

深度学习之旅

模型、算法与实践

张钟元

2018-01-19

内容

- 1 破冰之旅
- 2 前馈神经网络
- 3 卷积神经网络
- 4 循环神经网络
- 5 长短时记忆网络
- 6 TensorFlow 实现
- 7 致谢

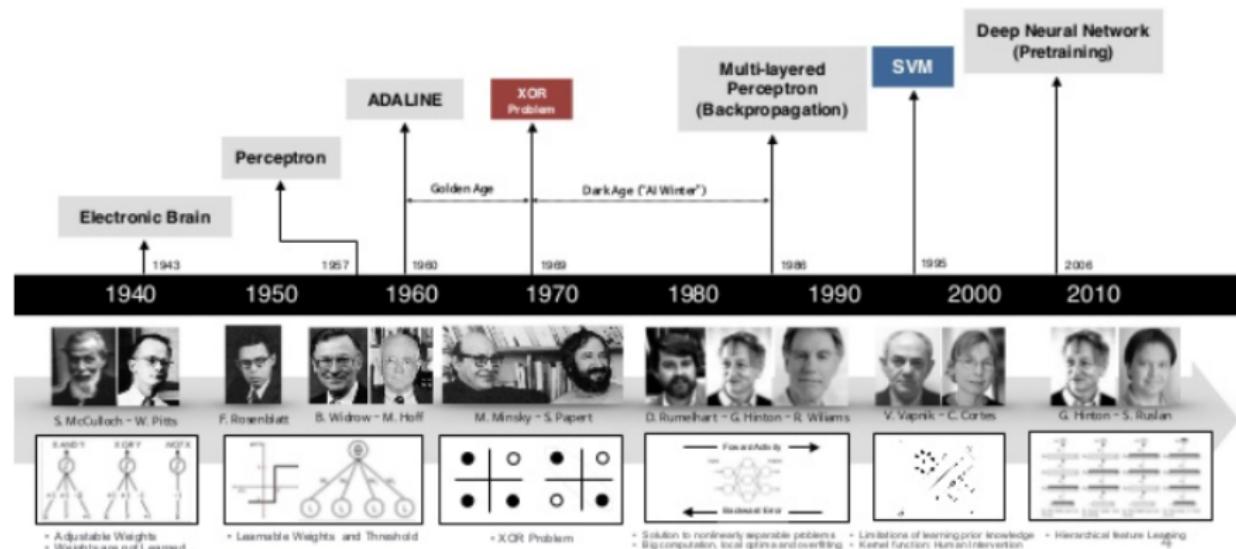


破冰之旅

- ① 历史
- ② 破冰
- ③ 架构

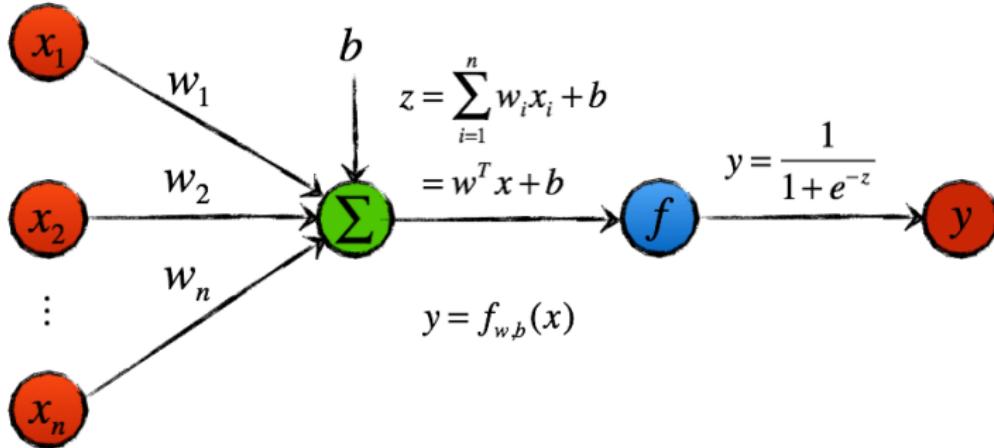
历史

历史



破冰

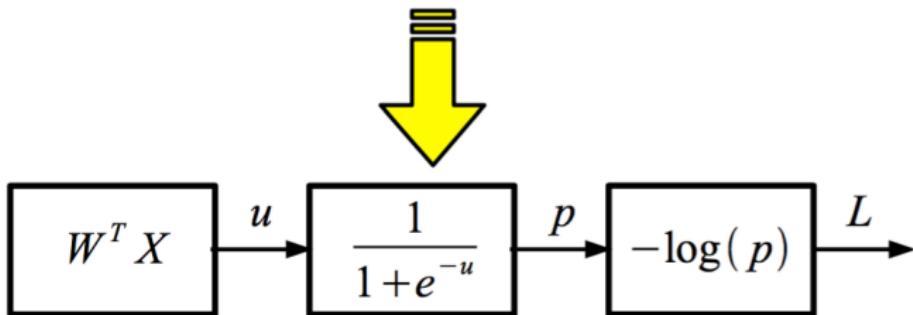
逻辑回归



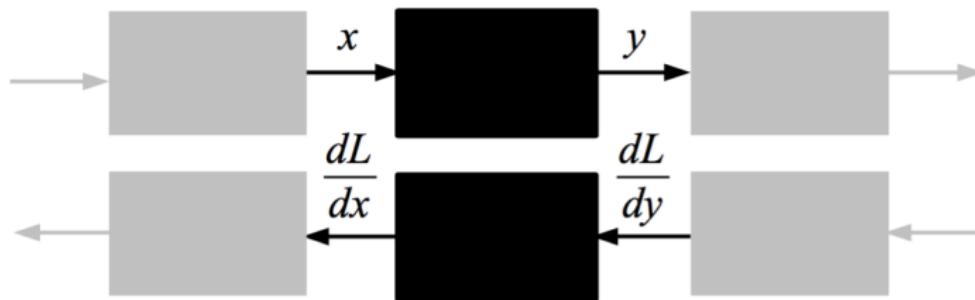
计算损失：函数组合

Complicated Function

$$-\log\left(\frac{1}{1+e^{-W^T X}}\right)$$



链式法则



Given $y(x)$ and dL/dy ,

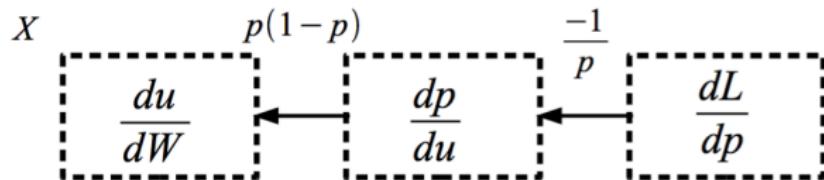
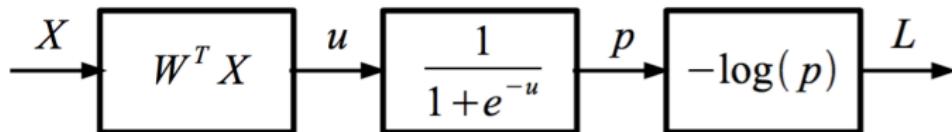
What is dL/dx ?

$$\rightarrow \frac{dL}{dx} = \frac{dL}{dy} \cdot \frac{dy}{dx}$$

All needed information is local!

破冰

计算梯度



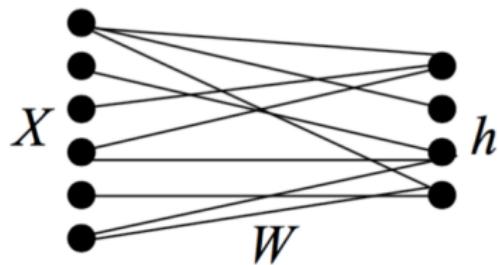
$$\frac{dL}{dW} = \frac{dL}{dp} \cdot \frac{dp}{du} \cdot \frac{du}{dW} = (p-1)X$$

网络架构

单层网络

$$\xrightarrow{} f(X; W) \xrightarrow{}$$

is equivalent to

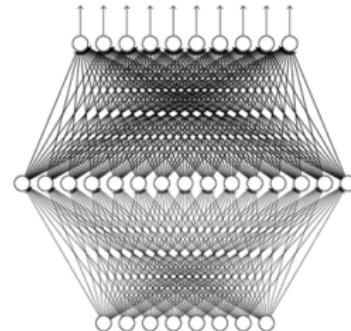


万能近似定理

Any continuous function

$$f : R^N \rightarrow R^M$$

Can be realized by a network
with one hidden layer, given
enough hidden neurons)



