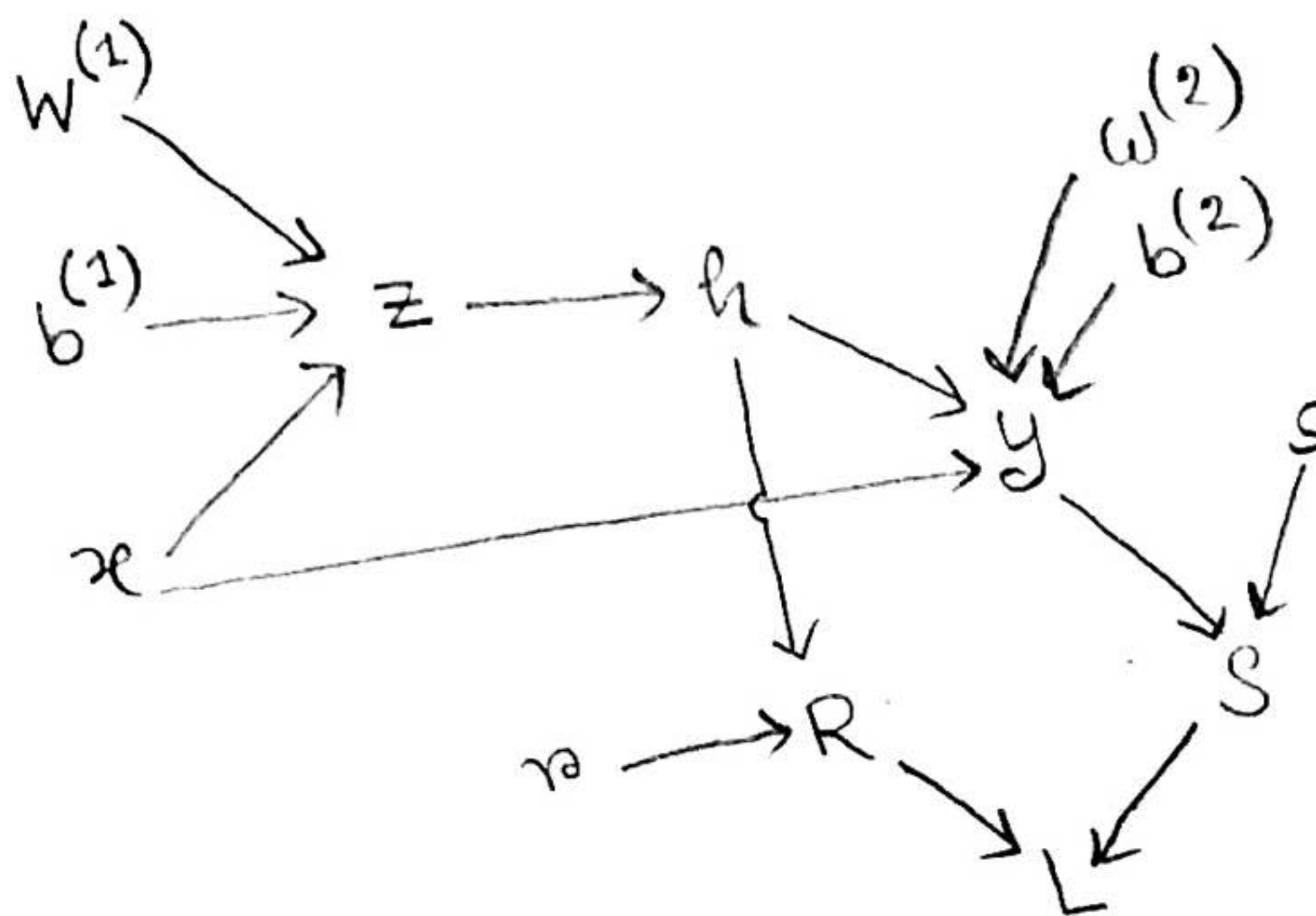Asadullah Hill Galib

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 memory are needed ]

