

Recap: FCNN



$$\text{hidden: } \begin{cases} h^{(i+1)} = \text{activation}(W^{(i+1)} h^{(i)} + b^{(i+1)}) \\ h^{(0)} = x \end{cases} \quad i = 0, \dots, T$$

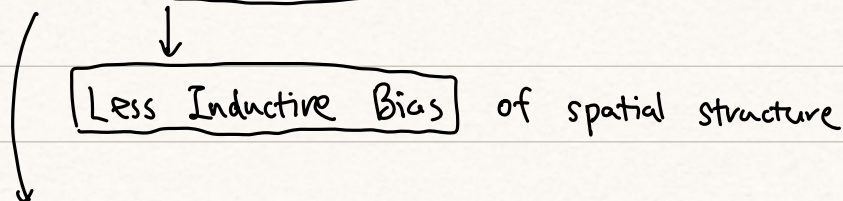
$$\text{output: } \hat{y} = w^T h^{(T)} + b$$

Observation: FCNN does not outperform other ML Algorithm

Today: CNN (Now also have Transformer)

(Weakness w.r.t Spatial / Temporal Data)

① Permutation Invariance of FCNN (not desirable)



Actually it is the Property of FCNN Hypothesis Space

$$\Rightarrow \forall f \in \mathcal{H}_{\text{FCNN}}, \text{ then } f_p \in \mathcal{H}_{\text{FCNN}}$$

where $f_p(x_1, \dots, x_d) = f(x_{p(1)}, \dots, x_{p(d)})$

Meaning: if we permute the feature index,

then we can re-train a FCNN to behave as good
as the previous unpermuted data → on permuted data

(w.r.t Hypothesis Space)

Not desirable for sequential (spatial) data
 Desirable for tabular feature (attributes)

② Conv. operation ("filter")

conv $\rightarrow x * w(t) = \int x(s) w(t-s) ds$

$$w(x) = \begin{cases} e^{-x} & , x \geq 0 \\ 0 & , x < 0 \end{cases} \rightarrow \underline{\text{Low-pass filter}}$$

\searrow Learnt by Data

Discrete Convolution

$$\begin{aligned} S(t) &:= x * w(t) = \sum_a x(a) w(t-a) \\ &= \sum_a x(a) \tilde{w}(a-t) \end{aligned} \quad \boxed{\tilde{w}(t) = w(-t)}$$

a) Circular convolution $\rightarrow x(-1) := x(n)$

3 2 5 1 4

\downarrow pad

4 3 2 5 1 4 3

\downarrow then do convolution

one way of padding

b) valid convolution \rightarrow no padding

\downarrow

output size will decrease (we do not want)

★

c) Zero-padding convolution \Rightarrow maintain the output size

2D-conv

$$S(i,j) := \sum_{m,n} X(m,n) w(i-m, j-n)$$

Convolution \iff Cross-correlation

$$\begin{aligned}\text{conv: } S(i,j) &= \sum_{m,n} X(m,n) W(i-m, j-n) \\ &= \sum_{m,n} X(m,n) \tilde{W}(m-i, n-j)\end{aligned}$$

[往右下 shift (i,j) unit]

$$\text{cross-correlation: } S(i,j) = \sum_{m,n} X(m,n) \underline{W(m+i, n+j)}$$

→ [implementation]

shift index

(since NN just needs an element-wise calculation)

Conv is just one Linear Operation

$$w * x = \boxed{A(w) x}$$

→ matrix multiplication

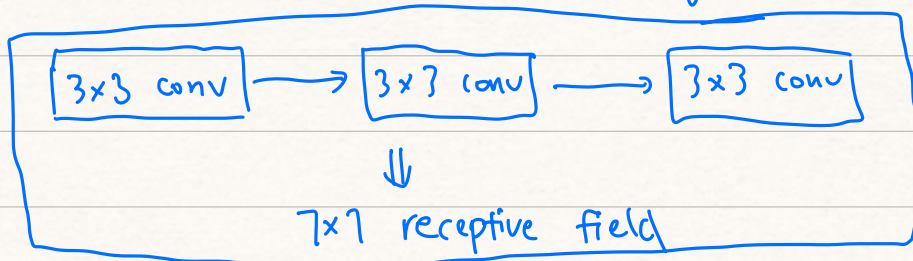
$$= \begin{pmatrix} w_2 & w_3 & 0 & 0 & 0 \\ w_1 & w_2 & w_3 & 0 & 0 \\ \vdots & & \vdots & & \vdots \\ 0 & 0 & 0 & w_1 & w_2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix}$$

