

Recap:

① FCNN $\begin{cases} \text{shallow} \\ \text{deep} \end{cases}$

→ formulation:

$$\begin{cases} f(x) = v^T f_T(x) \\ f_{t+1}(x) = \sigma(W_t f_t(x) + b_t) \quad t=0, 1, \dots, T-1. \end{cases}$$

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

 $\begin{cases} T: \text{depth of NN} \\ d_t: \text{width of NN} \end{cases}$

② "Richness" of Hypothesis Space w.r.t $\begin{cases} \text{width } M \\ \text{depth } T \end{cases}$

↖ model capacity

★

→ force NN to learn Hierarchical Feature Extractor

via forcing $f_\theta(x_i)$ matches y_i

↖ "learning"

③ Back-propagation Algorithm

→ forward: $x \rightarrow x_1 \rightarrow \dots \rightarrow x_T \rightarrow x_{T+1}$.

→ backward: $x_{T+1} \rightarrow p_{T+1} \rightarrow \nabla_{w_T} \leftarrow x_T$

\downarrow

$p_T \rightarrow \nabla_{w_{T-1}} \leftarrow x_{T-1}$

\vdots

\downarrow

$p_1 \rightarrow \nabla_{w_0} \leftarrow x_0$

Issue: when we use "mini-batch GD"



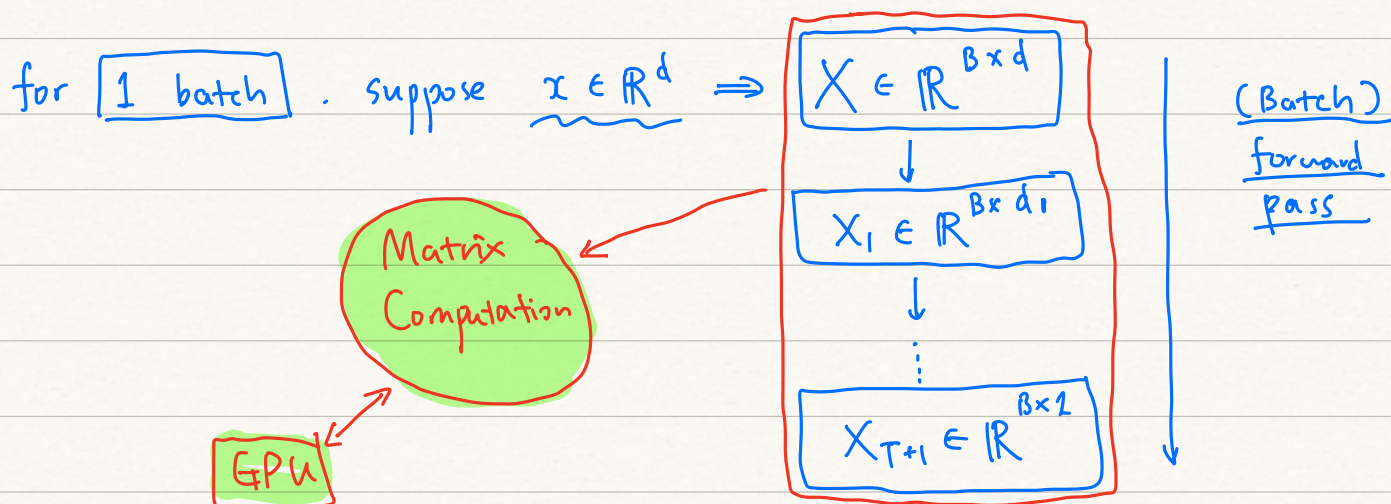
$$W_{t+1} = W_t - \alpha_t \frac{1}{|B|} \sum_{x \in B} \nabla_w f(f(x; w), y)$$

→ for each update, we need to go through $|B|$ times of forward pass and backward propagation.

Big number!

