



# 机器学习

苏州大学计算机科学与技术学院

自然语言处理实验室

主讲：周夏冰

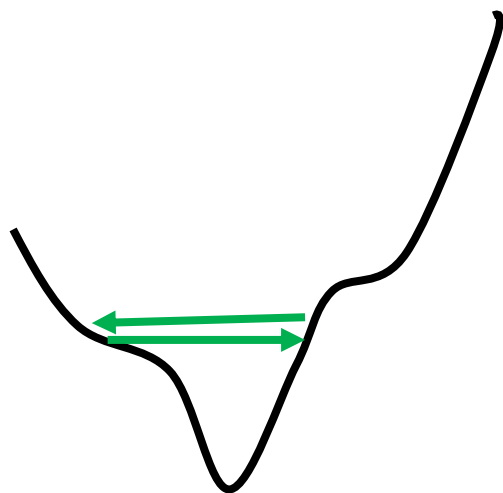
邮箱：[zhouxiabing@suda.edu.cn](mailto:zhouxiabing@suda.edu.cn)

The background of the slide is a soft-focus photograph of pink cherry blossoms. The flowers are in various stages of bloom, with some showing yellow centers. The overall tone is gentle and aesthetic.

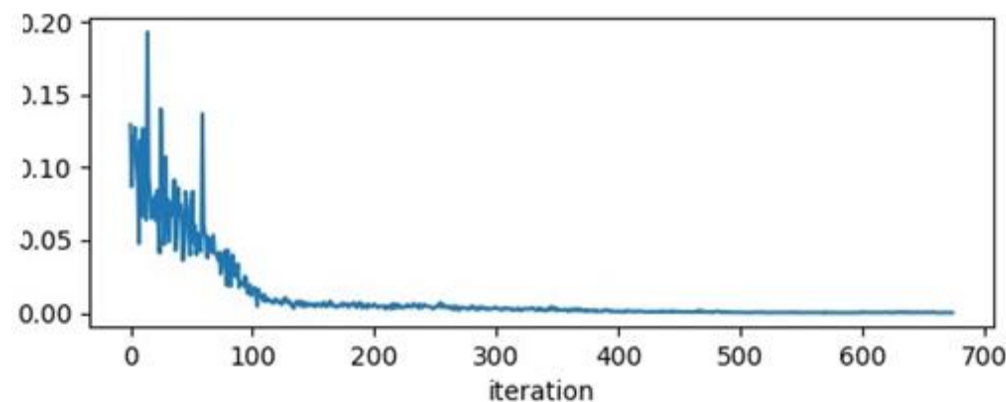
# Learning rate



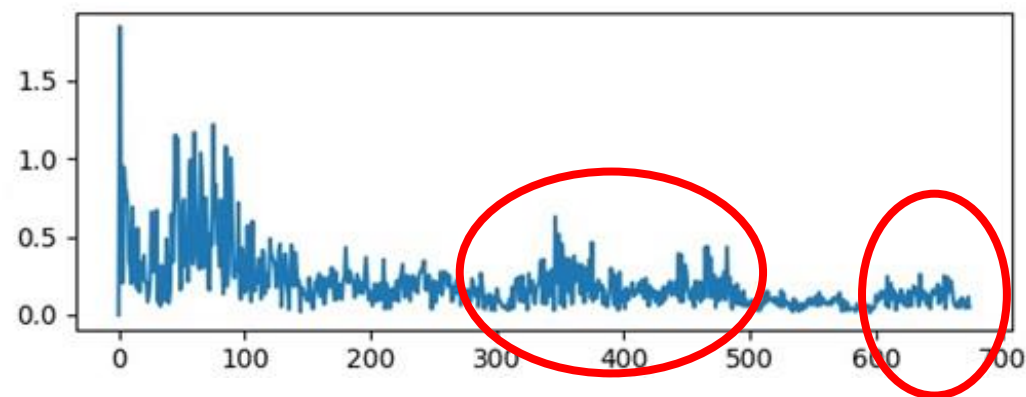
# Training stuck $\neq$ Small Gradient

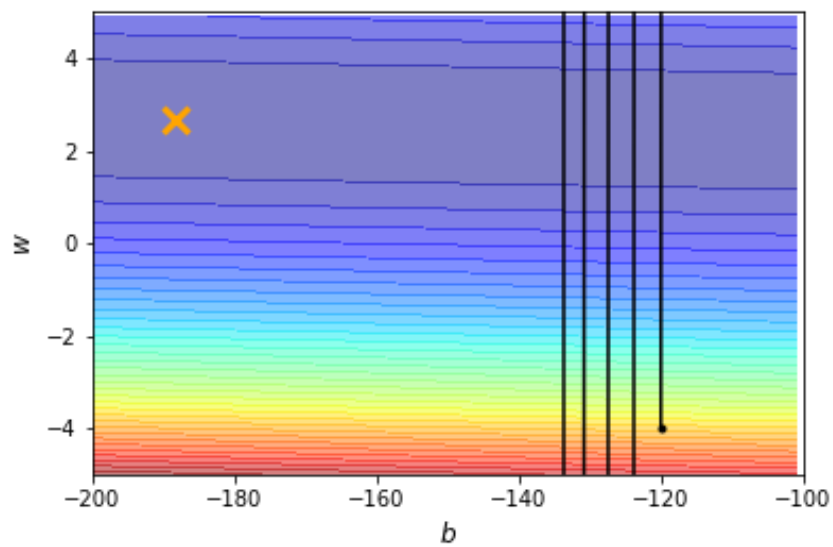
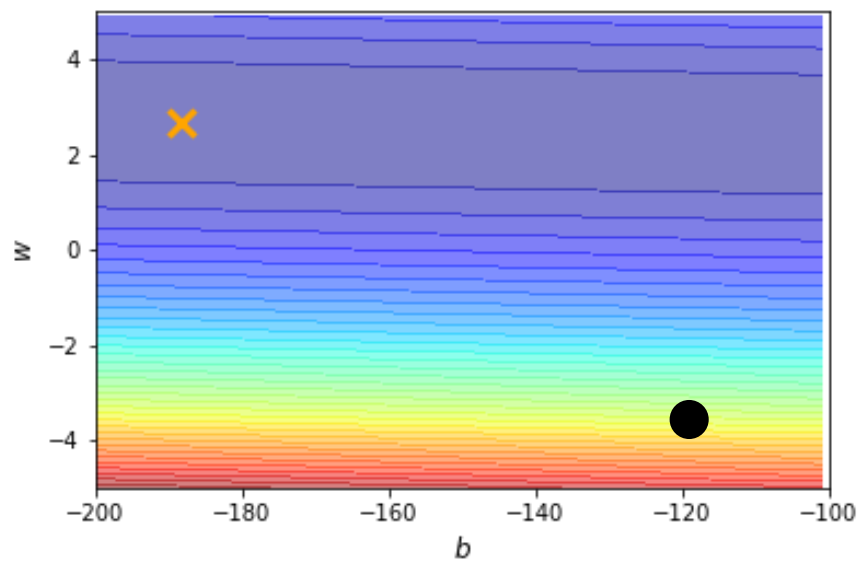


loss

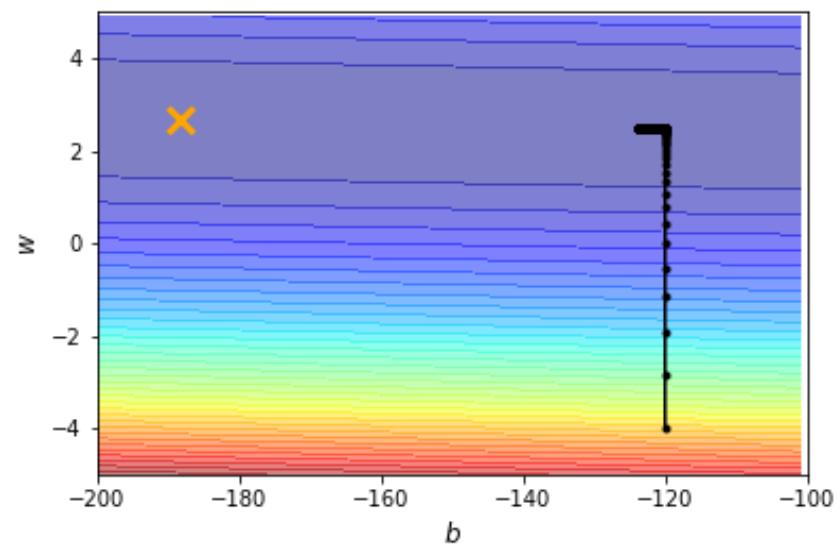


norm of gradient





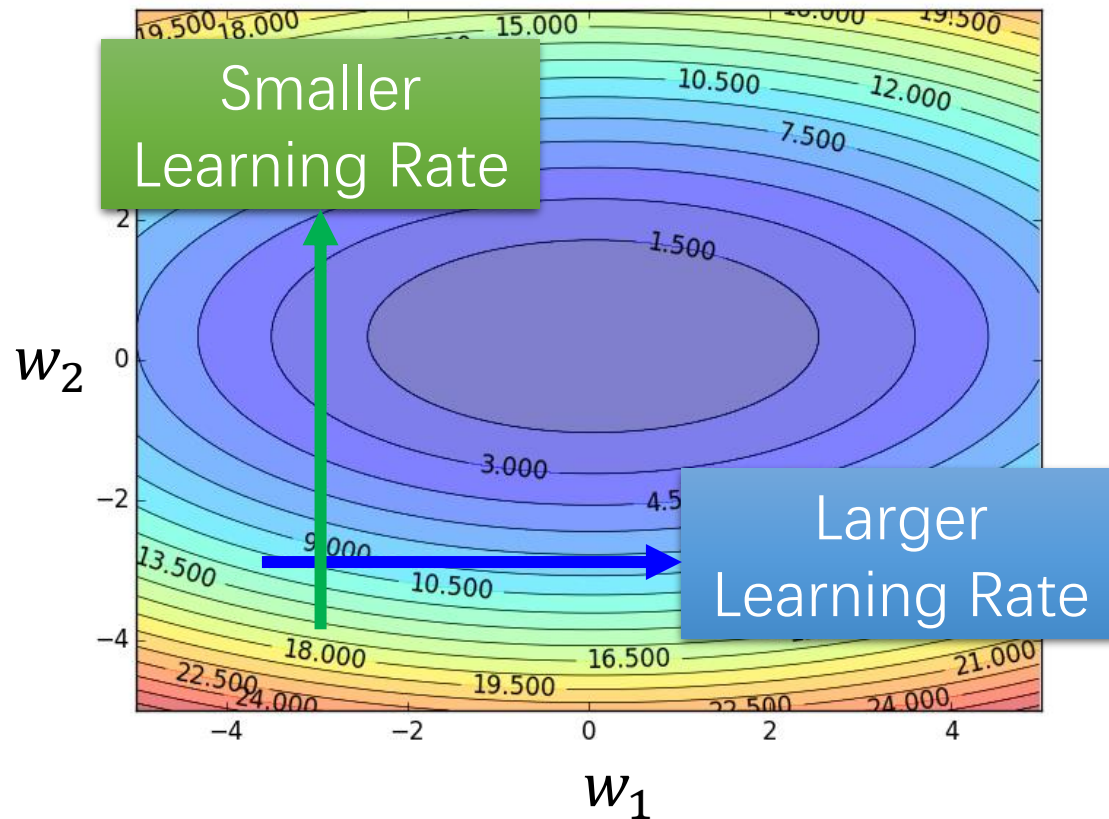
$$\eta = 10^{-2}$$



$$\eta = 10^{-7}$$



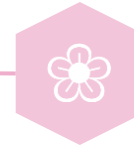
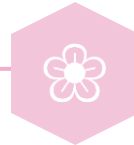
# Learning rate



$$\theta_i^{t+1} \leftarrow \theta_i^t - \boxed{\eta} g_i^t$$

$$g_i^t = \frac{\partial L}{\partial \theta_i} \bigg|_{\theta = \theta^t}$$

$$\theta_i^{t+1} \leftarrow \theta_i^t - \boxed{\frac{\eta}{\sigma_i^t}} g_i^t$$

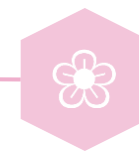
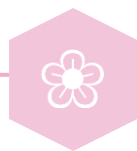
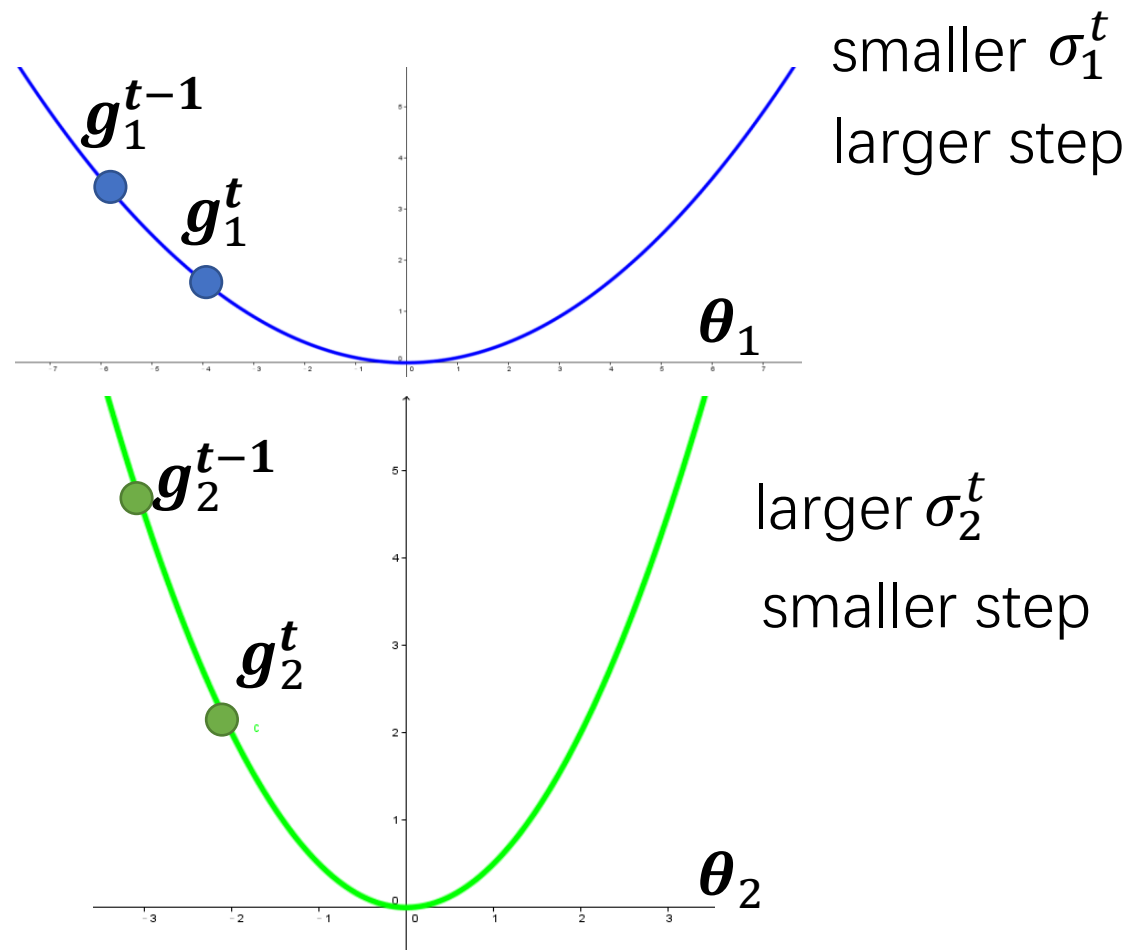


