# 机器学习

黄雷

人工智能研究院

huangleiAI@buaa.edu.cn
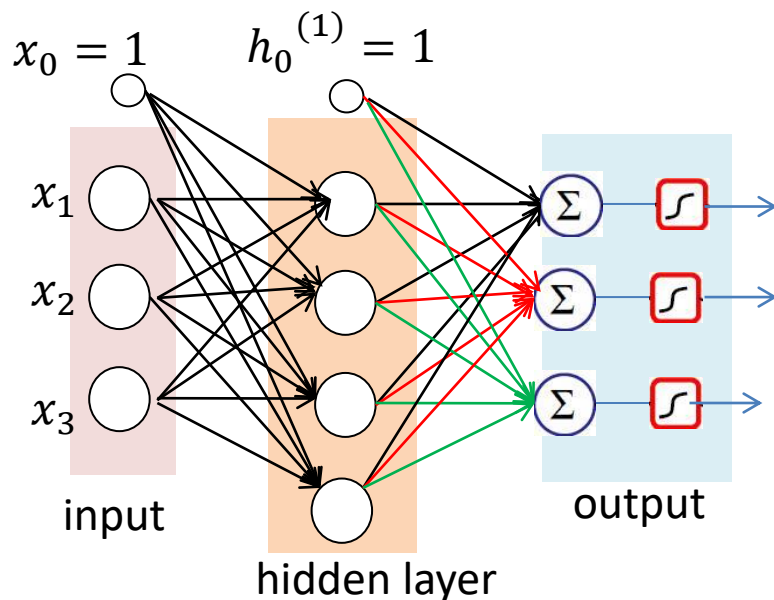
2021年10月25日

# 主要内容

- 多层感知机(Multi-layer Perceptron, MLP)
  - 反向传播算法(Back-propagation)
- 卷积神经网络 （Convolutional Neural Network, CNN）
  - 卷积操作和卷积层
  - 池化 （Pooling）
- 循环神经网络（Recurrent Neural Network, RNN）
  - 建模和训练
  - LSTM模型

# Multi-layer perceptron (多层感知机)

- **Multi-layer perceptron or feed-forward neural network**



$x_i$: 第$i$ 个输入节点

$h_i^{(k)}$: 第 $k$ 层隐藏层的第$i$个节点

$w_{ij}^{(k)}$: 第 $k$ 层隐藏层, 第 i 个输出神经元, 连接第$j$ 个输入神经元

$y_i$: 第$i$ 个输出节点

$x_0 = 1$  $h_0^{(1)} = 1$

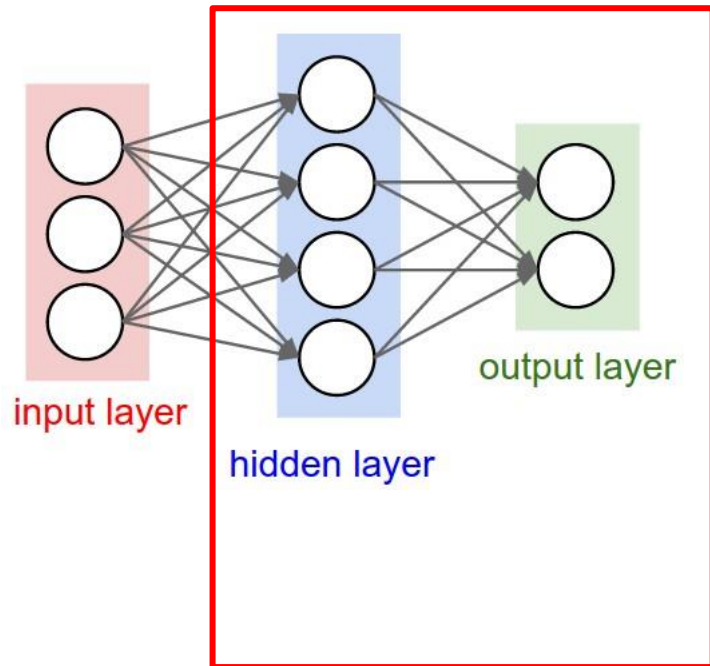input

hidden layer

output

预激活值：
（Pre-activation）

激活值
（Activation）

$$a^{(1)} = W^{(1)} \cdot x$$
$$h^{(1)} = \sigma(a^{(1)})$$

$$a^{(2)} = W^{(2)} \cdot h^{(1)}$$
$$y = \sigma(a^{(2)})$$

$$a^{(i)} = W^{(i)} \cdot h^{(i-1)}$$
$$h^{(i)} = \sigma(a)$$
$$(h^{(0)} = x, h^{(L)} = y)$$

# Neural Networks: Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

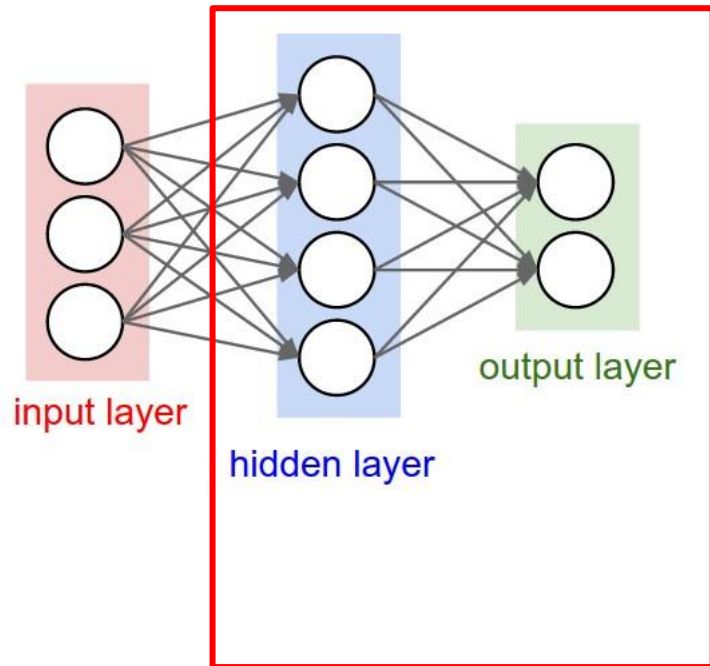Source: Andrej Karpathy & Fei-Fei Li

# Neural Networks: Architectures



Number of Neurons: ?
Number of Weights: ?
Number of Parameters: ?

# Neural Networks: Architectures



Number of Neurons: 4+2 = 6
Number of Weights: [4x3 + 2x4] = 20
Number of Parameters: 20 + 6 = 26 (biases!)

Number of Neurons: ?
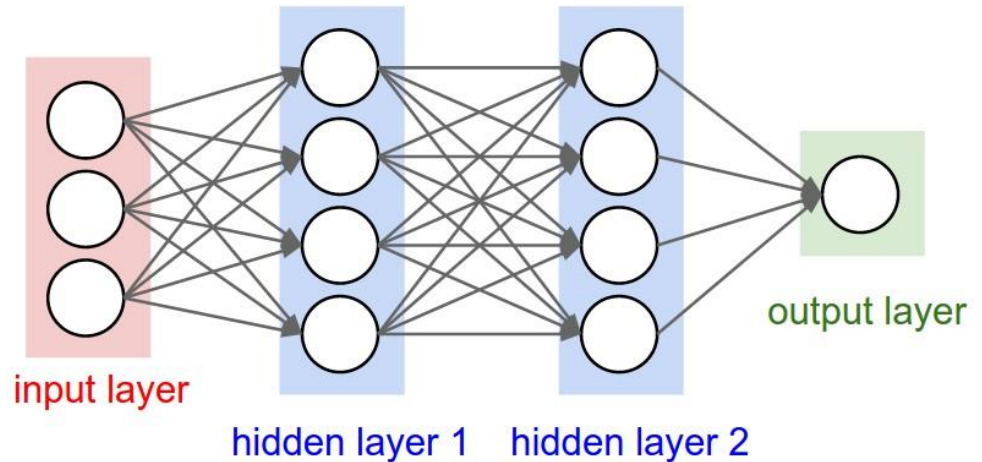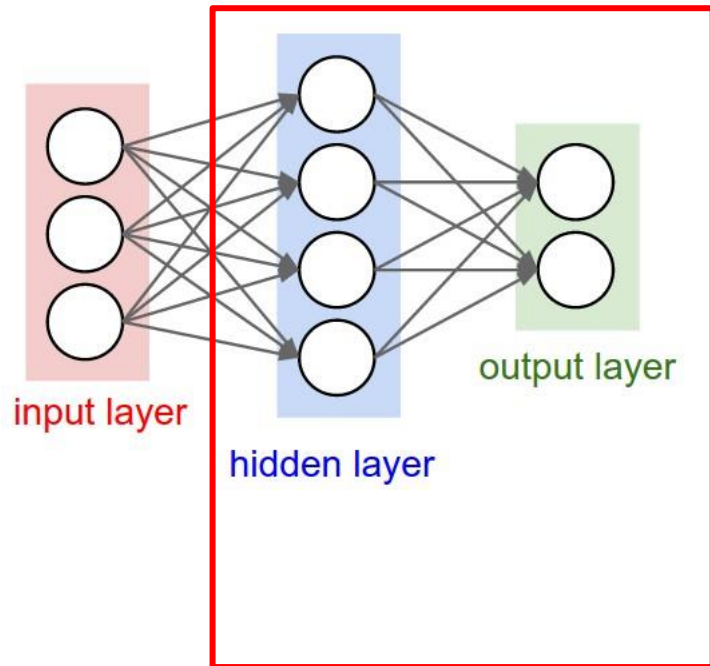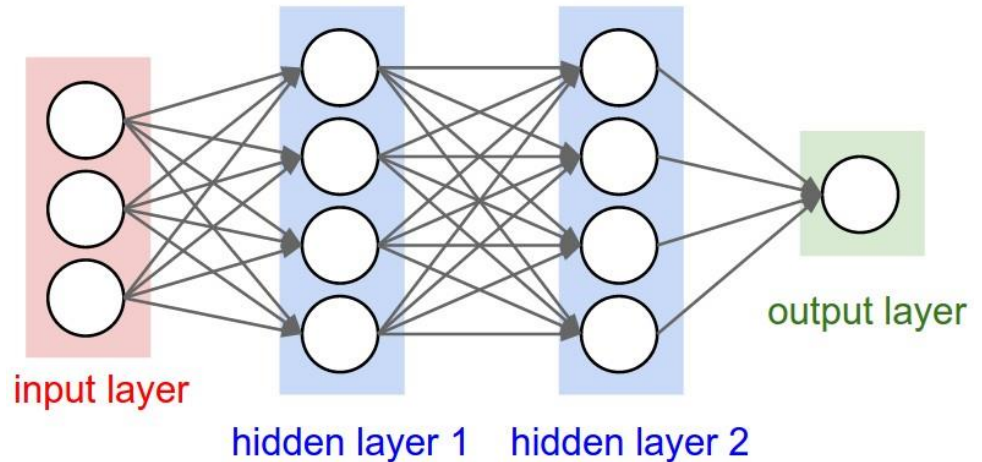Number of Weights: ?
Number of Parameters: ?

# Neural Networks: Architectures


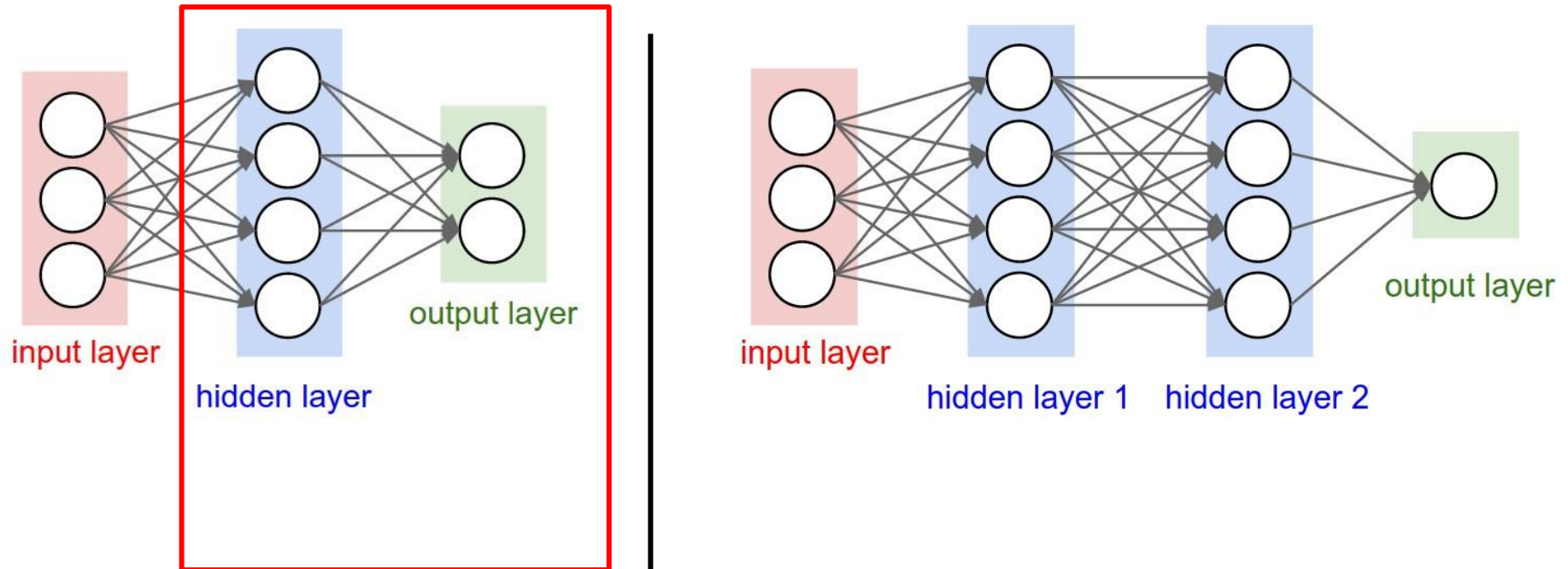
Number of Neurons: 4+2 = 6
Number of Weights: [4x3 + 2x4] = 20
Number of Parameters: 20 + 6 = 26 (biases!)

