 m_2 = -1 \Rightarrow$ perpendicular)

③ ③ ©

Using Pythagorean Theorem:

B is the gradient descent solution from previous part.

{
To show, B has the smallest

Euclidean norm, it is enough

to prove that OB perpendicularly

intersects $AD$ ($OB \perp AD$)
}



Here, $OA = 1$

$OC = 0.8$

$BC = 0.4$ [$BC \perp OA$]

So, $OB = \sqrt{OC^2 + BC^2} = \sqrt{\frac{2}{5}}$

$AB = \sqrt{CA^2 + BC^2} = \sqrt{(OA - OC)^2 + BC^2}$

$= \sqrt{(0.2)^2 + (0.4)^2} = \sqrt{\frac{1}{5}}$

$\triangle OBA$

In $\triangle OAB$,

$OB^2 + AB^2 = \frac{4}{5} + \frac{1}{5} = 1 = OA^2$

So, $\triangle OAB$ is a right-angled $\triangle$. $\Rightarrow OB \perp AB$
$\triangle OBA$

$\Rightarrow OB \perp AD$

[Proved]

—o—

Alternative Proof:

⊛ It can also be shown also using slopes:

$m_{AD} = -2$ and $m_{OB} = \frac{1}{2}$

So, $m_{AD} \cdot m_{OB} = -1 \Rightarrow OB \perp AD \Rightarrow OB$ has the smallest norm.

③ ④ ⓐ With, $\hat{w}(0) = 0$, we get gradient vector in the span of $X$.

As the gradient vector is always spanned by the rows of $X$, we can get $\hat{w}$ as a linear combination of $X$ and some other matrix. Let, $P$ is that matrix.

So, $\hat{w} = X^T P$

Then, $\frac{2}{n} X^T (X\hat{w} - t) = 0$

$\Rightarrow X^T(XX^TP - t) = 0$

$\Rightarrow XX^T(XX^TP - t) = X.0 = 0$

$\Rightarrow XX^TP - t = (XX^T)^{-1}.0 = 0$  [ for $d > n$, $XX^T$ is invertible]

$\Rightarrow XX^TP = t$

$\Rightarrow P = (XX^T)^{-1}t$

Thereby, $\hat{w} = X^T P$

$= X^T(XX^T)^{-1}t$

The solution is unique as it's just a linear transformation of $t$.

③ ④ ⑥

Zero-loss solution with $\hat{\omega}_1$.

So, $\hat{\omega}_1^T x - t = 0 \Rightarrow \hat{\omega}_1^T x = t \dots$ ①

$$(\hat{\omega} - \hat{\omega}_1)^T \hat{\omega} = (x^T(xx^T)^{-1}t - \hat{\omega}_1)^T \hat{\omega}$$

$$= (t^T(xx^T)^{-1}x - \hat{\omega}_1)\, \hat{\omega}\, (x^T(xx^T)^{-1}t)$$

$$= (t^T(xx^T)^{-1}xx^T(xx^T)^{-1}t - \hat{\omega}_1 x^T(xx^T)^{-1}t)$$

$$= t^T(xx^T)^{-1}t - t^T(xx^T)^{-1}t \quad [\text{using } ①]$$

$$= 0$$

So, $(\hat{\omega} - \hat{\omega}_1)$ and $\hat{\omega}$ are perpendicular to each other

$\Rightarrow \hat{\omega}_1$ and $\hat{\omega}$ are perpendicular to each other.

So, like $\hat{\omega}_1$, all other solutions are perpendicular to $\hat{\omega}$. And, this gradient descent solution, $\hat{\omega}$

has the smallest Euclidean norm similarwise we

proved before (using Pythagoream theorem).