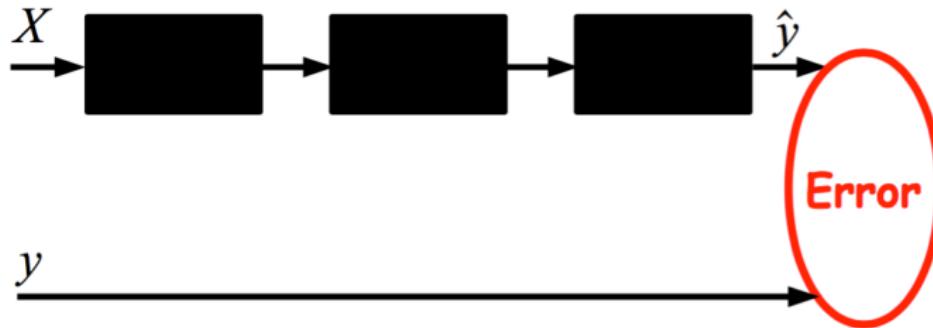
网络架构

多层网络：深度学习



网络架构

回归

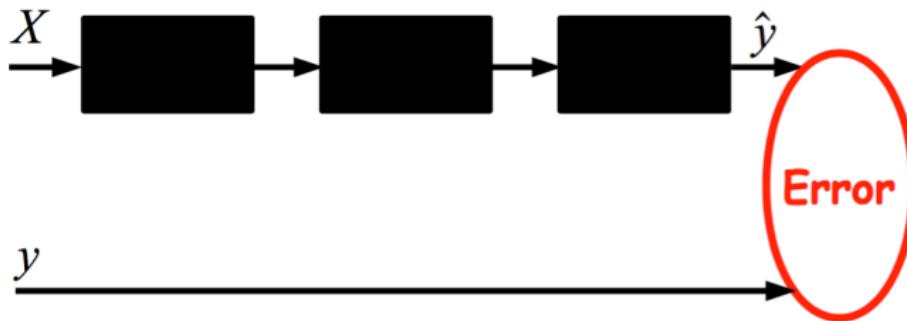


Square Euclidean Distance (regression):

$$y, \hat{y} \in R^N$$

$$L = \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

分类



Cross Entropy (classification):

$$y, \hat{y} \in [0,1]^N, \sum_{i=1}^N y_i = 1, \sum_{i=1}^N \hat{y}_i = 1$$

$$L = -\sum_{i=1}^N y_i \log \hat{y}_i$$

前向计算

A) Compute loss on small mini-batch

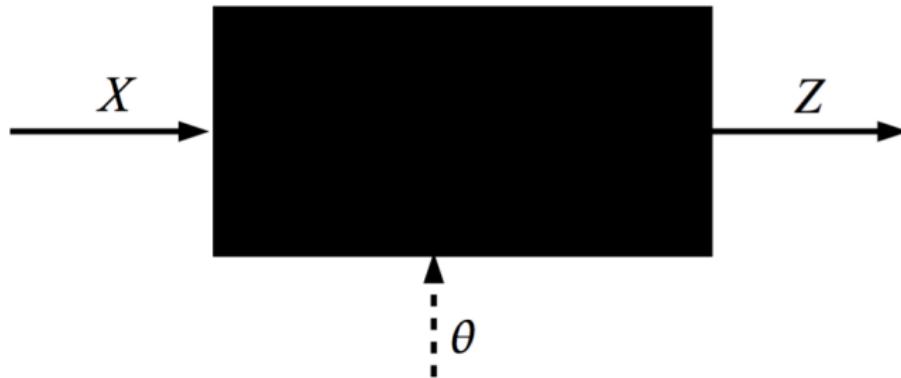
F-PROP



BP 算法

前向计算

F-PROP



反向传播

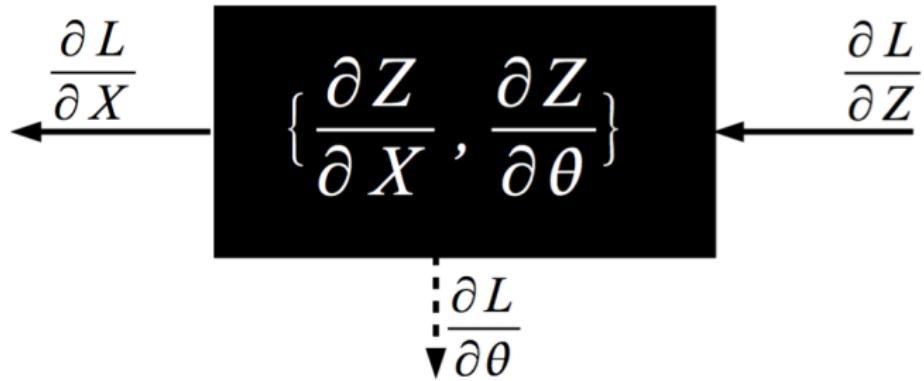
- A) Compute loss on small mini-batch
- B) Compute gradient w.r.t. parameters

B-PROP



反向传播

B-PROP



更新参数

- A) Compute loss on small mini-batch
- B) Compute gradient w.r.t. parameters
- C) Use gradient to update parameters $\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$



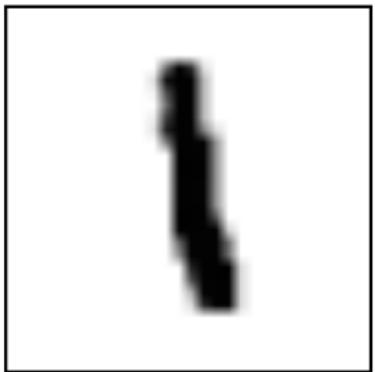


前馈神经网络

- ① 单层前馈网络
- ② 多层前馈网络

提出问题

图片： $28 \times 28 = 784$



~

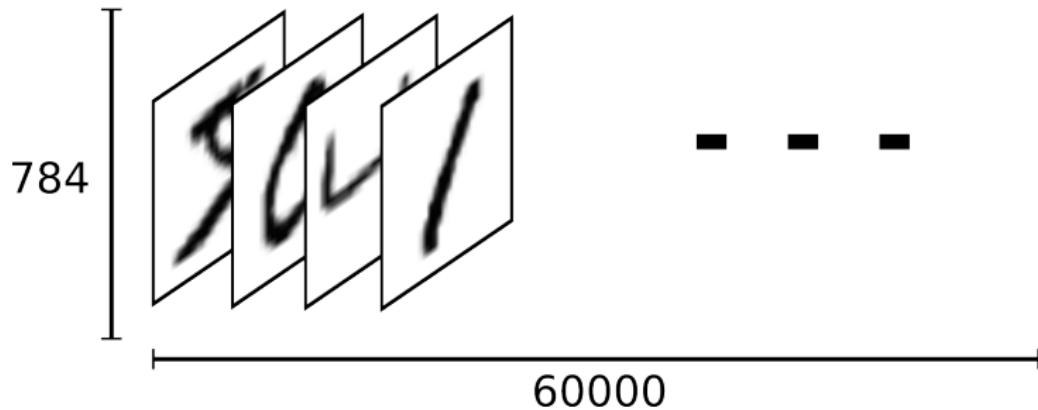
$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{6} & \textcolor{brown}{8} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{7} & \textcolor{brown}{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{7} & \textcolor{brown}{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{5} & \textcolor{brown}{1} & \textcolor{brown}{4} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{1} & \textcolor{brown}{4} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{1} & \textcolor{brown}{4} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{1} & \textcolor{brown}{7} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{1} & \textcolor{brown}{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{9} & \textcolor{brown}{1} & \textcolor{brown}{1} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcolor{brown}{3} & \textcolor{brown}{1} & \textcolor{brown}{1} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

A decorative horizontal border consisting of a repeating pattern of small, solid black circles.

提出问题

训练数据集输入: [100(*batch_size*), 784]

mnist.train.xs



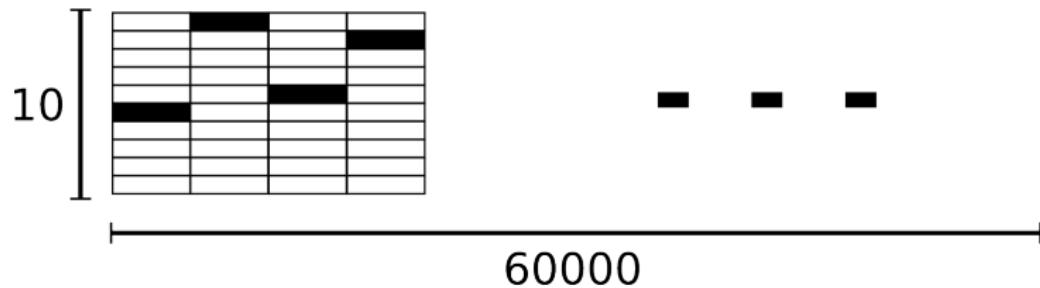
000000

提出问题

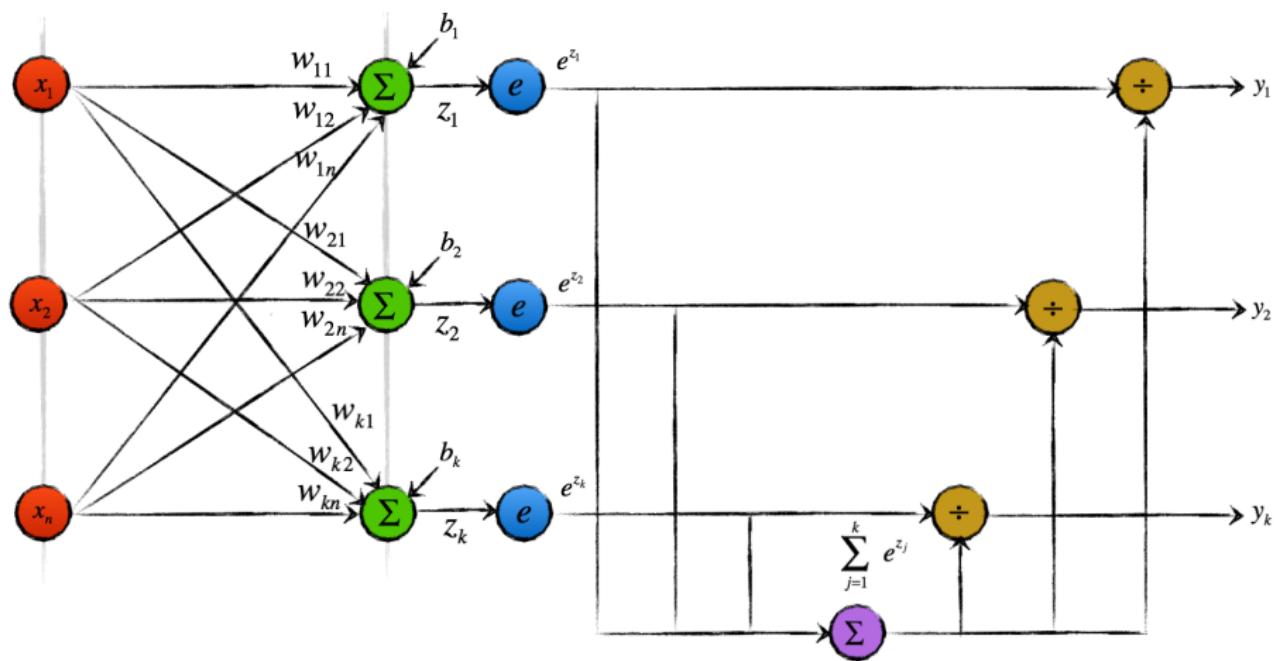
提出问题

训练数据集输出: [100(*batch_size*), 10]