Number of Neurons: 4 + 4 + 1 = 9
Number of Weights: [4x3+4x4+1x4]=32
Number of Parameters: 32+9 = 41

Source: Andrej Karpathy & Fei-Fei Li

# Neural Networks: Architectures



Modern CNNs: ~10 million neurons
Human visual cortex: ~5 billion neurons

Source: Andrej Karpathy & Fei-Fei Li

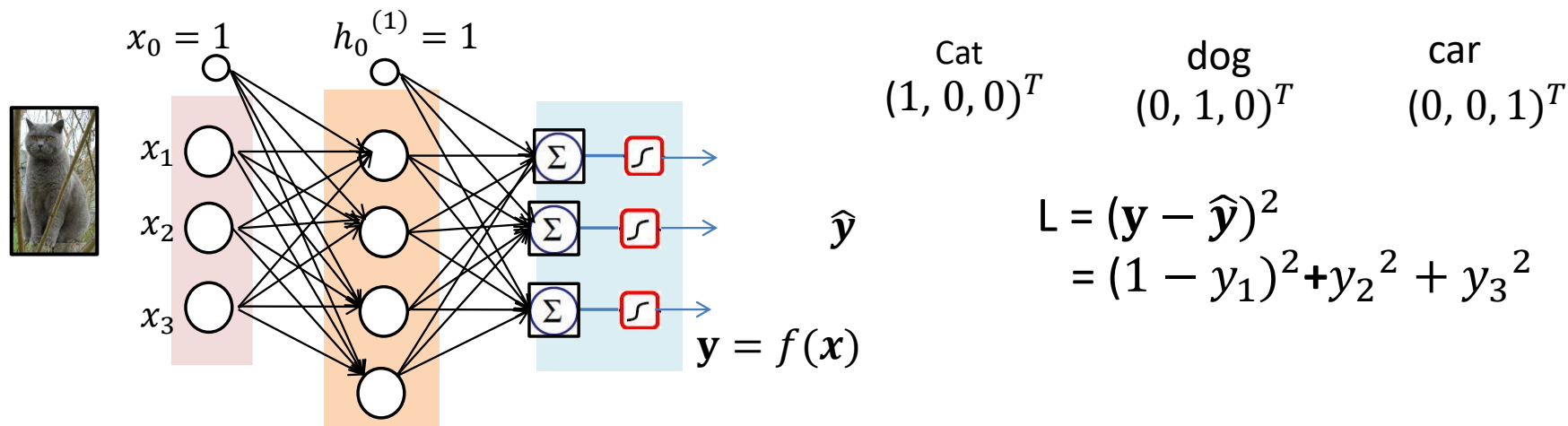# outline

- Multi-layer perceptron (多层感知机)
  - Model representation (模型表示)
  - <span style="color:red">Loss function: the goal for learning</span>
  - Training
    - Gradient based optimization
    - Back-propagation

# Target of learning: Loss function

- 损失函数（Loss function）



$x_0 = 1$   $h_0^{(1)} = 1$

$x_1$   $x_2$   $x_3$

$\hat{y}$

$\mathbf{y} = f(\mathbf{x})$

Cat
$(1, 0, 0)^T$

dog
$(0, 1, 0)^T$

car
$(0, 0, 1)^T$

$L = (\mathbf{y} - \hat{\mathbf{y}})^2$
$= (1 - y_1)^2 + y_2{}^2 + y_3{}^2$

均方误差（Mean Squared Error）：$L = (\mathbf{y} - \hat{\mathbf{y}})^2$，其中 $\mathbf{y} = f(\mathbf{x})$

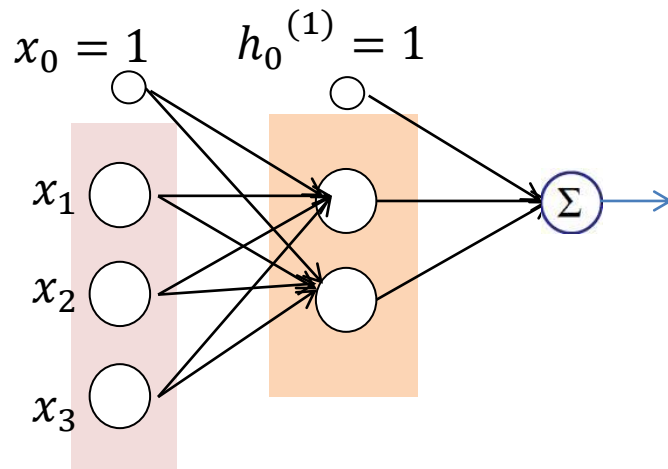- 优化目标函数: min $L = (\mathbf{y} - \hat{\mathbf{y}})^2$

# outline

- Multi-layer perceptron (多层感知机)
  - Model representation (模型表示)
  - Loss function: the goal for learning
  - Training
    - Gradient based optimization
    - Back-propagation

# 目标函数

- 求解目标函数

$$\min_{W} E_{(X,Y) \sim D} \ L(F(X;W), \hat{Y}))$$

- 方案一：令 $\dfrac{\partial L}{\partial W} = 0$，求解方程组

$x_0 = 1$     $h_0{}^{(1)} = 1$

$x_1$

$x_2$

$x_3$

$\Sigma$

$$L= (\frac{w_{21}^{(2)}}{1+e^{-(x_1 w_{11}^{(1)} + x_2 w_{12}^{(1)} + x_3 w_{13}^{(1)} + w_{10}^{(1)})}} + \frac{w_{22}^{(2)}}{1+e^{-(x_1 w_{11}^{(1)} + x_2 w_{12}^{(1)} + x_3 w_{13}^{(1)} + w_{10}^{(1)})}} + w_{20}^{(2)} - \hat{y})^2$$
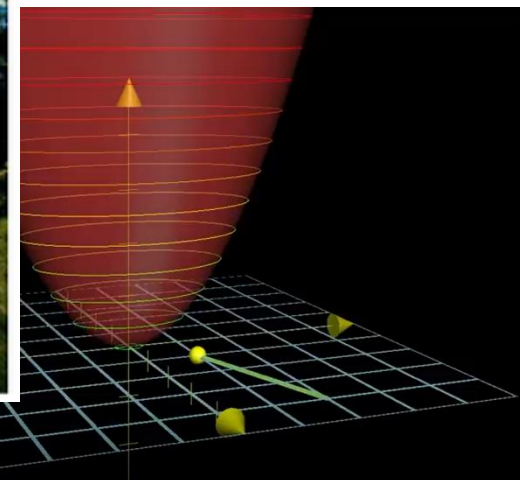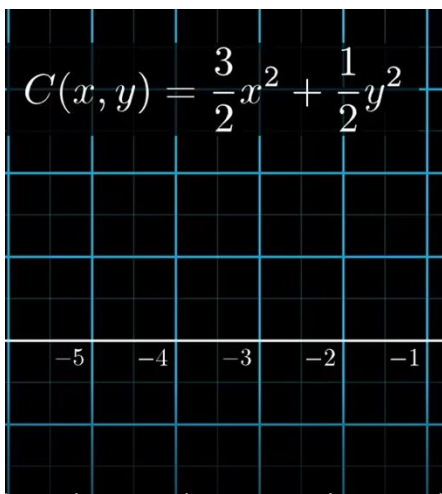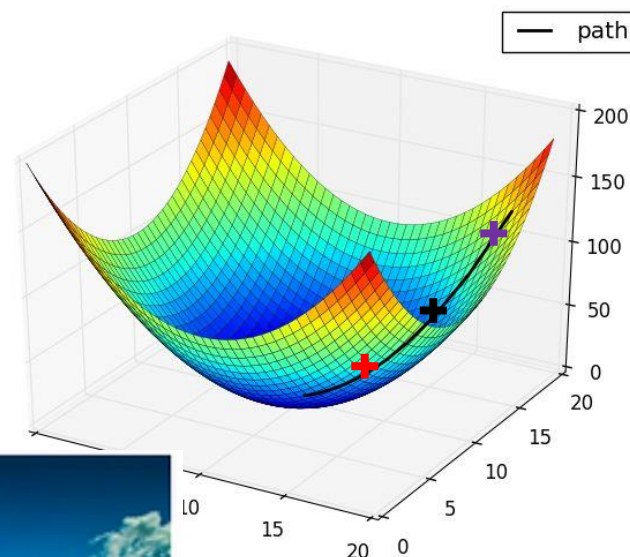
# 目标函数

- 求解目标函数

$$\min_{W} E_{(X,Y)\sim D} \ L(F(X;W),\hat{Y}))$$

# 基于迭代的训练方式

$$\min_{W} E_{(X,Y)\sim D} \ L(F(X;W),\hat{Y}))$$

- 局部下降搜索
  - 基于目前的参数$W^t$，给其多个扰动$\Delta W$ ，确保存在某个$\Delta W$，使得$L(W^t + \Delta W) < L(W^t),$
  - 更新 $W^{t+1} = W^t + \Delta W$

- 更高效的下降搜
  - 基于梯度的下

$$C(x,y) = \frac{3}{2}x^2 + \frac{1}{2}y^2$$

# 梯度下降算法

➤ 0.初始化权重　　$\boldsymbol{W}^{(0)}$
➤ 1. 前向过程：
　　➤ 1.1根据输入，计算输出值 $\boldsymbol{y}$
　　➤ 1.2.计算损失函数值$L(\boldsymbol{y}, \hat{\boldsymbol{y}})$。
➤ 2.计算梯度$\dfrac{d\,L}{d\,\boldsymbol{W}}$
➤ 3.更新梯度

$$\boldsymbol{W}^{(t+1)} = \boldsymbol{W}^{(t)} - \eta\,\dfrac{d\,\mathrm{L}}{d\,\boldsymbol{W}^{(t)}}$$



$L(w_1, w_2)$

$w_2$

$w_1$

gradient：$(\dfrac{\mathrm{dL}(w_1,w_2)}{\mathrm{w}_1}, \dfrac{\mathrm{dL}(w_1,w_2)}{w_2})$

# outline

- Multi-layer perceptron (多层感知机)
  - Model representation (模型表示)
  - Loss function: the goal for learning
  - Training
    - Gradient based optimization
    - Back-propagation

# 计算梯度：反向传播

- 求导基础知识回顾

  ➢实值函数对一维实值变量的导数:

  $$\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

  ➢实值函数对多维向量变量的梯度为向量（偏导数）:

  $$\frac{\partial \mathbf{f}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = (\frac{\partial f(\theta_0, \theta_1)}{\partial \theta_0}, \frac{\partial \mathbf{f}(\theta_0, \theta_1)}{\partial \theta_1}), \quad \boldsymbol{\theta} = (\theta_0, \theta_1)$$

# 计算梯度：反向传播

- 神经网络中的基本操作

加法： $f(x,y)=x+y$    $\dfrac{\partial f}{\partial x}=1$    $\dfrac{\partial f}{\partial y}=1$

乘法： $f(x,y)=xy$    $\dfrac{\partial f}{\partial x}=y$    $\dfrac{\partial f}{\partial y}=x$

非线性变换： $\sigma(x)=\dfrac{1}{1+e^{-x}}$

$$\frac{d\sigma(x)}{dx}=\frac{e^{-x}}{(1+e^{-x})^2}=\left(\frac{1+e^{-x}-1}{1+e^{-x}}\right)\left(\frac{1}{1+e^{-x}}\right)=(1-\sigma(x))\,\sigma(x)$$
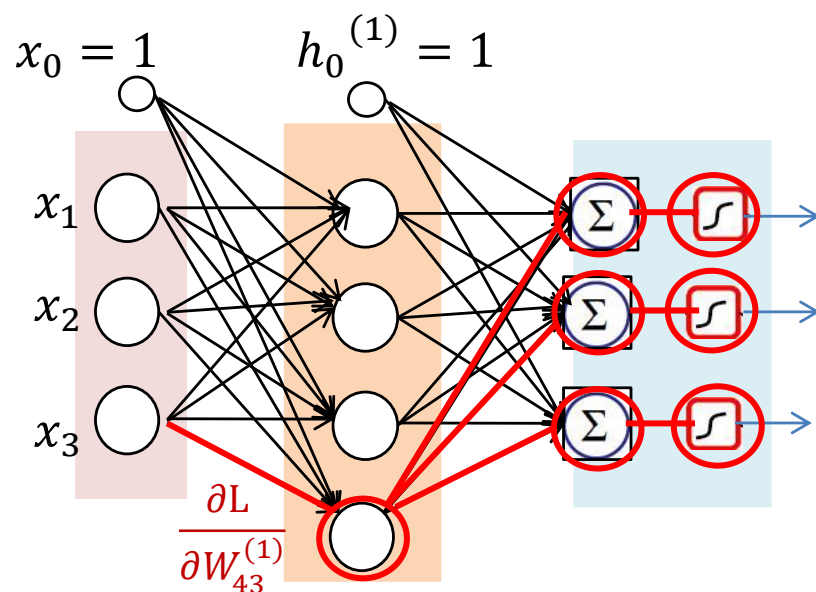
# 反向传播（Back-Propagation）

