

# Lec 7 $\rightarrow$ Gradient Vanishing $\leftrightarrow$ Model Architecture (stability)

Re-cap:

$$\textcircled{1} \begin{cases} \nabla_{W_t} \ell(f_\theta(x_{T+1}), y) = \nabla_{W_t} x_{t+1} \nabla_{x_{t+1}} \ell = \nabla_{W_t} g_t(x_t; W_t) \cdot p_{t+1} \\ p_t := \nabla_{x_t} \ell = \nabla_{x_t} x_{t+1} \cdot \nabla_{x_{t+1}} \ell = \nabla_{x_t} g_t(x_t; W_t) \cdot p_{t+1} \end{cases}$$

$\downarrow$   
for small  $t$ ,  $p_t \rightarrow 0$  since  $p_t = \prod_{i=t}^T \sigma_i p_{T+1} \rightarrow 0$

$\textcircled{2}$  For 1-dim FCNN without bias, we have

$$p_t = \frac{d}{dx_t} [g(W_t x_t)] \cdot p_{t+1}$$

$$= W_t \cdot \underbrace{g'(W_t x_t)}_{\text{choice of activation function}} \cdot p_{t+1}$$

$\underbrace{W_t}_{\text{weight init.}}$

choice of activation function

★ GV

$\rightarrow$  for activation function, we hope:  $g'(x) \approx 1 \leftrightarrow \text{ReLU}$

we don't want  $g'(x) \ll 1 \leftrightarrow \text{sigmoid}$

$\rightarrow$  for initialization of weight, we are interested in

$$r_t := \frac{p_t}{p_{t+1}} = \underbrace{W_t \cdot g'(W_t x_t)}_{\downarrow} x_0$$

Suppose  $W_t \sim N(0, \sigma^2)$  here,  $W_t$  &  $x_t$  are R.V.

Conclusion: a)  $\mathbb{E}_{W_t, x_t} [r_t] = 0$

$\rightarrow$  suggest  $\sigma^2 = 2$

$$b) \mathbb{E}_{W_t, x_t} [r_t^2] = \frac{\sigma^2}{2}$$

$\rightarrow$  a more general result: FCNN  $\left\{ \begin{array}{l} \text{depth } T \\ \text{width } d_t \text{ in } t\text{-th layer} \end{array} \right.$

★  
stabilize forward prop.  
& backward prop.

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

$$W_t \in \mathbb{R}^{d_{t+1} \times d_t}$$

$\rightarrow$  idea: stabilize  $\mathbb{E}_{W_t, x_t} [(x_{t+1}^i)^2]$

# Today's lecture :



1. All previous discussion is based on trivial FCNN / CNN

$$\Rightarrow \begin{cases} x_{t+1} = G(w_t \cdot x_t + b_t) & t=0, 1, \dots, T-1 \\ x_{T+1} = V^T x_T \end{cases}$$

Question: can we modify the ARCHITECTURE to solve Gradient Vanishing Problem?

2. Intuition: for 1-dim FCNN without bias example:

$$p_t = w_t \cdot G'(w_t \cdot x_t) \cdot p_{t+1}$$

may goes to 0

★

idea: can we achieve  $p_t = (1 + \star) \cdot p_{t+1}$ ?

goes to 1

ResNet Architecture:

$$\begin{cases} x_{t+1} = x_t + f_t(x_t) \\ x_{T+1} = V^T x_T \end{cases} \quad t=0, 1, \dots, T-1$$

t-th residual block

can be complicated

$$x_{t+1} - x_t = f_t(x_t)$$

"residual"

compute gradient for 1-dim case:

$$\rightarrow \frac{dp}{dx_t} = \frac{dp}{dx_{t+1}} \cdot \frac{dx_{t+1}}{dx_t} = \frac{dp}{dx_{t+1}} \cdot \left(1 + \frac{df_t}{dx_t}\right)$$

$$= \frac{dp}{dx_{T+1}} \prod_{k=t}^T \left(1 + \frac{df_k}{dx_k}\right)$$

not  $\approx 0$



Remark:

1. "1" plays crucial role in preserving the magnitude of gradient !!!  $\rightarrow$  why it will help solve GD

2. However, it may cause Gradient Exploding

$\rightarrow$  Solution:

① In practice, we can use Normalization Layer (BN....)



try to normalize the features after each layer

② modify  $\frac{dJ_k}{dx_k} \rightarrow \alpha_k \cdot \frac{dJ_k}{dx_k}$

$\Rightarrow$  scaling k-th residual blocks with factor  $\alpha_k$

3. Scaling k-th residual block to avoid Gradient Exploding

$\rightarrow$  Example: simple ResNet with width  $d$  (fixed)

$$\begin{cases} x_{t+1} = x_t + W_t \cdot \phi(x_t) \\ x_{T+1} = T_{fc} \cdot x_T \\ x_0 = W_{in} \cdot x \end{cases} \quad \begin{matrix} t=0, 1, \dots, T-1 \\ \rightarrow f_t(x_t) \end{matrix}$$