mnist.train.ys



Softmax 模型



k 分类问题

参数定义

$$W = \begin{pmatrix} w_{11} & \dots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{k1} & \dots & w_{kn} \end{pmatrix} \in \mathbb{R}^{k \times n}$$

$$b = (b_1, b_2, \dots, b_k)^T \in \mathbb{R}^k$$

$$x = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$$

模型定义

$$z = Wx + b \in \mathbb{R}^k$$

$$y = softmax(z) \in \mathbb{R}^k$$

$$softmax(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}, \quad i = 1, 2, \dots, k$$

损失函数

交叉熵损失函数

$$\begin{aligned} L(W, b; S) &= -\frac{1}{m} \sum_{i=1}^m t^{(i)} \odot \log(y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^k t_j^{(i)} \log(y_j^{(i)}) \end{aligned}$$

优化问题

$$W^*, b^* = \arg \min_{W, b} L(W, b; S)$$

oooooooooooooooooooo

ooooooo●oooooooooooo

oooooooooooooooooooo

oooooooooooooooooooo

单层前馈网络

优化算法

梯度下降

$$W \leftarrow W - \alpha \nabla_W L(W, b)$$

$$b \leftarrow b - \alpha \nabla_b L(W, b)$$

梯度计算

$$\nabla_W L(W, b; x, t) = (y - t) \otimes x$$

$$\nabla_b L(W, b; x, t) = (y - t)$$

SGD & BGD → SGD with MiniBatch

SGD

$$\begin{aligned} w &\leftarrow w - \alpha \left(y^{(i)} - t^{(i)} \right) x \\ b &\leftarrow b - \alpha \left(y^{(i)} - t^{(i)} \right) \end{aligned}$$

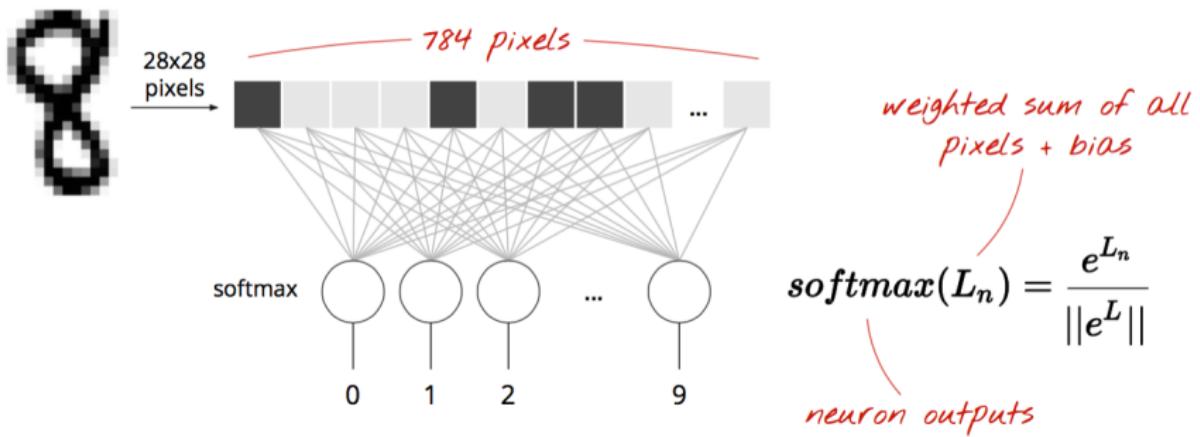
BGD

$$\begin{aligned} w &\leftarrow w - \alpha \left(\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - t^{(i)} \right) \right) x \\ b &\leftarrow b - \alpha \left(\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - t^{(i)} \right) \right) \end{aligned}$$



Softmax 实现

单层网络模型



占位符

```
# inputs  
x = tf.placeholder("float", [None, 784])  
  
# labels  
t = tf.placeholder("float", [None, 10])
```

- `tf.placeholder` 定义了一个占位 OP
- `None` 表示未确定的样本数目 (`batch_size`)
- `Session.run` 时提供 `feed_dict` 提供一个批次的样本数据

Softmax 实现

变量

```
# train parameters
W = tf.Variable(tf.zeros([784,10]))
b = tf.Variable(tf.zeros([10]))

# init_op
init_op = tf.global_variables_initializer()
```

- **变量**: 使用 `tf.Variable` 定义变量，常用于定义模型的权重和偏置
- **初始化**: 使用初始化 `init_op` 初始化所有全局变量

Softmax 实现

模型: Softmax

```
y = tf.nn.softmax(tf.matmul(x, W) + b)
```

$$Y = \text{softmax}(X \cdot W + b)$$

Predictions Images Weights Biases
 $Y[100, 10]$ $X[100, 784]$ $W[784, 10]$ $b[10]$

tensor shapes in []
 applied line by line matrix multiply broadcast on all lines

Softmax 实现

损失函数

```
cross_entropy = -tf.reduce_sum(t * tf.log(y))
```

0	1	2	3	4	5	6	7	8	9
0	0	0	0	0	0	1	0	0	0

actual probabilities, "one-hot" encoded

$$\text{Cross entropy: } - \sum Y'_i \cdot \log(Y_i)$$

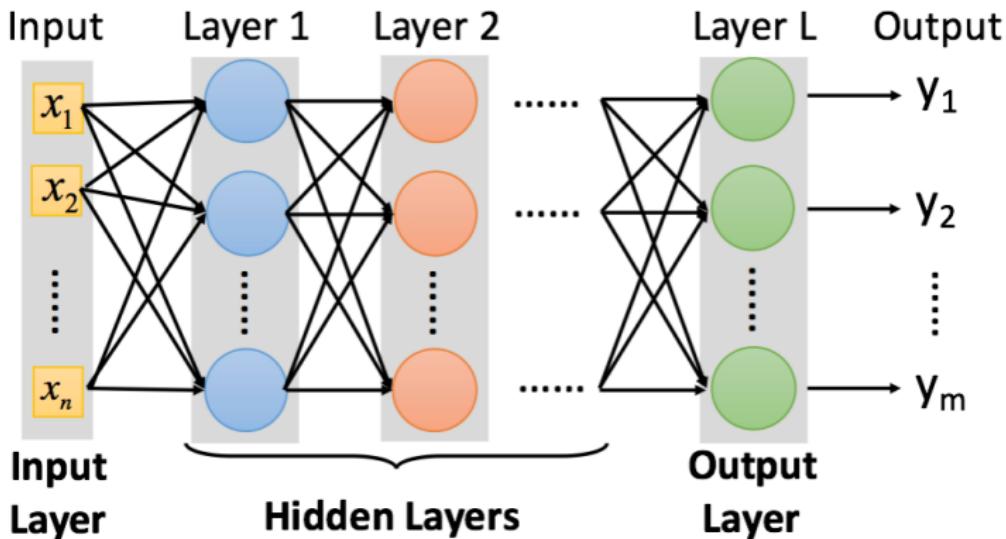
)
computed probabilities
this is a "6"

0.1	0.2	0.1	0.3	0.2	0.1	0.9	0.2	0.1	0.1
0	1	2	3	4	5	6	7	8	9

oooooooooooooooooo 前馈神经网络 oooooooooooooooo●ooooooooooooo 卷积神经网络 oooooooooooooooooo oooooooooooooooo 长短时记忆网络 oooooooooooooooooo

多层前馈网络

多层前馈网络

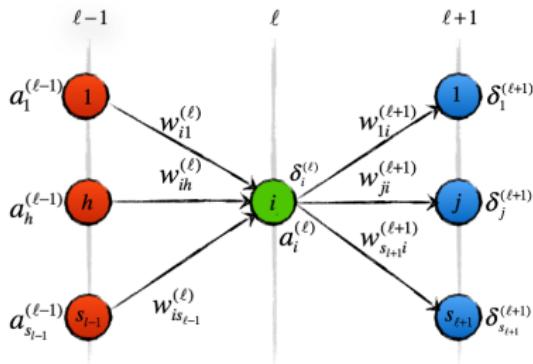


符号

- n_ℓ : 网络层数, 其中第 0 层为输入层, 第 n_ℓ 层为输出层
 - s_ℓ : 第 ℓ 层的节点数, $\ell = 0, 1, \dots, n_\ell$
 - $w_{ji}^{(\ell)}$: 第 $(\ell - 1)$ 层节点 i 与第 ℓ 层节点 j 之间的权重, $\ell = 1, \dots, n_\ell$
 - $b_i^{(\ell)}$: 第 ℓ 层节点 i 的偏置项, $\ell = 1, \dots, n_\ell$
 - $a_i^{(\ell)}$: 第 ℓ 层节点 i 的输出, $\ell = 1, \dots, n_\ell, x = a^{(0)}, y = a^{(n_\ell)}$
 - $z_i^{(\ell)}$: 第 ℓ 层节点 i 的权重和, $\ell = 1, \dots, n_\ell$
 - $\delta_i^{(\ell)}$: 第 ℓ 层节点 i 的误差项, $\ell = 1, \dots, n_\ell$
 - $S = \{(x^{(t)}, y^{(t)}); t = 1, 2, \dots, m\}$: 样本空间