- 链式法则（Chain rule）

$$f=q(x) \qquad \frac{\partial f}{\partial x} = \frac{\partial f}{\partial q}\frac{\partial q}{\partial x}$$
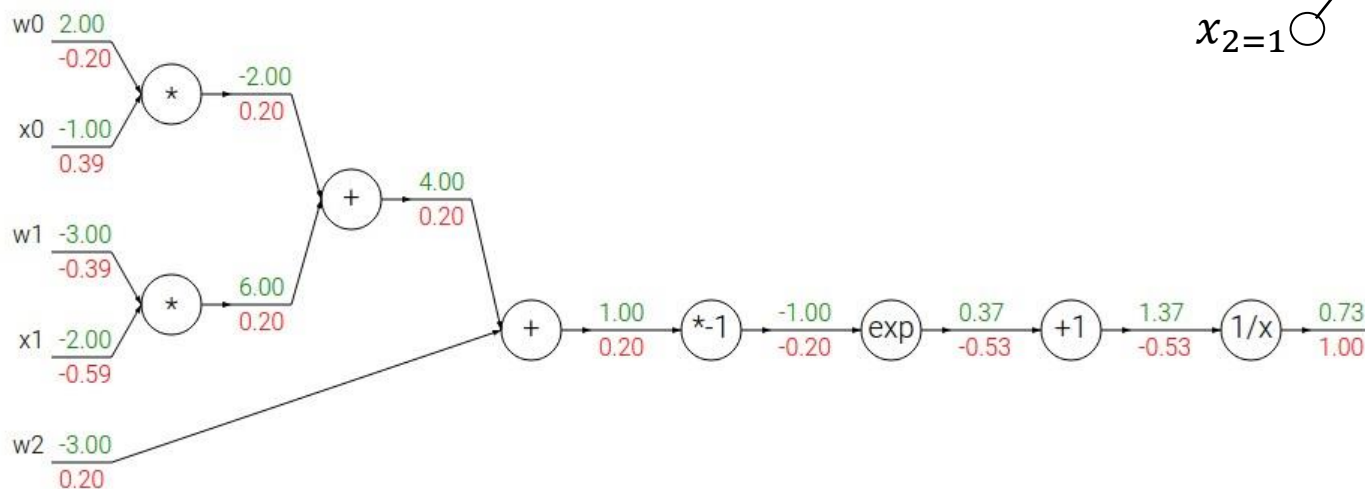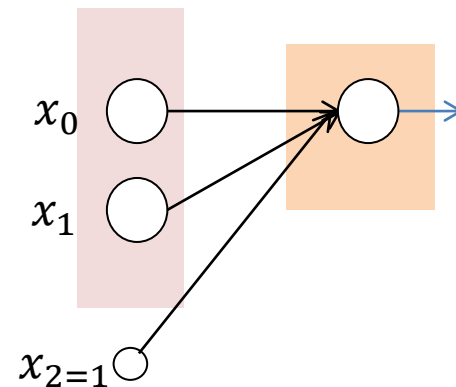
复合表达式: $f(x, y, z) = (x + y)z$

$$q=x+y \qquad \frac{\partial q}{\partial x} =1 \qquad \frac{\partial q}{\partial y} =1$$

$$f=qz \qquad \frac{\partial f}{\partial q} = z \qquad \frac{\partial f}{\partial z} = q$$
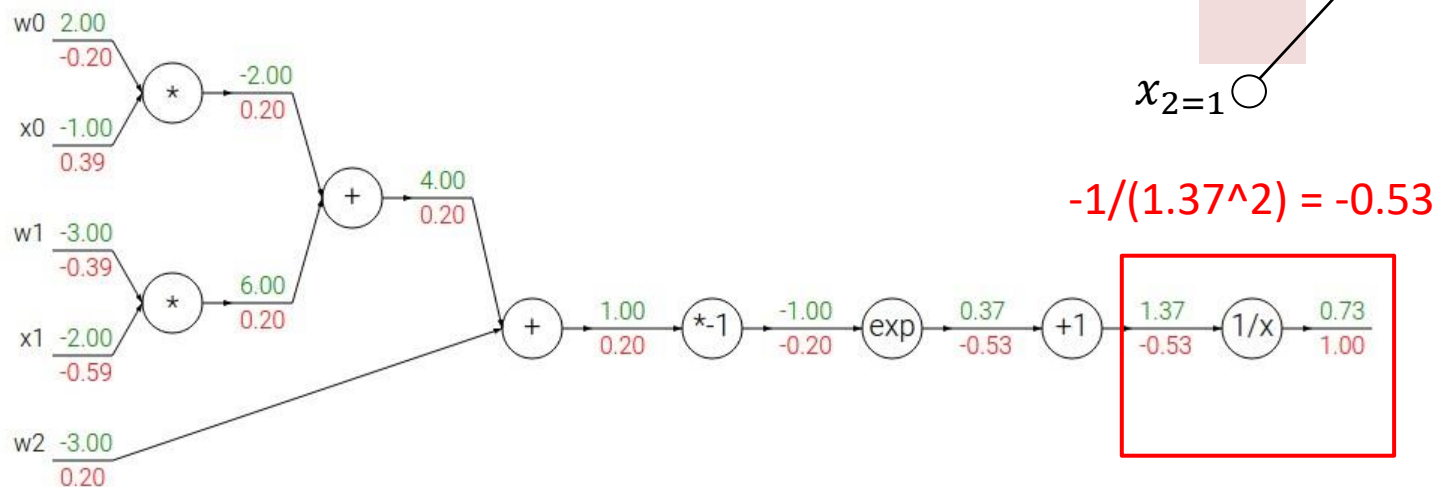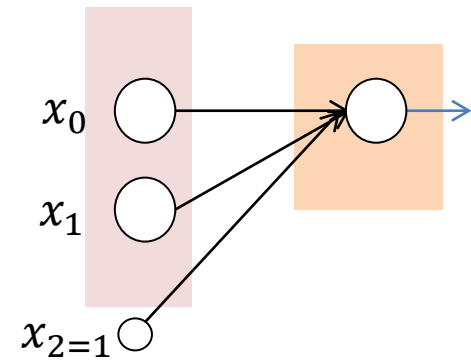
One example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Source: Andrej Karpathy & Fei-Fei Li

One example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$x_0$

$x_1$

$x_{2=1}$

-1/(1.37^2) = -0.53

w0 2.00 / -0.20

x0 -1.00 / 0.39

* → -2.00 / 0.20

w1 -3.00 / -0.39

x1 -2.00 / -0.59

* → 6.00 / 0.20

+ → 4.00 / 0.20

w2 -3.00 / 0.20

+ → 1.00 / 0.20 → *-1 → -1.00 / -0.20 → exp → 0.37 / -0.53 → +1 → 1.37 / -0.53 → 1/x → 0.73 / 1.00

$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x \qquad \Big| \qquad f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

Source: Andrej Karpathy & Fei-Fei Li

Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$x_0$

$x_1$

$x_{2=1}$

[local gradient] x [its gradient]

[1] x [-0.53] = -0.53

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Source: Andrej Karpathy & Fei-Fei Li

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[local gradient] x [its gradient]
[e^(-1)] x [-0.53] = -0.20

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Big| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Source: Andrej Karpathy & Fei-Fei Li

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[local gradient] x [its gradient]
[-1] x [-0.2] = 0.2

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Source: Andrej Karpathy & Fei-Fei Li

Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$x_0$

$x_1$

$x_{2=1}$

[local gradient] x [its gradient]
[1] x [0.2] = 0.2
[1] x [0.2] = 0.2  (both inputs!)

w0 2.00
-0.20

x0 -1.00
0.39

-2.00
0.20

w1 -3.00
-0.39

x1 -2.00
-0.59

4.00
0.20

6.00
0.20

w2 -3.00
0.20

1.00
0.20

*-1

-1.00
-0.20

exp

0.37
-0.53

+1

1.37
-0.53

1/x

0.73
1.00

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Source: Andrej Karpathy & Fei-Fei Li

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[local gradient] x [its gradient]
x0: [2] x [0.2] ~= 0.4
w0: [-1] x [0.2] = -0.2

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Source: Andrej Karpathy & Fei-Fei Li

# 反向传播

- 一层神经网络（线性模型）



输入: $\boldsymbol{x}$

输出: $\mathbf{y}$

$(1, 0, 0)^T$

➤ 1.给定输入，计算输出值：

$$a_i = \sum_{j=0}^{3} w_{ij} x_i = \boldsymbol{w}_{i.} \cdot \boldsymbol{x},$$
$$i = (1,2,3)$$
$$y_i = \sigma(a_i) = \frac{1}{1 + e^{-a_i}}$$

MSE Loss: L=$(\mathbf{y} - \widehat{\boldsymbol{y}})^2$ =$(1 - y_1)^2 + y_2^2 + y_3^2$

➤ 2.根据链规则，计算梯度$\frac{\partial L}{\partial \boldsymbol{w}}$：

$$\frac{\partial L}{\partial y_1} = 2(y_1\text{-}1)$$

$$\frac{\partial L}{\partial y_i} = 2y_i, \text{(i=2,3)}$$

$$\frac{\partial L}{\partial a_i} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial a_i} = \frac{\partial L}{\partial y_i} \sigma(a_i)(1\text{-}\sigma(a_i))$$

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial L}{\partial a_i} x_{ij}$$
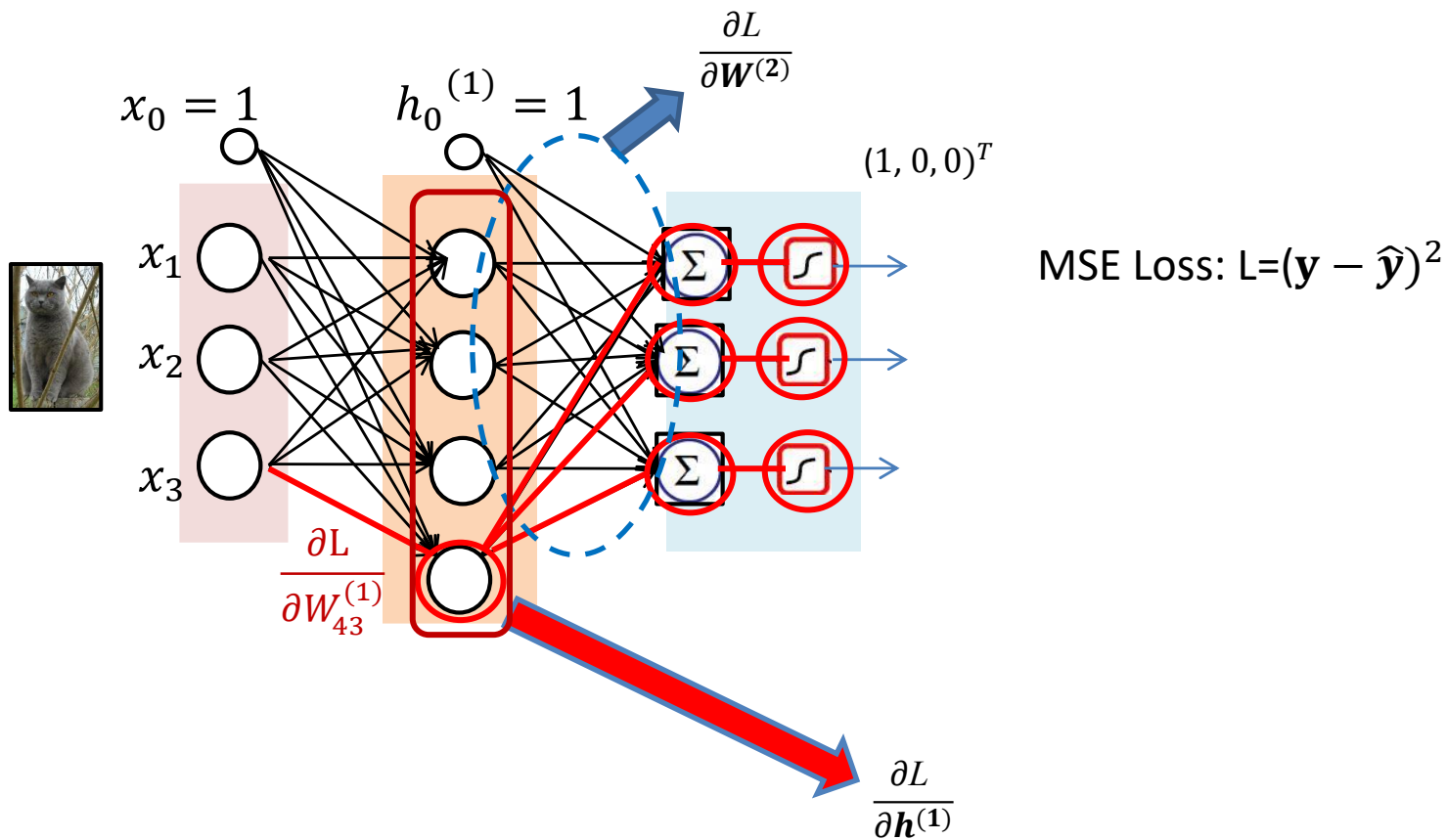$$= \frac{\partial L}{\partial y_i} \sigma(a_i)(1\text{-}\sigma(a_i))$$

$$\frac{\partial L}{\partial \mathbf{y}} = 2(\mathbf{y} - \widehat{\boldsymbol{y}})$$

$$\frac{\partial L}{\partial \boldsymbol{a}} = 2[(\mathbf{y} - \widehat{\boldsymbol{y}}) \cdot \sigma(\boldsymbol{a}) \cdot (1 - \sigma(\boldsymbol{a}))]^T$$

