DSA5202    Lec 6.

Recap:

① FCNN $\begin{cases} \text{shallow} \\ \text{deep} \to \begin{cases} f(x) = V^\top f_T(x) \\ f_{t+1}(x) = \sigma(W_t f_t(x) + b_t) \end{cases} \end{cases}$

$$t = 0, 1, \dots, T-1$$

② Empirical Risk Minimization (ERM)

Training set

$$\hat{\theta} = \underset{\theta \in \textcircled{D}}{\arg\min} \; \frac{1}{N} \sum_{i=1}^{N} \ell(f_\theta(x_i), y_i)$$

optimize $\longrightarrow \begin{cases} GD \\ SGD \\ \text{Momentum } GD \end{cases}$

③ BP Algorithm to calculate gradient $\begin{cases} \text{forward} \\ \text{backward} \end{cases}$

(Batch update in practice)

forward : $x \to x_1 \to \cdots \to x_T \to x_{T+1}$

backward : $x_{T+1} \to P_{T+1} := \nabla_{x_{T+1}} \ell \longrightarrow \nabla_{W_T} \ell \leftarrow x_T$

$\downarrow$     $x_T = g_{T-1}(x_{T-1}, W_{T-1})$

$P_T \xrightarrow{\hspace{2cm}} \nabla_{W_{T-1}} \ell \leftarrow x_{T-1}$

① $\underline{P_t := \nabla_{x_t} \ell(x_{T+1}, y)}$

$\qquad = \underline{\nabla_{x_t} \cdot g_t(x_t, W_t) \cdot P_{t+1}}$

$\downarrow$

$\vdots$

$\downarrow$

② $\underline{\nabla_{W_t} \ell(x_{T+1}, y)}$

$P_1 \xrightarrow{\hspace{2cm}} \nabla_{W_0} \ell \leftarrow x_0 = x$

$\qquad = \underline{\nabla_{W_t} g_t(x_t, W_t) \cdot P_{t+1}}$     depth

✶✶✶

$\longrightarrow$ Issue : when T is large , the gradient is potential to $\begin{cases} \text{vanish} \\ \text{explode} \end{cases}$

$\longrightarrow$ especially for RNN !

$\hookrightarrow$ comes from the nature of BP

$$P_t = \left( \prod_{i=t+1}^{T+1} \gamma_i \right) P_{T+1}$$

tend to $\begin{cases} \text{explode} \\ \text{vanish} \end{cases}$

④ <u>CNN</u> $\longrightarrow$ <u>convolution operation</u>

---

<u>Today's lecture</u>:

1. <u>Gradient Vanishing</u> / <u>Gradient Exploding</u>

E.g. (intuition)   (toy example)

$\downarrow$                                    $\rightarrow$ omit bias

<u>FCNN</u>:   $\underline{x_{t+1} = \sigma(w_t \cdot x_t)}$      $t = 0, 1, \ldots, T$

$w_t, x_t \in \mathbb{R}$

consider $\Longrightarrow P_t = \nabla_{x_t} g_t(x_t, w_t) \cdot P_{t+1}$

$= \underbrace{w_t} \cdot \underbrace{\sigma'(w_t \cdot x_t)} \cdot P_{t+1}$      $\underline{t = 0, 1, \ldots, T}$

initialization                                    $\rightarrow$ activation function

$P_{T+1} = \nabla_{x_{T+1}} \ell(x_{T+1}, y)$

<u>Insight</u>: $\boxed{\text{gradient vanishing / exploding}} \Longleftrightarrow \begin{cases} \text{activation function} \\ \text{weight initialization} \end{cases}$

$\searrow$ also depends on the $\boxed{\text{choice of architecture}}$

<u>In this example</u>: we use $\boxed{\text{FCNN}}$ ( do not utilize architecture

$\downarrow$                          like skip-connection )

$\boxed{\text{most naive one}}$

2. choice of [activation function]
   ↓
   [1d FCNN] (omit bias)

Recap last example : $P_t = w_t \cdot \boxed{6'(w_t \cdot x_t)} \cdot P_{t+1}$

                                    ↓

                            [ideally] , we want $|6'| \approx 1$

→ Pess likely for GV

[Relu] : $6(z) = \max\{0, z\} \implies 6'(z) = \begin{cases} 1 & , z > 0 \\ 0 & , z < 0 \end{cases}$