# Learning rate

- $\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t$

- $\sigma_i^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g_i^t)^2}$

- Adagrad



# Learning rate

- $\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t$

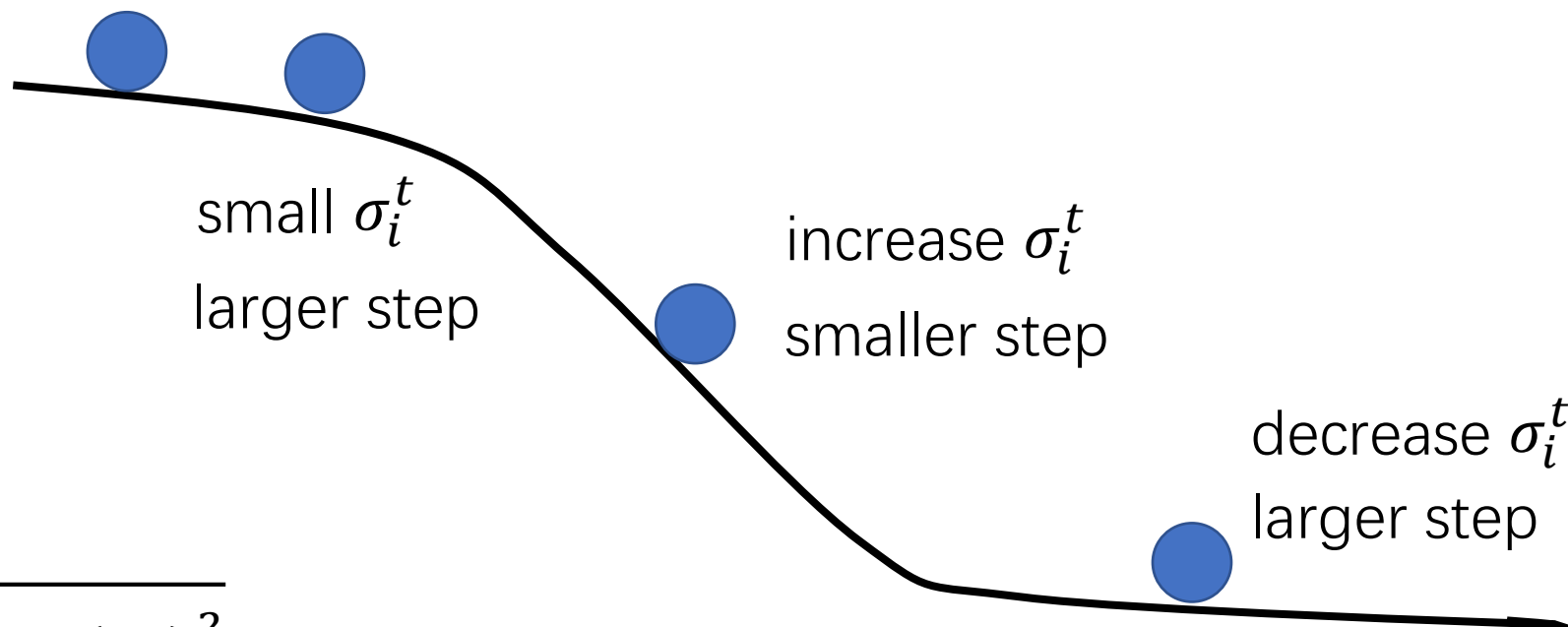
- $\sigma_i^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g_i^t)^2}$

- **Adagrad**

- $\sigma_i^t = \sqrt{\alpha(\sigma_i^{t-1})^2 + (1 - \alpha)(g_i^t)^2}$

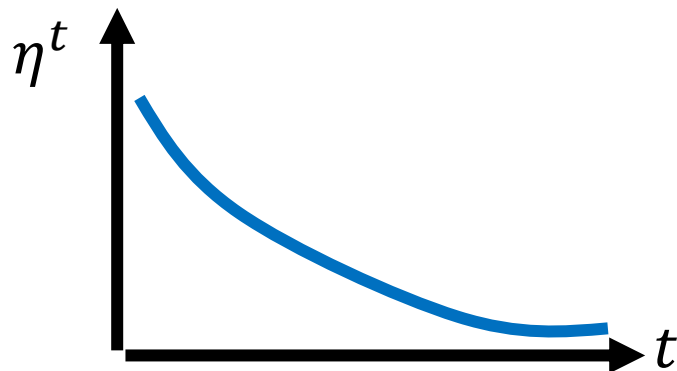
- **RMSProp**

**Adam:** 结合了动量梯度下降和RMSProp两种算法的优点

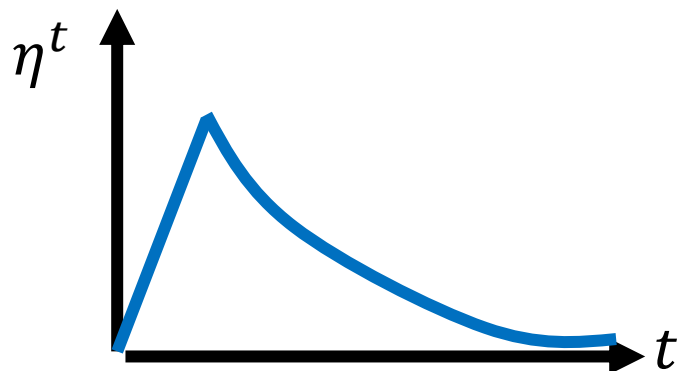


# Learning rate

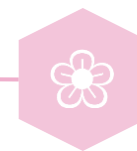
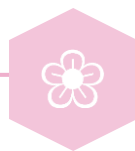
$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t$$