② ③

$z = J^T y$

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

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

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

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

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

So, $z = J^T y$

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

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

③① 

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

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

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

③② Underparameterized:

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

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

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

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

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

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

③③ Overparameterized Model: 2D

ⓐ

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

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

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

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

$$\Rightarrow \omega_2 = -2\omega_1 + 2 \quad \rightarrow \text{equation of line}$$

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 ⓑ

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

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

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

$$\frac{d\mathcal{L}}{d\hat{w}} = \frac{2}{n}(-x^T t) = \frac{1}{\cancel{n}}(-2x_1^T t_1)$$

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 ~~$\hat{w}$~~ $\hat{w}$ 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

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

Here,

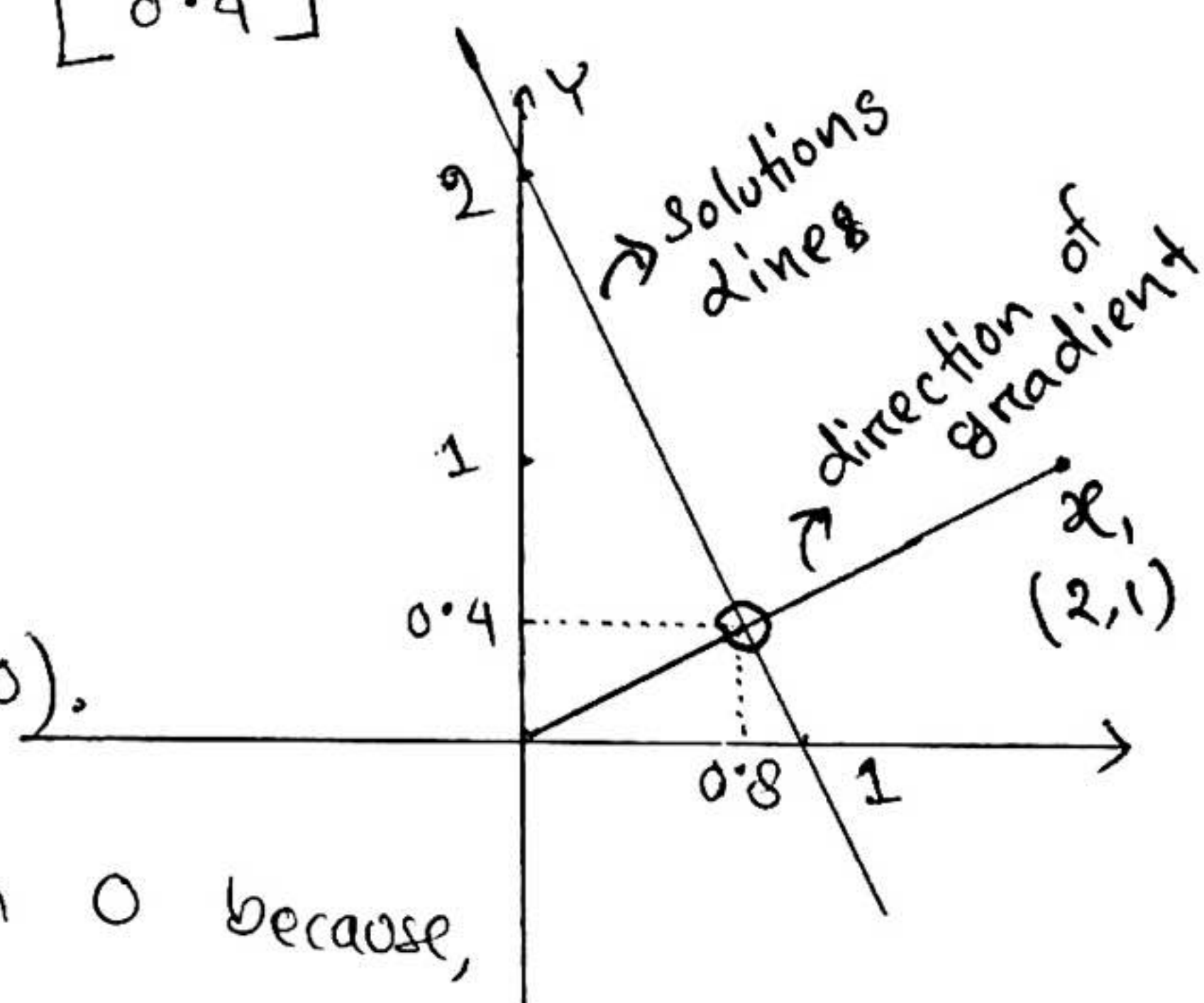the gradient descent finds the

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. ", $m_2 = \frac{1}{2}$
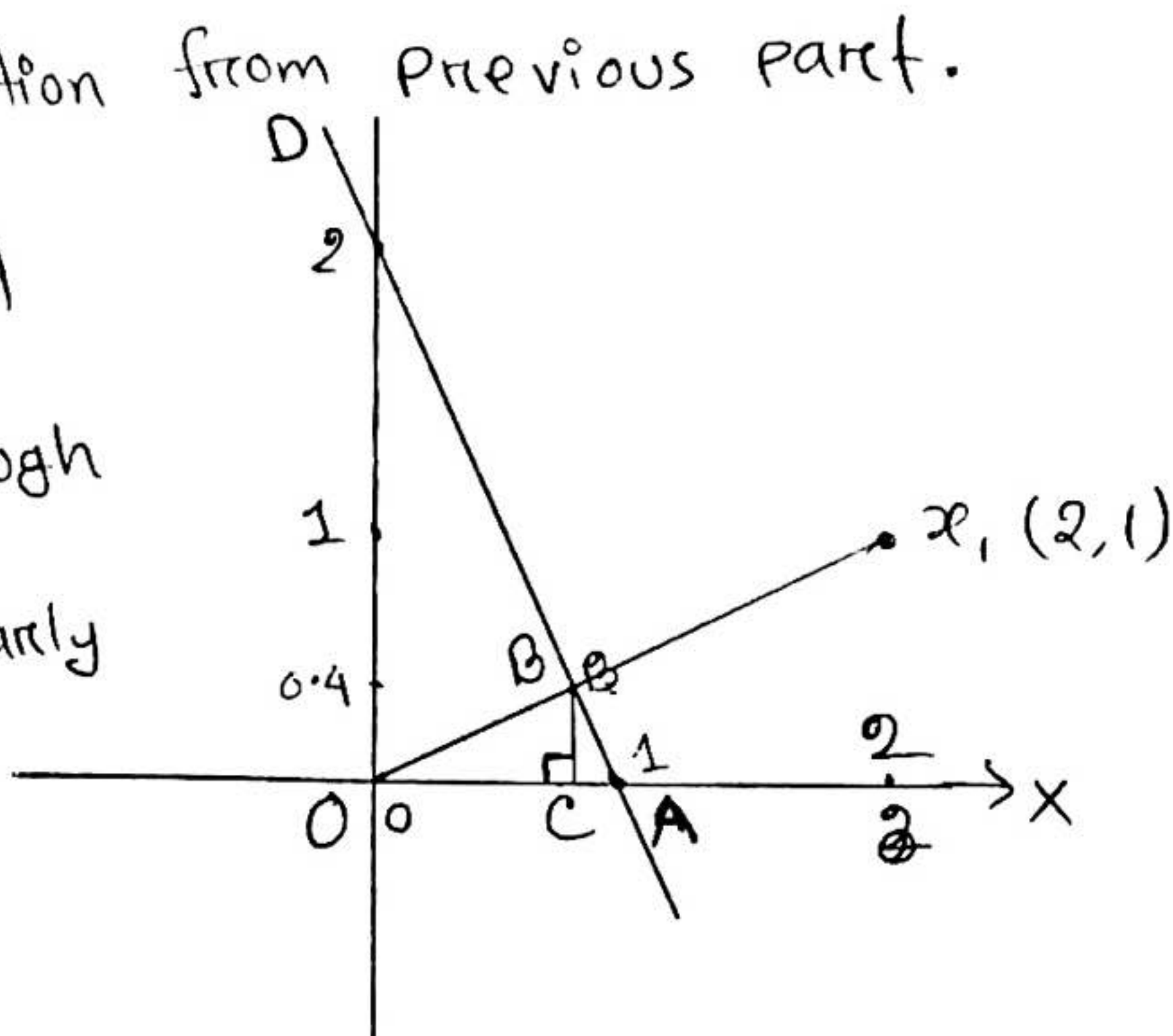
$m_1 m_2 = -1 \implies$ perpendicular)