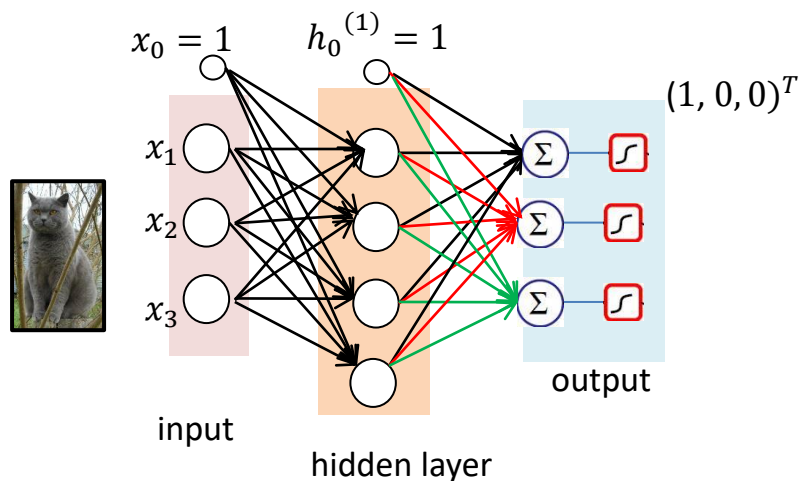$$\frac{\partial L}{\partial \boldsymbol{W}} = 2[(\mathbf{y} - \widehat{\boldsymbol{y}}) \cdot \sigma(\boldsymbol{a}) \cdot (1 - \sigma(\boldsymbol{a}))] \, \boldsymbol{x}^T$$

# 反向传播

- 两层的网络



$$\frac{\partial L}{\partial \boldsymbol{W^{(2)}}}$$

$x_0 = 1$

$h_0^{(1)} = 1$

$(1, 0, 0)^T$

MSE Loss: L=$(\boldsymbol{y} - \widehat{\boldsymbol{y}})^2$

$x_1$

$x_2$

$x_3$

$$\frac{\partial \text{L}}{\partial W_{43}^{(1)}}$$

$$\frac{\partial L}{\partial \boldsymbol{h^{(1)}}}$$

# 反向传播

- 一层神经网络（线性模型）



$x_0 = 1$

$x_1$

$x_2$

$x_3$

$(1, 0, 0)^T$

输入:$\boldsymbol{x}$　　　输出:$\mathbf{y}$

➤ 1.给定输入，计算输出值：

$$a_i = \sum_{j=0}^{3} w_{ij} x_i = \boldsymbol{w}_{i\cdot} \cdot \boldsymbol{x}, \qquad i = (1,2,3)$$

$$y_i = \sigma(a_i) = \frac{1}{1 + e^{-a_i}}$$

MSE Loss: L=$(\mathbf{y} - \widehat{\boldsymbol{y}})^2$ =$(1 - y_1)^2 + y_2{}^2 + y_3{}^2$

➤ 2.根据链规则，计算梯度$\frac{\partial L}{\partial \boldsymbol{x}}$：

$$\frac{\partial L}{\partial y_1} = 2(1 - y_1)$$

$$\frac{\partial L}{\partial y_i} = 2y_i, \text{(i=2,3)}$$

$$\frac{\partial L}{\partial a_i} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial a_i} = \frac{\partial L}{\partial y_i} \sigma(a_i)(1-\sigma(a_i))$$

$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial a_i} \frac{\partial a_i}{\partial x_i} = \frac{\partial L}{\partial a_i} \sum_{j=0}^{3} w_{ij}$$

$$\frac{\partial L}{\partial \boldsymbol{x}} = \frac{\partial L}{\partial \boldsymbol{a}} \boldsymbol{W}$$

# 反向传播

- 两层的网络



$x_0 = 1$    $h_0^{(1)} = 1$

$(1, 0, 0)^T$

$x_1$

$x_2$

$x_3$

output

input

hidden layer

➢1给定输入，计算输出值：

$$a^{(1)} = W^{(1)} \cdot x$$
$$h^{(1)} = \sigma\big(a^{(1)}\big)$$
$$a^{(2)} = W^{(2)} \cdot h^{(1)}$$
$$y = \sigma\big(a^{(2)}\big)$$

➢MSE Loss: L=$(y - \widehat{y})^2$

➢2根据链规则，计算梯度$\frac{\partial L}{\partial a^{(i)}}$ ,$\frac{\partial L}{\partial x}$ ：

$$\frac{\partial L}{\partial y} = 2(y - \widehat{y})$$

$$\frac{\partial L}{\partial a^{(2)}} = \frac{\partial L}{\partial y} \cdot \sigma\big(a^{(2)}\big) \cdot \big(1 - \sigma\big(a^{(2)}\big)\big)$$

$$\frac{\partial L}{\partial h^{(1)}} = \frac{\partial L}{\partial a^{(2)}} W^{(2)}$$

$$\frac{\partial L}{\partial a^{(1)}} = \frac{\partial L}{\partial h^{(1)}} \cdot \sigma\big(a^{(1)}\big) \cdot \big(1 - \sigma\big(a^{(1)}\big)\big)$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial a^{(1)}} W^{(1)}$$

➢3.根据链规则，计算梯度$\frac{\partial L}{\partial W^{(i)}}$ ：

$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial a^{(2)}} h^{(1)}$$

$$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial a^{(1)}} x$$

# 反向传播

- 利用链式法则计算梯度
- 利用了动态规划的思想

前向计算：

$x \rightarrow$ $\boxed{\boldsymbol{h}_1 = W_1 \boldsymbol{x}_1}$ $\xrightarrow{\boldsymbol{x}_2 = f_1(\boldsymbol{h}_1)}$ $\cdots$ $\boxed{\boldsymbol{h}_l = W_l \boldsymbol{x}_l}$ $\cdots$ $\rightarrow$ $y$

反向传播：

$\boxed{\frac{\partial L}{\partial x_1} = W_1^T \frac{\partial L}{\partial h_1}}$ $\leftarrow$ $\cdots$ $\boxed{\frac{\partial L}{\partial \boldsymbol{x}_l} = W_l^T \frac{\partial L}{\partial \boldsymbol{h}_l}}$ $\cdots$ $\leftarrow$ $\hat{y}$

$$\frac{\partial L}{\partial W_l} = \sum (\boldsymbol{x}_l \frac{\partial L}{\partial \boldsymbol{h}_l})^T$$

Loss

# 前馈神经网络梯度下降训练算法

➢0.初始化权重 $W^{(0)}$

➢1. 前向过程：
  ➢1.1根据输入，计算输出值 $y$
  ➢1.2.计算损失函数值L$(\mathbf{y}, \hat{\mathbf{y}})$。

➢2.反向传播

  ➢计算$\dfrac{\partial \mathrm{L}}{\partial \mathbf{y}}$

  ➢后向传播直到计算$\dfrac{\partial \mathrm{L}}{\partial \mathbf{x}}$

➢3.计算梯度$\dfrac{\partial L}{\partial W}$

➢4.更新梯度

$$W_{t+1} = W_t - \eta \frac{\partial \mathrm{L}}{\partial w_t}$$



$x_0 = 1$   $h_0^{(1)} = 1$

$x_1$

$x_2$

$x_3$

$(1, 0, 0)^T$

Input:

hidden layer

output

# 课外拓展研究

- 对抗样例（Adversarial example）



$f(X + \Delta X) \neq f(X)$, $\Delta X$ should be imperceptible by human

X ΔX $X + \Delta X$

+ =

Panda Car

# 主要内容

- 多层感知机(Multi-layer Perceptron, MLP)
  - 反向传播算法(Back-propagation)
- 卷积神经网络 （Convolutional Neural Network, CNN）
  - 卷积操作和卷积层
  - 池化 （Pooling）
- 循环神经网络（Recurrent Neural Network, RNN）
  - 建模和训练
  - LSTM模型

# Feature extraction

- Feature extraction
  - Pixel-wise input
  - Correlation between features



$x_0 = 1$

$h_0^{(1)} = 1$

$(1, 0, 0)^T$

Input:

hidden layer

output

$x_1$

$x_2$

$x_3$

$(x_1, x_2, x_3)$

Convolutional Neural Network(CNN),卷积神经网络

What the computer sees

# Convolution Neural Network

- Lenet-5



Convolution related layers

全连接层

# outline

- Convolutional layer (module)
  - <span style="color:red">Convolution operation</span>
  - Filters
  - Convolution module in a network
- Pooling layer (module)

# Convolution (卷积)

➢ 一维相关操作（correlation）例子

$y$= | $y_1$ | $y_2$ |

$$y_1 = w_1 x_1 + w_2 x_2$$

$w$= | $w_1$ | $w_2$ |

$x$= | $x_1$ | $x_2$ | $x_3$ |

$$y_2 = w_1 x_2 + w_2 x_3$$

Correlation operator (similarity)

$$y_{i'} = \sum_{i=1}^{M_f=2} w_i x_{i'+i-1}$$

➢ 一维卷积（correlation）例子

➢ Flip

$\overline{w}$= | $w_2$ | $w_1$ |

$$y_{i'} = \sum_{i=1}^{M_f=2} w_{M_f+1-i} x_{i'+i-1}$$

Convolution

# Convolution (卷积)

- 连续空间的卷积:

$$y(t) = x(t) * h(t) = \int_{-\infty}^{+\infty} x(s)h(t-s)\, ds$$

- 离散空间卷积：

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

- 图像卷积是二维离散卷积

$$g(i,j) = \sum_{k,l} f(k,l)\, w(i-k, j-l)$$



kernel

input

output

# Convolution （卷积）

- ## 图像卷积，二维, 离散
  - ### Correlation Operator(相关算子)

定义：g = f⊗$w$

$$g(i,j) = \sum_{k,l} f(i+k, j+l)\, w(k,l)$$

image

Kernel(filters)

**f:**

| 17 | 24 | 1 | 8 | 15 |
|----|----|----|----|----|
| 23 | 5 | 7 | 14 | 16 |
| 4 | 6 | 13 | 20 | 22 |
| 10 | 12 | 19 | 21 | 3 |
| 11 | 18 | 25 | 2 | 9 |

**w:**

| 8 | 1 | 6 |
|---|---|---|
| 3 | 5 | 7 |
| 4 | 9 | 2 |



Values of correlation kernel

Image pixel values

Center of kernel

| 17 | 24 | $1^8$ | $8^1$ | $15^6$ |
|----|----|-------|-------|--------|
| 23 | 5 | $7^3$ | $14^5$ | $16^7$ |
| 4 | 6 | $13^4$ | $20^9$ | $22^2$ |
| 10 | 12 | 19 | 21 | 3 |
| 11 | 18 | 25 | 2 | 9 |

# Convolution （卷积）

- ## 图像卷积，二维，离散
  - ### Convolution operator (卷积算子)

定义：$g = f * w$

$$g(i,j) = \sum_{k,l} f(k,l)\, w(i-k, j-l)$$

image

Kernel(filters)

旋转180

Values of rotated convolution kernel

f:

| 17 | 24 | 1 | 8 | 15 |
|----|----|----|----|----|
| 23 | 5 | 7 | 14 | 16 |
| 4 | 6 | 13 | 20 | 22 |
| 10 | 12 | 19 | 21 | 3 |
| 11 | 18 | 25 | 2 | 9 |

w:

| 8 | 1 | 6 |
|---|---|---|
| 3 | 5 | 7 |
| 4 | 9 | 2 |

Image pixel values

Center of kernel

| 17 | 24 | $1^2$ | $8^9$ | $15^4$ |
|----|----|----|----|----|
| 23 | 5 | $7^7$ | $14^5$ | $16^3$ |
| 4 | 6 | $13^6$ | $20^1$ | $22^8$ |
| 10 | 12 | 19 | 21 | 3 |
| 11 | 18 | 25 | 2 | 9 |

# outline

- Convolutional layer (module)
  - Convolution operation
  - <span style="color:red">Filters</span>
  - Convolution module in a network
- Pooling layer (module)

# Practice with linear filters(线性滤波器)



Original

| 0 | 0 | 0 |
|---|---|---|
| 0 | 1 | 0 |
| 0 | 0 | 0 |

Filter

Filtered
(no change)

# Practice with linear filters



Original

| 0 | 0 | 0 |
|---|---|---|
| 0 | 0 | 1 |
| 0 | 0 | 0 |

Filter



Shifted *left*
By 1 pixel

# Practice with linear filters



$$\frac{1}{9}$$

| 1 | 1 | 1 |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 1 |

Original

Filter

Blur (with a box filter)

# Practice with linear filters

Original

| 0 | 1 | 0 |
|---|----|---|
| 1 | -4 | 1 |
| 0 | 1 | 0 |

Filter

Output Image

**Edge detect** （边缘检测）

# Filters in practise



Input image

filter

output image

- ## Size of output image
  - How to move? stride
  - How about the border? padding

# filters: stride（步幅）



7x7 input
assume 3x3 connectivity, stride 1

# filters: stride （步幅）



7x7 input
assume 3x3 connectivity, stride 1

# filters: stride （步幅）



7x7 input
assume 3x3 connectivity, stride 1

# filters: stride （步幅）

7x7 input
assume 3x3 connectivity, stride 1

# filters: stride （步幅）



7x7 input
assume 3x3 connectivity, stride 1

# filters: stride （步幅）

7x7 input
assume 3x3 connectivity, stride 1
=> **5x5 output**

# filters: stride （步幅）

7x7 input
assume 3x3 connectivity, stride 1
=> **5x5 output**

what about stride 2?