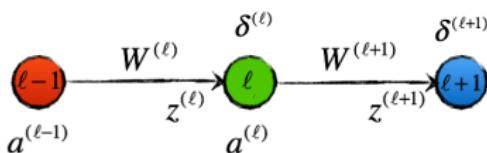
多层前馈网络

计算梯度 (标量)



$$\begin{aligned}
 \nabla_{w_{ih}^{(\ell)}} L(w, b) &= \frac{\partial L}{\partial w_{ih}^{(\ell)}} = \frac{\partial L}{\partial z_i^{(\ell)}} \frac{\partial z_i^{(\ell)}}{\partial w_{ih}^{(\ell)}} = \delta_i^{(\ell)} a_h^{(\ell-1)} \\
 z_i^{(\ell)} &= \sum_{h=1}^{s_{\ell-1}} (W_{ih}^{(\ell)} a_h^{(\ell-1)} + b_i^{(\ell)}) \\
 \frac{\partial L}{\partial z_i^{(\ell)}} &= \sum_{j=1}^{s_{\ell+1}} \frac{\partial L}{\partial z_j^{(\ell+1)}} \frac{\partial z_j^{(\ell+1)}}{\partial z_i^{(\ell)}} \frac{\partial a_i^{(\ell)}}{\partial z_i^{(\ell)}} \\
 &= f'(z_i^{(\ell)}) \sum_{j=1}^{s_{\ell+1}} \frac{\partial L}{\partial z_j^{(\ell+1)}} W_{ji}^{(\ell+1)} \\
 \delta_i^{(\ell)} &= f'(z_i^{(\ell)}) \sum_{j=1}^{s_{\ell+1}} \delta_j^{(\ell+1)} W_{ji}^{(\ell+1)}
 \end{aligned}$$

计算梯度(张量)



$$\nabla_{w^{(\ell)}} L(w, b) = \frac{\partial L}{\partial W^{(\ell)}} = \cancel{\frac{\partial L}{\partial z^{(\ell)}}} \cancel{\frac{\partial z^{(\ell)}}{\partial W^{(\ell)}}} = \delta^{(\ell)} (a^{(\ell-1)})^T$$

$$\frac{\partial L}{\partial z^{(\ell)}} = \frac{\partial a^{(\ell)}}{\partial z^{(\ell)}} \frac{\partial z^{(\ell+1)}}{\partial a^{(\ell)}} \frac{\partial L}{\partial z^{(\ell+1)}}$$

$$= f'(\zeta^{(\ell)}) \odot (W^{(\ell+1)})^T \frac{\partial L}{\partial z^{(\ell+1)}}$$

$$\delta^{(\ell)} = f'(z^{(\ell)}) \odot (W^{(\ell+1)})^T \delta^{(\ell+1)}$$

$$\delta^{(n_t)} = f'(z^{(\ell)}) \odot \nabla_y L(W, b)$$

$$a^{(\ell)} = f(z^{(\ell)})$$

$$z^{(\ell)} = W^{(\ell)} a^{(\ell-1)} + b^{(\ell)}$$

$$\frac{\partial z^{(\ell)}}{\partial W^{(\ell)}} = \left(a^{(\ell-1)} \right)^T$$

oooooooooooooooooooo oooooooooooooooooooo●ooooooo oooooooooooooooooooo oooooooooooooooo oooooooooooooooooooo

多层前馈网络

前向输出

$$z^{(\ell)} = w^{(\ell)} a^{(\ell-1)} + b^{(\ell)}$$

$$a^{(\ell)} = f(z^{(\ell)})$$

$$a^{(0)} = x$$

$$y = a^{(n_\ell)}$$

oooooooooooooooo

oooooooooooooooooooo●ooooooo

oooooooooooooooooooo

oooooooooooooooooooo

多层前馈网络

反向传播

$$\delta^{(\ell)} = \begin{cases} f'(z^{(\ell)}) \odot (w^{(\ell+1)})^T \delta^{(\ell+1)}; & \ell \neq n_\ell \\ f'(z^{(\ell)}) \odot \nabla_y L(W, b); & \ell = n_\ell \end{cases}$$

梯度计算

$$\nabla_{w^{(\ell)}} L(w, b; x, y) = \delta^{(\ell)} \otimes a^{(\ell-1)}$$

$$\nabla_{b^{(\ell)}} L(w, b; x, y) = \delta^{(\ell)}$$

$$\ell = 1, 2, \dots, n_\ell$$

oooooooooooooooo

oooooooooooooooooooo●ooooo ooooooooooooooooooooo ooooooooooooo ooooooooooooo

多层前馈网络

变化量

$$\Delta w^{(\ell)} \leftarrow \Delta w^{(\ell)} + \nabla_{w^{(\ell)}} L(w, b; x^{(k)}, t^{(k)})$$

$$\Delta b^{(\ell)} \leftarrow \Delta b^{(\ell)} + \nabla_{b^{(\ell)}} L(w, b; x^{(k)}, t^{(k)})$$

$$k = 1, 2, \dots, m; \ell = 1, 2, \dots, n_\ell$$

权重更新

$$w^{(\ell)} \leftarrow w^{(\ell)} - \alpha \left(\frac{\Delta w^{(\ell)}}{m} \right)$$

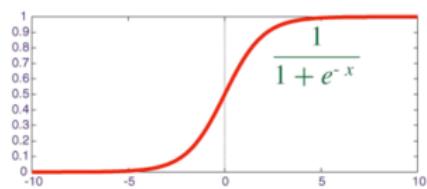
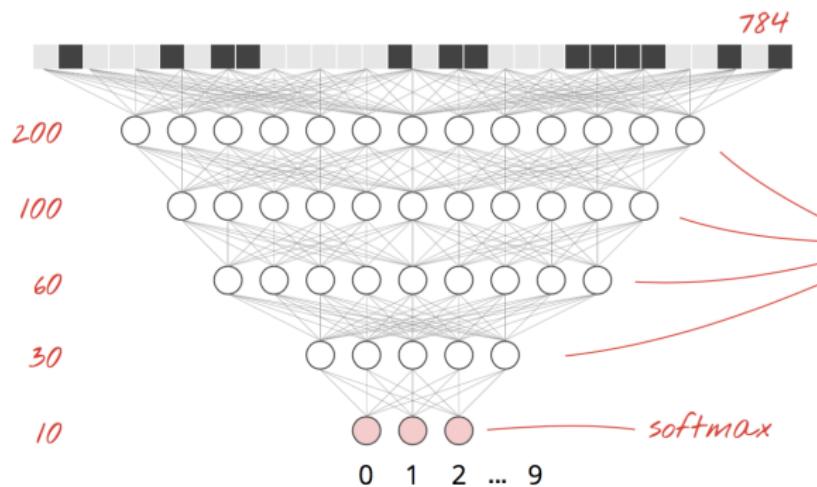
$$b^{(\ell)} \leftarrow b^{(\ell)} - \alpha \frac{\Delta b^{(\ell)}}{m}$$

$$\ell = 1, 2, \dots, n_\ell$$



实现网络

多层前馈网络



sigmoid function

softmax

实现网络

定义参数

```
K = 200  
L = 100  
M = 60  
N = 30
```

```
W1 = tf.Variable(tf.truncated_normal([28*28, K], stddev=0.1))  
B1 = tf.Variable(tf.zeros([K]))
```

```
W2 = tf.Variable(tf.truncated_normal([K, L], stddev=0.1))  
B2 = tf.Variable(tf.zeros([L]))
```

```
W3 = tf.Variable(tf.truncated_normal([L, M], stddev=0.1))  
B3 = tf.Variable(tf.zeros([M]))  
W4 = tf.Variable(tf.truncated_normal([M, N], stddev=0.1))  
B4 = tf.Variable(tf.zeros([N]))  
W5 = tf.Variable(tf.truncated_normal([N, 10], stddev=0.1))  
B5 = tf.Variable(tf.zeros([10]))
```

*weights initialised
with random values*



实现网络

定义模型

```
X = tf.reshape(X, [-1, 28*28])
```

weights and biases



```
Y1 = tf.nn.sigmoid(tf.matmul(X, W1) + B1)
```

```
Y2 = tf.nn.sigmoid(tf.matmul(Y1, W2) + B2)
```

```
Y3 = tf.nn.sigmoid(tf.matmul(Y2, W3) + B3)
```

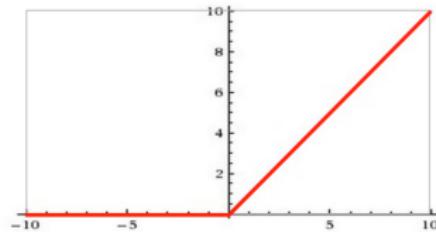
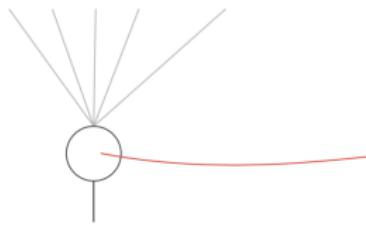
```
Y4 = tf.nn.sigmoid(tf.matmul(Y3, W4) + B4)
```

```
Y = tf.nn.softmax(tf.matmul(Y4, W5) + B5)
```