(3) (3) (c)

### Using Pythagorean Theorem:

B is the gradient descent solution from previous part.

$\begin{cases} \text{To show, } B \text{ has the smallest} \\ \text{Euclidean norm, it is enough} \\ \text{to prove that } OB \text{ perpendicularly} \\ \text{intersects } AD \; (OB \perp AD) \end{cases}$



Hence, $OA = 1$

$OC = 0.8$

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

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

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

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

$\triangle OBA$

In $\cancel{\triangle OAB}$,

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

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

$\Rightarrow OB \perp AD$

[Proved]

— o —

**Alternative Proof:** ⊛ It can also be shown also using slopes:

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

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

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

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

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

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

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

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

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

$\Rightarrow XX^T P = t$

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

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

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

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

③ ④ ⑥

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

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

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

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

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

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

$$= 0$$

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

$\implies \hat{\omega}_1$ 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).
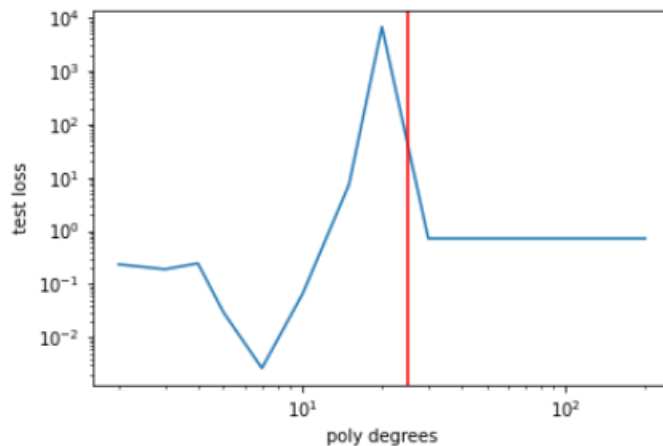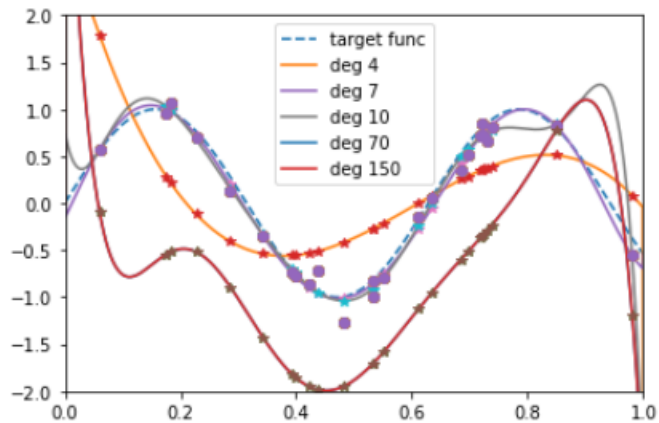
Question 3 -5:

```
1 # to be implemented; fill in the derived solution for the und
2
3 def fit_poly(X, d, t):
4     X_expand = poly_expand(X, d=d, poly_type=poly_type)
5     if d > n:
6         W = X_expand.T@np.linalg.inv(X_expand@X_expand.T)@t
7     else:
8         W = np.linalg.inv(X_expand.T@X_expand)@X_expand.T@t
9     return W
```

```
2 0.23473638175555447
3 0.19020505096352716
4 0.24537346900180165
5 0.02963124207539845
7 0.0026501350788807432
10 0.06635344938407685
15 7.3835492062768155
20 6706.8570800068255
30 0.7170381150111599
50 0.7170381150111599
70 0.7170381150111599
100 0.7170381150111599
150 0.7170381150111599
200 0.7170381150111599
```