Learning Rate Decay



Warm Up





# Summary of Optimization

## (Vanilla) Gradient Descent

$$\theta_i^{t+1} \leftarrow \theta_i^t - \eta g_i^t$$

## Various Improvements

$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta^t}{\sigma_i^t} m_i^t$$

Learning rate scheduling

Momentum: weighted sum of the previous gradients

root mean square of the gradients

确定方向

确定大小

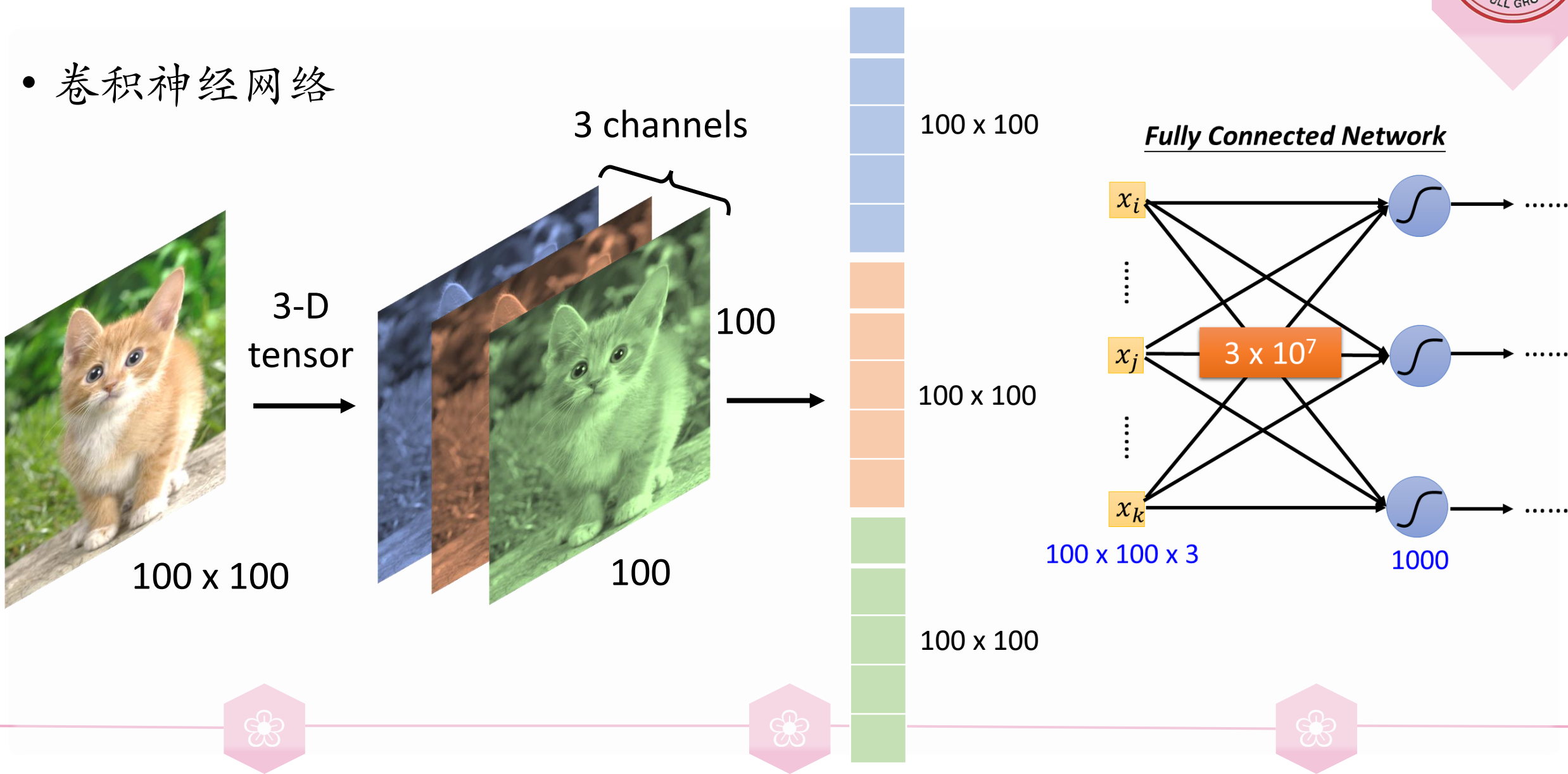
# 02

## 卷积神经网络 CNN



# CNN

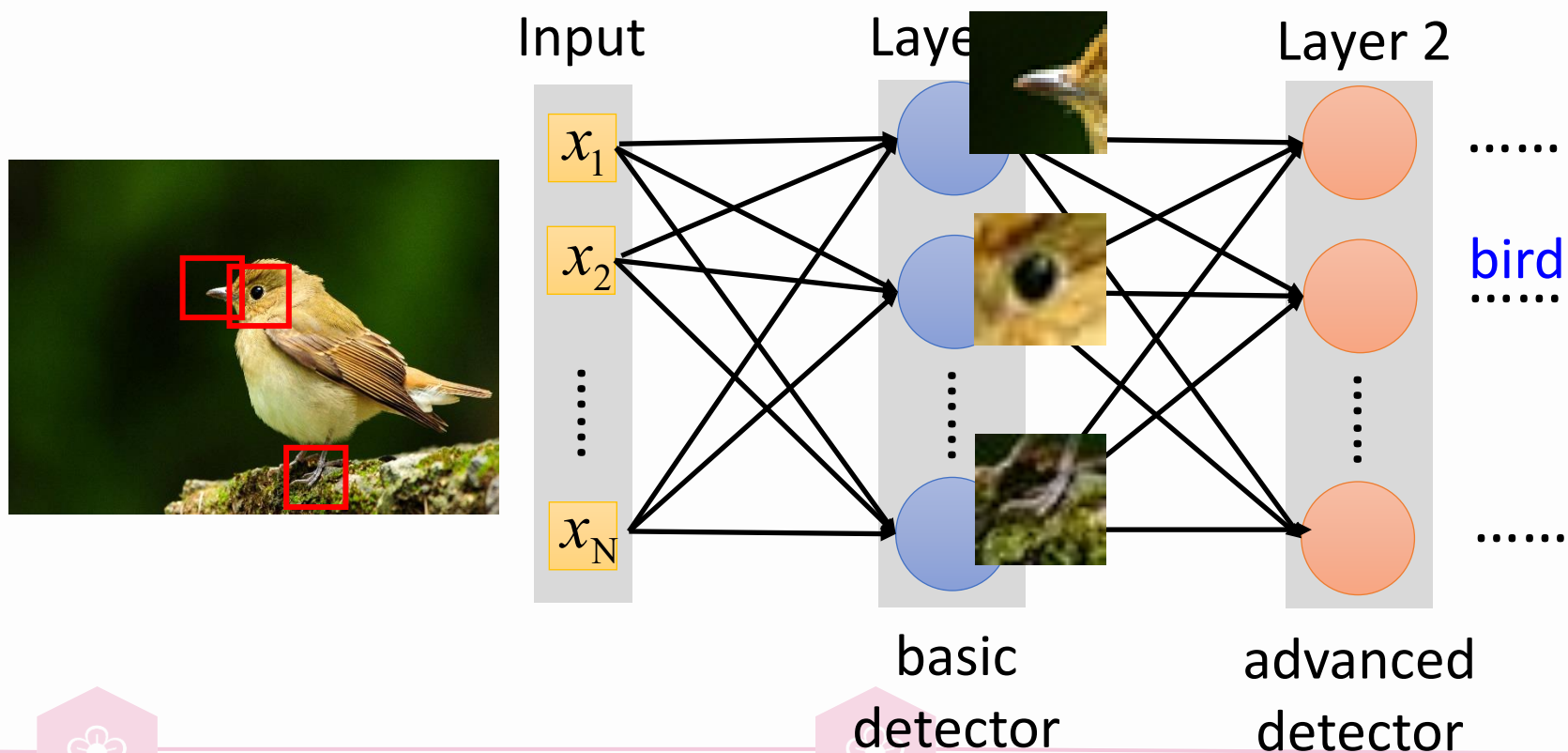
- 卷积神经网络



# CNN

- 关键的模式只占整个图片的很小的区域

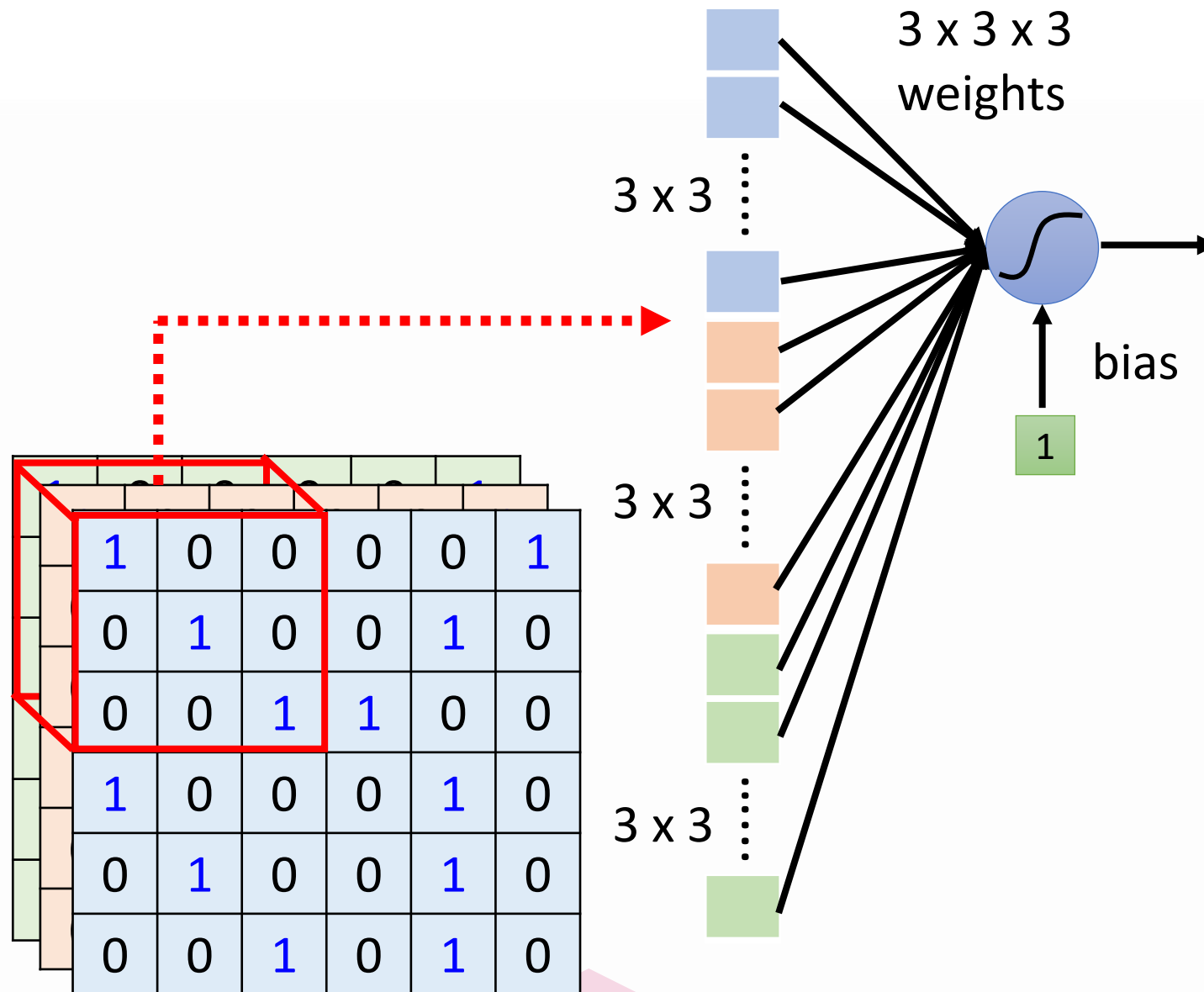
对于神经元，它不需要观测整个图像去识别pattern



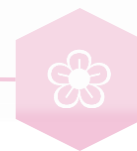
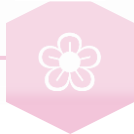
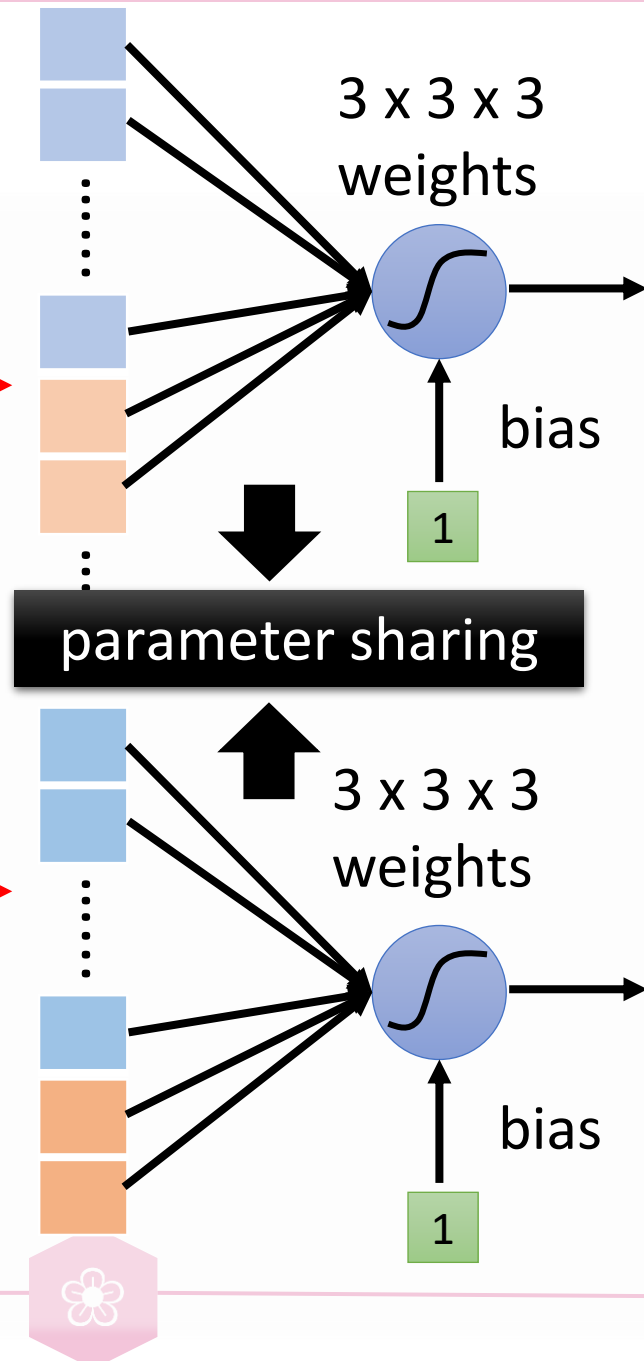
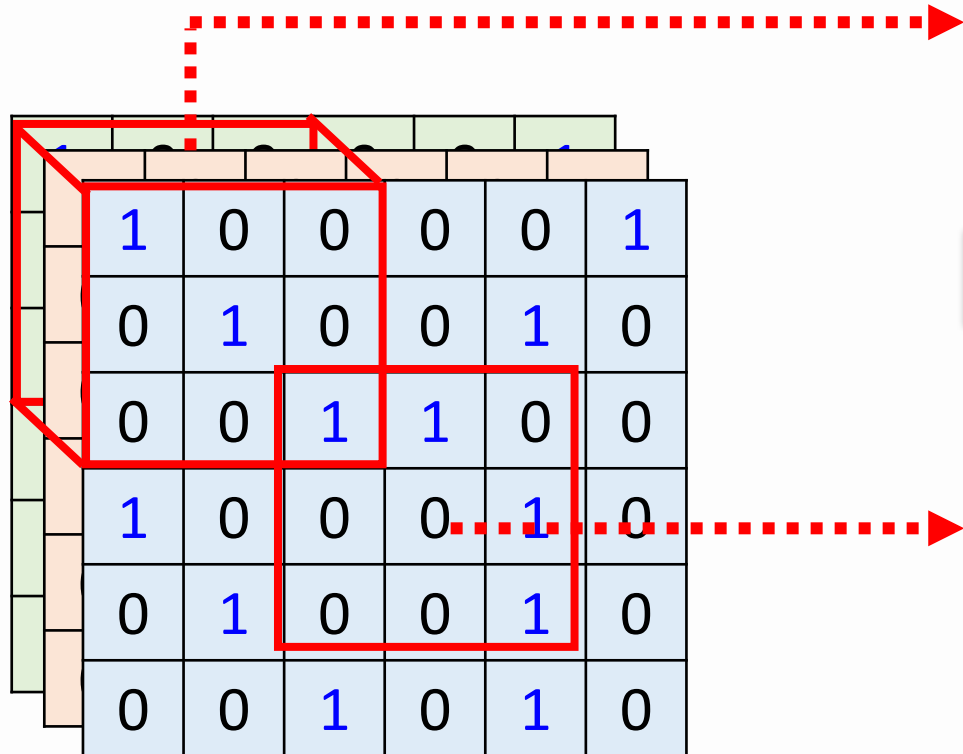
# CNN



感受野  
Receptive field

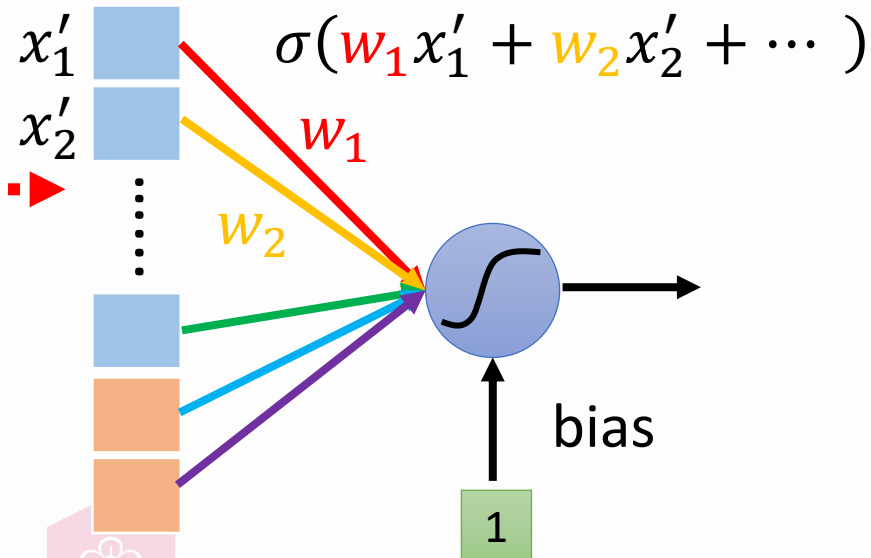
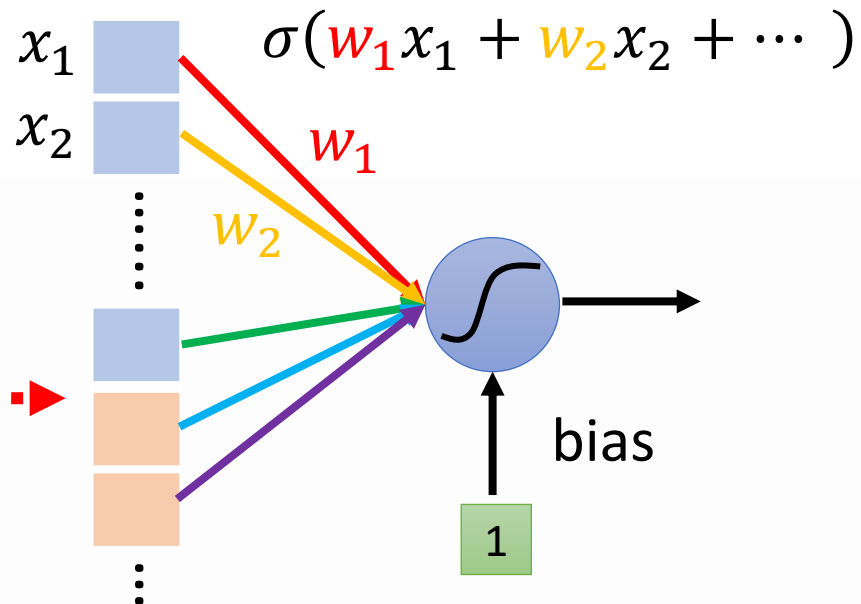
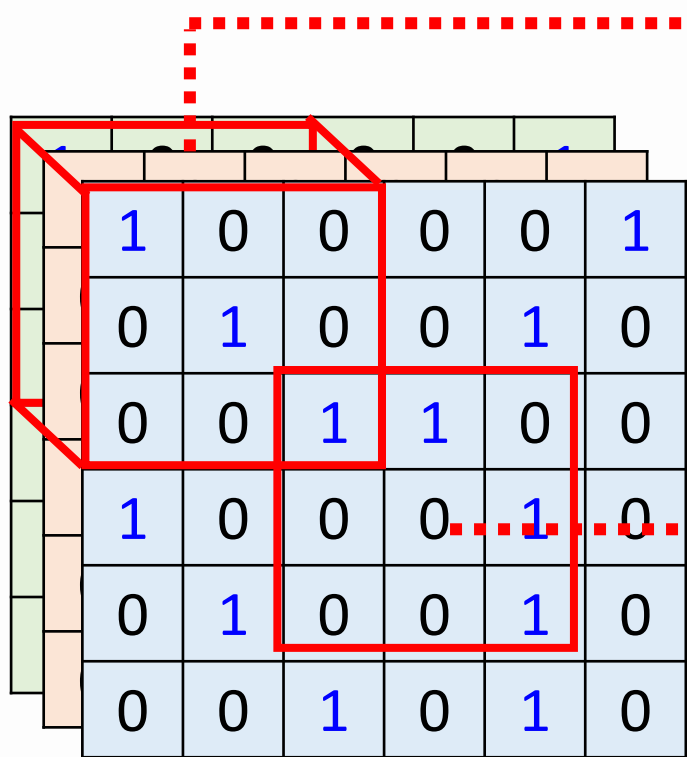