[Sigmoid] : $6(z) = \dfrac{1}{1 + \exp(-z)} \implies 6'(z) = \dfrac{\exp(-z)}{(1 + \exp(-z))^2}$

→ potential GV

                                    ⟱

                    $0 \leq 6'(z) \leq \dfrac{1}{1 + \exp(-z)} < 1$

                                    ↘ may lead to [Gradient Vanish]

3. choice of [initialization]

    Recap : $P_t = w_t \cdot 6'(w_t \cdot x_t) \cdot P_{t+1}$

                            ⟱                          [idea] :

    $r_t := \dfrac{P_t}{P_{t+1}} = w_t \cdot 6'(w_t \cdot x_t) \implies$ we want to control this rate

→ Random Initialization : $w_t \sim \text{Gaussian}(0, \gamma^2)$

    ↳ Here, we choose (assume) $6 \longrightarrow$ ReLu

⟹ Question : How to choose $\gamma^2$ (variance) ?

**Answer:** we want to analyze $r_t := \frac{P_t}{P_{t+1}} = w_t \cdot 6'(w_t \cdot x_t)$

$$6'(z) = \begin{cases} 1, & z > 0 \\ 0, & z < 0 \end{cases}$$

⊛

$\boxed{x_t \text{ is also random}}$

**Conclusion:** $\begin{cases} \mathbb{E}_{w_t, x_t}[r_t] = 0 \\ \mathbb{E}_{w_t, x_t}[r_t^2] = \frac{\gamma^2}{2} \end{cases}$ $\Rightarrow$ suggest us to choose $\boxed{\gamma^2 = 2}$

**Derivation:** b) $\mathbb{E}_{w_t, x_t}[w_t^2 \, 6'(w_t \cdot x_t)]$

symmetric.

if $w_t \sim p(\cdot)$

$= \mathbb{E}_{w_t, x_t}[w_t^2 \, \mathbb{1}\{w_t \cdot x_t > 0\}]$

then: $\mathbb{E}_{w_t}[f(w_t)]$

$= \mathbb{E}_{w_t}[f(-w_t)]$ $= \mathbb{E}_{w_t, x_t}[(-w_t)^2 \, \mathbb{1}\{(-w_t) \cdot x_t > 0\}]$

$\Rightarrow 2\,\mathbb{E}_{w_t, x_t}[r_t^2] = \boxed{\mathbb{E}_{w_t, x_t}[w_t^2]} = \boxed{\gamma^2}$

$\|$

$\Rightarrow \underline{\mathbb{E}_{w_t, x_t}[r_t^2] = \frac{\gamma^2}{2}}$

$\mathbb{E}_{w_t}[w_t^2]$

$\longrightarrow$ Generally, the initialization scheme is:

In DNNs with width $d$, this becomes $\gamma_d^2 = 2/d$. This is known as Kaiming (or He) initialization scheme. More generally, with different widths $d_t$, we have

$$W_t^{ij} \sim \mathcal{N}(0, \frac{2}{d_t}).$$

$\rightarrow$ This also solves a problem of vanishing/exploding during forward propagation!!

$\rightarrow$ not only for the gradient back-propagation

for general width $d_t$ (previously, $d_t = 1$)

Consider a **FCNN**:
$$x_{t+1} = \sigma(W_t \cdot x_t), \quad \underline{W_t \in \mathbb{R}^{d_{t+1} \times d_t}}$$

$\boxed{\text{Goal}} \longrightarrow$

Here, $\underline{W_t^{ij} \sim N(0, \gamma_t^2)} \rightarrow$ we want to determine $\underline{\gamma_t^2}$

☆

$\color{red}\boxed{x_t \text{ is also random!}}$

$\longrightarrow x_t \in \mathbb{R}^{d_t} \longrightarrow x_t = \begin{pmatrix} x_t^1 \\ \vdots \\ x_t^{d_t} \end{pmatrix} \in \mathbb{R}^{d_t}$

$$\Rightarrow \boxed{x_{t+1}^i = \sigma\left( \sum_{\hat{j}=1}^{d_t} W_t^{i\hat{j}} x_t^{\hat{j}} \right)} \rightarrow \underline{\text{scaler form}}$$

$$\mathbb{E}_{W_t^{i\cdot}, x_t}\left[ (x_{t+1}^i)^2 \right] = \mathbb{E}_{W_t^{i\cdot}, x_t}\left[ \left\{ \sigma\left( \sum_{\hat{j}=1}^{d_t} W_t^{i\hat{j}} x_t^{\hat{j}} \right) \right\}^2 \right]$$