No, overparameterization does not always lead to overfitting. Here, overparameterization give stable and better performance than the medium range of parameters (9-35). Implicit regularization induced by gradient descent is reason for this trend.

④ ① SGD:

$$X\hat{\omega} = t \quad\Rightarrow\quad X\hat{\omega} - t = 0 \quad\cdots\text{①}$$

In SGD, all $x_i$ is contained in the span o $X$. And

the SGD update do steps don't ever leave the span of

$X$. Because, $\dfrac{d}{d\hat{\omega}_p}\left(x_i\hat{\omega}_p - t_i\right)^2 = 0$ will ~~git~~ update

$\hat{\omega}_p$ as ~~so~~ some combination of $x_i$, and $t_i$.

Thereby, we can assume the SGD solution is spanned

by $X$ : $\hat{\omega} = X^T S$, where $S$ is a arbitrary matrix.

From ①, $X\hat{\omega} - t = 0$

$$\Rightarrow X X^T S - t = 0$$

$$\Rightarrow S = (X X^T)^{-1} t$$

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

$$= \omega^*$$

[ Showed ]

④ ② Mini-batch SGD:

Yes, mini-batch SGD also obtains minimum norm solution on convergence.

Because the batch $B$ is taken from the rows of $X$. So, the solution $\hat{\omega}$ is spanned by the rows of $X$.

$$\hat{\omega} = BS = X\tilde{S}$$

$\downarrow$
Batch

So, $X\hat{\omega} - t = XX^T\tilde{S} - t = 0$

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

$$\therefore \hat{\omega} = X^T(XX^T)^{-1} t = \omega^*$$

④ ③ Adaptive Methods: Adagrad

$x_1 = [2, 1]$    $\omega_0 = [0, 0]$    $\ell = [2]$

Using minimum norm solution with GD,

we got    $\omega^* = \begin{bmatrix} 0.8 \\ 0.4 \end{bmatrix}$ . and

$$\nabla_{\omega^*} \mathcal{L}(\omega) = -2x_1 \ell_1$$

Using Adagrad,

$\omega_0 = [0, 0]$

$$\omega_1 = \omega_0 - \frac{\eta}{\sqrt{G_{1}} + \epsilon} \nabla_{\hat{\omega}_0} \mathcal{L}(\omega)$$

$$G_1 = 0 + \left(\nabla_{\hat{\omega}_0} \mathcal{L}(\omega)\right)^2$$

Let, assume,    $\nabla_{\hat{\omega}_0} \mathcal{L}(\omega) = -2x_1 \ell_1$ [similar to the GD]

then,    $\omega_1 = \omega_0 - \dfrac{\eta}{(-2x_1 \ell_1) + \epsilon} \cdot (-2x_1 \ell_1)$

as $\epsilon$ is small, $\omega_1$ looses $x_1$ term almost, ~~that means~~ the direct because numerator and

denominator both contains $(-2x_1 \ell_1)$, So,

$\omega_1$ has a little impact from $x_1$, which indicates the direction of the gradient is no longer along $x_1$ as much as the $\omega^*$ case.
(GD with minimum norm sol.)

Thereby, AdaGrad doesn't always obtain the minimum norm solution.

This
~~Same~~ results. holds
true for other adaptive
~~models~~ methods (RMSProp, Adam) in general.

gradient of the AdaGrad

gradient of the minimum norm sol.

Because the scaling part in the weight update may divert solution gradient from the span of X and it may get outside of the span of X sometimes.