

1. Consider stochastic gradient descent method to learn the house price model

$$h(x_1, x_2) = \sigma(b + w_1x_1 + w_2x_2),$$

where  $\sigma$  is the sigmoid function.

Given one single data point  $(x_1, x_2, y) = (1, 2, 3)$ , and assuming that the current parameter is  $\theta^0 = (b, w_1, w_2) = (4, 5, 6)$ , evaluate  $\theta^1$ .

Just write the expression and substitute the numbers; no need to simplify or evaluate.

$$\begin{aligned} \text{Loss}(\theta) &= \frac{1}{N} \sum_{i=1}^N |y^i - h(x^i; \theta)| \\ &= |y^1 - h(x^1; \theta)| \quad (N=1) \\ &= 3 - h(1, 2; \theta) \\ &= 3 - \sigma(b + w_1 + 2w_2) \end{aligned}$$

$$\theta^1 = \theta^0 - \alpha \nabla_{\theta} \text{Loss}(\theta^0) \quad , \quad \alpha > 0 .$$

$$= \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial \text{Loss}}{\partial b}(4, 5, 6) \\ \frac{\partial \text{Loss}}{\partial w_1}(4, 5, 6) \\ \frac{\partial \text{Loss}}{\partial w_2}(4, 5, 6) \end{bmatrix}$$

2. (a) Find the expression of  $\frac{d^k}{dx^k} \sigma$  in terms of  $\sigma(x)$  for  $k = 1, \dots, 3$  where  $\sigma$  is the sigmoid function.

(b) Find the relation between sigmoid function and hyperbolic function.

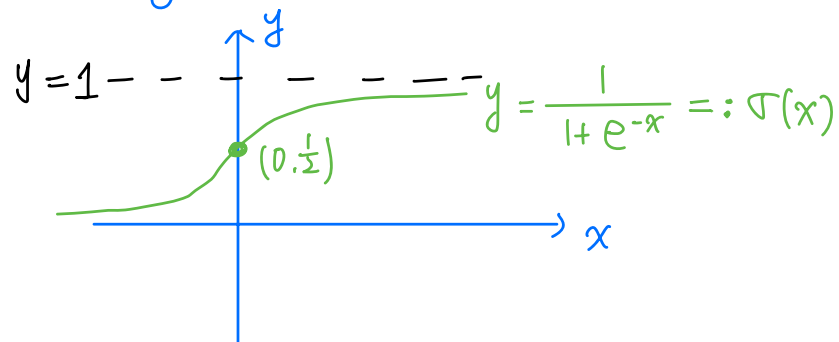
(a)  $k=1$  : 
$$\frac{d}{dx} \sigma(x) = \frac{d}{dx} \left( \frac{1}{1+e^{-x}} \right) = \frac{1}{(1+e^{-x})^2} \cdot e^{-x} = \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}+1}{1+e^{-x}} = \sigma(x)[1-\sigma(x)]$$

$k=2$  : 
$$\begin{aligned} \frac{d^2}{dx^2} \sigma(x) &= \frac{d}{dx} \sigma = \frac{d}{dx} [\sigma(1-\sigma)] = \sigma'(1-\sigma) + \sigma(-\sigma') = \sigma(1-\sigma)(1-\sigma) - \sigma \cdot \sigma(1-\sigma) \\ &= \sigma(1-\sigma)(1-\sigma-\sigma) \\ &= \sigma(1-\sigma)(1-2\sigma) \end{aligned}$$

