# L10
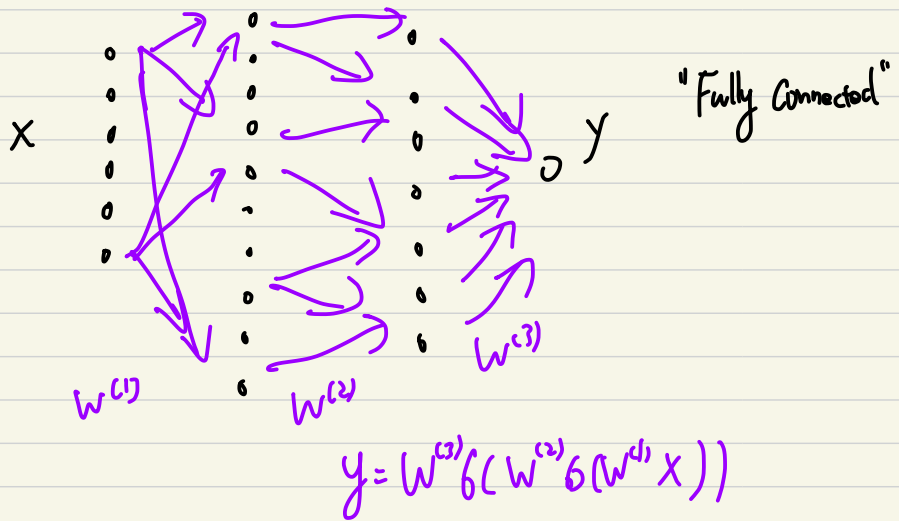
midterm solution will be posted tomorrow.

* Recap
* Convolutional Neural Network (CNNs)
* Recurrent Neural Network (RNNs)

---



$X$

"Fully Connected"

$y$

$W^{(1)}$    $W^{(2)}$    $W^{(3)}$

$$y = W^{(3)} 6 ( W^{(2)} 6 ( W^{(1)} X ))$$

problems of this picture:
* Computational Cost

e.x.   large processing
       Input :    RGB images of 400x400
       1 layer of neurons : 1000 neurons
       Output : 10 classes

# edges (weights)

$$= 3 \times 400 \times 600 \times 1000 + 1000 \times 10$$

RGB ↗            ⌣ bias

$$= 700 \text{ million weights}$$

Generally:    L layer network

$$O(L \cdot d^2)$$

☆ Sample complexity

- Require a very large training dataset

★ Loss of context

- all input is flattened into vector

- e.g images have special information

- e.g NLP/speech

have long-range temporal dependencies

混淆的

☆ Confounding features:

—

- Background cluster

- illumination

—

)
.

# Convolution Neural Network (CNNs) 最著名 image classification

Two main principles:

    ① Locality :    The model should look at local regions of

          your input image    如图脉识别，直接只看物品那一部分

    ② shift - invariance : The model should make the

    same prediction to the same object no matter where

    it appears

        | Convolutional layers |

---

## Convolution

$$X [t] \qquad W [t] \text{ called}$$

input ⟶          ⟶ filter / kernel

$$Y(t) = (x * w)[t] \qquad \text{卷积表达式}$$

$$= \sum_{\tau} X[t - \tau] w(\tau)$$

"Translate - and - scale"

Not hard to prove:    $y(t) = \sum_{\tau} X(\tau) w(t - \tau)$

    "Flip and dot"

* linear, associative, commu ~~ ?

* shift-invariance ||

2D - images       $X(s,t)$     $w(s,t)$

$2\frac{4}{5}\frac{6}{7}$       $(x * w)[s,t] = \sum_{6} \sum_{\tau} x[s-6, t-\tau] w[6, \tau]$

---

Convolution layers

Same    as    convolution, exept no flipping.

