

Deep learning HK4. Participants 1° Hanzhi Zhuang 2° Rui Song 3° Tengyinhao Yang
 5578347 5576534 5587382

1. Pen and Paper tasks.

$$1) \hat{y}_0 = w_0^T X = [1 \ 0] \cdot \begin{bmatrix} 2 \\ -1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = [2 \ -1]$$

First update step:

$$1^\circ w_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \theta_0$$

$$2^\circ \text{ Sample from: } \left\{ \begin{bmatrix} 2 \\ -1 \end{bmatrix}, 2 \right\}, \left\{ \begin{bmatrix} -1 \\ 3 \end{bmatrix}, -1 \right\}$$

$$3^\circ \hat{g} = \frac{1}{m} \nabla_0 \sum (\hat{y} - y)^2 = \frac{1}{m} \nabla_0 \sum (w^T x - y)^2 = \frac{1}{m} \sum 2(w^T x - y) \cdot x = \frac{\sum (w^T x - y) x}{m}$$

$$= \frac{(2-3) \begin{bmatrix} 2 \\ -1 \end{bmatrix} + (-1-1) \begin{bmatrix} -1 \\ 3 \end{bmatrix}}{2}$$

$$= \begin{bmatrix} -2+2 \\ 1-6 \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ -5 \end{bmatrix}$$

$$4^\circ V_1 = \beta \cdot V_0 - \alpha \hat{g} = 0.8 \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.2 \begin{bmatrix} 0 \\ -5 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$5^\circ W_1 = W_0 + V_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \underline{\underline{\begin{bmatrix} 1 \\ 1 \end{bmatrix}}}$$

Second update step:

$$\hat{y}_1 = w_1^T X = [1 \ 1] \cdot \begin{bmatrix} 2 \\ -1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} = [1 \ 2]$$

$$1^\circ w_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \theta_1$$

$$2^\circ \text{ Sample from } \left\{ \begin{bmatrix} 2 \\ -1 \end{bmatrix}, 1 \right\}, \left\{ \begin{bmatrix} -1 \\ 3 \end{bmatrix}, 2 \right\}$$

$$3^\circ \hat{g}_1 = \sum (w^T x - y) x = (1-3) \begin{bmatrix} 2 \\ -1 \end{bmatrix} + (2-1) \begin{bmatrix} -1 \\ 3 \end{bmatrix} = \begin{bmatrix} -5 \\ 5 \end{bmatrix}$$

$$4^\circ V_2 = \beta V_1 - \alpha \hat{g}_1 = 0.8 \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 0.2 \begin{bmatrix} -5 \\ 5 \end{bmatrix} = \begin{bmatrix} 1 \\ -0.2 \end{bmatrix}$$

$$5^\circ W_2 = W_1 + V_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ -0.2 \end{bmatrix} = \underline{\underline{\begin{bmatrix} 2 \\ 0.8 \end{bmatrix}}}$$

Third update step:

$$1^{\circ} \hat{y}_2 = W_2^T X = [2 \ 0.8] \cdot \begin{bmatrix} -2 \\ -3 \end{bmatrix} = [3.2 \ 0.4]$$

$$1^{\circ} W_2 = \begin{bmatrix} 2 \\ 0.8 \end{bmatrix} = \theta_2.$$

2^o Sample from $\left\{ \begin{pmatrix} -2 \\ -1 \end{pmatrix}, 3.2 \right\}, \left\{ \begin{pmatrix} -1 \\ 3 \end{pmatrix}, 0.4 \right\} \right\}$.

$$3^{\circ} \hat{g}_2 = \sum (w^T x - y) X = (3.2 - 3) \begin{bmatrix} -2 \\ -1 \end{bmatrix} + (0.4 - 1) \begin{bmatrix} -1 \\ 3 \end{bmatrix} = \begin{pmatrix} 1 \\ -2 \end{pmatrix}$$

$$4^{\circ} V_3 = \beta V_2 - \alpha \hat{g}_2 = 0.8 \begin{bmatrix} 1 \\ -2 \end{bmatrix} - 0.2 \begin{pmatrix} 1 \\ -2 \end{pmatrix} = \begin{pmatrix} 0.6 \\ 0.24 \end{pmatrix}$$

$$5^{\circ} W_3 = W_2 + V_3 = \begin{bmatrix} 2 \\ 0.8 \end{bmatrix} + \begin{bmatrix} 0.6 \\ 0.24 \end{bmatrix} = \begin{bmatrix} 2.6 \\ 1.04 \end{bmatrix}$$

$$2) \hat{S} = \frac{S}{1 - \rho^t}, \quad S_t = \rho S + (1 - \rho) \hat{g}.$$

proof: for $t=1$: $S_1 = \rho S_0 + (1 - \rho) \hat{g} = (1 - \rho) \hat{g}.$

$$\hat{S}_1 = \frac{S_1}{1 - \rho} = \frac{(1 - \rho) \hat{g}}{1 - \rho} = \hat{g}.$$

assume it holds true for $t=n$, where $t \geq 1$, such as $\hat{S}_n = \hat{g}$, when $t=n+1$:

$$\therefore S_{n+1}^{\wedge} = \hat{g} \quad (\text{under assumption}).$$

$$\text{also } S_{n+1}^{\wedge} = \frac{S_n}{1 - \rho^n} = \hat{g}$$

$$\therefore \text{we can get } S_n = (1 - \rho^n) \hat{g}.$$

$$\begin{aligned}
 \text{for } s_{n+1} &= p \cdot s_n + (1-p) \cdot \hat{g} \\
 &= p \cdot (1-p^n) \hat{g} + (1-p) \hat{g} \\
 &= (p - p^{n+1} + 1 - p) \hat{g} \\
 &= (1 - p^{n+1}) \hat{g}.
 \end{aligned}$$

$$\therefore \hat{s}_{n+1} = \frac{s_{n+1}}{(1-p^{n+1})} = \frac{(1-p^{n+1}) \hat{g}}{1-p^{n+1}} = \hat{g}$$

\therefore The assumption also holds for \hat{s}_{n+1} , when \hat{s}_n holds.

\therefore if $\hat{s}_n = \hat{g}$ holds, then $\hat{s}_{n+1} = \hat{g}$ holds for $n \geq 1$.

\therefore By mathematical induction, $\hat{s}_n = \hat{g}$ for every step.