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

$$h(x_1, x_2) = \sigma(b + w_1 x_1 + w_2 x_2),$$

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.

* data $\left\{ (x_1^i, x_2^i, y^i) \right\}_{i=1}^{N}$

* Loss function :

$$\text{Loss}(b, w_1, w_2) = \frac{1}{M} \sum_{i=1}^{M} \left( y^i - h(x_1^i, x_2^i) \right)^2$$

* Gradient descent algorithm

$$\theta^{h+1} = \theta^h - \alpha \nabla_\theta \text{Loss}$$

$1^0$  $h^0(x_1, x_2) = 6(4 + 3 - 1 + 6 \cdot 2)$

$\qquad\qquad = 6(21)$

$2^0$  $\text{Loss} = \frac{1}{2}(h^0 - 3)^2$

$3^{\circ}$

let $b + w_1 x_1 + w_2 x_2 = l$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial l} \frac{\partial l}{\partial b}$$

$$= (h^{\circ} - y) \cdot \overset{\circ}{h}(1 - h^{\circ}) \cdot 1$$

$$= (h^{\circ} - 3) h^{\circ} (1 - h^{\circ})$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial l} \frac{\partial l}{\partial w_1}$$

$$= (h^{\circ} - y) \cdot h^{\circ}(1 - h^{\circ}) \cdot x_1$$

$$= (h^{\circ} - 3) h^{\circ} (1 - h^{\circ})$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial l} \frac{\partial l}{\partial w_2}$$

$$= (h^{\circ} - y) \cdot \overset{\circ}{h}(1 - \overset{\circ}{h}) \cdot x_2$$

$$= (h^{\circ} - 3) h^{\circ} (1 - h^{\circ}) \cdot 2$$

$u^0$

$b^1 = 4 - \alpha (h^0 - 3) h^0 (1 - h^0)$

$w_1^1 = 5 - \alpha (h^0 - 3) h^0 (1 - h^0)$

$w_2^1 = 6 - 2\alpha (h^0 - 3) h^0 (1 - h^0)$

(learning rate: $\alpha$, $h^0 \approx 6$ (≥1)

$\theta^1 = (b^1, w_1^1, w_2^1)$  #1

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

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

(a)

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{d}{dx}\sigma(x)\ \sigma'(x) = \frac{0 \cdot (1 + e^{-x}) - 1 \cdot (0 - e^{-x})}{(1 + e^{-x})^2}$$

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

$$= \frac{1}{1 + e^{-x}} \cdot \left(1 - \frac{1}{1 + e^{-x}}\right)$$

$$= \sigma(x) \cdot (1 - \sigma(x))$$

$$\frac{d^2}{dx^2}\sigma(x)\ \sigma''(x) = \sigma'(x) \cdot (1 - \sigma(x)) + \sigma(x) \cdot (-\sigma'(x))$$

$$= \sigma(x) \cdot (1 - \sigma(x))^2 - \sigma^2(x)(1 - \sigma(x))$$

$$= \sigma(x)(1 - \sigma(x))(1 - \sigma(x)) - \sigma(x))$$

$$= \sigma(x)(1 - \sigma(x))(1 - 2\sigma(x))$$

$$\frac{d^3}{dx^3}\sigma = \sigma'''(x)$$

$$= \quad \sigma'(x)\,(1-\sigma(x))\,(1-2\sigma(x))$$
$$+\ \sigma(x)\cdot(-\sigma'(x))\,(1-2\sigma(x))$$
$$+\ \sigma(x)\,(1-\sigma(x))\,(-2\sigma'(x))$$

$$= \quad \sigma(x)\,(1-\sigma(x))^2\,(1-2\sigma(x))$$
$$-\ \sigma^2(x)\,(1-\sigma(x))\,(1-2\sigma(x))$$
$$-\ 2\sigma^2(x)\,(1-\sigma(x))^2$$

$$= \quad \sigma(x)\,(1-\sigma(x))\cdot$$

$$\Big((1-\sigma(x))(1-2\sigma(x)) - \sigma(x)\,(1-2\sigma(x))$$
$$-\ 2\sigma(x)\,(1-\sigma(x))\Big)$$

$$= \sigma(x)\,(1-\sigma(x))\left(1 - 6\,\sigma(x) - 6\sigma^2(x)\right)$$

$\mathcal{E}1$

(b)

$$i^\circ \quad \sigma(x) = \frac{1}{1+e^{-x}}$$

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

$$= \frac{e^{\frac{x}{2}}}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}} \cdot \left(2 \cdot \frac{1}{2}\right)$$

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

$$= \frac{1}{2} \cdot \frac{e^{\frac{x}{2}} + e^{\frac{x}{2}} + \left(e^{-\frac{x}{2}} - e^{-\frac{x}{2}}\right)}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}}$$

$$= \frac{1}{2} \cdot \left(1 + \frac{e^{\frac{x}{2}} - e^{-\frac{x}{2}}}{e^{\frac{x}{2}} + e^{-\frac{x}{2}}}\right)$$

$$= \frac{1}{2}\left(1 + \tanh\left(\frac{x}{2}\right)\right)$$

$$= \frac{1 + \tanh\left(\frac{x}{2}\right)}{2} \qquad \#1$$

$$\Rightarrow \tanh\left(\tfrac{x}{2}\right) = 2\sigma(x) - 1 \qquad ✠$$

$$2^\circ \quad \sigma'(x) = \sigma(x)\left(1 - \sigma(x)\right) \qquad \text{by } 2(a)$$

$$= \frac{1 + \tanh\left(\tfrac{x}{2}\right)}{2} \cdot \frac{1 - \tanh\left(\tfrac{x}{2}\right)}{2}$$

$$= \frac{1}{4} \cdot \left(1 - \tanh^2\left(\tfrac{x}{2}\right)\right)$$

$$= \frac{1}{4} \operatorname{sech}^2\left(\tfrac{x}{2}\right) \qquad ✠$$

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.

In "derivation of MSE loss" and "gradient descent", the sample index often starts from $i = 1 \cdots M$, while in implementations often $i = 0 \cdots M-1$, reminds me of statistic, where degree of freedom often reduced ',

Are these differences mainly notation conventions, or do they reflect mathematical or practical implications?

sorry for didn't notice assignment earlier : c