# filters: stride （步幅）



7x7 input
assume 3x3 connectivity, stride 1
=> **5x5 output**

what about stride 2?

# filters: stride （步幅）

7x7 input
assume 3x3 connectivity, stride 1
=> **5x5 output**

what about stride 2?

# filters: stride （步幅）



7x7 input
assume 3x3 connectivity, stride 1
=> **5x5 output**

what about stride 2?
=> **3x3 output**

# filters: stride （步幅）



Output size:
**(N - F) / stride + 1**

e.g. N = 7, F = 3:
stride 1 => (7 - 3)/1 + 1 = 5
stride 2 => (7 - 3)/2 + 1 = 3

# filters: padding

- In practice: Common to zero pad the border



e.g. input 7x7
neuron with receptive field 3x3, stride 1
pad with 1 pixel border => what is the
output?

7x7 => preserved size!

# Filters in practise

- **"Same convolution" (preserves size)**

Input [9x9]

3x3 neurons, stride 1, pad **1** =>
[9x9]

- No headaches when sizing architectures
- Works well

- **"Valid convolution" (shrinks size)**

Input [9x9]

3x3 neurons, stride 1, pad **0** =>
[7x7]

- **Headaches** with sizing the full architecture
- Works Worse!

# outline

- Convolutional layer (module)
  - Convolution operation
  - Filters
  - <span style="color:red">Convolution module in a network</span>
- Pooling layer (module)

# Feature detection (特征检测)

- Learning filters (weights)



one layer

input Image

I to NI convolutions
or correlation
or template matching

NI to NI
subsampling
or pooling
or L norm pool

non-linearity
tanh
etc...

# Convolution Layer (卷积层)

**Input: X**$\in$
$R^{d_{in}\times h\times w}$

**weight: W**$\in$
$R^{d_{out}\times d_{in}\times F_h\times F_w}$

**output: Y**$\in$
$R^{d_{out}\times h\times w}$



channels

Feature maps

---

**Input: x**$\in R^{d_{in}}$

**weight: W**$\in R^{d_{out}\times d_{in}}$

**output: y**$\in R^{d_{out}}$

# Forward (前向过程)

**Input: X**∈
$R^{d_{in} \times h \times w}$

**weight: W**∈
$R^{d_{out} \times d_{in} \times F_h \times F_w}$

**output: Y**∈
$R^{d_{out} \times h \times w}$



$$y_{f',i',j'} = \sum_{f=1}^{d_{in}} \sum_{i=-\frac{F_h}{2}}^{\frac{F_h}{2}} \sum_{j=-\frac{F_w}{2}}^{\frac{F_w}{2}} x_{f,i'+i-1,j'+j-1} \, w_{f',f,i,j} + b_{f'}$$

Feature maps

# example

$$y_{4,10,10} = \sum_{i=-\frac{F_h}{2}}^{\frac{F_h}{2}} \sum_{j=-\frac{F_w}{2}}^{\frac{F_w}{2}} x_{1,10+i-1,10+j-1} \, w_{4,1,i,j} + b_4$$



Input layer    (S1) 4 feature maps    (C1) 4 feature maps    (S2) 6 feature maps    (C2) 6 feature maps

convolution layer    sub-sampling layer    convolution layer    sub-sampling layer    fully connected MLP

$$y_{4,10,100} = \sum_{i=-\frac{F_h}{2}}^{\frac{F_h}{2}} \sum_{j=-\frac{F_w}{2}}^{\frac{F_w}{2}} x_{1,10+i-1,100+j-1} \, w_{4,1,i,j} + \sum_{i=-\frac{F_h}{2}}^{\frac{F_h}{2}} \sum_{j=-\frac{F_w}{2}}^{\frac{F_w}{2}} x_{2,10+i-1,100+j-1} \, w_{4,2,i,j} +$$

$$\sum_{i=-\frac{F_h}{2}}^{\frac{F_h}{2}} \sum_{j=-\frac{F_w}{2}}^{\frac{F_w}{2}} x_{3,10+i-1,100+j-1} \, w_{4,3,i,j} + \sum_{i=-\frac{F_h}{2}}^{\frac{F_h}{2}} \sum_{j=-\frac{F_w}{2}}^{\frac{F_w}{2}} x_{4,10+i-1,100+j-1} \, w_{4,4,i,j} + b_4$$

# A toy ConvNet: X's and O's

Says whether a picture is of an X or an O

A two-dimensional
array of pixels

CNN → **X** or **O**

# For example

# Trickier cases



translation     scaling     rotation     weight

# Deciding is hard

# What computers see

# What computers see

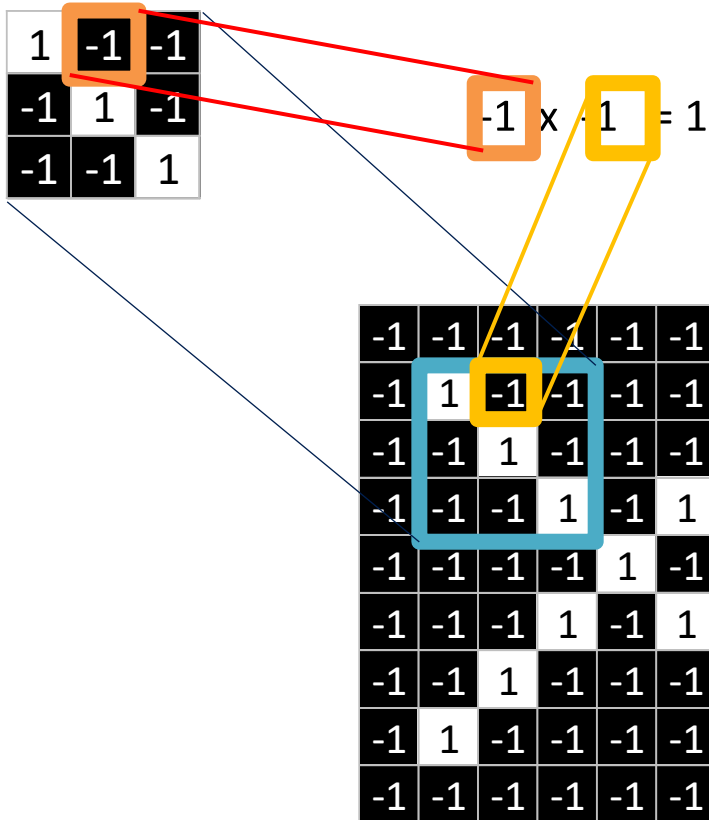# Computers are literal

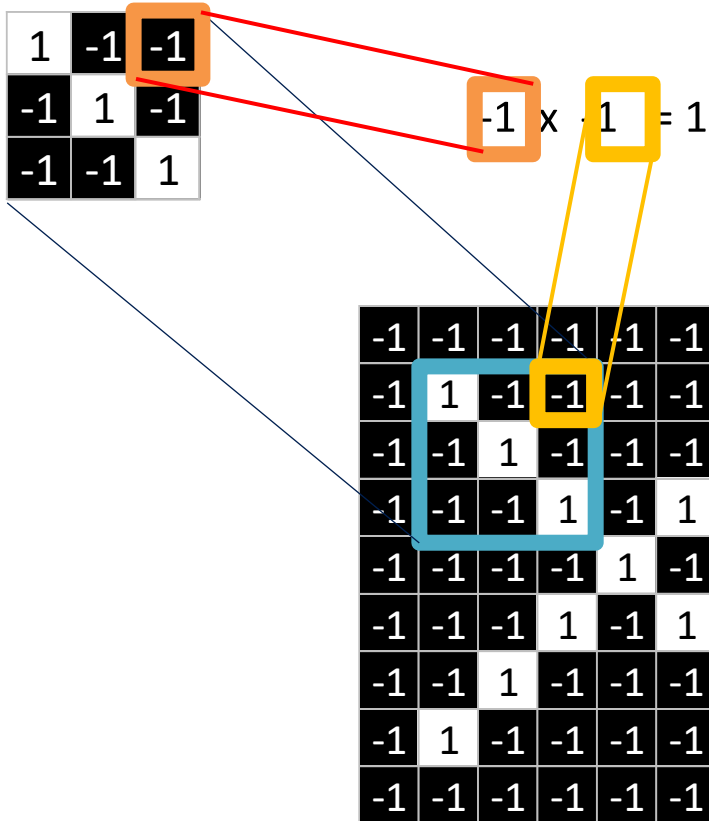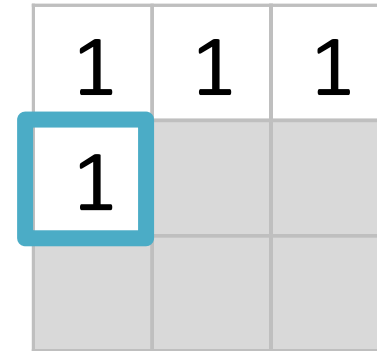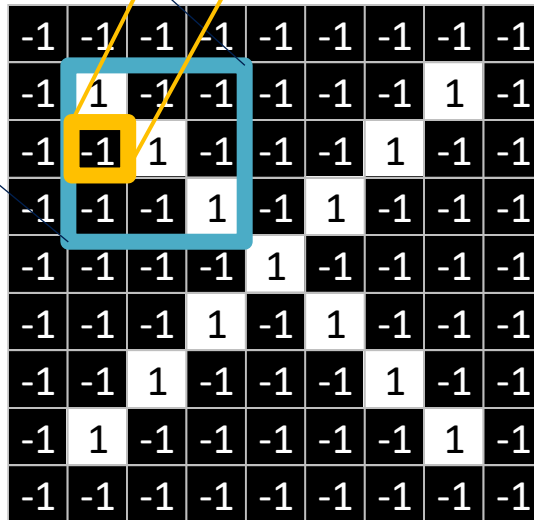# ConvNets match pieces of the image

# Features match pieces of the image

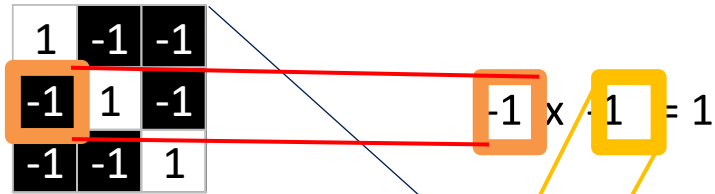# Filtering: The math behind the match
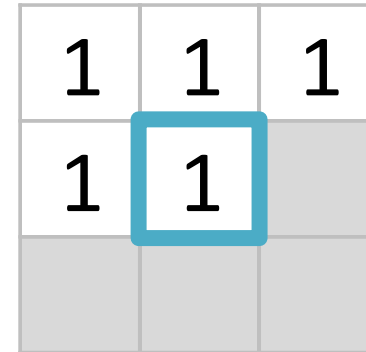
# Filtering: The math behind the match

1. Line up the feature and the image patch.
2. Multiply each image pixel by the corresponding feature pixel.
3. Add them up.
4. Divide by the total number of pixels in the feature.

