

# LEC 6

## Recap

1. **NN** → universal approximator
  - Shallow NN ↔ **MLP**
  - Deep NN

2. **Optimization** → **GD**

SGD with mini-batch

learning rate ↔ Trade-off

batch size

convergence rate  
quality of solution

3. (Also, we talk about NN will not suffer curse of dimensionality)

## Architecture of NN

(FCNN)

→ Last Week → fully-connected NN architecture

simulate the neuron structure

Issue

(naturally)

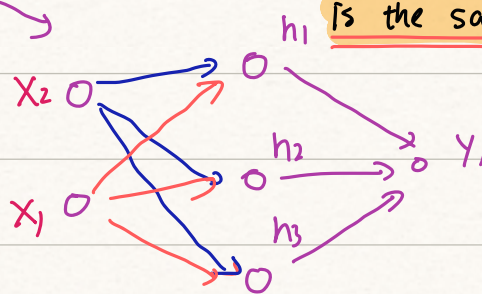
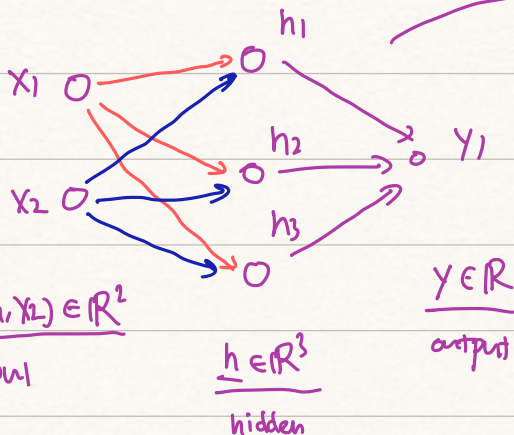
permutation Invariance

$\{1, 2, 3\} \rightarrow \{3, 1, 2\}$

① **本质** ⇒  $w^T x$  is permutation invariant

FCNN Invariance

② after permutation, the hypothesis space is the same!!!



## Procedure

$$\begin{cases} X_i = (x_i^1, \dots, x_i^d) \in \mathbb{R}^d \\ y_i \in \mathbb{R} \end{cases}$$

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$$

$\hat{\theta}_1$

① we train the NN by minimizing  $R_{\text{emp}}(\underline{\theta}) = \frac{1}{N} \sum_{i=1}^N \ell(f_{\underline{\theta}}(x_i), y_i) \quad \underline{\theta} \in \mathcal{H}$

② if we apply PERMUTATION to  $x_i \rightsquigarrow \boxed{x_i^p}$

$$\min_{\underline{\theta} \in \mathcal{H}}: R_{\text{emp}}^p(\underline{\theta}) = \frac{1}{N} \sum_{i=1}^N \ell(f_{\underline{\theta}}(x_i^p), y_i)$$

$\hat{\theta}_2$

then  $\hat{\theta}_1$  &  $\hat{\theta}_2$  is different only under permutation

Mathematically, FCNN invariance is:

$\Rightarrow f \in \mathcal{H}$ , then for any permutation  $p$ ,

$f_p(x_1, \dots, x_d) := f(x_{p(1)}, \dots, x_{p(d)})$ , then  $f_p \in \mathcal{H}$

$\downarrow$   
this is very bad when inputs have natural order (sequence) (spatial/temporal structure)

