

Deep learning

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



$$\frac{\Delta C}{\Delta b_1} \approx \frac{\Delta C}{\Delta b_1}$$

$$\Delta a_1 \approx \frac{\partial a_1}{\partial z} \frac{\partial z}{\partial b} \Delta b_1$$

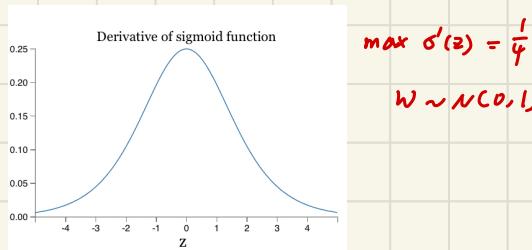
$$= \sigma'(z_1) \Delta b_1$$

$$\Delta a_2 \approx \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial a_1} \Delta a_1$$

$$= \sigma'(z_2) - w_2 \sigma'(z_1) \Delta b_1$$

...

$$\frac{\Delta C}{\Delta b_1} = \sigma'(z_4) w_4 \underbrace{\sigma'(z_3)}_{< \frac{1}{4}} w_3 \underbrace{\sigma'(z_2)}_{< \frac{1}{4}} w_2 \underbrace{\sigma'(z_1)}_{< \frac{1}{4}}$$



Avoid 1) $w_j = 100$ $z_j = 0$

$$b_j = -w_j x$$

P

1) $|w\sigma'(wz+b)| \geq 1$
 $|\sigma'| \leq \frac{1}{4} \quad \therefore |w| \geq 4$

2) $|w| \frac{e^{wz+b}}{(e^{wz+b} + 1)^2} \geq 1$
 $\frac{e^{wz+b}}{(e^{wz+b} + 1)^2} \geq \frac{1}{|w|}$

$$(e^{wz+b})^2 + (2 - |w|) e^{wz+b} + 1 \leq 0$$

$$\frac{-2 - |w|}{2} \leq e^{wz+b} \leq \frac{|w| - 2 + \sqrt{|w|(|w| - 4)}}{2}$$

$$-2 - |w| \leq \frac{1}{|w|} \ln \left(\frac{|w| - 2 + \sqrt{|w|(|w| - 4)}}{2} \right) - 1 + \frac{b}{|w|}$$

$$\text{so band } \leq \frac{2}{|w|} \ln \left(\frac{|w| + \sqrt{|w|^2 - 4|w|}}{2} \right) - 1$$

3) greatest value

$$\text{band}' = -\frac{2}{|w|^2} \ln(\cdot) + \frac{2}{|w|} \frac{1}{(\cdot)} - \frac{1}{2} \left(1 + \frac{1}{2} \frac{1}{\sqrt{(\cdot)}} \right) = 0$$

$$-2 \ln(\cdot) \cdot \frac{1}{\sqrt{(\cdot)}} + |w| \left(\frac{1}{\sqrt{(\cdot)}} + \frac{1}{2} \right) = 0$$

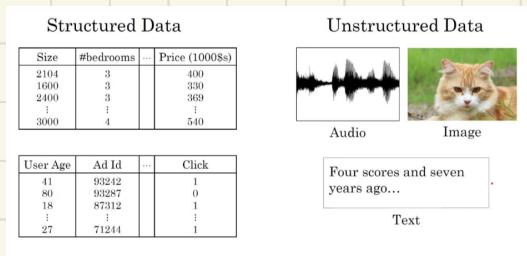
$$\sqrt{(\cdot)} |w| - 2 \ln(\cdot) + \frac{1}{2} |w| = 0$$

Deep Learning

Supervised Learning

Input(x)	Output (y)	Application
Home features	Price	Real Estate
Ad, user info	Click on ad? (0/1)	Online Advertising
Image	Object (1,...,1000)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