# Filtering: The math behind the match



$1 \times 1 = 1$

# Filtering: The math behind the match

# Filtering: The math behind the match

# Filtering: The math behind the match

# Filtering: The math behind the match

# Filtering: The math behind the match

# Filtering: The math behind the match

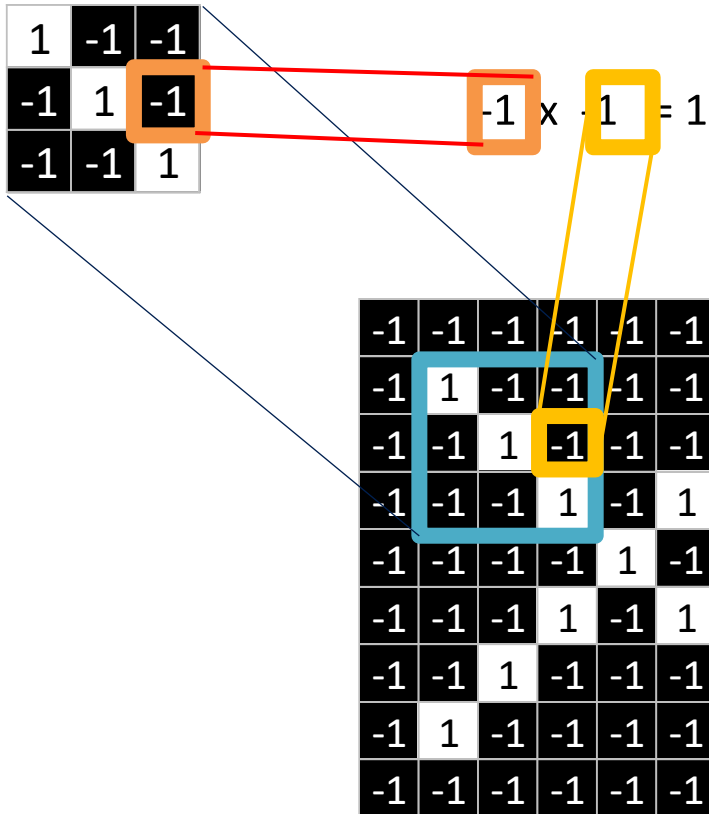# Filtering: The math behind the match



-1 x -1 = 1

# Filtering: The math behind the match
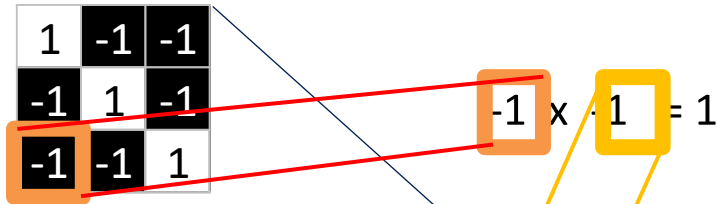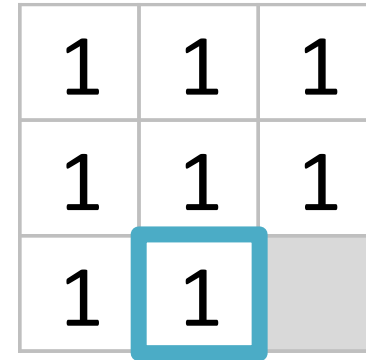
# Filtering: The math behind the match

# Filtering: The math behind the match



| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

| 1 | 1 | 1 |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 1 |

$$\frac{1+1+1+1+1+1+1+1+1}{9} = 1$$

# Filtering: The math behind the match



$1 \times 1 = 1$

# Filtering: The math behind the match

# Filtering: The math behind the match

# Filtering: The math behind the match

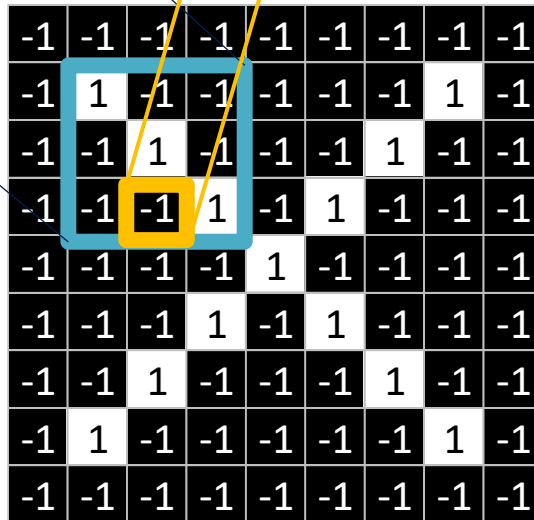| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

| 1 | 1 | -1 |
|---|---|----|
| 1 | 1 | 1 |
| -1 | 1 | 1 |

$$\frac{1 + 1 - 1 + 1 + 1 + 1 - 1 + 1 + 1}{9} = .55$$
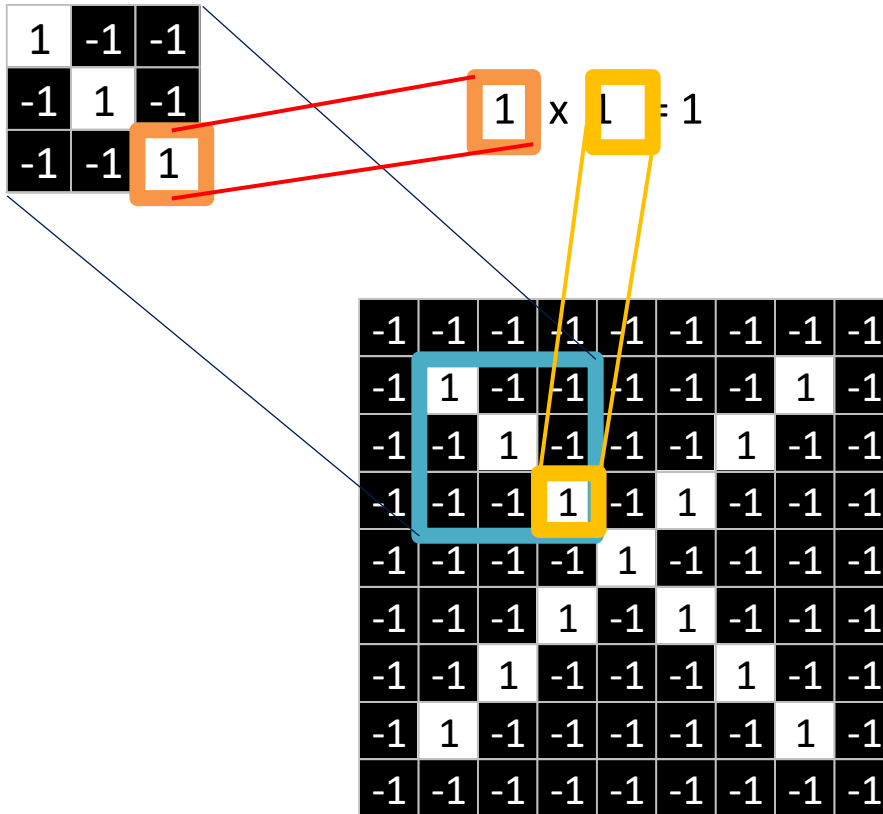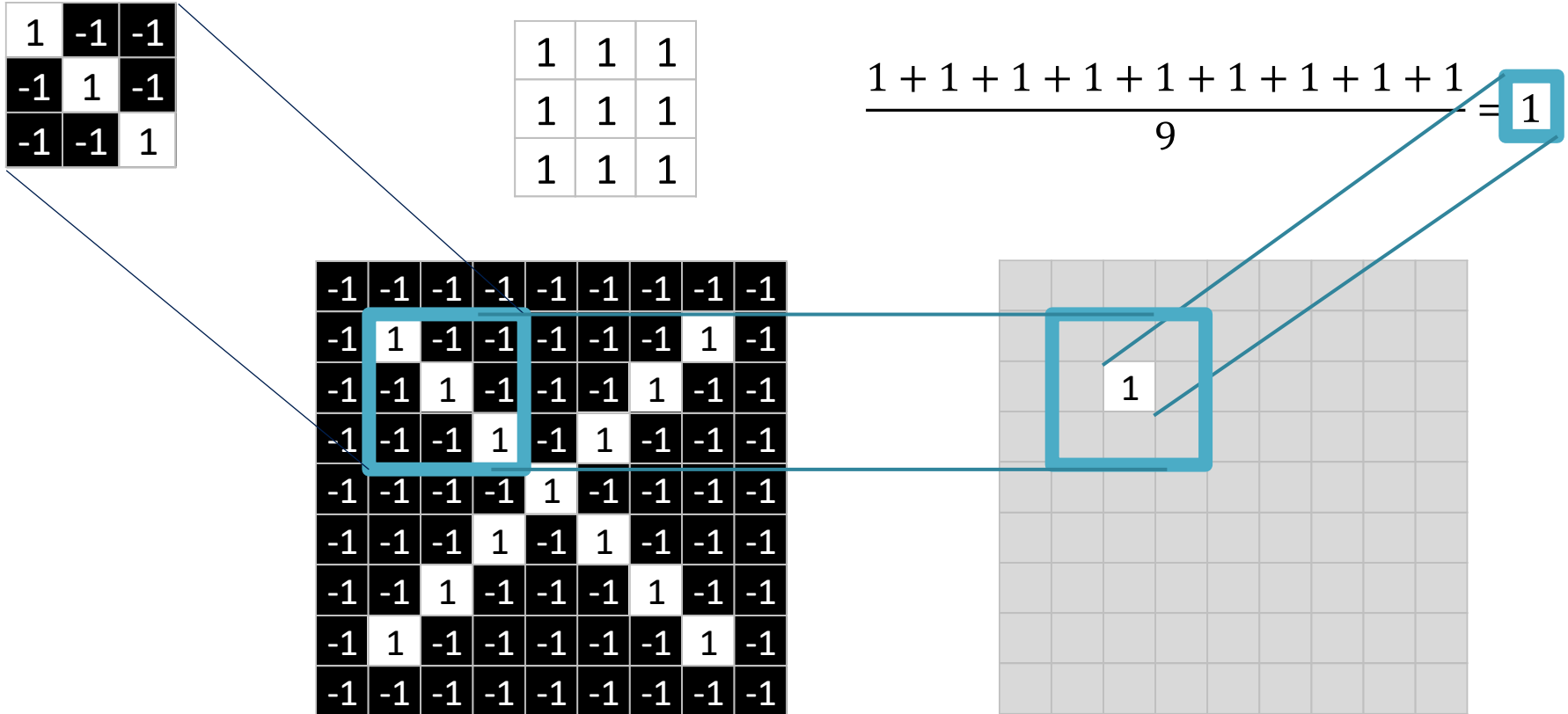
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
|----|----|----|----|----|----|----|----|----|
| -1 | 1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 |
| -1 | -1 | 1 | -1 | -1 | -1 | 1 | -1 | -1 |
| -1 | -1 | -1 | 1 | -1 | 1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 | 1 | -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 1 | -1 | 1 | -1 | -1 | -1 |
| -1 | -1 | 1 | -1 | -1 | -1 | 1 | -1 | -1 |
| -1 | 1 | -1 | -1 | -1 | -1 | -1 | 1 | -1 |
| -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |

1

.55

# Convolution: Trying every possible match

# Convolution: Trying every possible match

# Convolution layer

One image becomes a stack of filtered images

# outline

- Convolutional layer (module)
  - Convolution operation
  - Filters
  - Convolution module in a network
- Pooling layer (module)

# POOLING Layer

- In ConvNet architectures, **Conv** layers are often followed by **Pooling** layers
  - makes the representations smaller and more manageable without losing too much information.
  - Invariant in region.

32

downsampling

16

16

32

Source: Andrej Karpathy & Fei-Fei Li

# MAX POOLING

Single depth slice

| | | | |
|---|---|---|---|
| 1 | 1 | 2 | 4 |
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

x

y

max pool with 2x2 filters and stride 2

| | |
|---|---|
| 6 | 8 |
| 3 | 4 |

Source: Andrej Karpathy & Fei-Fei Li

# Max Pooling

# Max Pooling

maximum

# Max Pooling

# Max Pooling

# Max Pooling

maximum

| 0.77 | -0.11 | 0.11 | 0.33 | 0.55 | -0.11 | 0.33 |
| -0.11 | 1.00 | -0.11 | 0.33 | -0.11 | 0.11 | -0.11 |
| 0.11 | -0.11 | 1.00 | -0.33 | 0.11 | -0.11 | 0.55 |
| 0.33 | 0.33 | -0.33 | 0.55 | -0.33 | 0.33 | 0.33 |
| 0.55 | -0.11 | 0.11 | -0.33 | 1.00 | -0.11 | 0.11 |
| -0.11 | 0.11 | -0.11 | 0.33 | -0.11 | 1.00 | -0.11 |
| 0.33 | -0.11 | 0.55 | 0.33 | 0.11 | -0.11 | 0.77 |

| 1.00 | 0.33 | 0.55 | 0.33 |
| 0.33 | | | |
| | | | |
| | | | |

# Max Pooling

| 0.77 | -0.11 | 0.11 | 0.33 | 0.55 | -0.11 | 0.33 |
|---|---|---|---|---|---|---|
| -0.11 | 1.00 | -0.11 | 0.33 | -0.11 | 0.11 | -0.11 |
| 0.11 | -0.11 | 1.00 | -0.33 | 0.11 | -0.11 | 0.55 |
| 0.33 | 0.33 | -0.33 | 0.55 | -0.33 | 0.33 | 0.33 |
| 0.55 | -0.11 | 0.11 | -0.33 | 1.00 | -0.11 | 0.11 |
| -0.11 | 0.11 | -0.11 | 0.33 | -0.11 | 1.00 | -0.11 |
| 0.33 | -0.11 | 0.55 | 0.33 | 0.11 | -0.11 | 0.77 |

| 1.00 | 0.33 | 0.55 | 0.33 |
|---|---|---|---|
| 0.33 | 1.00 | 0.33 | 0.55 |
| 0.55 | 0.33 | 1.00 | 0.11 |
| 0.33 | 0.55 | 0.11 | 0.77 |

| 0.33 | -0.55 | 0.11 | -0.11 | 0.11 | -0.55 | 0.33 |
|---|---|---|---|---|---|---|
| -0.55 | 0.55 | -0.55 | 0.33 | -0.55 | 0.55 | -0.55 |
| 0.11 | -0.55 | 0.55 | -0.77 | 0.55 | -0.55 | 0.11 |
| -0.11 | 0.33 | -0.77 | 1.00 | -0.77 | 0.33 | -0.11 |
| 0.11 | -0.55 | 0.55 | -0.77 | 0.55 | -0.55 | 0.11 |
| -0.55 | 0.55 | -0.55 | 0.33 | -0.55 | 0.55 | -0.55 |
| 0.33 | -0.55 | 0.11 | -0.11 | 0.11 | -0.55 | 0.33 |

| 0.55 | 0.33 | 0.55 | 0.33 |
|---|---|---|---|
| 0.33 | 1.00 | 0.55 | 0.11 |
| 0.55 | 0.55 | 0.55 | 0.11 |
| 0.33 | 0.11 | 0.11 | 0.33 |

| 0.33 | -0.11 | 0.55 | 0.33 | 0.11 | -0.11 | 0.77 |
|---|---|---|---|---|---|---|
| -0.11 | 0.11 | -0.11 | 0.33 | -0.11 | 1.00 | -0.11 |
| 0.55 | -0.11 | 0.11 | -0.33 | 1.00 | -0.11 | 0.11 |
| 0.33 | 0.33 | -0.33 | 0.55 | -0.33 | 0.33 | 0.33 |
| 0.11 | -0.11 | 1.00 | -0.33 | 0.11 | -0.11 | 0.55 |
| -0.11 | 1.00 | -0.11 | 0.33 | -0.11 | 0.11 | -0.11 |
| 0.77 | -0.11 | 0.11 | 0.33 | 0.55 | -0.11 | 0.33 |

| 0.33 | 0.55 | 1.00 | 0.77 |
|---|---|---|---|
| 0.55 | 0.55 | 1.00 | 0.33 |
| 1.00 | 1.00 | 0.11 | 0.55 |
| 0.77 | 0.33 | 0.55 | 0.33 |

# Max Pooling layer

A stack of images becomes a stack of smaller images.