not a good idea to use FCNN

Images  
stock prices  
text

## CNN

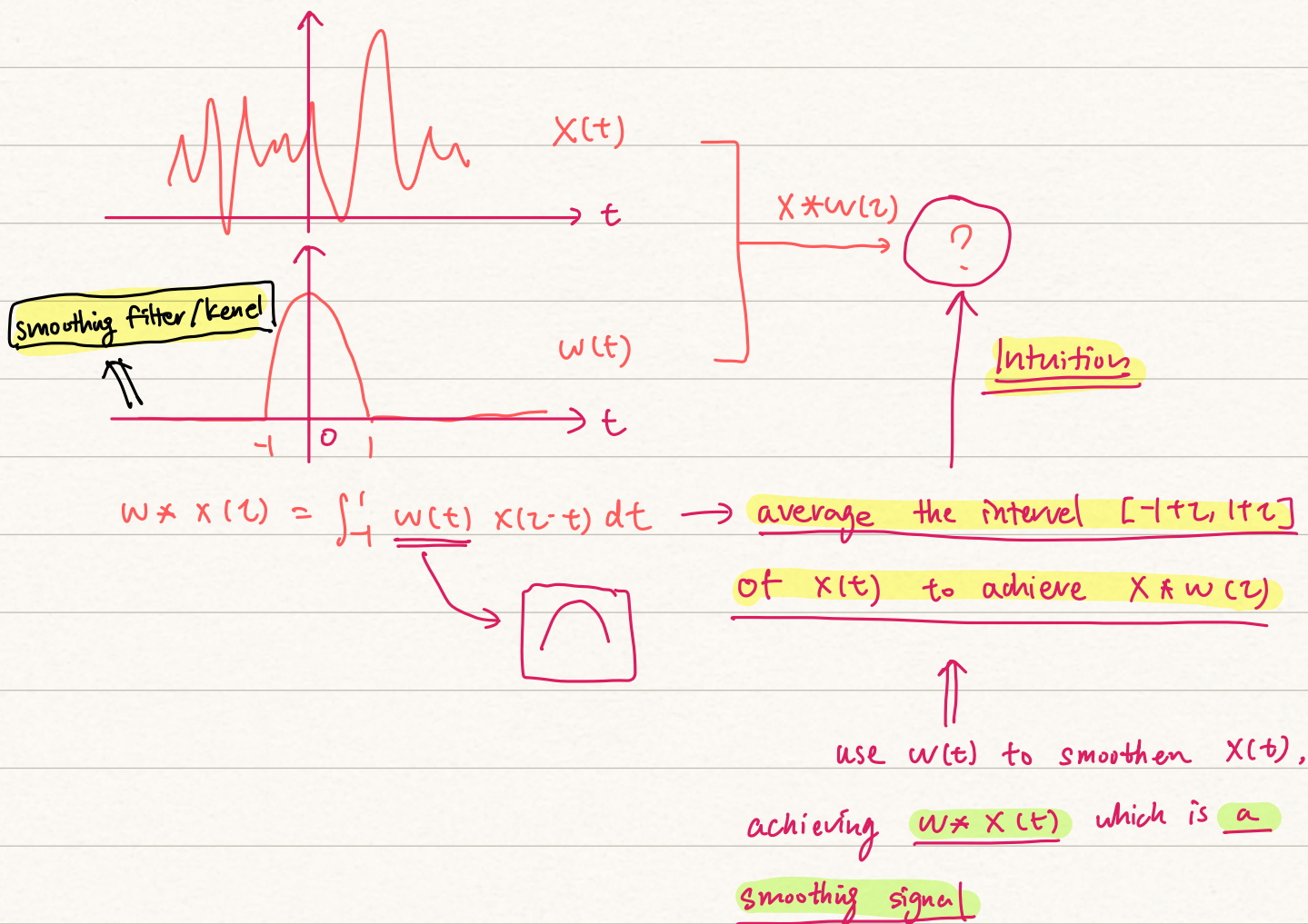
① convolution operation

$$a) w * x(\tau) = \int_{-\infty}^{+\infty} w(t) x(\tau-t) dt = \int_{-\infty}^{+\infty} x(t) w(\tau-t) dt = x * w(\tau)$$

$$\rightarrow \begin{cases} \text{commutative} \Rightarrow w * x(\tau) = x * w(\tau) \\ \text{linearity} \Rightarrow w * (\lambda_1 x_1 + \lambda_2 x_2)(\tau) = \lambda_1 w * x_1(\tau) + \lambda_2 w * x_2(\tau) \end{cases}$$



b) what can conv. do?



c) Discrete Convolution  $\rightarrow$  real-life application

$$\Rightarrow (w * x)(k) = \sum_{i=-\infty}^{+\infty} w(i) x(k-i) \quad (\text{discrete convolution})$$

(discrete Cross-Correlation)  $\rightarrow$  applied in CNN

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

truncate

$$(w * x)(k) = \sum_{i=1}^m w(i) x(k+i) \quad k = 0, 1, 2, \dots, d-m$$



where  $w \in \mathbb{R}^m$   $x \in \mathbb{R}^d$   $m < d$   $w * x \in \mathbb{R}^{d-m+1}$

kernel signal

dimension is decreasing

BAD

cannot have large number of conv.

solved by padding (zero padding)

$$\rightarrow x \in \mathbb{R}^d \xrightarrow{\text{pad}} x_{\text{pad}} \in \mathbb{R}^{d+m-1} \xrightarrow{w(\cdot)} w * x_{\text{pad}} \in \mathbb{R}^d$$

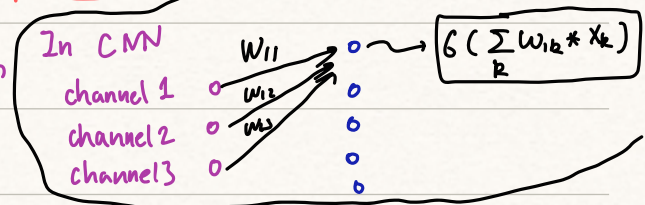
填充  $m-1$  个 0

dimension preserve



d) multi-dimension conv.