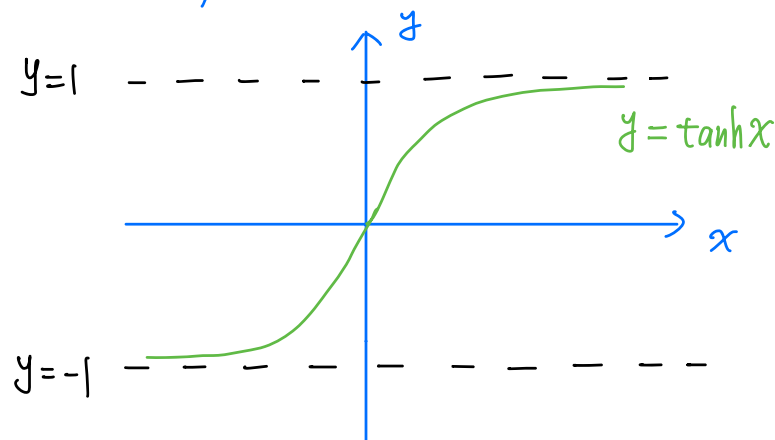
$k=3$  : 
$$\begin{aligned} \frac{d^3}{dx^3} \sigma(x) &= \frac{d}{dx} \left( \frac{d^2}{dx^2} \sigma \right) = \frac{d}{dx} [\sigma(1-\sigma)(1-2\sigma)] = \sigma'(1-\sigma)(1-2\sigma) + \sigma(-\sigma')(1-2\sigma) + \sigma(1-\sigma)(-2\sigma') \\ &= \sigma(1-\sigma)(1-\sigma)(1-2\sigma) - \sigma \cdot \sigma(1-\sigma)(1-2\sigma) - 2\sigma(1-\sigma)\sigma(1-\sigma) \\ &= \sigma(1-\sigma) [(1-\sigma)(1-2\sigma) - \sigma(1-2\sigma) - 2\sigma(1-\sigma)] \\ &= \sigma(1-\sigma) (1-3\sigma+2\sigma^2-\sigma+2\sigma^2-2\sigma+2\sigma^2) \\ &= \sigma(1-\sigma) (6\sigma^2-6\sigma+1) \end{aligned}$$

(b)

Sigmoid function



hyperbolic function



$$\text{Let } \tanh x = C_1 \sigma(C_2 x + C_3) + C_4 \quad C_i \in \mathbb{R}, i=1, \dots, 4$$

根據觀察圖形, guess  $C_3 = 0$

$$\text{Let } x=0, \text{ then } 0 = C_1 \sigma(C_2 \cdot 0 + 0) + C_4 = C_1 \sigma(0) + C_4 = \frac{1}{2} C_1 + C_4 \Rightarrow C_4 = -\frac{1}{2} C_1$$

$$\therefore \text{Let } \tanh x = C_1 \sigma(C_2 x) - \frac{1}{2} C_1 = \frac{C_1}{1+e^{-C_2 x}} - \frac{1}{2} C_1$$

$$\frac{d}{dx} \Rightarrow 1 - \tanh^2 x = \frac{C_1 C_2 e^{-C_2 x}}{(1+e^{-C_2 x})^2}$$

$$\Rightarrow 1 - \left[ \frac{C_1^2}{(1+e^{-C_2 x})^2} - \frac{C_1^2}{1+e^{-C_2 x}} + \frac{1}{4} C_1^2 \right] = \frac{C_1 C_2 e^{-C_2 x}}{(1+e^{-C_2 x})^2}$$

$$\Rightarrow 1 - \frac{1}{4} C_1^2 + \frac{-C_1^2 + C_1^2 + C_1^2 e^{-C_2 x}}{(1+e^{-C_2 x})^2} = \frac{C_1 C_2 e^{-C_2 x}}{(1+e^{-C_2 x})^2} \Rightarrow C_1 = \pm 2, C_2 = C_1 = \pm 2$$

(C<sub>1</sub>=C<sub>2</sub>=2)

$$\therefore \tanh x = \frac{2}{1+e^{-2x}} - 1$$

$$= \frac{-2}{1+e^{2x}} + 1$$

(C<sub>1</sub>=C<sub>2</sub>=-2)

3. There are unanswered questions during the lecture, and there are likely more questions we haven't covered. Take a moment to think about them and write them down here.

Gradient descent algorithm :  $\theta^{n+1} = \theta^n - \alpha \nabla_{\theta} \text{Loss}$  .  $\alpha > 0$  : learning rate

為什麼會收斂？

會收斂到同個值嗎？

收斂的值是最大/最小值嗎？

$\alpha$  要怎麼選？