Standard NN
 CNN
 RNN
 Custom Hyper



Good Results on Training Data

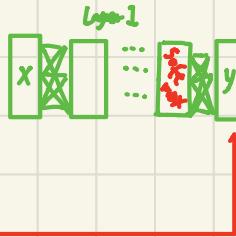
(0.) softmax :

$$a_j^e = \frac{e^{z_j^e}}{\sum e^{z_i^e}}$$

1. choosing proper loss : cross entropy

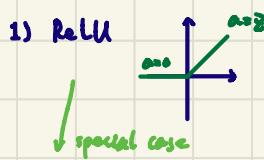
$-y_i \ln y_i$

\times square error $(\hat{y}_i - y_i)^2$



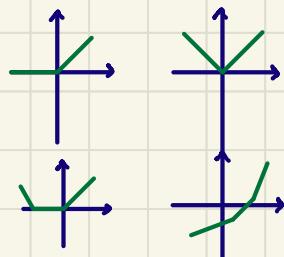
2. Mini batch

3. New activation func



- fast
- infinite sigmoid with biases
- vanishing gradient problem

2) Tanh



4. Learning Rate

- 1) small \rightarrow $y \leftarrow$ large
fast \rightarrow result \leftarrow converge

2) Adagrad $g_w(y, \frac{\partial L}{\partial w})$

$$w \leftarrow w - g_w \frac{\partial L}{\partial w}$$

$$g_w = \frac{y}{\sqrt{\sum (\frac{\partial L}{\partial w})^2}}$$

- $g_w \propto$ always
- if $\frac{\partial L}{\partial w}$ small, g_w big

5. Momentum

$$1) \text{Movement} = -\frac{\partial L}{\partial w} + \text{Momentum}$$

$$\downarrow +$$

2) Adam RMSprop (Advanced Adagrad) + Momentum

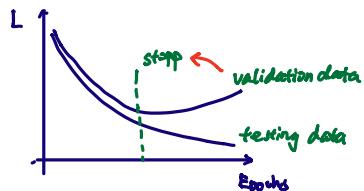
Good Results on Testing Data

(Overfitting)

more training $\xleftarrow[\text{create}]{\text{get}}$

1. Early Stopping

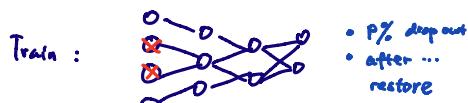
testing data \rightarrow validation data



2. Weight Decay (Regularization)

$$J = \sum (f - g)^2 + \lambda \cdot \sum w^2$$

3. Dropout



- $p\%$ dropout
- after ... restore

- Test:
- No drop out
 - $w \leftarrow w \cdot (1-p)\%$

4. NN Structure

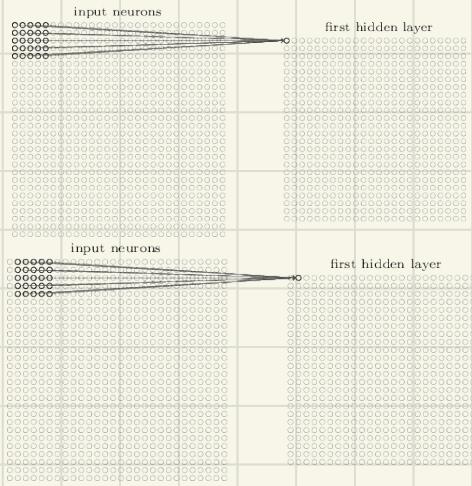
CNN Convolutional : convolution $\} \times h \rightarrow$ Flatten \rightarrow Fully Connected
 RNN Recurrent : max pooling Feedforward Network

Convolutional Neural Network

local receptive fields

shared weights

pooling



$$\text{input } 5 \times 5 \rightarrow 1 \times 1$$

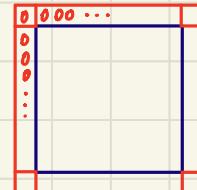
$$28 \times 28 \rightarrow 24 \times 24$$

$$\sigma(b + \sum_{l=0}^{4 \times 4} \sum_{m=0}^{4 \times 4} w_{l,m} a_{j+l, k+m})$$

shared bias L shared weights ($n=25$)