2-D conv:  $(w * x)(k, l) = \sum_{i=-\infty}^{+\infty} \sum_{j=-\infty}^{+\infty} w(i, j) x(k+i, l+j)$



multi-channels

image with channels (RGB)  $\rightarrow$  3 channels

$$G\left(\sum_k w_{ek} * x_k + b_e\right)$$

the feature map (dx) given by  $e$ -th kernel



FCNN:  $G\left(\sum_k w_{ek} \cdot x_k + b_e\right) \Rightarrow$  hidden unit  $h_e$

e) weight sharing

Matrix computation

$$w * x \text{ in 1D} \leftrightarrow w \cdot x$$

$$= \begin{pmatrix} w_1 & w_2 & w_3 & 0 & 0 \\ 0 & w_1 & w_2 & w_3 & 0 \\ 0 & 0 & w_1 & w_2 & w_3 \\ w_3 & 0 & 0 & w_1 & w_2 \\ w_2 & w_3 & 0 & 0 & w_1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix}$$

Recall: FCNN Invariance  
(permutation invariance issues)



★  
 $\Rightarrow$  CNN solves  
the Invariance Issue  
 (permutation)

permutation variant !!!

Since convolution has its own structure  
 of Matrix  $W$

not so much freedoms

Compare: FCNN :  $Wx = \begin{pmatrix} w_{11} & & \\ & \ddots & \\ & & w_{ss} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_s \end{pmatrix}$

all is free to choose

Permutation Invariant !

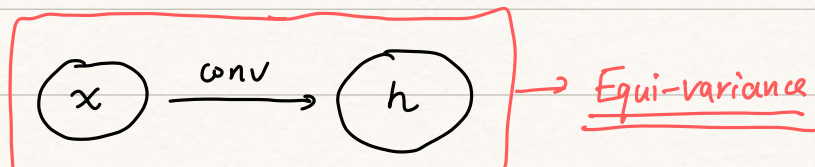
f) feature extractor  $\rightarrow$  we can design different  
kernels/convolutions to achieve some function

- blur
- vertical edges
- ...

Defn: Equi-variance (  $\hat{f}$  is equi-var with Translation  $T$  )

$$\underline{f(T(x)) = T(f(x))}$$

Conclusion: Conv. is Equi-variance with Shift-Translation



Defn: Invariance:

$$f(T(x)) = f(x)$$

[e.g.]  $f(x) = \sum_{i,j}^n x(i,j)$        $X \in \mathbb{R}^{n \times n}$

$$f(x) = \prod_{i,j}^n x(i,j)$$

☆

Conclusion: for some translation  $T$ , if  $\begin{cases} f \rightarrow \text{invariant} \\ g \rightarrow \text{equivariant} \end{cases}$

then  $f \circ g$  is invariant

→ firstly apply equivariant function  $g$ .

secondly apply invariant function  $f$ .

$$f \circ g_1 \circ \dots \circ g_j \rightarrow \text{invariant!}$$

[e.g.] ~~FCNN~~

FCNN is permutation invariant !!!

$$f_{FC}(x) = \sum_{i=1}^m v_i \sigma(w_i^T x + b_i)$$

$$\mathcal{D} = \{ (x_i, y_i) \}_{i=1}^N$$

$$\downarrow \quad \tilde{x}_i = x_i + a = T(x_i)$$

$$\mathcal{D}_s = \{ (\tilde{x}_i, y_i) \}_{i=1}^N$$

$$f_{FC}(x) = \sum_{i=1}^M v_i \sigma(w_i^T x + b_i)$$



$$= \sum_{i=1}^M v_i b(w_i^T \tilde{x} + \tilde{b}_i) \quad \text{where } \tilde{b}_i = b_i - w_i^T a$$

$$= f_{FC}(\tilde{x})$$

$$= f_{FC}(T(x))$$

Shrink hypothesis class  $\implies \mathcal{H}_{\text{conv}} = \{f: f(T(x)) = f(x)\}$

↑  
invariant function class

Pooling Layer  $\rightarrow$  Max Pooling  $\rightarrow$  can shrink dimensionality  
(with a stride p)