# CNN





# CNN



# CNN



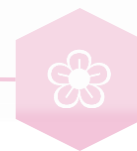
- 一定的压缩不影响识别

bird



subsampling

bird

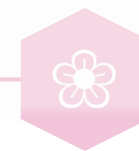
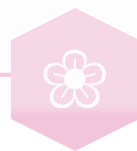
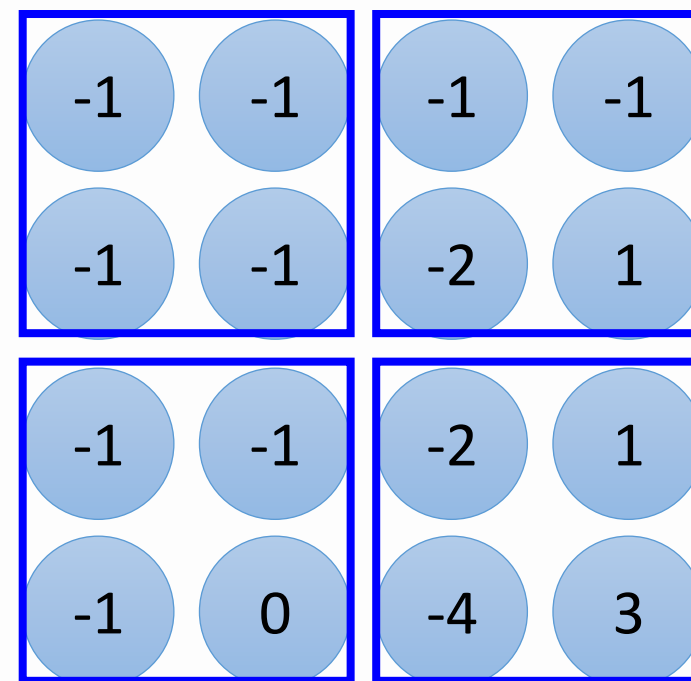
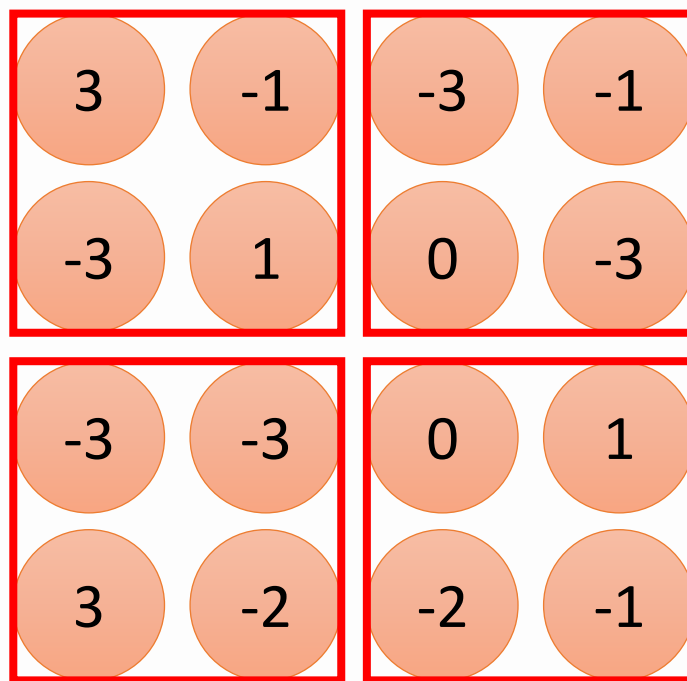


# CNN



- 池化层

- 最大池化
- 平均池化



# CNN



## Property 1

- Some patterns are much smaller than the whole image

## Property 2

- The same patterns appear in different regions.

## Property 3

- Subsampling the pixels will not change the object

Convolution

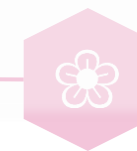
Max Pooling

Convolution

Max Pooling

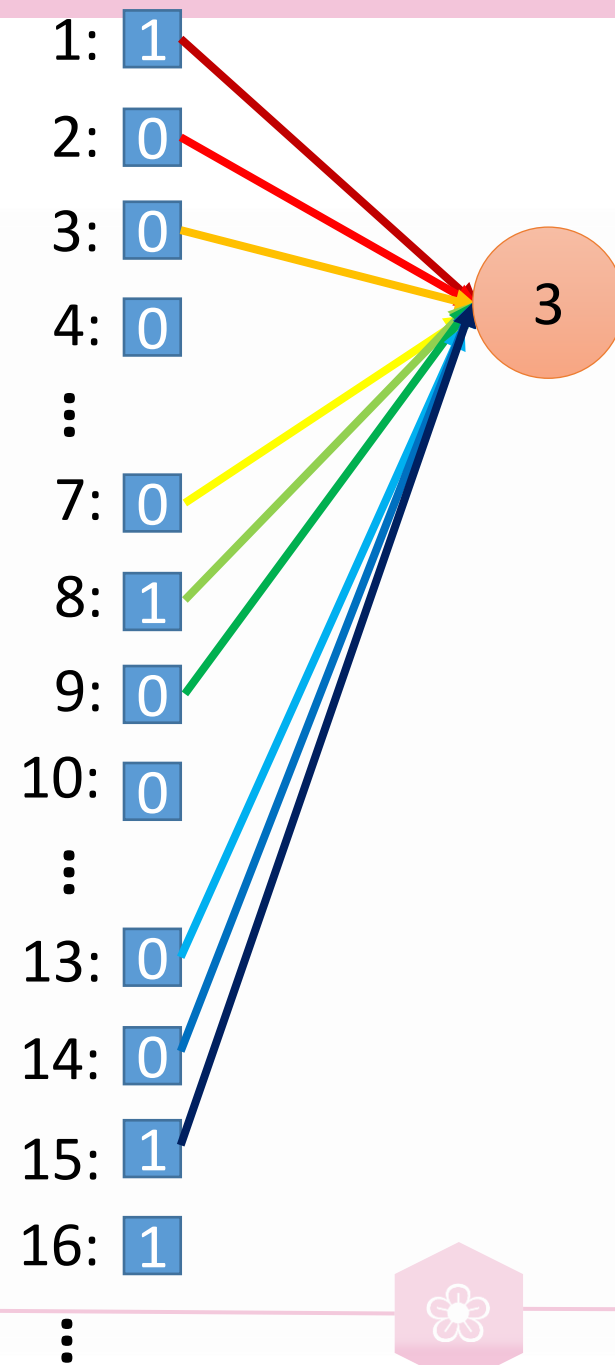
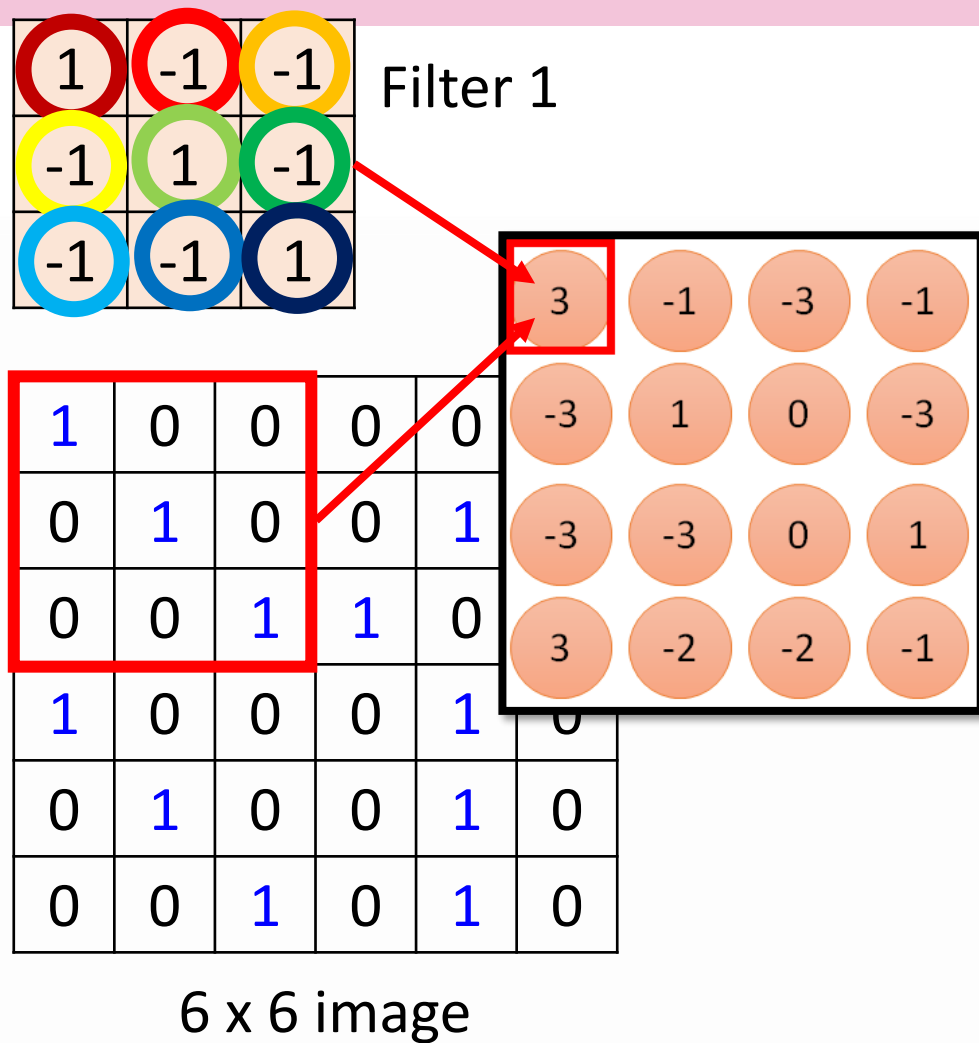
Flatten