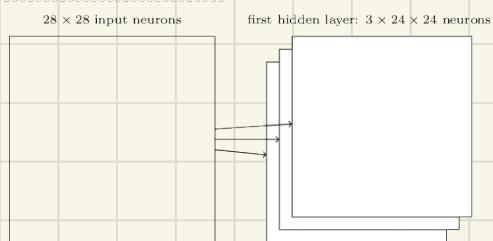
$\begin{matrix} (i,j) \\ (l,m) \end{matrix}$

→ Zero Padding

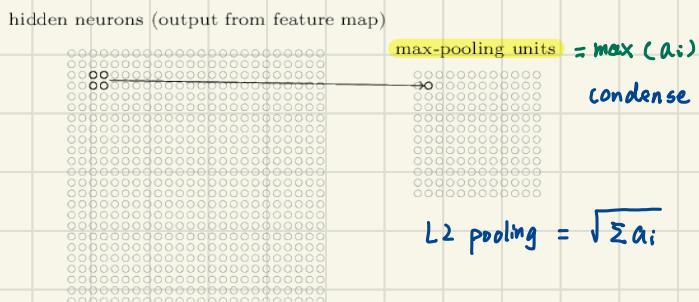


① Convolutional

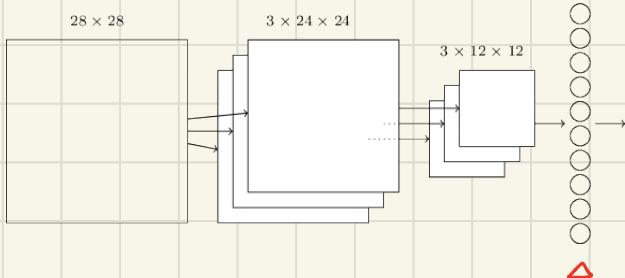
layer



3 Filter



full connected



②

Flatten

① Initial $w \sim N(0, 1)$
 $b \sim N(0, 1)$

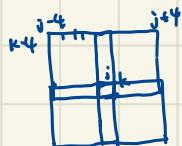
② forward
1) $L_{i,i} : a_{jk} = \sigma \left(\sum_{l=0}^m w_{l,m} a_{j+l, k+m} + b \right) = z_{jk}$

2) $L_{2i} : a_{jk} = \max (a_{2j-1, 2k-1} / a_{2j-1, 2k}, a_{2j, 2k-1} / a_{2j, 2k})$

3) $L_{3i} : a_j = \sum_w a_{kj}$

③ Backpropagation

$L_0 \rightarrow L_1 : \Delta a_{jk}^0 :$



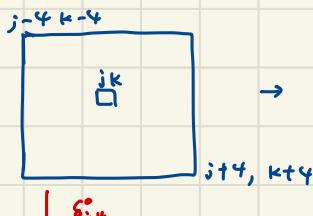
$$a_{j-4, k-4}^1 = \sigma \left[b + \sum_{l=0}^m w_{l,m} (a_l^0 + \Delta a_l^0) \right] \Rightarrow \frac{\partial a_{j-4, k-4}^1}{\partial a_{jk}^0} = \sigma' w_{4,4}$$

$$\begin{aligned} a_{j-3, k-4}^1 \\ a_{j-2, k-4}^1 \\ a_{j-1, k-4}^1 \\ \vdots \\ a_{j, k}^1 \end{aligned}$$