实现网络

梯度消失：引入 ReLU

RELU = Rectified Linear Unit



$Y = \text{tf.nn.relu}(\text{tf.matmul}(X, W) + b)$

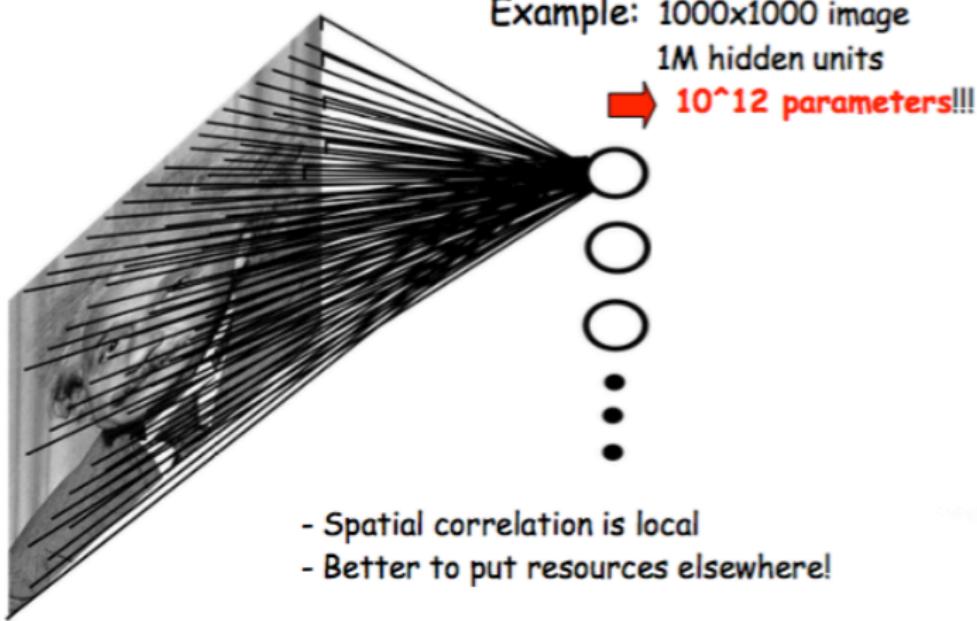


卷积神经网络

- ① 动机
- ② 卷积层
- ③ 池化层
- ④ 计算

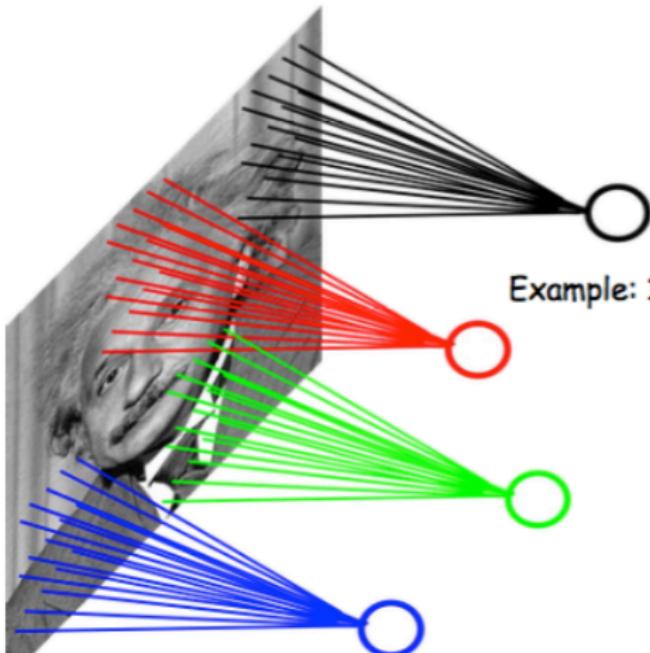
全连接网络

FULLY CONNECTED NEURAL NET



局部连接网络

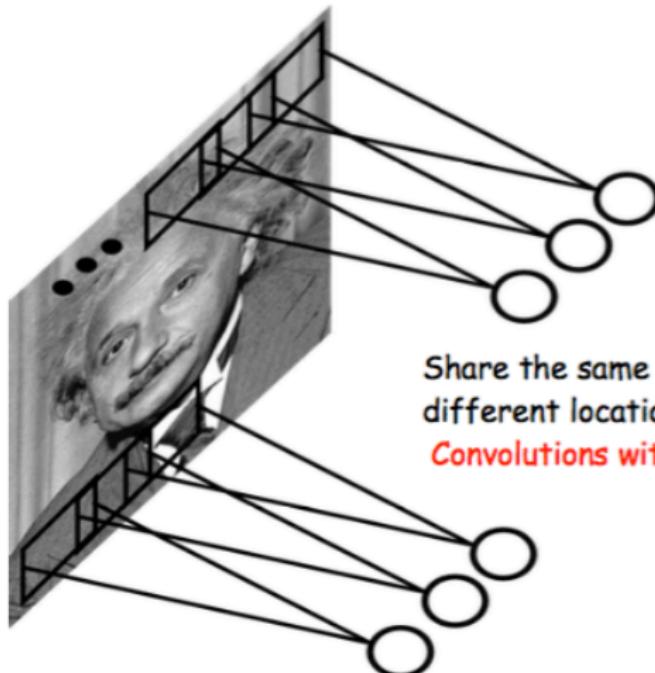
LOCALLY CONNECTED NEURAL NET



Example: 1000x1000 image
1M hidden units
Filter size: 10x10
100M parameters

权值共享：一个卷积核

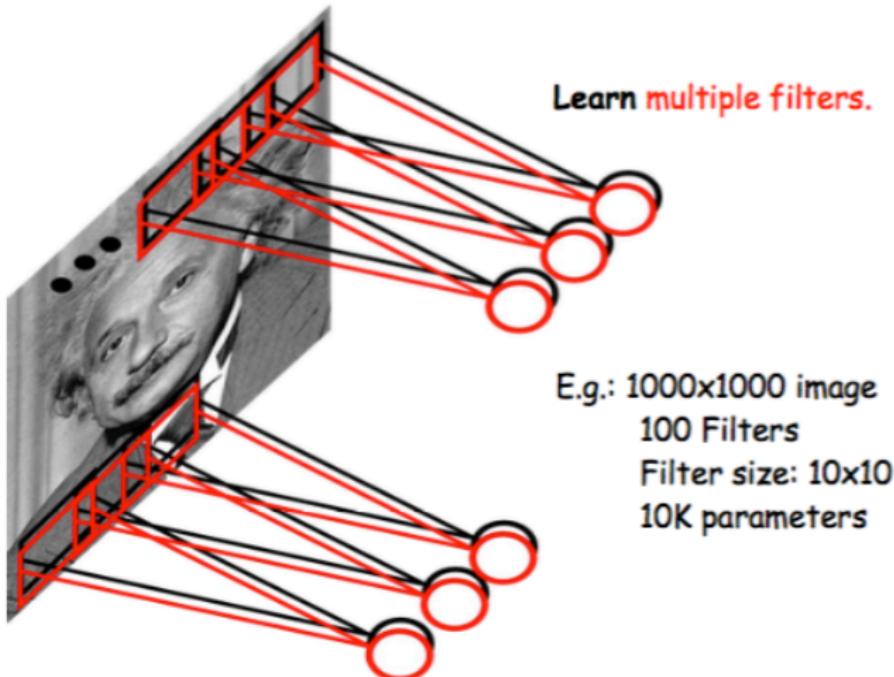
CONVOLUTIONAL NET



Share the same parameters across different locations:
Convolutions with learned kernels