since $\underline{W_t^{i1}, \ldots, W_t^{id_t} \sim N(0, \gamma_t^2)}$

$$\Rightarrow \sum_{\hat{j}=1}^{d_t} x_t^{\hat{j}} \cdot W_t^{i\hat{j}} \Big| x_t \sim N\left( 0, \sum_{\hat{j}=1}^{d_t} (x_t^{\hat{j}})^2 \cdot \gamma_t^2 \right)$$

$$\Rightarrow \mathbb{E}_{W_t^{i\cdot}, x_t}\left[ (x_{t+1}^i)^2 \right] = \mathbb{E}_{x_t}\left[ \mathbb{E}_{W_t^{i\cdot}}\left[ \left\{ \sigma\left( \sum_{\hat{j}=1}^{d_t} W_t^{i\hat{j}} x_t^{\hat{j}} \right) \right\}^2 \Big| x_t \right] \right]$$

$\boxed{\begin{array}{l} \text{Lemma: if } \underline{Z \sim N(0, \alpha^2)} \\ \\ \text{then } \mathbb{E}[\sigma^2(Z)] = \dfrac{\alpha^2}{2} \\ \qquad\qquad \Downarrow \\ \mathbb{E}[\sigma^2(Z)] \\ = \mathbb{E}[Z^2 \cdot \mathbb{1}\{Z>0\}] \\ = \dfrac{\alpha^2}{2} \end{array}}$

$$= \mathbb{E}_{x_t}\left[ \sum_{\hat{j}=1}^{d_t} (x_t^{\hat{j}})^2 \gamma_t^2 \cdot \frac{1}{2} \right]$$

$$= \color{blue}\boxed{\frac{\gamma_t^2}{2} \cdot \mathbb{E}_{x_t} \cdot \left[ \| x_t \|_2^2 \right]}$$

$\color{blue}\rightarrow$ independent with respect to $i$

$$\color{red}= \frac{\gamma_t^2}{2} \cdot d_t \, \mathbb{E}_{x_t^i}\left[ (x_t^i)^2 \right]$$ $\color{blue}\text{choice of node in}$

$\color{blue}\underline{\text{next layer}}$

We want $\boxed{\dfrac{\gamma_f^2 \, dt}{2} = 1}$ $\rightsquigarrow$ stablize

$$\Updownarrow$$

$$\gamma_f^2 = \frac{2}{dt}$$

---

Recap: $\qquad P_t = W_t \cdot \sigma'(W_t x_t) \cdot P_{t+1}$ $\qquad \boxed{x_{t+1} = \sigma(W_t x_t)}$

$$\Downarrow$$

$$r_t := \frac{P_t}{P_{t+1}} = W_t \cdot \sigma'(W_t \cdot x_t)$$

① 
$$\Rightarrow \mathbb{E}_{W_t, x_t}[r_t] = \mathbb{E}_{x_t}\left\{\mathbb{E}_{W_t}[r_t \mid x_t]\right\}$$

$$= \mathbb{E}_{x_t} \cdot \left\{ \mathbb{E}_{W_t}\left[ W_t \cdot \mathbb{1}\{W_t x_t > 0\} \mid x_t \right] \right\}$$

$$= \mathbb{E}_{x_t}\left[ f(x_t) \right]$$

$$f(x_t) = \begin{cases} \mathbb{E}_{W_t}\left[ W_t \, \mathbb{1}\{W_t > 0\} \right] & \text{if } \underline{x_t > 0} \\[2mm] \mathbb{E}_{W_t}\left[ W_t \, \mathbb{1}\{W_t < 0\} \right] & \text{if } \underline{x_t < 0} \end{cases}$$

$$\left( \text{Assume: } \begin{bmatrix} \mathbb{P}(x_t > 0) = \frac{1}{2} \\ \mathbb{P}(x_t < 0) = \frac{1}{2} \end{bmatrix} \right) \qquad = \frac{1}{2}\left( \mathbb{E}_{W_t}\left[ W_t \, \mathbb{1}\{W_t > 0\} \right] \right.$$

$$\left. + \mathbb{E}_{W_t}\left[ W_t \, \mathbb{1}\{W_t < 0\} \right] \right)$$

$$= \frac{1}{2} \mathbb{E}_{W_t}\left[ W_t \right] = \underline{0}$$