$W_t \in \mathbb{R}^{d \times d}$ , we assume  $W_t^{ij} \sim \mathcal{N}(0, \frac{\gamma_t^2}{d})$

$\rightarrow$  Our interest is:  $\mathbb{E}_{W_t, x_t} [(x_{t+1}^i)^2]$

Calculation:  $(x_{t+1}^i)^2 = (x_t^i + \sum_{j=1}^d W_t^{ij} \cdot \phi(x_t^j))^2$

Therefore,  $\mathbb{E}_{W_t, x_t} [(x_{t+1}^i)^2]$   $x_t = f(W_{t+1}, \dots, W_0, W_{in}; x)$

$$= \mathbb{E}_{w_t, x_t} [(x_t^i)^2] + \mathbb{E}_{w_t, x_t} \left[ \left( \sum_{j=1}^d w_t^{ij} b(x_t^j) \right)^2 \right]$$

Second term

$$+ 2 \cdot \mathbb{E}_{w_t, x_t} \left[ x_t^i \cdot \sum_{j=1}^d w_t^{ij} b(x_t^j) \right]$$

$$\mathbb{E}_{w_t, x_t} \left[ \left( \sum_{j=1}^d w_t^{ij} b(x_t^j) \right)^2 \right]$$

$$= \mathbb{E}_{x_t} \left\{ \mathbb{E}_{w_t} \left[ \sum_{1 \leq j, j' \leq d} w_t^{ij} w_t^{ij'} b(x_t^j) \cdot b(x_t^{j'}) \mid x_t \right] \right\}$$

$$= \mathbb{E}_{x_t} \left\{ \mathbb{E}_{w_t} \left[ \sum_{j=1}^d (w_t^{ij})^2 \cdot b(x_t^j)^2 \mid x_t \right] \right\}$$

$$= \mathbb{E}_{x_t} \left\{ \sum_{j=1}^d \mathbb{E}_{w_t} \left[ (w_t^{ij})^2 \cdot b(x_t^j)^2 \mid x_t \right] \right\}$$

$$= \sum_{j=1}^d \mathbb{E}_{x_t, w_t} \left[ (w_t^{ij})^2 \cdot b(x_t^j)^2 \right] \quad \leftarrow \begin{array}{l} \text{independence} \\ x_t = f(w_{t+1}, \dots, w_0) \perp w_t \end{array}$$

$$= \sum_{j=1}^d \mathbb{E}_{x_t^j} [b(x_t^j)^2] \cdot \mathbb{E}_{w_t^{ij}} [(w_t^{ij})^2] \quad \begin{array}{l} \text{Note:} \\ \{x_t^j\}_{j=1}^d \rightarrow \text{same distribution} \\ \text{(not independent)} \end{array}$$

$$= \sum_{j=1}^d \mathbb{E}_{x_t^j} [b(x_t^j)^2] \cdot \frac{\sigma_t^2}{d}$$

$$= \sigma_t^2 \cdot \mathbb{E}_{x_t^1} [b(x_t^1)^2]$$

Third term

$$\mathbb{E}_{w_t, x_t} \left[ x_t^i \sum_{j=1}^d w_t^{ij} b(x_t^j) \right]$$

$$= \sum_{j=1}^d \mathbb{E}_{w_t^{ij}} [w_t^{ij}] \cdot \mathbb{E}_{x_t} [x_t^i \cdot b(x_t^j)] \quad (\text{independence similarly})$$

$$= 0$$



Therefore:  $\mathbb{E}_{w_t, x_t} [(\chi_{t+1}^i)^2]$

$$= \mathbb{E}_{x_t^i} [(\chi_t^i)^2] + \gamma_t^2 \mathbb{E}_{x_t^i} [6(\chi_t^i)^2]$$

Note: given  $x$ ,  $x_0 = W_{in} x \Rightarrow x_0^i$  is symmetric w.r.t pdf

since  $x_0^i = \sum_{j=1}^d x^j w_{in}^{ij}$  Gaussian(0,  $\frac{\gamma^2}{d}$ )

Now,  $x_0 = W_{in} x \rightarrow$  symmetric

$x_1 = x_0 + W_0 \cdot 6(x_0) \rightarrow$  symmetric

$$= \mathbb{E}_{x_t^i} [(\chi_t^i)^2] + \frac{\gamma_t^2}{2} \cdot \mathbb{E}_{x_t^i} [(\chi_t^i)^2]$$

$$= (1 + \frac{\gamma_t^2}{2}) \mathbb{E}_{x_t^i} [(\chi_t^i)^2]$$

$$= \dots = (1 + \frac{\gamma_t^2}{2})^t \mathbb{E}_{x_1^i} [(\chi_1^i)^2]$$

exploding behavior!!!!