# Average POOLING

Single depth slice

| | | | |
|---|---|---|---|
| 1 | 1 | 2 | 4 |
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

x

y

average pool with 2x2 filters and stride 2

| | |
|---|---|
| 4.25 | 5.25 |
| 2 | 2 |

# CNN: Intuitive example



CONV    CONV    POOL CONV    CONV    POOL CONV    CONV    POOL    FC (Fully-connected)

ReLU   ReLU   ReLU   ReLU   ReLU   ReLU   ReLU

truck
car
airplane
ship
horse

Source: Andrej Karpathy & Fei-Fei Li

# Famous Net Architecture



**Year 2012**

SuperVision

[Krizhevsky NIPS 2012]

**Year 2014**

GoogLeNet    VGG

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096
FC-1000
softmax

Convolution
Pooling
Softmax
Other

[Szegedy arxiv 2014]    [Simonyan arxiv 2014]

**Year 2015**

MSRA

34-layer residual

# 主要内容

- 多层感知机(Multi-layer Perceptron, MLP)
  - 反向传播算法(Back-propagation)
- 卷积神经网络 （Convolutional Neural Network, CNN）
  - 卷积操作和卷积层
  - 池化 （Pooling）
- 循环神经网络（Recurrent Neural Network, RNN）
  - 建模和训练
  - LSTM模型

# outline

- Recurrent Neural Network
  - Modeling
  - Training
- Long Short Term Memory (LSTM)
  - Motivation
  - Modeling
- Application
  - Generate article

# Classification: MLP and CNN



Convolution (卷积)



$x_0 = 1$    $h_0{}^{(1)} = 1$

$x_1$

$x_2$

$x_3$

Input:

hidden layer

output

$(1, 0, 0)^T$

$(x_1, x_2, x_3)$



What the computer sees

I go to cinema, and I book a ticket

# One example-modeling: motivation

- **Task**: **Character-level language model**
  - **example**
    - Vocabulary [h,e,l,o]
    - Training sequence "hello"

- Data representation：
  - X: {h, e, l, l}
  - Y: {e, l , l, o}

Input: **x**   hidden layer   output layer: Y

# Modeling

- MLP → RNN

$$y = W_{yh}h$$

$$y = \sigma(W_{yh}h) \qquad \sigma(a) = \frac{1}{1 + e^{-a}}$$



$$h = \sigma(W_{hx}x) \qquad \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

hyperbolic tangent function

# Modeling

- MLP $\rightarrow$ RNN

$$\mathbf{y} = \mathbf{W}_{yh}\mathbf{h}$$

$$\mathbf{h} = tanh(\mathbf{W}_{hx}\mathbf{x})$$

$$\mathbf{y}_t = \mathbf{W}_{yh}\mathbf{h}_t$$

$$\mathbf{h}_t = tanh(\mathbf{W}_{hx}\mathbf{x} + \mathbf{W}_{hh}\mathbf{h}_{t-1})$$

# Modeling

- RNN-unrolling(展开)

$$y_t = W_{yh}h_t$$

$$h_t = tanh(W_{hx}x + W_{hh}h_{t-1})$$

# Example

- Character-level language model

  - Training sequence:
    - "Hello"
  - Presentation:

  X: {h, e, l, l}

  Y: {e, l, l, o}



One-hot aka one-of-K encoding

# Example

- Examples:

  – Training sequence:
    - "Hello"
  – Presentation:

  X: {h, e, l, l}

  Y: {e,  l , l, o}



| | "h" | "e" | "l" | "l" |
|---|---|---|---|---|
| input layer | 1 0 0 0 | 0 1 0 0 | 0 0 1 0 | 0 0 1 0 |
| input chars: | "h" | "e" | "l" | "l" |

# Example

- Examples:

  - Training sequence:
    - "Hello"
  - Presentation:

  X: {h, e, l, l}

  Y: {e,  l , l, o}

# Example

- Examples:

  – Training sequence:
    - "Hello"
  – Presentation:
  X: {h, e, l, l}
  Y: {e, l , l, o}
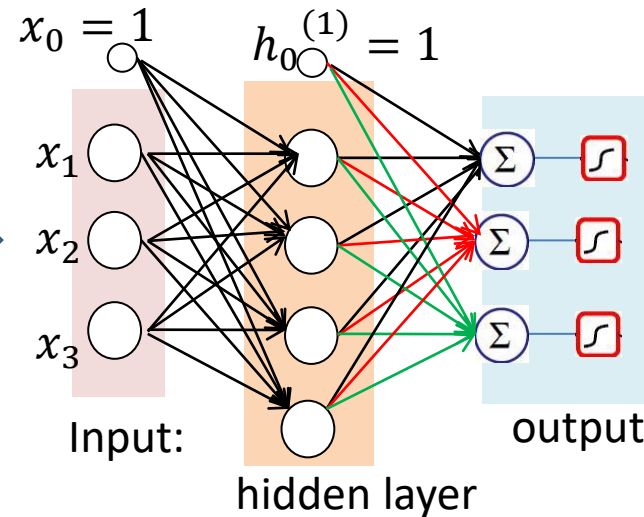
# outline

- Recurrent Neural Network
  - Modeling
  - <span style="color:red">Training</span>
- Long Short Term Memory (LSTM)
  - Motivation
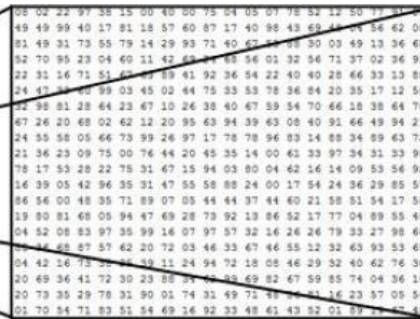  - Modeling
- Application
  - Generate article

# Neural Network

- Training Algorithm



➤ 0.初始化权重　$\boldsymbol{W}^{(0)}$
➤ 1. 前向过程：
   ➤ 1.1根据输入$\boldsymbol{x}$，计算输出值 $\boldsymbol{y}$
   ➤ 1.2.计算损失函数值$L(\boldsymbol{y}, \hat{\boldsymbol{y}})$。
➤ 2.后向传播
   ➤ 计算$\dfrac{d\,L}{\boldsymbol{y}}$
   ➤ 后向传播直到计算$\dfrac{d\,L}{\boldsymbol{x}}$
➤ 3.计算梯度$\dfrac{d\,L}{d\,\boldsymbol{W}}$
➤ 4.更新梯度
$$\boldsymbol{W}^{(t+1)} = \boldsymbol{W}^{(t)} - \eta\,\frac{d\,L}{d\,\boldsymbol{W}^{(t)}}$$

# Training

- learning
  - Sequence length=3

- Back-propagation

$$-\{\frac{\mathrm{d}\,\mathrm{L}}{d\,y_i} \Rightarrow \frac{\mathrm{d}\,\mathrm{L}}{d\,h_i} \Rightarrow \frac{\mathrm{d}\,\mathrm{L}}{d\,x_i}\}$$



**Target:** $\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$

# Training

- Back-propagation  **Target:** $\widehat{y}_3$

$$\frac{\mathrm{d}\,L}{d\,h_3} = \frac{\mathrm{d}\,L}{d\,y_3}\,\frac{\mathrm{d}y_3}{d h_3}$$

# Training

- Back-propagation **Target:**

$$\frac{\mathrm{d}\,\mathrm{L}}{d\,\boldsymbol{h}_3} = \frac{\mathrm{d}\,\mathrm{L}}{\boldsymbol{d}\,\boldsymbol{y}_3}\,\frac{\mathrm{d}\boldsymbol{y}_3}{d\boldsymbol{h}_3}$$

$$\frac{\mathrm{d}\,\mathrm{L}}{d\,\boldsymbol{h}_2} = \frac{\mathrm{d}\,\mathrm{L}}{\boldsymbol{d}\,\boldsymbol{y}_2}\,\frac{\mathrm{d}\boldsymbol{y}_2}{d\boldsymbol{h}_2} + \frac{\mathrm{d}\,\mathrm{L}}{\boldsymbol{d}\,\boldsymbol{y}_3}\,\frac{\mathrm{d}\boldsymbol{y}_3}{d\boldsymbol{h}_3}\frac{\mathrm{d}\boldsymbol{h}_3}{d\boldsymbol{h}_2}$$

# Training

- Back-propagation

**Target:** $\hat{y}_1$      $\hat{y}_2$      $\hat{y}_3$

$$\frac{d\,L}{d\,h_3} = \frac{d\,L}{d\,y_3}\frac{dy_3}{dh_3}$$

$$\frac{d\,L}{d\,h_2} = \frac{d\,L}{d\,y_2}\frac{dy_2}{dh_2} + \frac{d\,L}{d\,y_3}\frac{dy_3}{dh_3}\frac{dh_3}{dh_2}$$

$$\frac{d\,L}{d\,h_1} = \frac{d\,L}{d\,y_1}\frac{dy_1}{dh_1} + \frac{d\,L}{d\,y_2}\frac{dy_2}{dh_2}\frac{dh_2}{dh_1} +$$

$$\frac{d\,L}{d\,y_3}\frac{dy_3}{dh_3}\frac{dh_3}{dh_2}\frac{dh_2}{dh_1}$$

$$\frac{d\,L}{d\,h_t} = \sum_{s=t}^{T=3} \frac{d\,L}{d\,y_s}\frac{dy_s}{dh_s}\frac{dh_s}{dh_t}$$

# Training

- Gradient r.t Weight

**Target:** $\hat{y}_1$     $\hat{y}_2$     $\hat{y}_3$

$$\frac{d\,L}{d\,\boldsymbol{h}_t} = \sum_{s=t}^{T} \frac{d\,L}{d\,\boldsymbol{y}_s}\,\frac{d\boldsymbol{y}_s}{d\boldsymbol{h}_s}\frac{d\boldsymbol{h}_s}{d\boldsymbol{h}_t}$$

➢ Calculate gradient respect to weight

$$\frac{d\,L}{d\mathbf{W}_{\text{yh}}{}^{(i)}} \qquad \frac{d\,L}{d\mathbf{W}_{\text{hh}}{}^{(i)}} \qquad \frac{d\,L}{d\mathbf{W}_{\text{hx}}{}^{(i)}}$$

➢ Update weight

$$\boldsymbol{W}^{(t+1)} = \boldsymbol{W}^{(t)} - \eta \sum_{s=1}^{T} \frac{d\,L}{d\,\boldsymbol{W}^{(s)}}$$

# RNN

- longer
- deeper

**Target:** $\hat{y}_1$ $\qquad$ $\hat{y}_2$ $\qquad$ $\hat{y}_3$ $\qquad$ $\hat{y}_T$

$y_1$ $\qquad$ $y_2$ $\qquad$ $y_3$ $\qquad$ ... $\qquad$ $y_T$

$W_{yh}$ $\qquad$ $W_{yh}$ $\qquad$ $W_{yh}$ $\qquad$ $W_{yh}$

$h_1$ $\;\xrightarrow{W_{hh}}\;$ $h_2$ $\;\xrightarrow{W_{hh}}\;$ $h_3$ $\;...\xrightarrow{W_{hh}}\;$ $h_T$

$W_{hx}$ $\qquad$ $W_{hx}$ $\qquad$ $W_{hx}$ $\qquad$ $W_{hx}$

$x_1$ $\qquad$ $x_2$ $\qquad$ $x_3$ $\qquad$ ... $\qquad$ $x_T$

# Relation: MLP, CNN and RNN



$x \in R^{d_{in}}$   $W \in R^{d_{out} \times d_{in}}$   $y \in R^{d_{out}}$

$W_{xh} \in R^{d_{out} \times d_{in}}$

$W_{hh}$

$t_1$

$t_2$

...