Can repeat many times



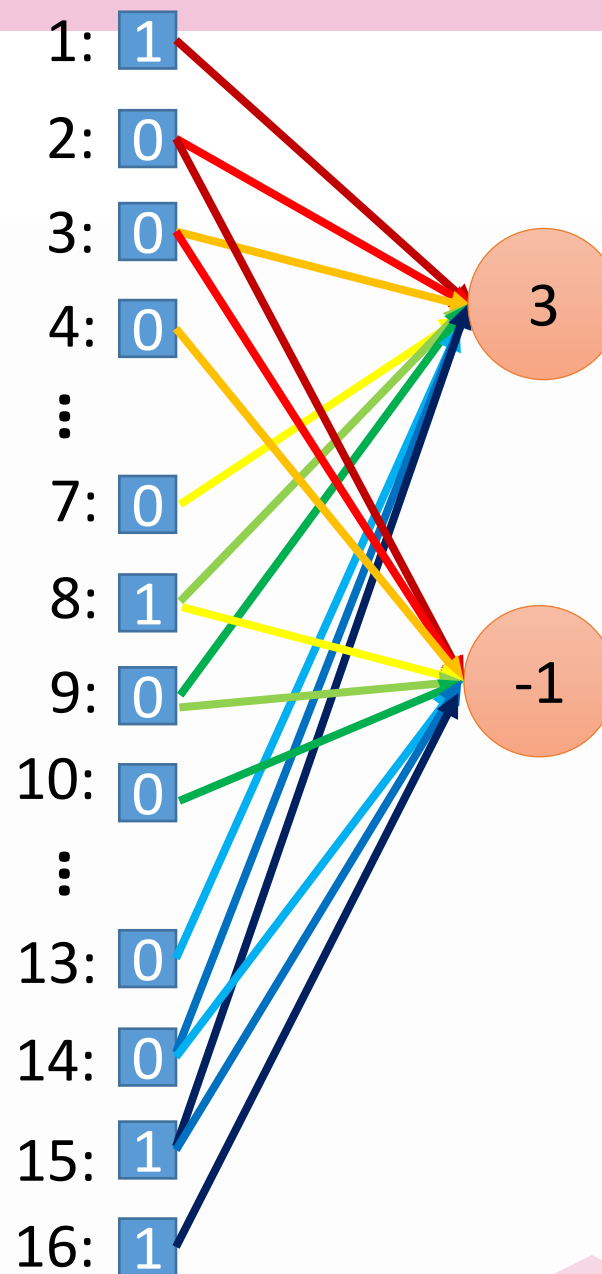
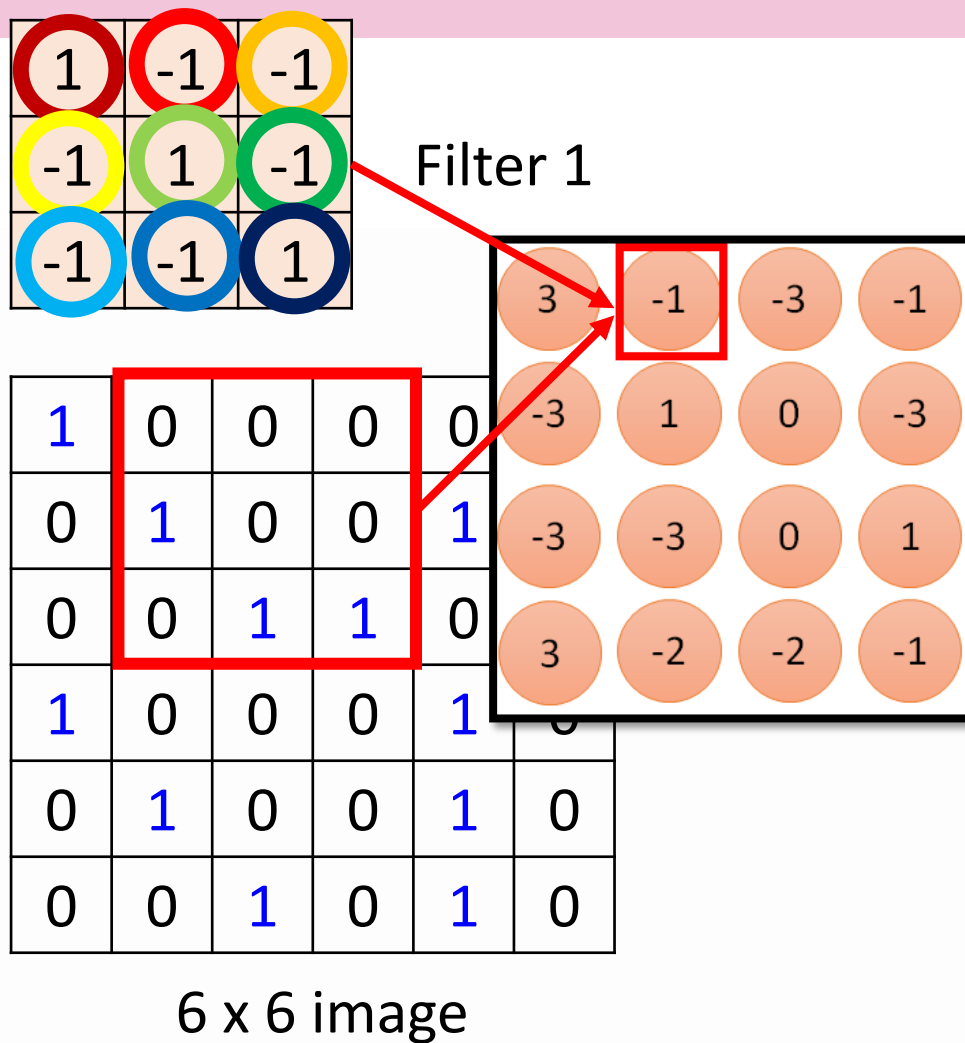
# CNN

- 卷积层



# CNN

- 卷积层



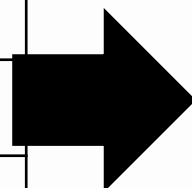


# CNN



1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

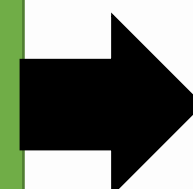
6 x 6 image



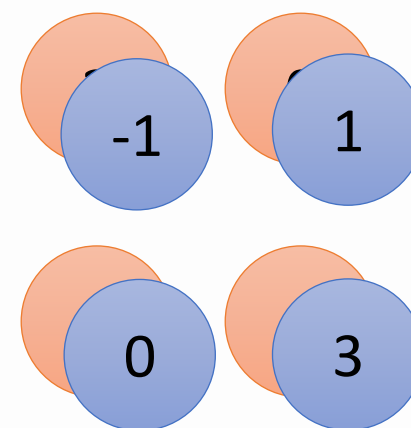
Conv



Max  
Pooling

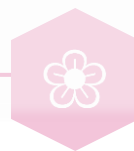


新的图像



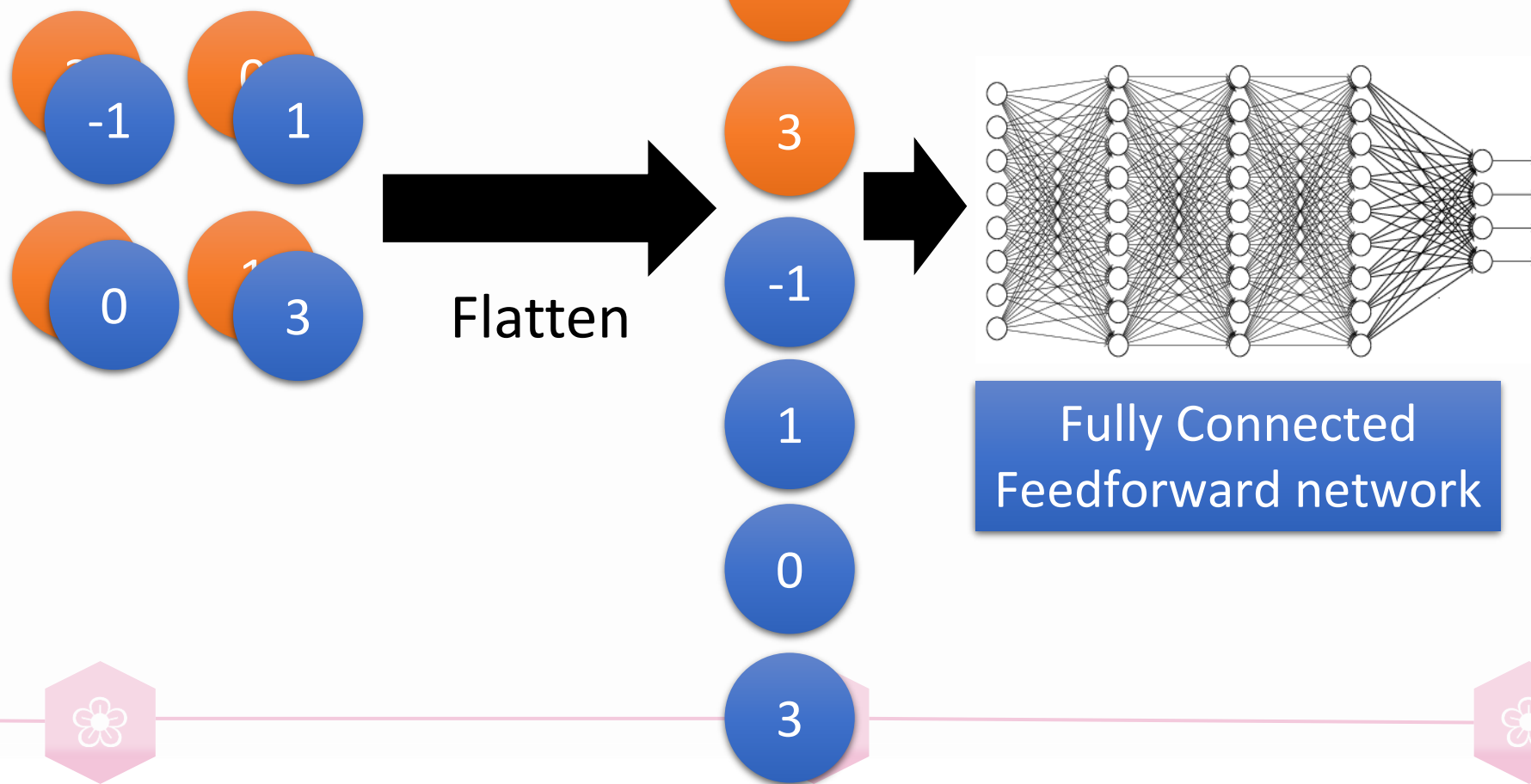
2 x 2 image

Each filter  
is a channel

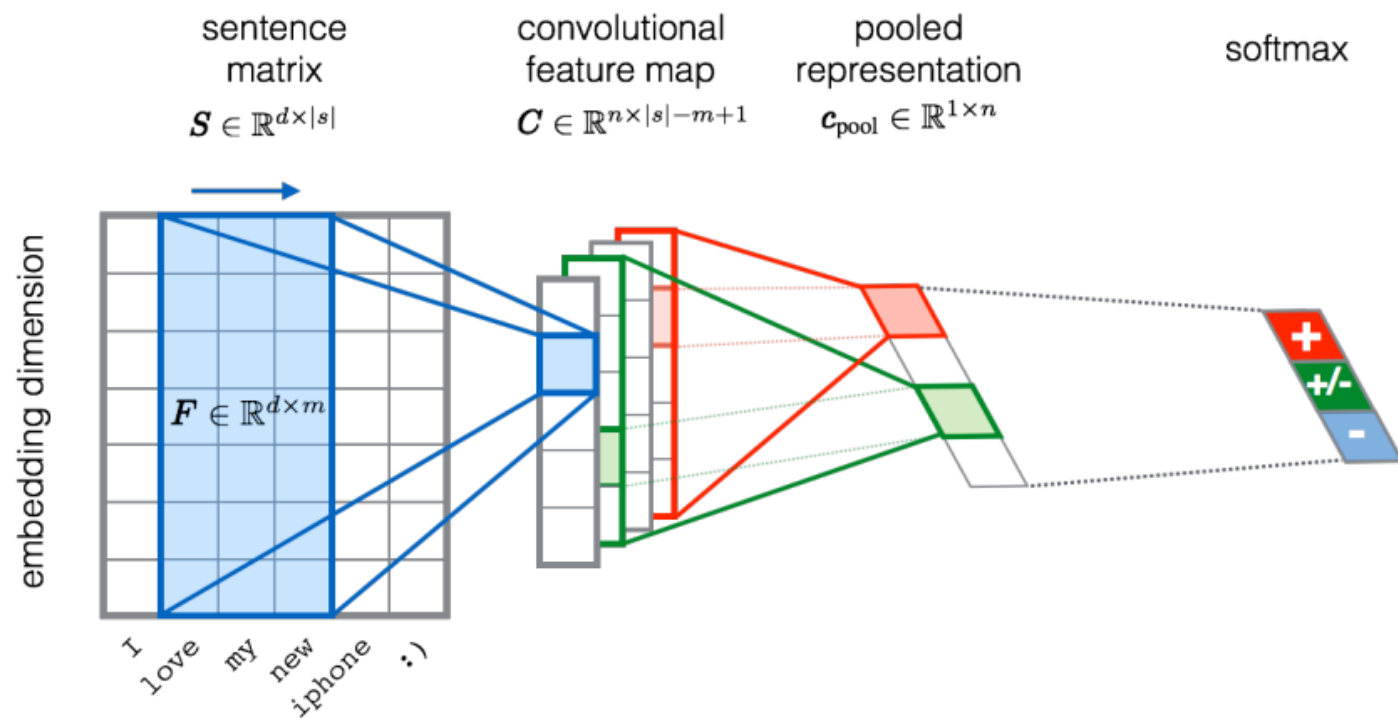
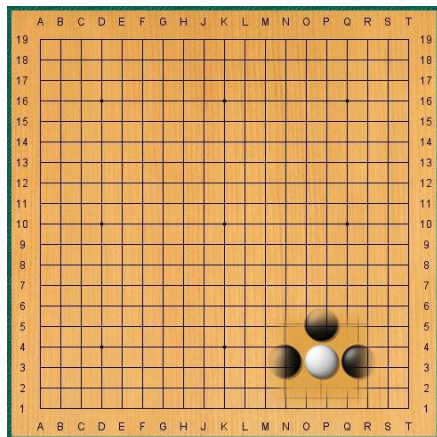
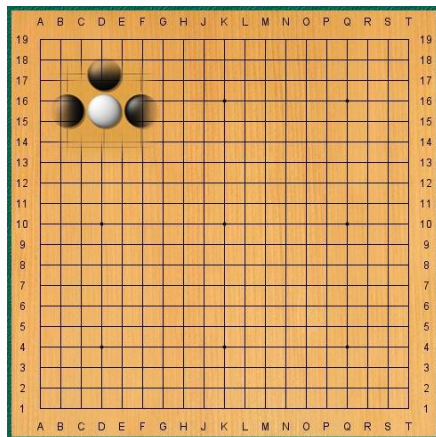