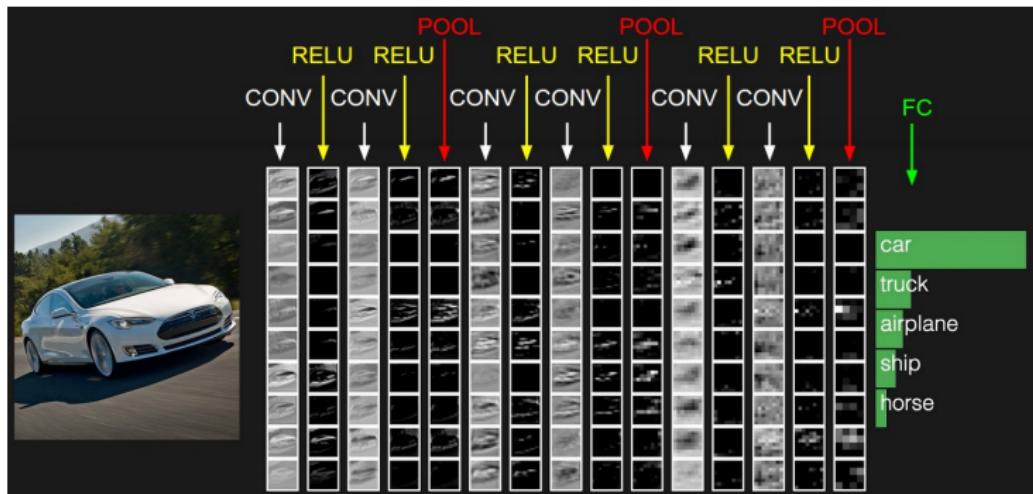
权值共享：多个卷积核

CONVOLUTIONAL NET



架构模式

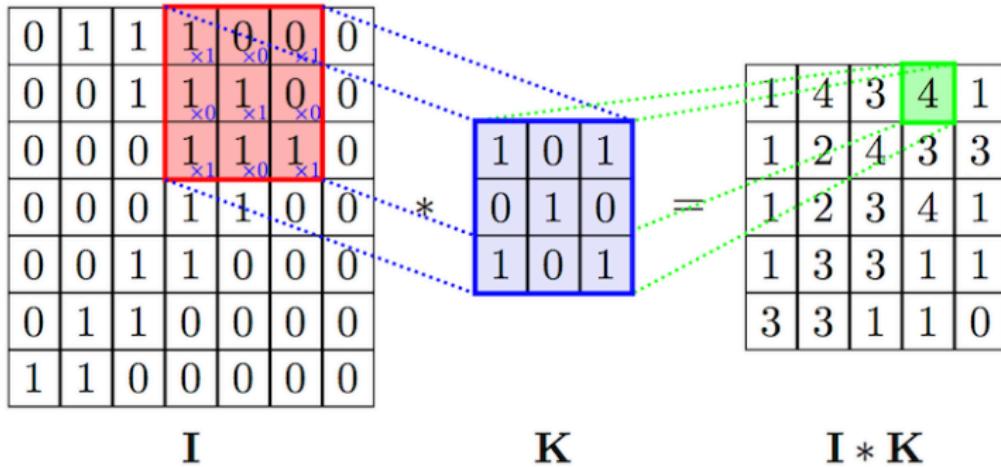
架构模式



INPUT -> [[CONV -> RELU]*N -> POOL?] *M -> [FC -> RELU]*K -> FC

卷积层

卷积运算



ooooooooooooooo

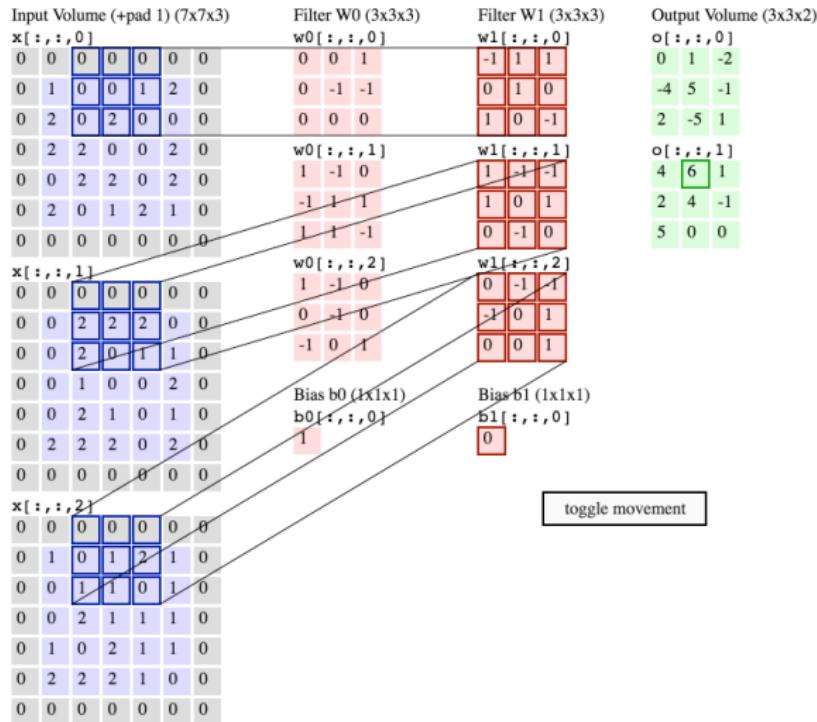
oooooooooooooooooooo

ooooooooooooooo

ooooooooooooooo

卷积层

卷积运算



oooooooooooooooo

oooooooooooooooooooo

oooooooooooo●oooo

oooooooooooo

oooooooooooooooo

卷积层

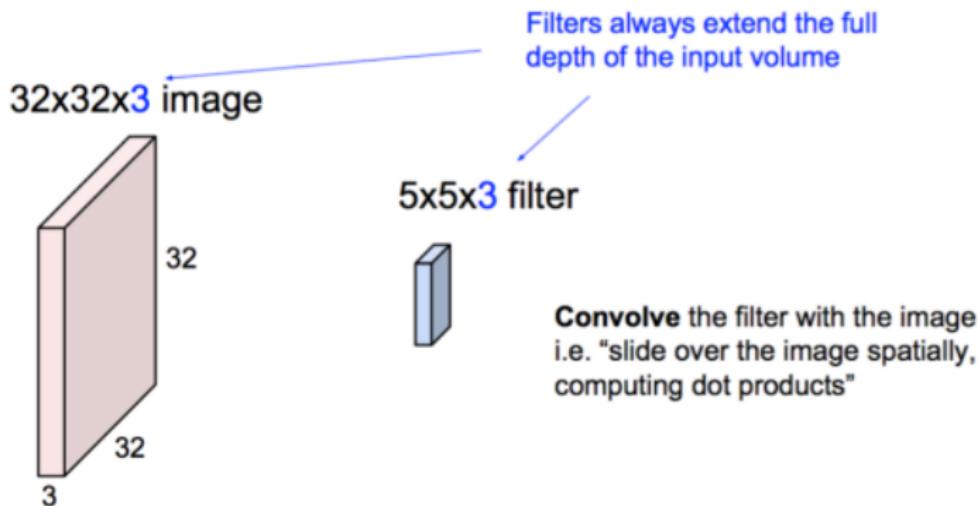
卷积运算

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.



卷积层

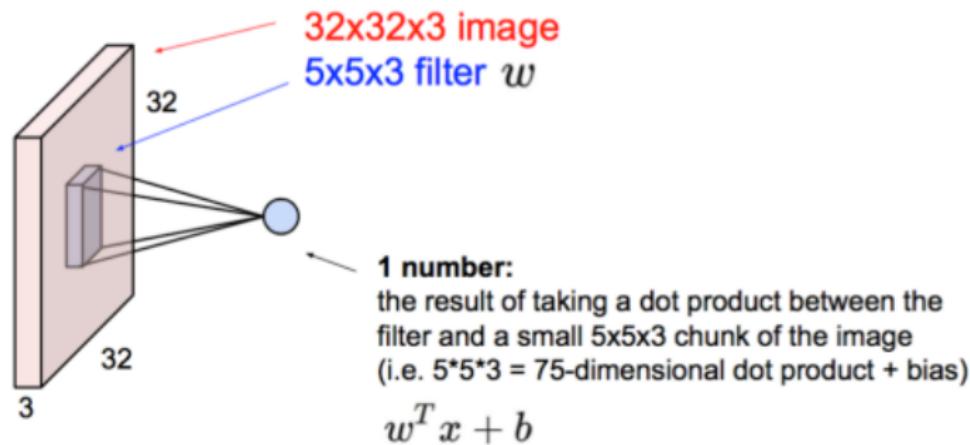
卷积运算





卷积层

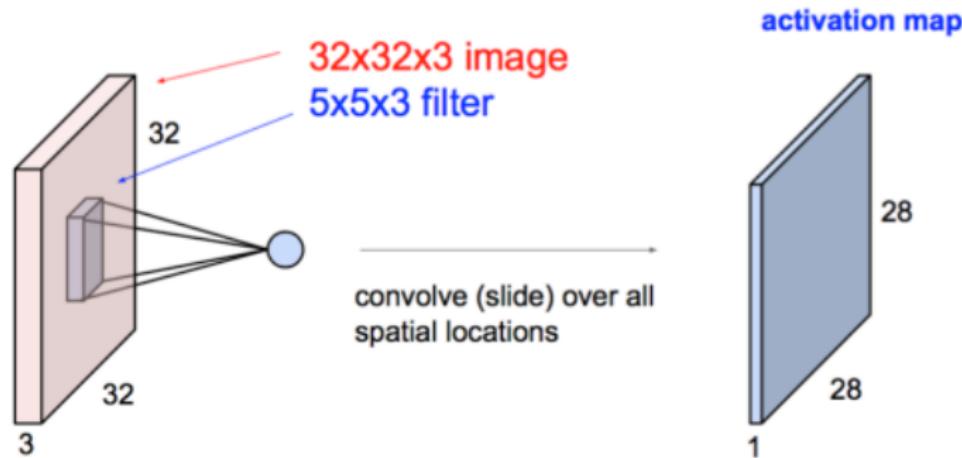
卷积运算





卷积层

卷积运算





卷积层

卷积运算



ooooooooooooooo

oooooooooooooooooooo

oooooooooooo●oooo

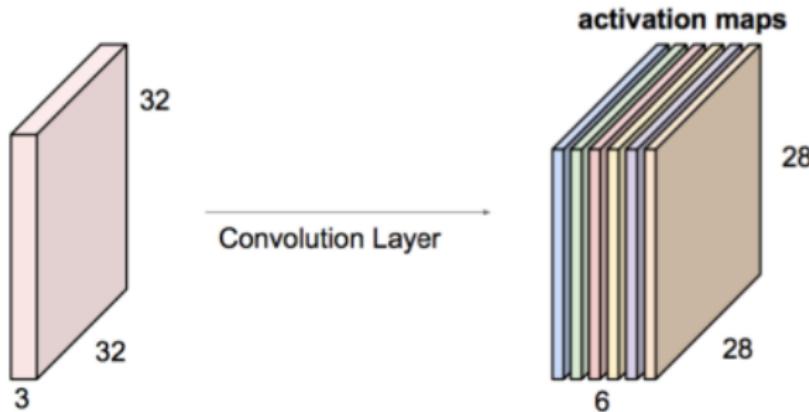
oooooooooooo

oooooooooooooooo

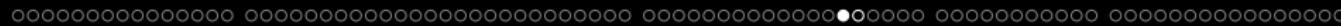
卷积层

卷积运算

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

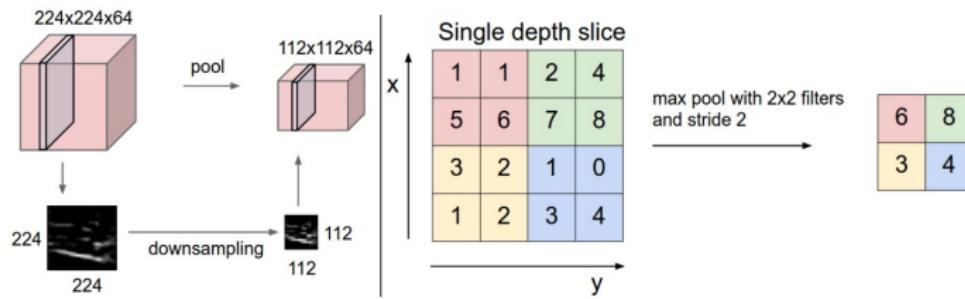


We stack these up to get a “new image” of size $28 \times 28 \times 6$!