Feature maps

$X \in R^{d_{in} \times h \times w}$

$W \in R^{d_{out} \times d_{in} \times F_h \times F_w}$

$Y \in R^{d_{out} \times h \times w}$

# outline

- Recurrent Neural Network
  - Modeling
  - Training

- Long Short Term Memory (LSTM)
  - Motivation
  - Modeling

- Application
  - Generate article

# Long Short Term Memory (LSTM)

- Motivation

"the clouds are in the *sky*,"

# Long Short Term Memory (LSTM)

• Motivation

"I grew up in France… I speak fluent *French*."

# LSTM

- Motivation
  - Gradient explosion (梯度爆炸)
  - Gradient vanishing(梯度弥散)

$$\frac{\mathrm{d}\, \mathrm{L}}{\mathrm{d}\, \boldsymbol{h}_t} = \sum_{s=t}^{T} \frac{\mathrm{d}\, \mathrm{L}}{\boldsymbol{d}\, \boldsymbol{y}_s} \frac{\mathrm{d}\boldsymbol{y}_s}{\mathrm{d}\boldsymbol{h}_s} \frac{\mathrm{d}\boldsymbol{h}_s}{\mathrm{d}\boldsymbol{h}_t}$$

**Target:** $\widehat{\boldsymbol{y}}_1$ $\qquad$ $\widehat{\boldsymbol{y}}_2$ $\qquad$ $\widehat{\boldsymbol{y}}_3$
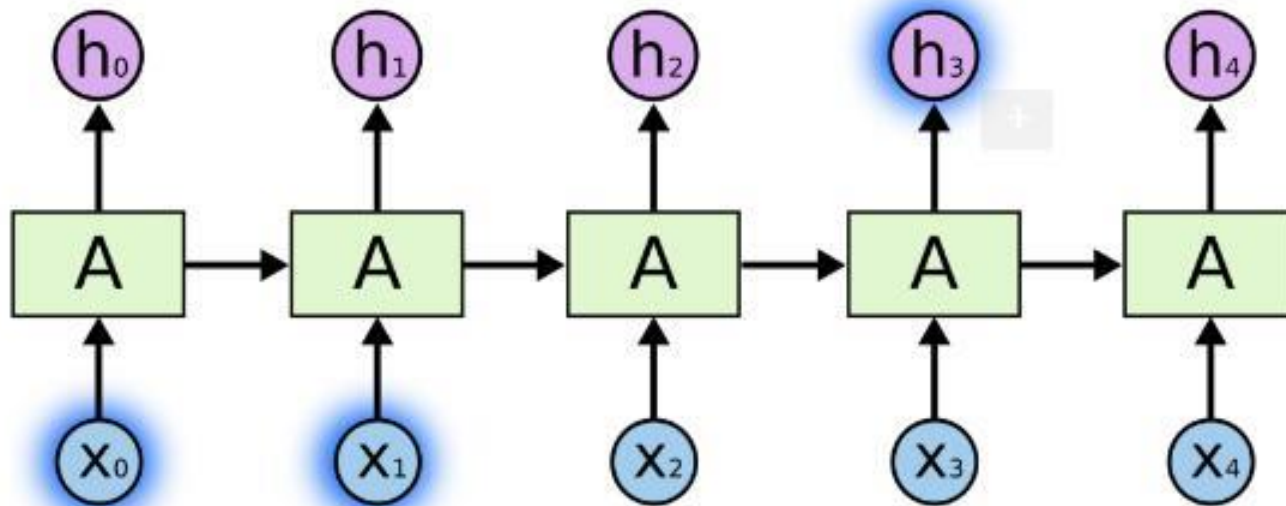
# outline

- Recurrent Neural Network
  - Modeling
  - Training

- Long Short Term Memory (LSTM)
  - Motivation
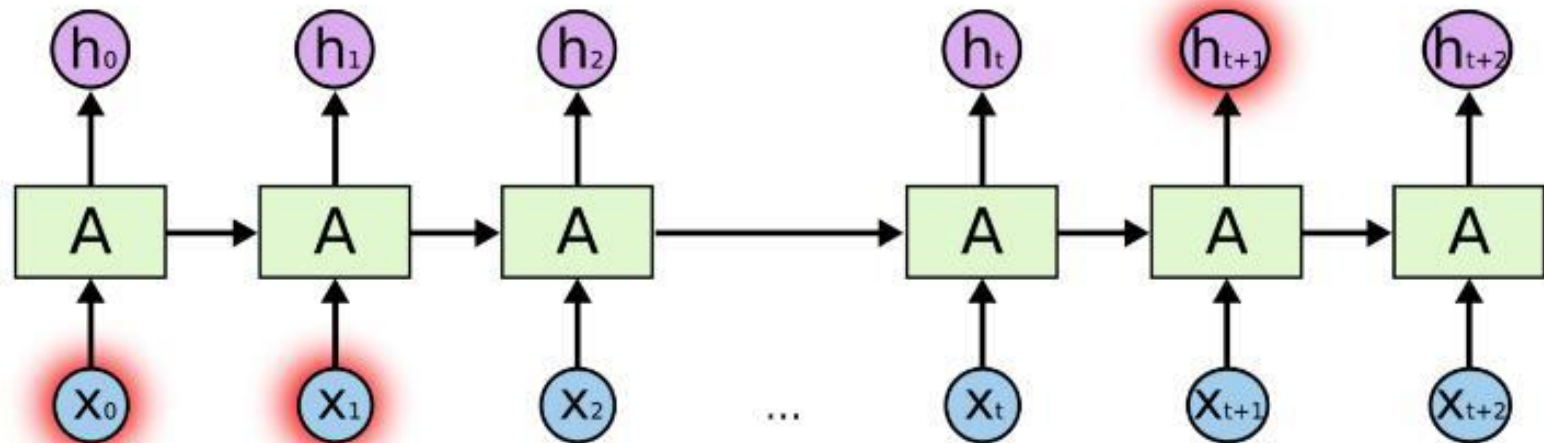  - Modeling

- Application
  - Generate article

# LSTM

- Modeling

$$h_t = tanh(W_{hx}x + W_{hh}h_{t-1})$$

# LSTM

- Modeling
  - Input gate
  - Input information



$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# LSTM

- Modeling

$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

Forget Gate    Input Gate

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Previous Info    Input Info

# LSTM

- Modeling
  - Output gate
  - Output information



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

# LSTM

- Modeling



$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$

$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

$$o_t = \sigma\left(W_o\,[h_{t-1}, x_t] + b_o\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

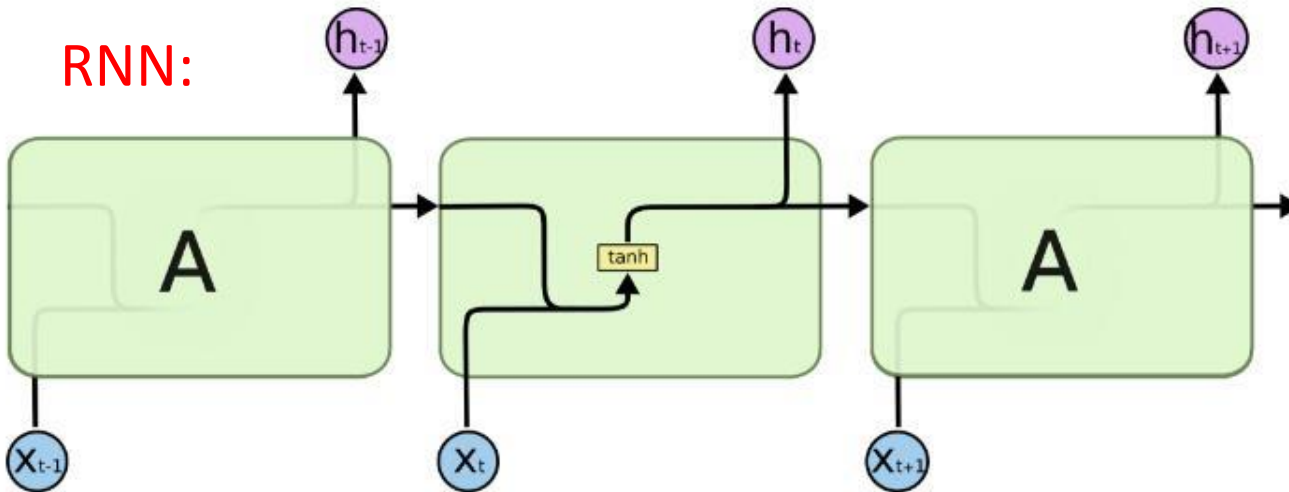$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$h_t = o_t * \tanh(C_t)$$

# outline

- Recurrent Neural Network
  - Modeling
  - Training

- Long Short Term Memory (LSTM)
  - Motivation
  - Modeling

- <span style="color:red">Application</span>
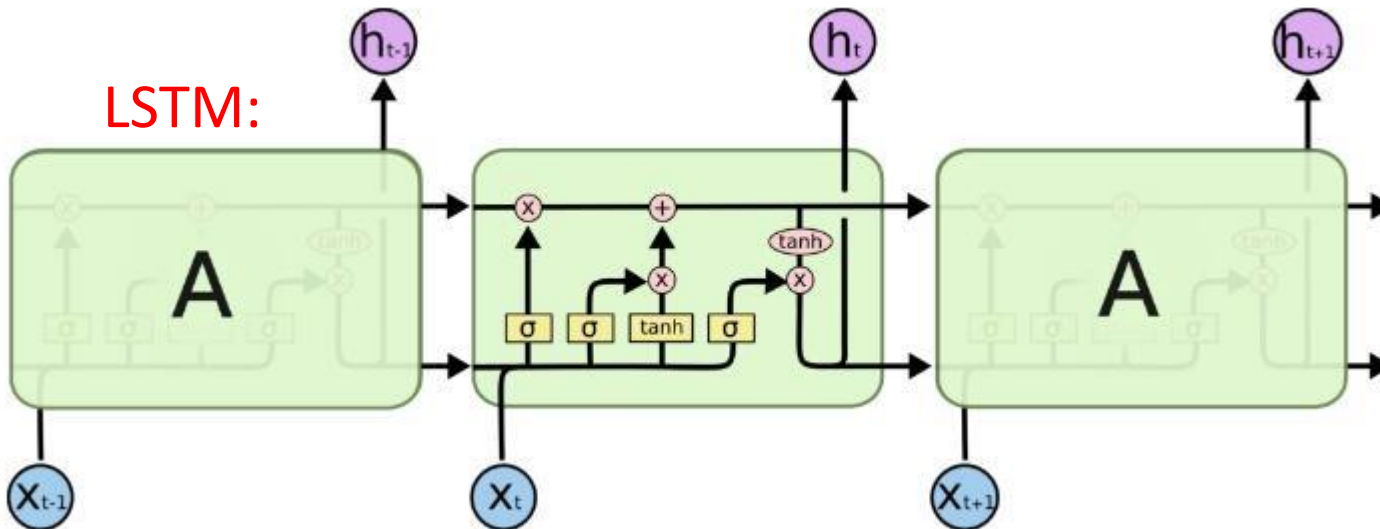  - Generate article

# RNN-application

- Generate article

# RNN-application

- Generate article



Sonnet 116 – Let me not …

by William Shakespeare

Let me not to the marriage of true minds
    Admit impediments. Love is not love
Which alters when it alteration finds,
    Or bends with the remover to remove:
O no! it is an ever–fixed mark
    That looks on tempests and is never shaken;
It is the star to every wandering bark,
    Whose worth's unknown, although his height be taken.
Love's not Time's fool, though rosy lips and cheeks
    Within his bending sickle's compass come:
Love alters not with his brief hours and weeks,
    But bears it out even to the edge of doom.
If this be error and upon me proved,
    I never writ, nor no man ever loved.

# RNN-application

- Generate article

at first:

tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd t o idoe ns,smtt   h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

↓ train more

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.

# RNN-application

- Generate article

PANDARUS:

Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

# RNN-application

- Generate article

VIOLA:
Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:
O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.

# RNN-application

- Generate C code

# RNN-application

- Generated C code

```c
static void do_command(struct seq_file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x00000000ffffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
  }
  /* Free our user pages pointer to place camera if all dash */
  subsystem_info = &of_changes[PAGE_SIZE];
  rek_controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control_check_polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)
    seq_puts(s, "policy ");
}
```

Generated
C code

# RNN-application

- Generated C code

```c
#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/setew.h>
#include <asm/pgproto.h>

#define REG_PG     vesa_slot_addr_pack
#define PFM_NOCOMP  AFSR(0, load)
#define STACK_DDR(type)     (func)

#define SWAP_ALLOCATE(nr)      (e)
#define emulate_sigs()  arch_get_unaligned_child()
#define access_rw(TST)  asm volatile("movd %%esp, %0, %3" : : "r" (0));   \
  if (__type & DO_READ)

static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
          pC>[1]);

static void
os_prefix(unsigned long sys)
{
#ifdef CONFIG_PREEMPT
  PUT_PARAM_RAID(2, sel) = get_state_state();
  set_pid_sum((unsigned long)state, current_state_str(),
          (unsigned long)-1->lr_full; low;
}
```