# CNN

- faltten



# CNN





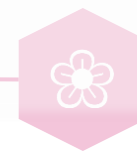
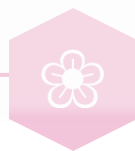
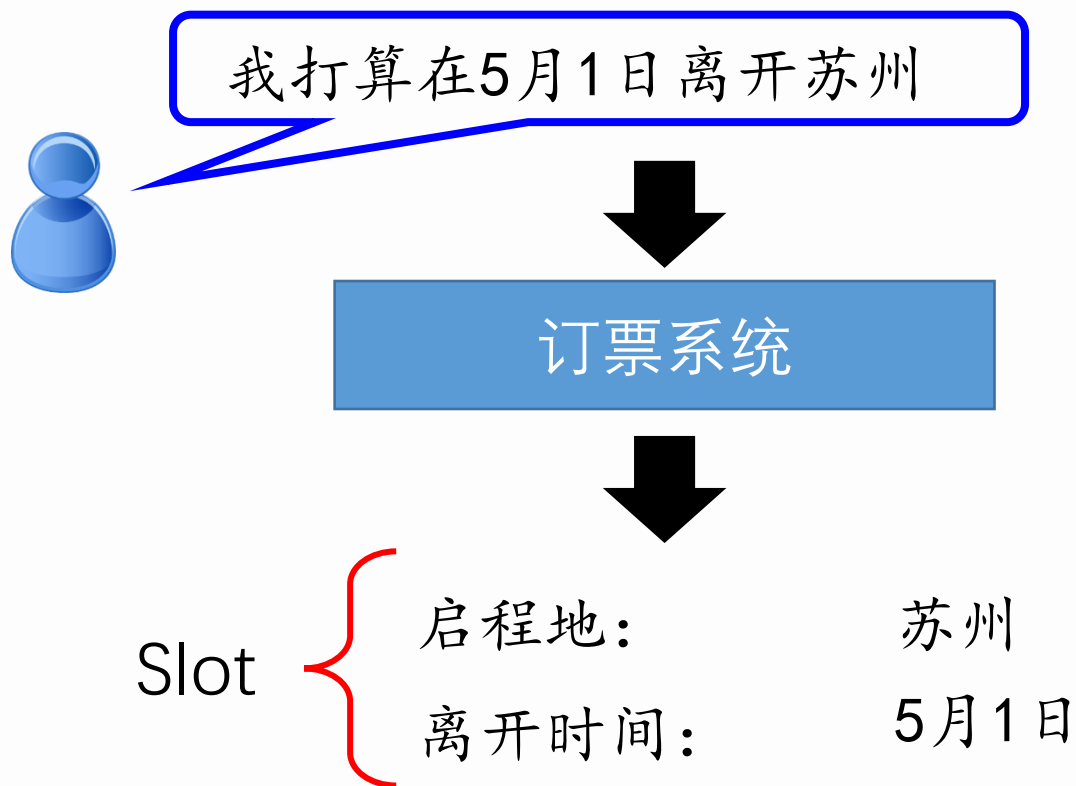
# 03

## 循环神经网络 RNN&LSTM



# 一个小例子

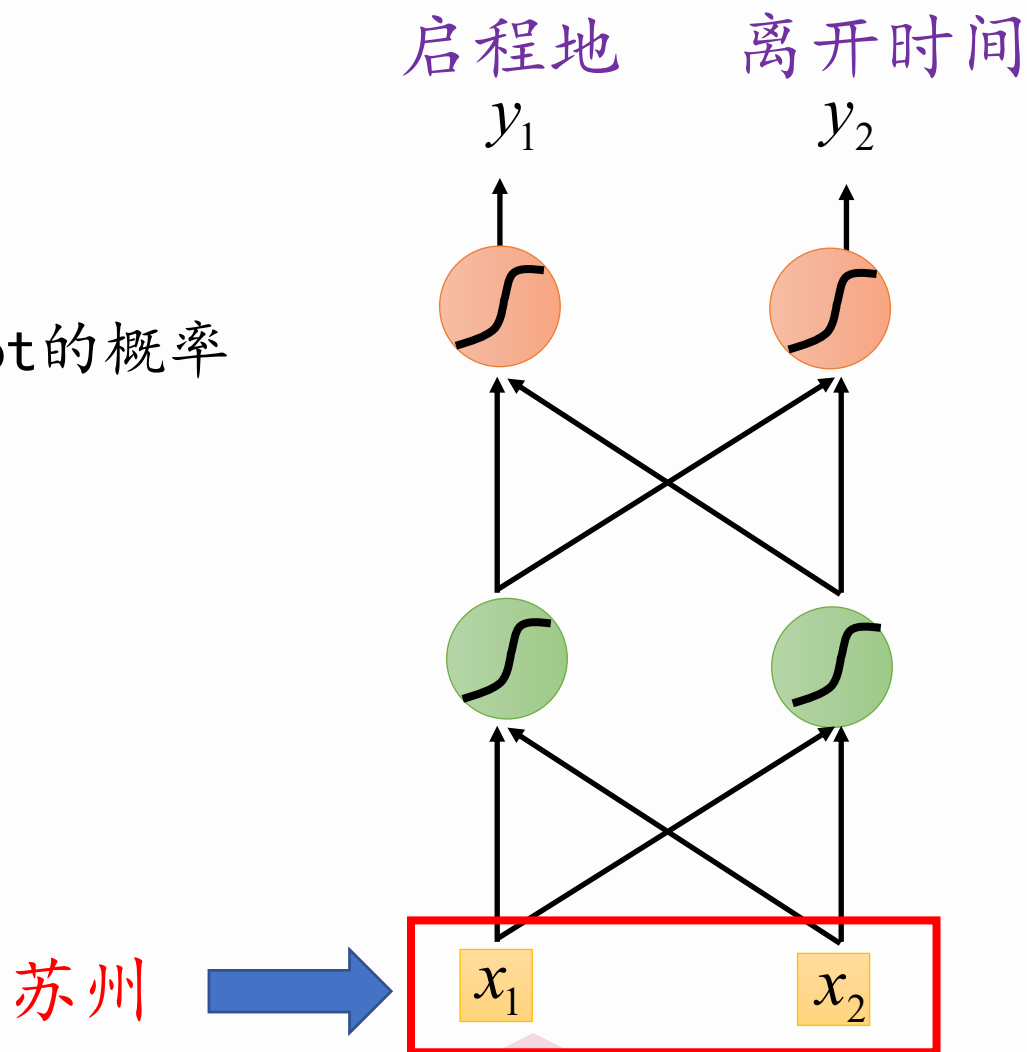
- Slot Filling



# 一个小例子

- 采用前馈神经网络

- 输入词向量
- 输出：隶属不同slot的概率





# 一个小例子

我打算在5月1日离开苏州

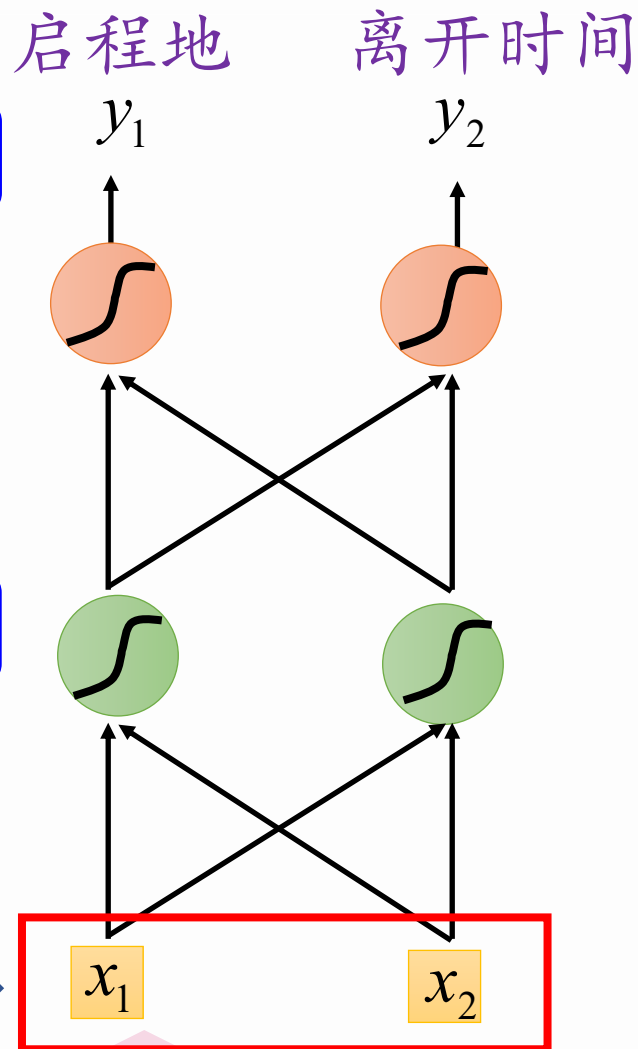
other other other 时间 other 启程地

我打算在5月1日到达苏州

目的地

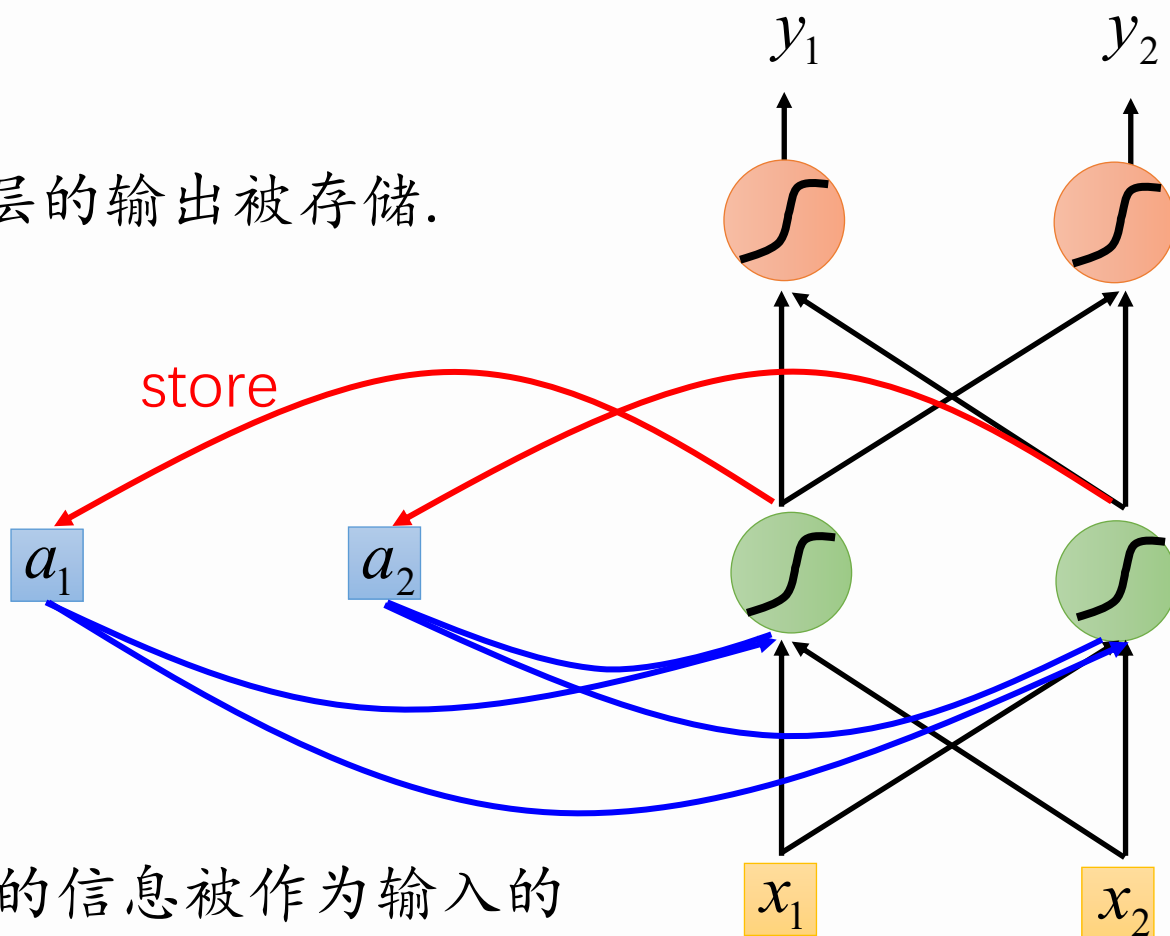
Neural network  
needs memory!

苏州

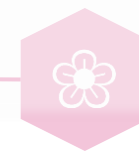
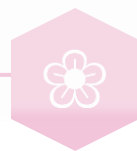


# 循环神经网络RNN

隐层的输出被存储.