☆☆☆

Remedy: In practice, we do (Batch) forward pass & (Batch) backward propagation



Similarly, we also calculate gradient for the total batch

$$\begin{array}{ccccc} X_{T+1} \in \mathbb{R}^{B \times 1} & \longrightarrow & P_{T+1} & \longrightarrow & \nabla_{w_T} \in \mathbb{R}^{B \times d_{T+1} \times d_T} \longleftarrow X_T \\ & & \downarrow & & \\ & & P_T & \longrightarrow & \nabla_{w_{T-1}} \in \mathbb{R}^{B \times d_T \times d_{T-1}} \longleftarrow X_{T-1} \\ & & \downarrow & & \\ & & \vdots & & \\ & & P_1 & \longrightarrow & \nabla_{w_0} \in \mathbb{R}^{B \times d_1 \times d} \longleftarrow X_0 \end{array}$$

Today's Lecture:

1. Issue of FCNN → structure-agnostic
- ↓ formulate
- idea: for FCNN, there is too much freedom of choosing weight
- if $f \in \mathcal{H}_{\text{FCNN}}$, then $f_T \in \mathcal{H}_{\text{FCNN}}$
- Here, $f_T(x) := f(T(x))$, $T(\cdot)$ is translation
- for FCNN, it will treat original input and translated input equivalently \Rightarrow FCNN is structure agnostic, which means it will ignore the natural structure (temporal, spatial) of input.
- ↓
- Dis-advantage if our inputs do have some structure
- ↓ { image, time series

2. CNN: → Convolutional Neural Network

① Convolution operation on infinitely-long vecs

$$(w * x)(k) = \sum_{i=-\infty}^{+\infty} w(i) x(k+i)$$

$\begin{cases} w = \{w(i) : i \in \mathbb{Z}\} \rightarrow \text{filter} \\ x = \{x(i) : i \in \mathbb{Z}\} \rightarrow \text{signal} \end{cases}$

smoothing \leftrightarrow low-pass
edge \leftrightarrow high-pass

avg

② conv. on finite-long vecs \leftrightarrow padding

a) circular padding: $x \in \mathbb{R}^n$ $w \in \mathbb{R}^m$ $m \leq n$

$$(w * x)(k) = \sum_{i=0}^{m-1} w(i) \cdot x(k+i) \quad k=0, 1, \dots, n-1$$

$x \rightarrow$ extend so that $x(\bar{j}) = x(\bar{j}-n)$ for $\bar{j} \geq n$

e.g. $x(n+3) = x(3)$

b) zero padding : $x \in \mathbb{R}^n$ $w \in \mathbb{R}^m$ $m < n$

center around $x(k)$

$$(w * x)(k) = \sum_{i=0}^{m-1} w(i) \cdot x(k+i - \lfloor \frac{m}{2} \rfloor) \quad k=0,1,\dots,n-1$$

$$x(j) = 0 \text{ for } j \notin \{0,1,2,\dots,n-1\}$$



$x(j) = \text{constant}$ \Rightarrow c) constant padding

③ Compare between MatMul and Convolution

FCNN : $x_{t+1} = b(W_t \cdot x_t + b_t) \rightarrow$ matrix multiplication

CNN : $x_{t+1} = b(W_t * x_t + b_t) \rightarrow$ convolution



include multiple w (kernel)

④ Advantage of CNN

$w * x$ =
$$\begin{bmatrix} w_0 & w_1 & w_2 & 0 & 0 \\ 0 & w_0 & w_1 & w_2 & 0 \\ 0 & 0 & w_0 & w_1 & w_2 \\ w_2 & 0 & 0 & w_0 & w_1 \\ w_1 & w_2 & 0 & 0 & w_0 \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix}$$

1d case

circular padding

C_w \rightarrow less freedom (3 dof)

Recap : In FCNN, $\text{dof} = 5 \times 5 = 25$

$f_T \in \mathcal{H}_{\text{FCNN}}$ if $f \in \mathcal{H}_{\text{FCNN}}$

- a) weight sharing \rightarrow break the invariant of translation
- b) sparsity \rightarrow storage
- c) translation equivariant

for FCNN

$f_T \in \mathcal{H}_{\text{CNN}}$ if $f \in \mathcal{H}_{\text{CNN}}$

↓
Structure - in agnostic!

3. Pooling { down-sampling
approximately translation invariant

4. E.g. of Deep CNN architecture

$$\begin{array}{l} \text{Body} \\ \text{head} \end{array} \left\{ \begin{array}{l} x_0 = x \\ x_{t+1} = T_{mp} \circ T_{conv} \circ x_t \quad t = 0, 1, \dots, T-1 \\ f(x) = T_{fcnn} \circ \text{flatten} \circ x_T \end{array} \right.$$