$(x * w)[s,t] = \sum_{\tau=-\Delta}^{\Delta} \sum_{6=-\Delta}^{\Delta} x[s+6, t+\tau] w[6, \tau]$

$\Delta \rightarrow$ locality ||

example $\Delta: 2 \rightarrow$ restrict to   $X[-2:2, -2:2]$     $5 \times 5 = 25$ indices
                                                                    of input

$h[s,t] = \phi(z[s,t])$

$h[s,t] = \phi((x * w)[s,t])$

$x_1$    0 ————→ 6  $z_1$
$x_2$    0 ———→ 0  $z_2$
$x_3$    0 ———→  0  $z_3$
$x_4$    0
         0
         0

✱ import to understand.

Consequences: 1) sparsely connected network

2) weights sharing

↓

detect some part satisfy feature we want.

it is like a certain kernel determined to fine one feature can be used at any location of the image

In practice, we have $D$ input channels

$F$ output channels.

Input $X_1 \cdots X_2 \cdots X_D$

Output $Z_1, Z_2 \cdots Z_F$

$$Z_1 = \sum_j X_j * W_{ij}$$

$$h_i = \phi(Z_1)$$

Computational benefits:

$$(2\sigma+1)^2 \cdot D \cdot F = O(\sigma^2 \cdot D \cdot F)$$

输入由于卷积核涉 $(2\sigma+1)$ ... W

Eg. $(20+1)^2 \cdot bf = 25 \cdot 3 \cdot 6 = \leq 600$

↑    ↑    ↑

⬜  RGB  output 6 channel

conv.

different  CNNs :

o Le Net — 5

o AlexNet (2012)   similar but deeper ⎤
                                        ⎥
• VGG Net  (2014)                       ⎥
                                        ⎥
• Google Net (2014)                    ⎦

⎡ o Residual Networks ⎤

Training :  Back propagation

# Recurrent Neural Networks ( Useful for Time Series/NLP analysis )

NLP
- Document retrieval
- Speech to text
- Language translation
- Sentiment prediction

— short range

— long range      dependencies   eg. 句子 长短

Classical:   Markov Models.
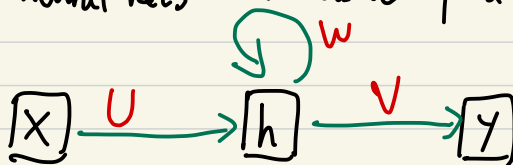
$$P ( \{ w_1, w_2, \dots, w_d \} )$$

往前
dependency

$$= P(w_1) \; P(w_2|w_1) + P(w_3|w_2) \dots P(w_d|w_{d-1})$$

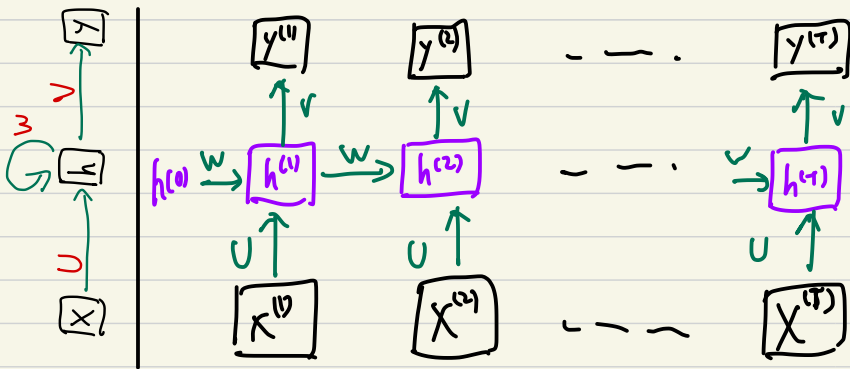$P(w_3|(w_1,w_2))$  需要9 word 依赖 $w_1, w_2 \rightarrow$ 更复杂]

long range
dependency

Recurrent neural nets : introduce feed back into neural nets

$$\boxed{X} \xrightarrow{U} \boxed{h} \xrightarrow{V} \boxed{Y}$$

$W$ (feedback loop on $h$)

$X^t$ indexed with time

$$h^t = \phi(Ux^t + Wh^{t-1})$$

$$y^t = Vh^t$$



U. V. W <u>don't</u> change over time

∴ <u>weight sharing</u>

Examples:
- Basic RNN  (h is single layer)
- GRU   ( Gated recurrent units)
- LSTM  ( Long -short term memory)
- Attention networks