存储的信息被作为输入的一部分



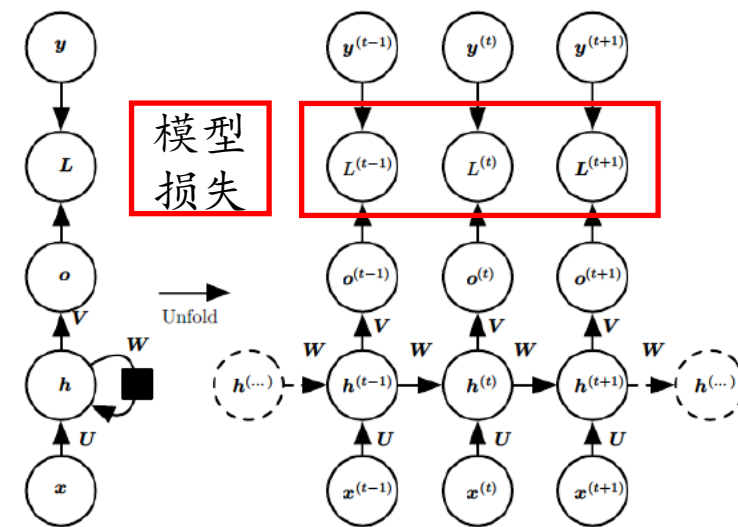
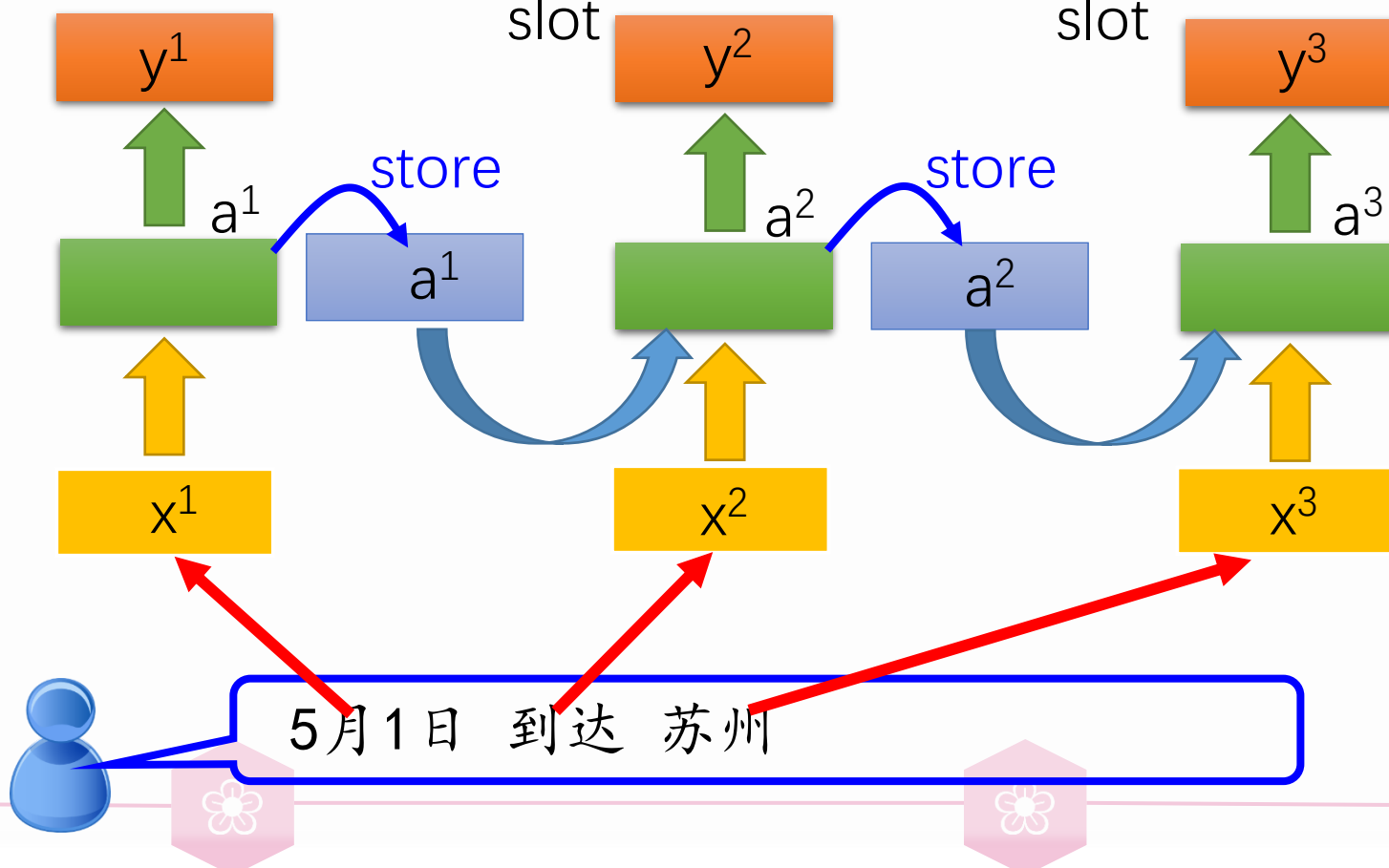
# RNN



Probability of “time”  
in each slot

Probability of  
“Arrive” in each  
slot

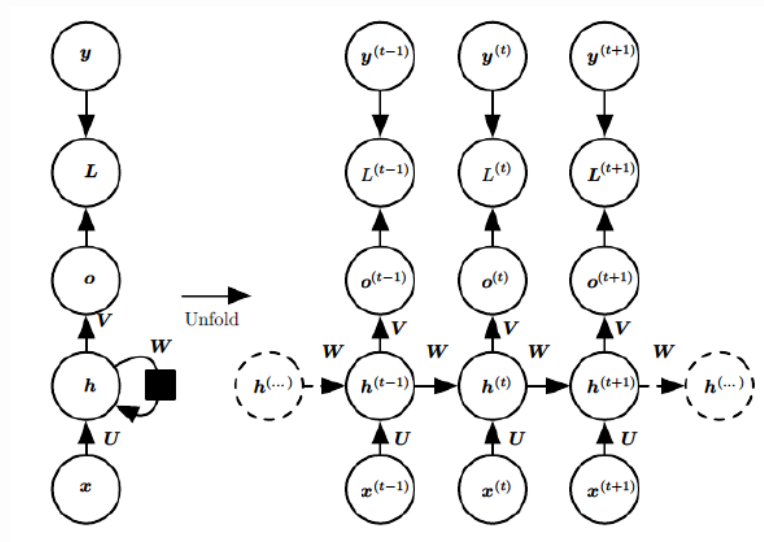
Probability of  
“Suzhou” in each  
slot



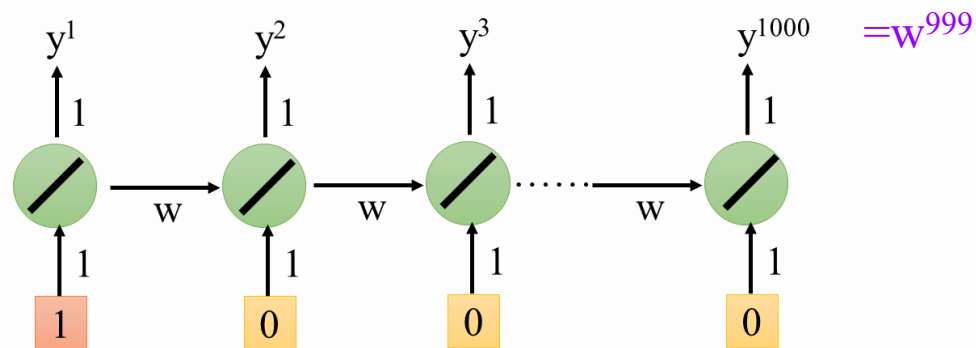
# RNN

## • 前向传播

- 任一序列索引号  $t$
- $h^t = \sigma(z^t) = \sigma(\mathbf{U}x^t + \mathbf{W}h^{t-1})$
- $\hat{y}^t = \sigma(o^t) = \sigma(\mathbf{V}h^t)$



$$\Rightarrow h^t = \sigma(\mathbf{U}x^t + \mathbf{W}(\sigma(\mathbf{U}x^{t-1} + \mathbf{W}h^{t-2})))$$



$$\begin{aligned} w = 1 &\Rightarrow y^{1000} = 1 \\ w = 1.01 &\Rightarrow y^{1000} \approx 20000 \end{aligned}$$

Large  $\partial L / \partial w$

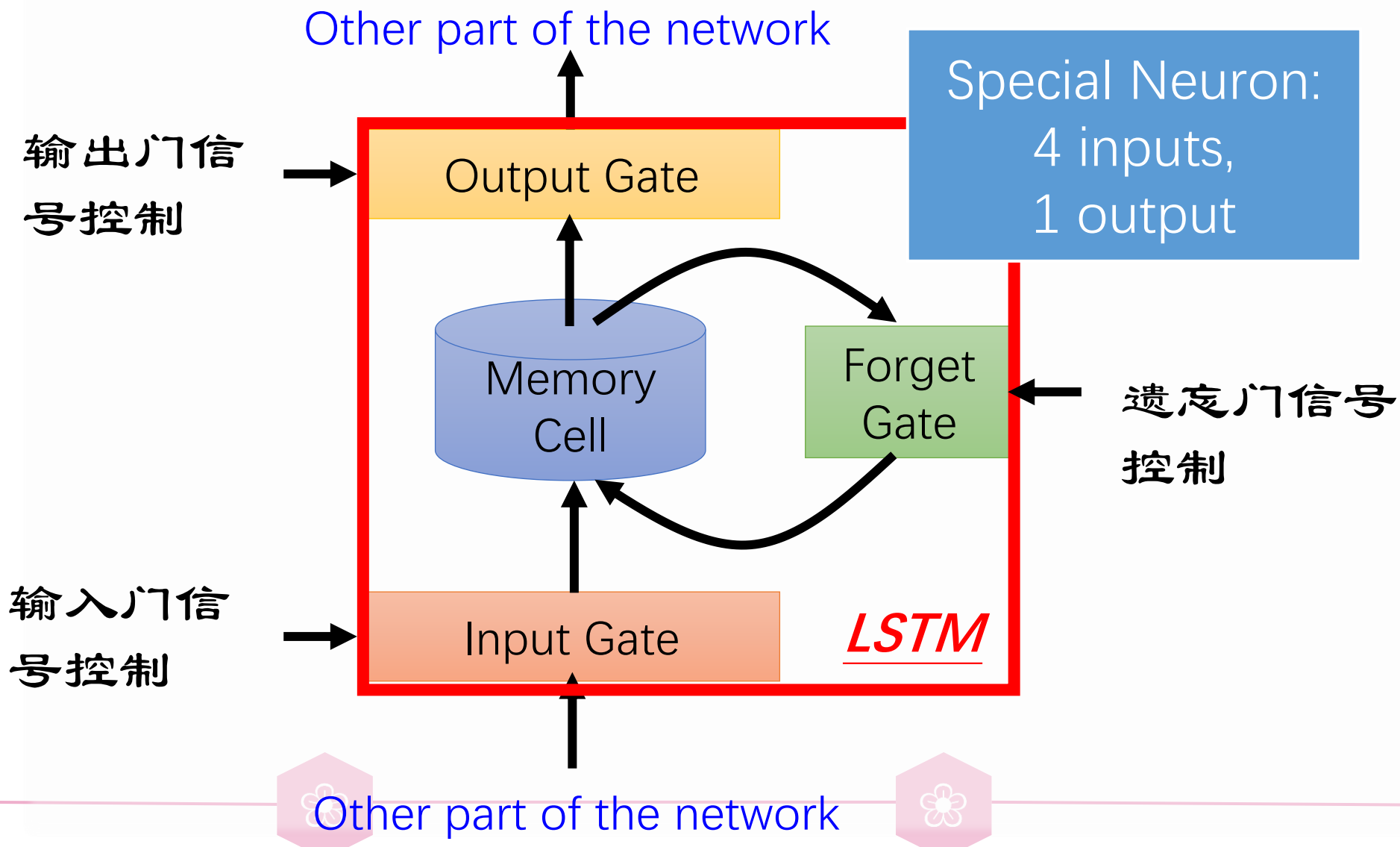
Small Learning rate?

$$\begin{aligned} w = 0.99 &\Rightarrow y^{1000} \approx 0 \\ w = 0.01 &\Rightarrow y^{1000} \approx 0 \end{aligned}$$

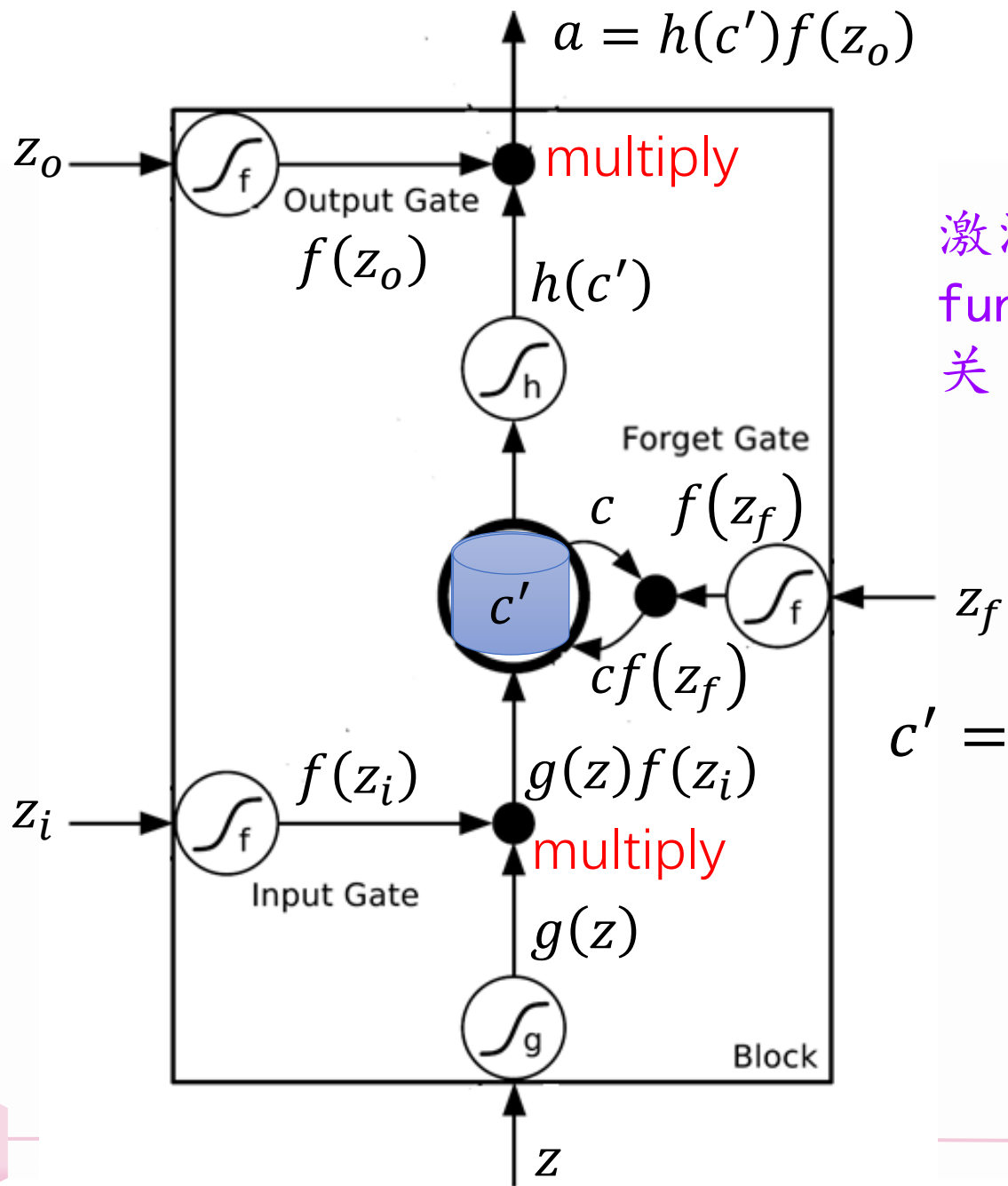
small  $\partial L / \partial w$

Large Learning rate?