$$w_{in}^{ij} \sim \mathcal{N}(0, \frac{\gamma^2}{d})$$

$\rightarrow$  Solution: intuition: initialize the weights whose variance is decreasing with depth

① We can choose  $\gamma_t^2 = \frac{1}{t}$   $\rightarrow$  decrease with depth

then  $\mathbb{E}_{w_t, x_t} [(\chi_{t+1}^i)^2] = (1 + \frac{1}{2t})^t \mathbb{E}[(\chi_1^i)^2]$

$\rightarrow e^{\frac{1}{2}}$  as  $t \rightarrow \infty$

no longer exploding

② Moreover, we can modify the architecture to solve.

$$\begin{cases} x_{t+1} = x_t + \underline{\lambda_t} \cdot w_t \cdot \phi(x_t) \\ w_t^{ij} \sim \mathcal{N}(0, \frac{\sigma^2}{d}) \end{cases}$$

then  $\mathbb{E}_{x_t, w_t} [(x_{t+1}^i)^2] = \left(1 + \frac{\sigma^2 \cdot \lambda_t^2}{2}\right) \cdot \mathbb{E}[(x_1^i)^2]$

we can choose  $\lambda_t^2 = \frac{1}{T}$  to avoid exploding

$$\Rightarrow x_{t+1} = x_t + \frac{1}{\sqrt{T}} \cdot w_t \cdot \phi(x_t)$$

→ modification on architecture

Show the symmetry of  $x_t$

①  $x_0 = W_{in} x$   $\Rightarrow x_0^i = \sum_{j=1}^d x^j \cdot \underline{w_{in}^{ij}}$  → symmetric (Gaussian)

$$\begin{aligned} & \mathbb{P}(aX + bY \leq c) \\ &= \mathbb{E}_{x,y} [\mathbb{1}\{aX + bY \leq c\}] \\ &= \mathbb{E}_{x,y} [\mathbb{1}\{-aX - bY \leq -c\}] \\ &= \mathbb{P}(aX + bY \geq -c) \end{aligned}$$

$\Rightarrow x_0^i$  is symmetric RV  $i=1, 2, 3, \dots, d$

naturally symmetric

$$(\cdot)^i | x_t = \sum_{j=1}^d \phi(x_t^j) w_t^{ij} \mid x_t \sim \mathcal{N}(\cdot)$$

②  $x_{t+1} = x_t + \underline{w_t \cdot \phi(x_t)}$  → symmetric

$$\mathbb{P}\left(\sum_{j=1}^d \phi(x_t^j) w_t^{ij} \leq x\right)$$



$$= \mathbb{E} \left[ \mathbb{1} \left\{ \sum_j 6(x_t^j) w_{t+1}^{ij} \leq x \right\} \right]$$

$$= \mathbb{E}_{x_t} \left[ \mathbb{E}_{w_t} \left[ \mathbb{1} \left\{ \sum_j 6(x_t^j) w_{t+1}^{ij} \leq x \right\} \mid x_t \right] \right]$$

$$= \mathbb{E}_{x_t} \left[ \mathbb{E}_{w_t} \left[ \mathbb{1} \left\{ \sum_j 6(x_t^j) w_{t+1}^{ij} \geq -x \right\} \mid x_t \right] \right]$$

$$= \mathbb{P} \left( \sum_{j=1}^d 6(x_t^j) w_{t+1}^{ij} \geq -x \right) \quad \text{since } \sum_j 6(x_t^j) w_{t+1}^{ij} \mid x_t \text{ is symmetric}$$

Recap:  $x_{t+1} = \underbrace{x_t}_{\text{symmetric}} + \underbrace{w_{t+1} \cdot 6(x_t)}_{\text{symmetric}} \Rightarrow \underline{x_{t+1} \rightarrow \text{symmetric}}$

☆☆☆

→ Remark:

Up to now, we have explored recipes for Gradient Vanishing

① Activation function →  $g'(z) \approx 1$   $\Leftrightarrow$  Relu

② Architecture → skip-connection (Residual Block)

③ Initialization:

a) FCNN:  $w_{t+1}^{ij} \sim N(0, \frac{2}{dt})$

the first update  
★ for the first time of  
calculating forward  
propagation

→ maintain  $\mathbb{E}[(x_{t+1}^i)^2] = \mathbb{E}[(x_1^i)^2]$  in forward prop.

b) ResNet:  $w_{t+1}^{ij} \sim N(0, \frac{1}{dt} \cdot \frac{1}{t})$  → decay for later layer

→ maintain  $\mathbb{E}[(x_{t+1}^i)^2] = C \cdot \mathbb{E}[(x_1^i)^2]$  → not explode  
as  $t \rightarrow \infty$

Lemma: if  $z \rightarrow$  symmetric R.V.

$$\text{then } \mathbb{E}_z [G^2(z)]$$

$$= \mathbb{E}_z [z^2 \cdot \mathbb{1}\{z \geq 0\}]$$

$$= \frac{1}{2} \mathbb{E}_z [z^2]$$