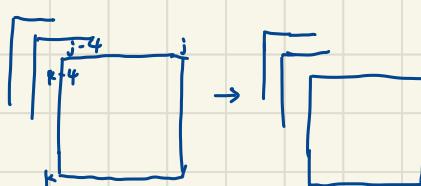


$$\frac{\partial a_{j+1, k+m}^1}{\partial a_{jk}^0} = \sigma' \cdot w_{l,m}^{1,i} \quad \text{Layer 1 convolution}_i$$

$l \in [-4, 0]$
 $m \in [-4, 0]$



$$\begin{aligned} \delta_{jk}^0 \\ \frac{\partial C}{\partial z_{jk}^0} = \frac{1}{2} \sum_{l=0}^4 \frac{\partial C}{\partial z_{j+l, k+m}^1} \cdot \frac{\partial z_{j+l, k+m}^1}{\partial z_{jk}^0} \\ = \sum_{l=0}^4 \sum_{m=0}^4 \delta_{j+l, k+m}^1 \cdot w_{l,m} \\ = \sum_{l=0}^4 [\delta^l] \cdot [w_{l,m}^{1,i}] \end{aligned}$$



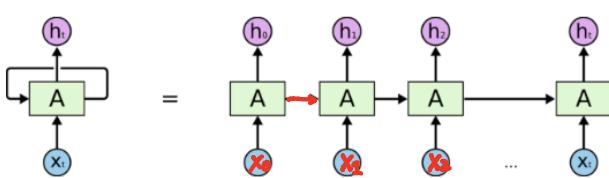
$$\begin{aligned} \frac{\partial C}{\partial z_{jk}^1} &= \frac{\partial C}{\partial z_{j,k}^1} \cdot \frac{\partial z_{j,k}^1}{\partial z_{jk}^0} \\ &= \delta^1 \cdot 1 \cdot f(\text{forward}) \cdot 0 \\ \frac{\partial C}{\partial z_{jk}^L} &= \sum_{m=0}^L \frac{\partial C}{\partial z_{j,k}^m} \cdot \frac{\partial z_{j,k}^m}{\partial z_{jk}^0} \\ &= [\delta^L] \cdot [w_{j,k,m}^{1,i}] \end{aligned}$$



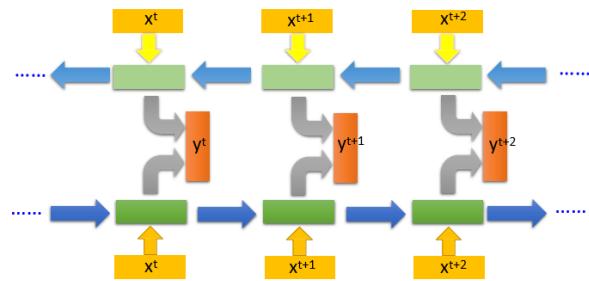
$$\begin{aligned} \frac{\partial C}{\partial z_{jk}^L} &= \frac{\partial C}{\partial a_{jk}^0} \cdot \frac{\partial a_{jk}^0}{\partial z_{jk}^L} \\ &= \nabla_a C \cdot \sigma'(z_{jk}^L) \end{aligned}$$

$$\begin{aligned} \frac{\partial C}{\partial w_{l,m}^{1,i}} &= \sum_{j,k} \frac{\partial C}{\partial z_{jk}^1} \cdot \frac{\partial z_{jk}^1}{\partial w_{l,m}^{1,i}} \\ &= \sum_{j,k} \delta_{j,k}^1 \cdot \theta_{j-l, k-m}^0 \end{aligned}$$