# LSTM



# LSTM



激活函数常用sigmoid function (0~1) 来模拟门开关

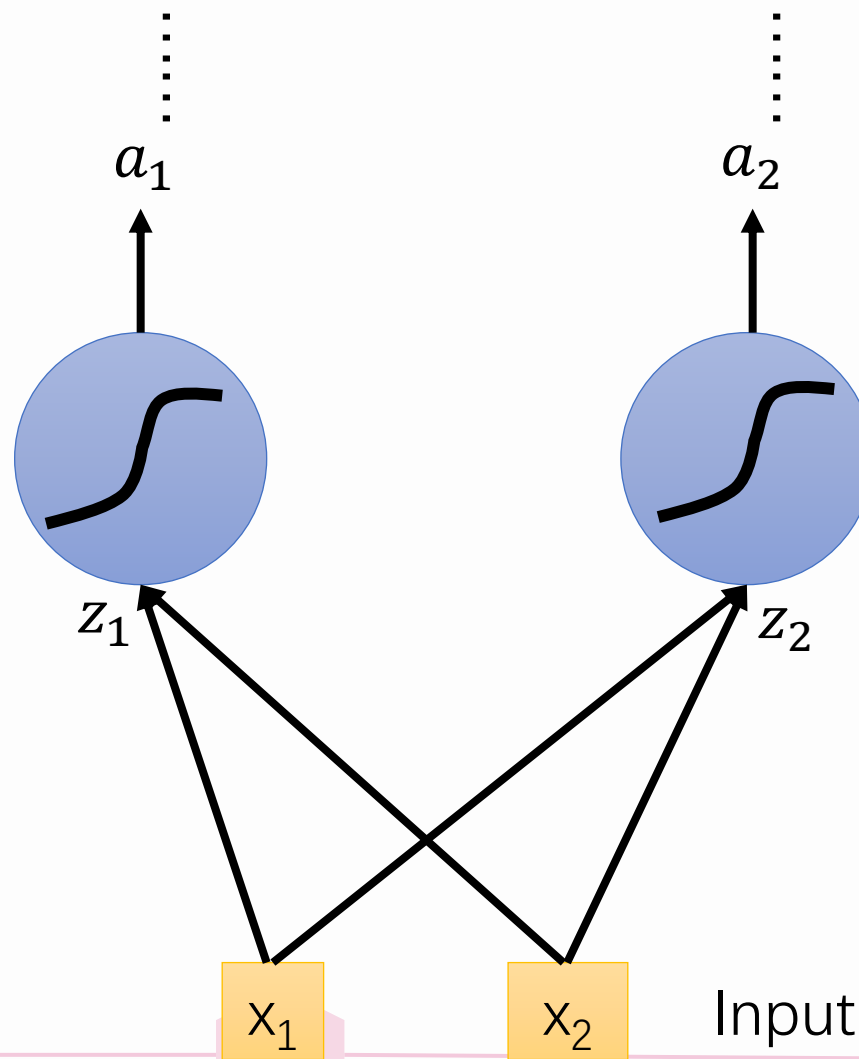
$$c' = g(z)f(z_i) + cf(z_f)$$



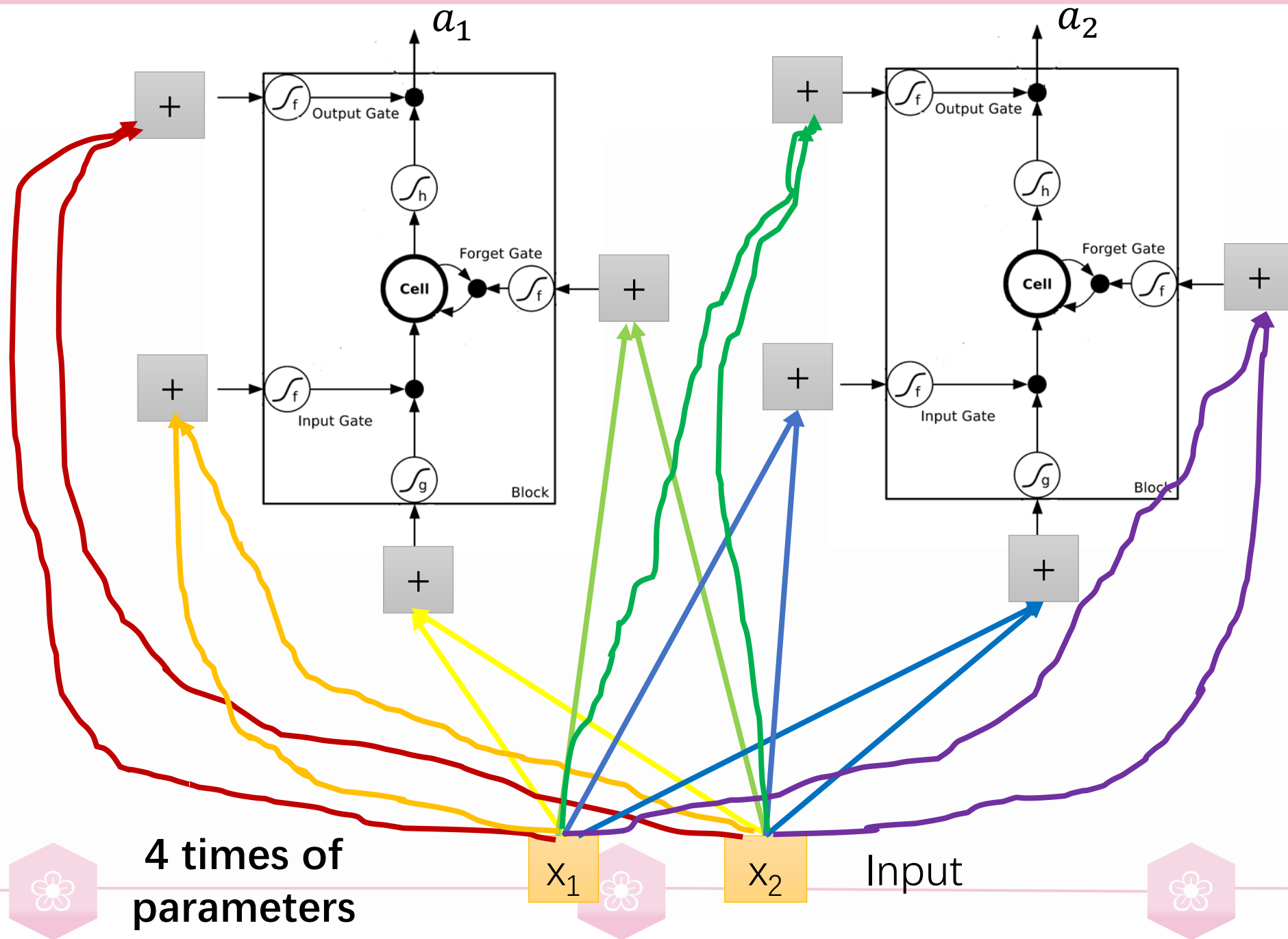
# LSTM



- 常规形式



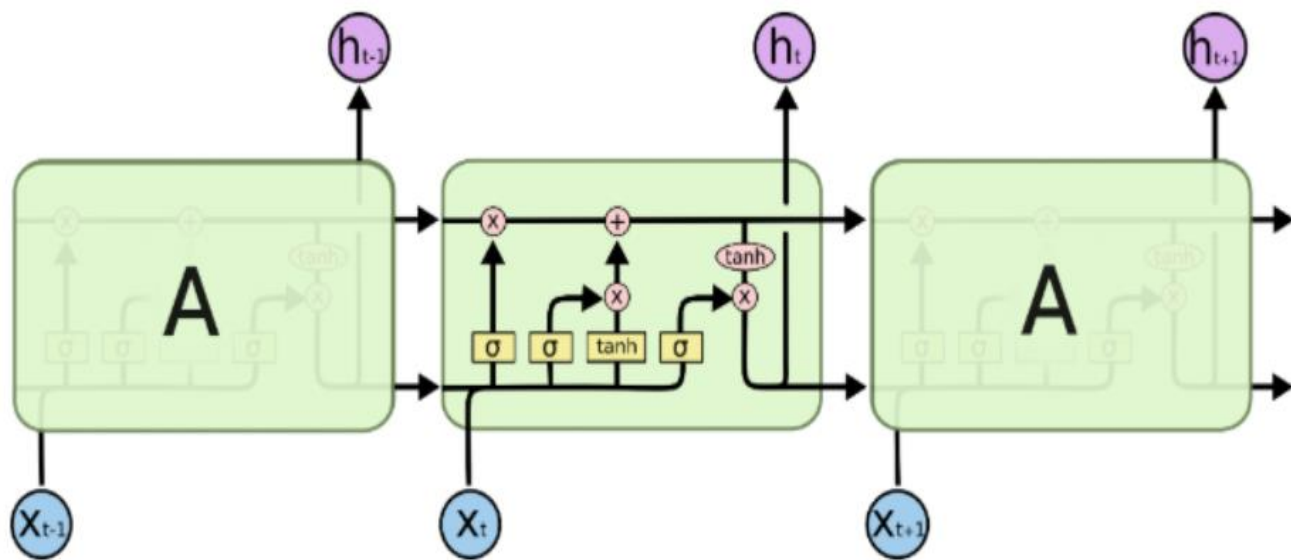
# LSTM



# LSTM

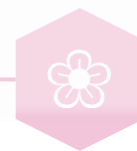
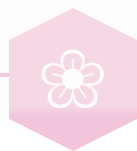
## • 前向传播

- **遗忘门**  $f^t = \sigma(W_f h^{t-1} + U_f x^t + b_f)$
- 输入部分:
  - **输入门**  $i^t = \sigma(W_i h^{t-1} + U_i x^t + b_i)$
  - **输入**  $a^t = \tanh(W_a h^{t-1} + U_a x^t + b_a)$
- 记忆存储部分
  - **cell更新**  $C^t = C^{t-1} \odot f^t + i^t \odot a^t$
- 输出部分:
  - **输出门**  $o^t = \sigma(W_o h^{t-1} + U_o x^t + b_o)$
  - **输出**  $h^t = o^t \odot \tanh(C^t)$
- $y^t = \sigma(Vh^t + c)$



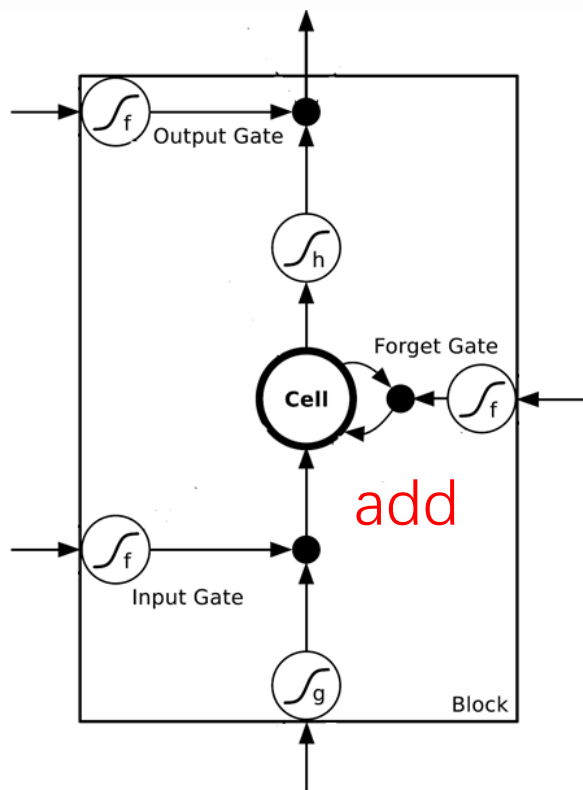
# RNN

- $h^t = \sigma(z^t) = \sigma(\mathbf{U}x^t + \mathbf{W}h^{t-1})$
- $\frac{\partial L}{\partial W} = \sum_{t=1}^{\tau} \sum_{k=1}^t \frac{\partial L_t}{\partial \hat{y}^t} \frac{\partial \hat{y}^t}{\partial h^t} \prod_{j=k+1}^t \frac{\partial h^j}{\partial h^{j-1}} \frac{\partial h^k}{\partial W}$
- $\frac{\partial h^j}{\partial h^{j-1}} = \sigma' W$ 
  - 小于1, j和k距离过大, 梯度消失
  - 大于1, j和k距离过大, 梯度爆炸
  - 长程依赖问题



# LSTM

- 记忆部分是相加
- 处理梯度消失，但不能处理梯度爆炸



$$C^t = C^{t-1} \odot f^t + i^t \odot a^t$$

$$\frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial C_t}{\partial f^t} \frac{\partial f_t}{\partial h^{t-1}} \frac{\partial h_{t-1}}{\partial C_{t-1}} + \frac{\partial C_t}{\partial i^t} \frac{\partial i_t}{\partial h^{t-1}} \frac{\partial h_{t-1}}{\partial C_{t-1}} + \frac{\partial C_t}{\partial a^t} \frac{\partial a_t}{\partial h^{t-1}} \frac{\partial h_{t-1}}{\partial C_{t-1}} + f^t$$