池化层

下采样：Max Pooling

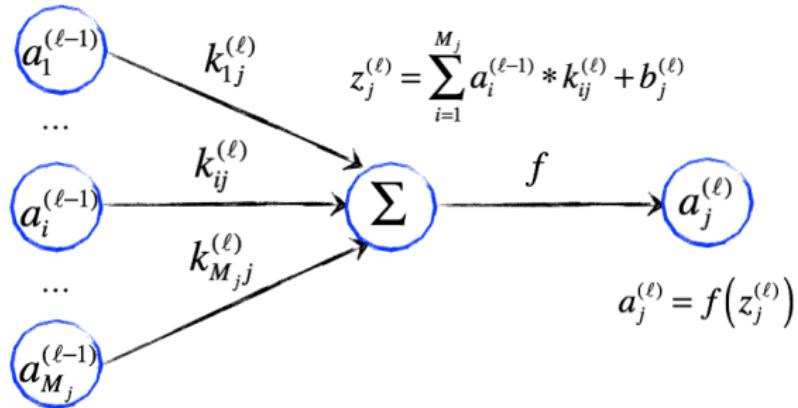


下采样

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires two hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

计算

前向计算：卷积



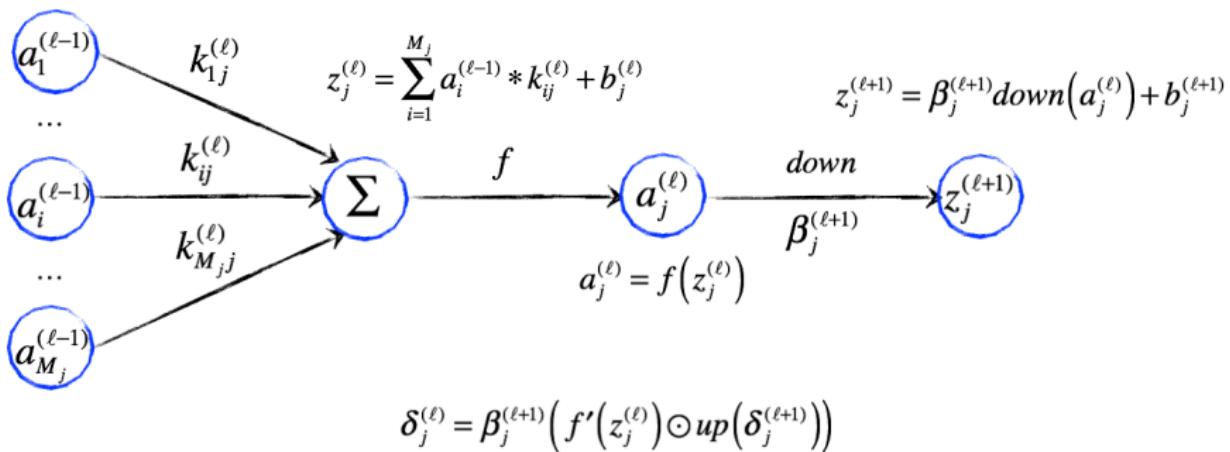
$$\frac{\partial L}{\partial k_{ij}^{(\ell)}} = a_i^{(\ell-1)} * \delta_j^{(\ell)}$$

$$\frac{\partial L}{\partial k_{ij}^{(\ell)}} = rot180(conv2(a_i^{(\ell-1)}, rot180(\delta_j^{(\ell)}), 'valid'))$$

$$\frac{\partial L}{\partial b_j^{(\ell)}} = \sum_{u,v} (\delta_j^{(\ell)})_{u,v}$$

计算

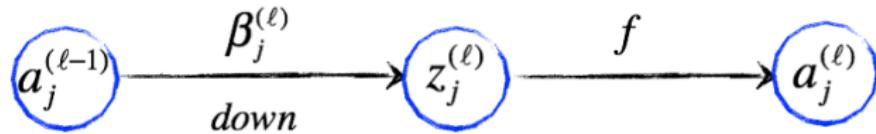
反向传播：下采样



计算

前向计算：下采样

$$z_j^{(\ell)} = \beta_j^{(\ell)} \text{down}\left(a_j^{(\ell-1)}\right) + b_j^{(\ell)} \quad a_j^{(\ell)} = f\left(z_j^{(\ell)}\right)$$

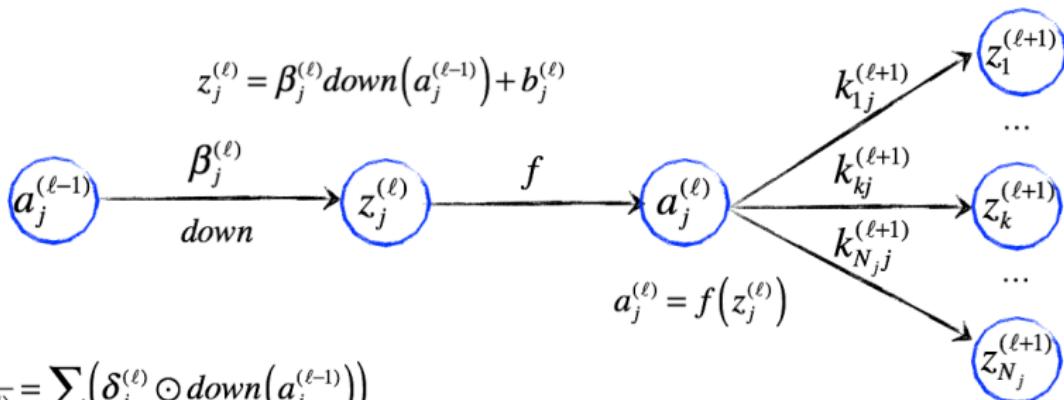


$$\frac{\partial L}{\partial \beta_j^{(\ell)}} = \sum_{u,v} \left(\delta_j^{(\ell)} \odot \text{down}\left(a_j^{(\ell-1)}\right) \right)_{u,v}$$

$$\frac{\partial L}{\partial b_j^{(\ell)}} = \sum_{u,v} \left(\delta_j^{(\ell)} \right)_{u,v}$$

计算

反向传播：卷积

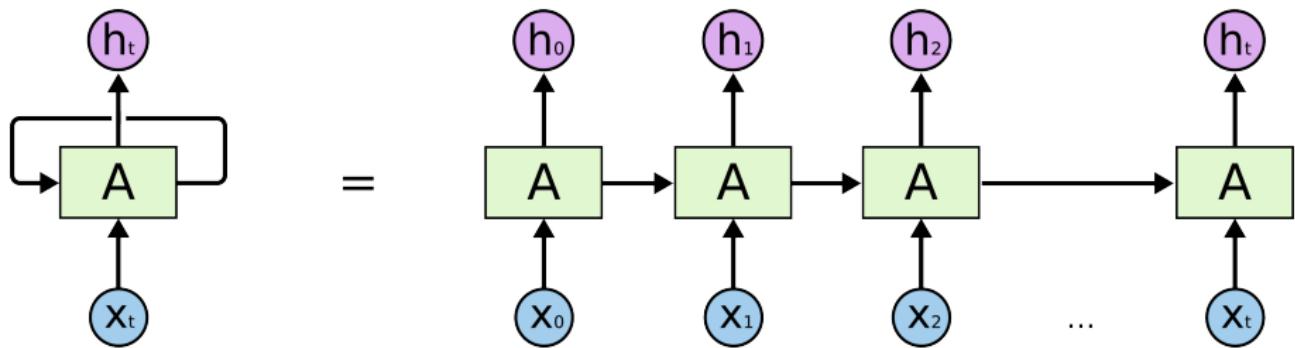


$$\begin{aligned}\delta_j^{(\ell)} &= f'(z_j^{(\ell)}) \odot \sum_{k=1}^{N_j} \text{padding}(\delta_k^{(\ell+1)}) * \text{rot180}(k_{kj}^{(\ell+1)}) \\ &= f'(z_j^{(\ell)}) \odot \sum_{k=1}^{N_j} \text{conv2}(\delta_k^{(\ell+1)}, \text{rot180}(k_{kj}^{(\ell+1)}), 'full')\end{aligned}$$