RNN

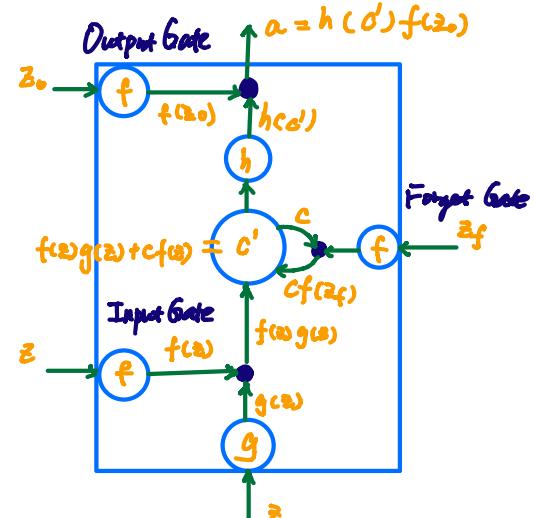
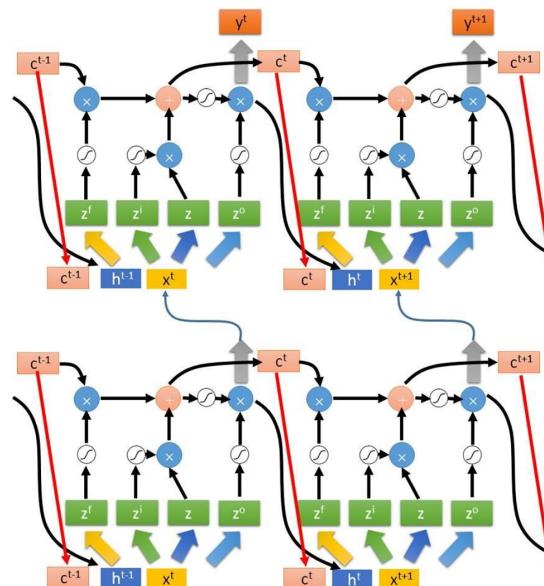
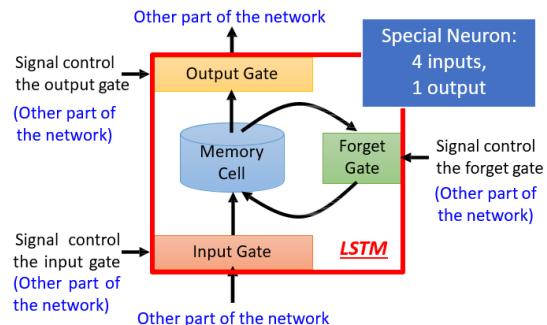


1. Bi-directional RNN

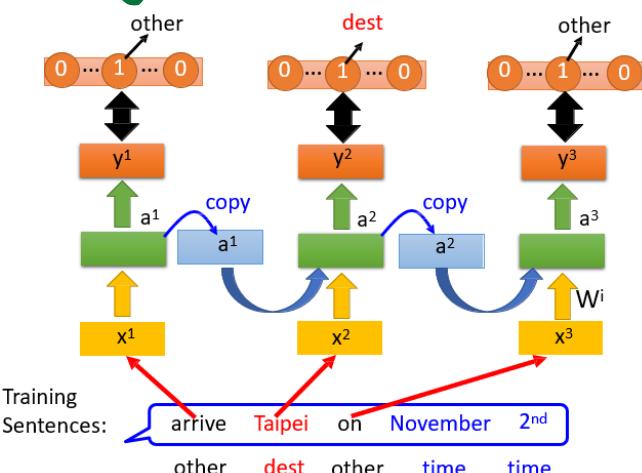


2. Long Short-Term Memory

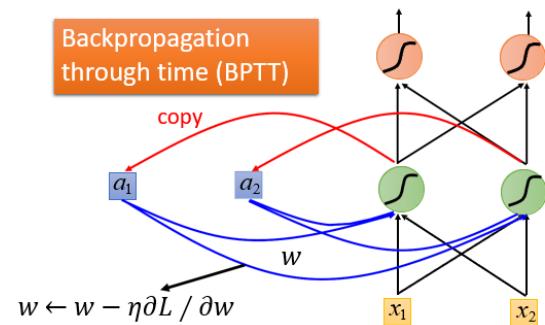
LSTM



Training



Learning



RNN Learning is very difficult in practice.