③ Meaning of conv (filter) \Rightarrow [feature extractor]

④ Property of conv → [inductive bias]

- a) [Sparsity] → [against over-fitting] (only focus on neighbour input)
- b) [weight sharing] → { [extract features]
[save memory storage space]

Note: [Receptive Field] → when we stack MORE CNN Layers
then it will enlarge Receptive Fields



c) translation - equivariance (of convolution operation)

→ g-equivariance : $f(g(x)) = g(f(x))$

[e.g.] ReLU func is equivariant w.r.t scaling func

→ g - invariance : $f(g(x)) = f(x)$

[e.g.] $\left\{ \begin{array}{l} \text{norm func. is invariant w.r.t orthogonal projection} \\ \text{max / sum is invariant w.r.t translation} \end{array} \right.$

↓
special affine transformation

1. T_c : translation of signal x

then we have $\omega * T_c(x) = T_c(\omega * x)$

Conclusion: ① for circular conv, it is translation-equivariant

② for other convolution (zero-padding etc.),

It is not strictly translation-equivariance for boundary.

But it is "almost" translation-equivariance

2. f_1, f_2 is g -equivariant. then $f_1 \circ f_2$ is also g -equivariance

Pf.: $f_1 \circ f_2 (g(x)) = f_1 (g(f_2(x))) = g(f_1 \circ f_2(x))$

3. F is g -invariant. f is g -equivariant.

then $F \circ f$ is g -invariant

Pf: $F \circ f(g(x)) = F(g(f(x))) = F(f(x)) = F \circ f(x)$

★

Corollary: $f_i \rightarrow$ convolution layer i

$F \rightarrow$ pooling layer

then $F \circ (f_1 \circ \dots \circ f_i)$ is translation-invariant

↓
CNN stack

⑤ Infinitely Strong Prior

↓

$\mathcal{H} := \{ f : f(T(x)) \approx f(x) \} \rightarrow$ T -invariant Hypothesis

↓↓

Make sense for image classification

⑥ Pooling Layer \rightarrow "almost" translation-invariant

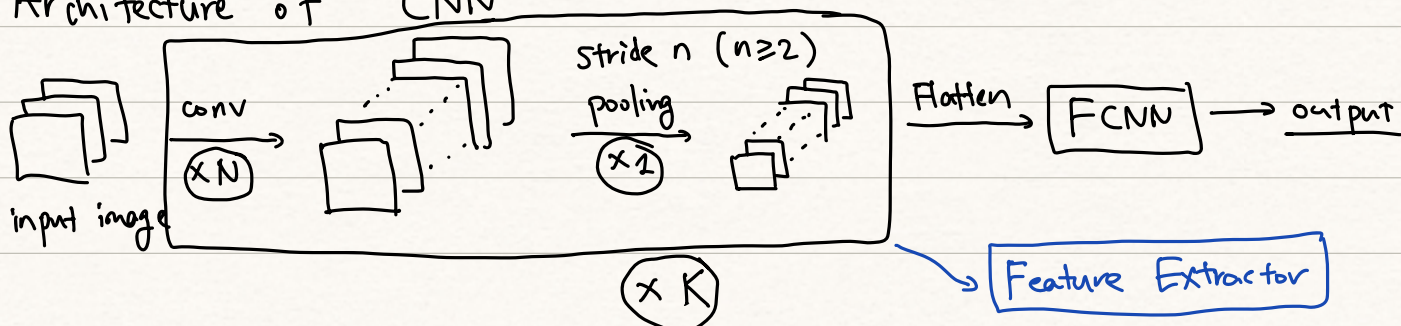
↓

for small translation

type of pooling $\left\{ \begin{array}{l} \text{average} \\ \text{max} \\ \vdots \end{array} \right.$

stride

⑦ Architecture of CNN



⇒ This is why we can use FCNN in the last layer
since those inputs can be viewed as "attributes of original image"