循环神经网络

- ① BPTT 算法
- ② 双向 RNN
- ③ 深层 RNN

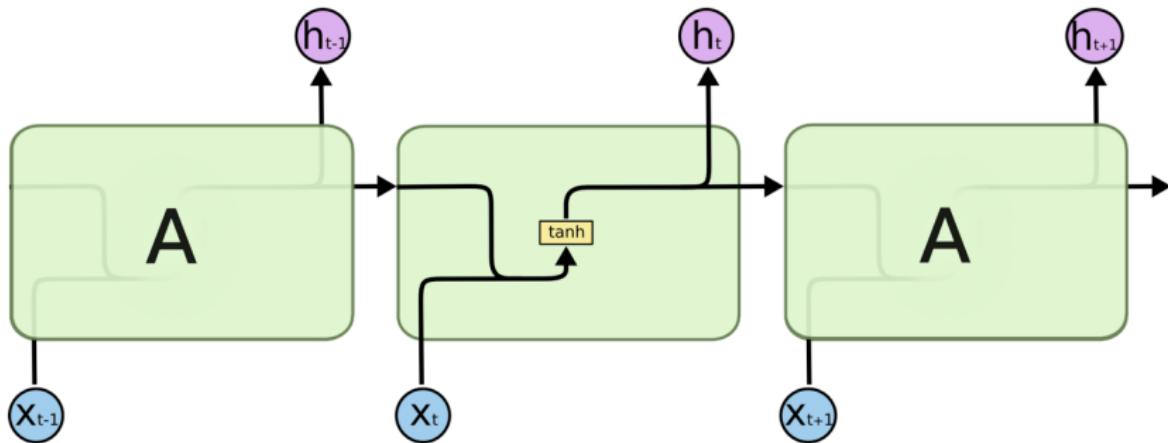
网络架构

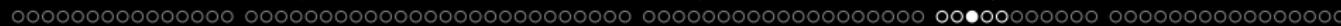


网络架构

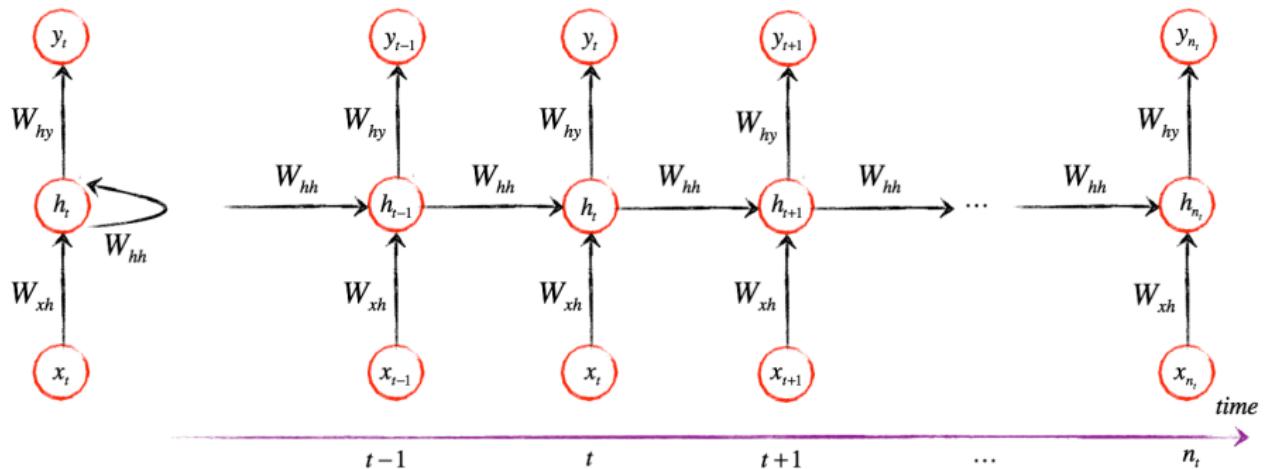
$$u_t = W[x_t, h_{t-1}]$$

$$h_t = \tanh(u_t)$$

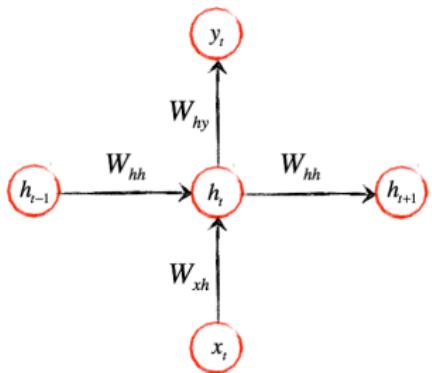




网络架构



前向计算



$$u_t = W_{xh}x_t + W_{hh}h_{t-1} \quad u_t = W_{xh}x_t + W_{hh}h_{t-1}$$

$$h_t = f(u_t) \quad = W[x_t, h_{t-1}]$$

$$v_t = W_{hy}h_t$$

$$y_t = g(v_t)$$

$$\begin{aligned} &= \begin{pmatrix} W_{xh} & W_{hh} \end{pmatrix} \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix} \end{aligned}$$

$$\nabla_{W_{hy}} L = \sum_{t=1}^{n_t} \nabla_{W_{hy}} L_t$$

$$\nabla_{W_{hy}} L_t = \frac{\partial L_t}{\partial W_{hy}} = \frac{\partial L_t}{\partial v_t} \frac{\partial v_t}{\partial W_{hy}} = \delta_t^y \otimes h_t$$

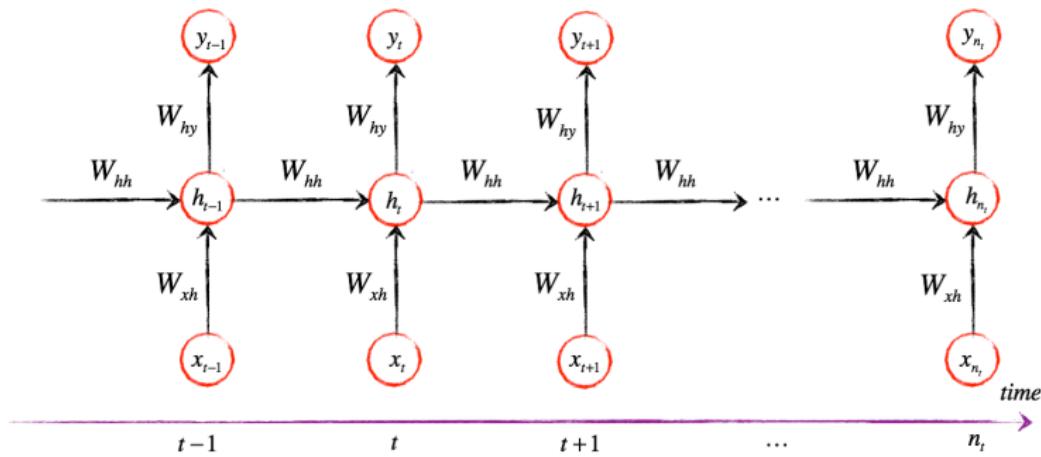
$$\nabla_{W_{hh}} L = \sum_{t=1}^{n_t} \nabla_{W_{hh}} L_t$$

$$\nabla_{W_{hh}} L_t = \frac{\partial L_t}{\partial W_{hh}} = \frac{\partial L_t}{\partial u_t} \frac{\partial u_t}{\partial W_{hh}} = \delta_t^h \otimes h_{t-1}$$

$$\nabla_{W_{xh}} L = \sum_{t=1}^{n_t} \nabla_{W_{xh}} L_t$$

$$\nabla_{W_{xh}} L_t = \frac{\partial L_t}{\partial W_{xh}} = \frac{\partial L_t}{\partial u_t} \frac{\partial u_t}{\partial W_{xh}} = \delta_t^h \otimes x_t$$

反向传播：BPTT



$$t = n_t - 1, n_t - 2, \dots, 1$$

$$\delta_t^y = (\nabla_{y_t} L_t) L \odot g'(v_t)$$

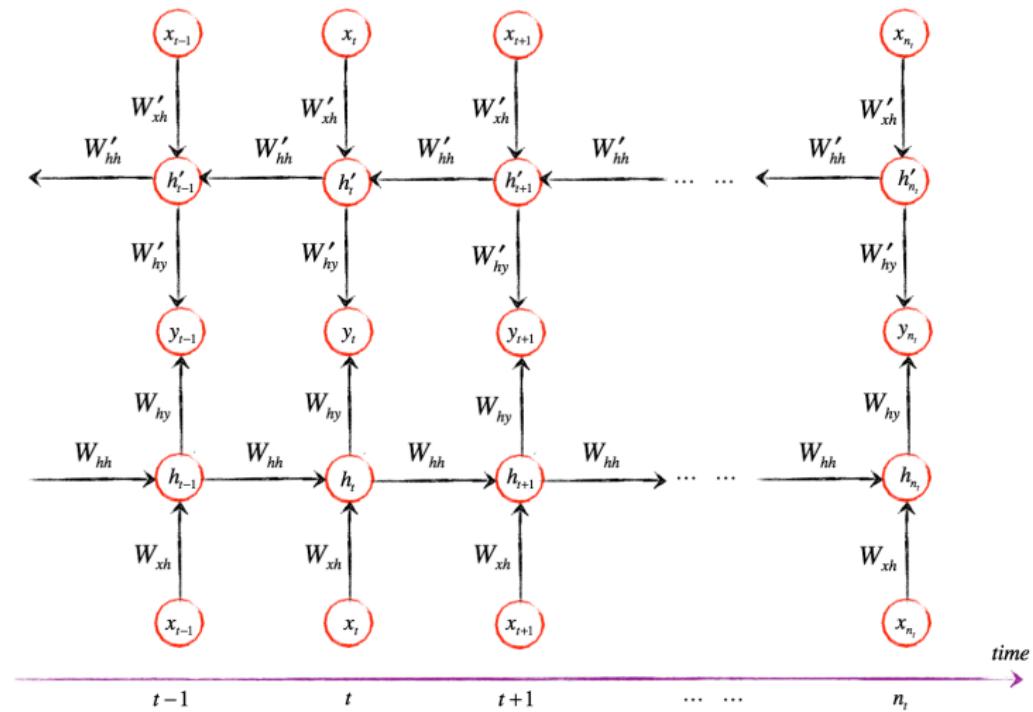
$$\delta_t^h = \left[(W_{hy})^T \delta_t^y + (W_{hh})^T \delta_{t+1}^h \right] \odot f'(u_t)$$

$$t = n_t$$

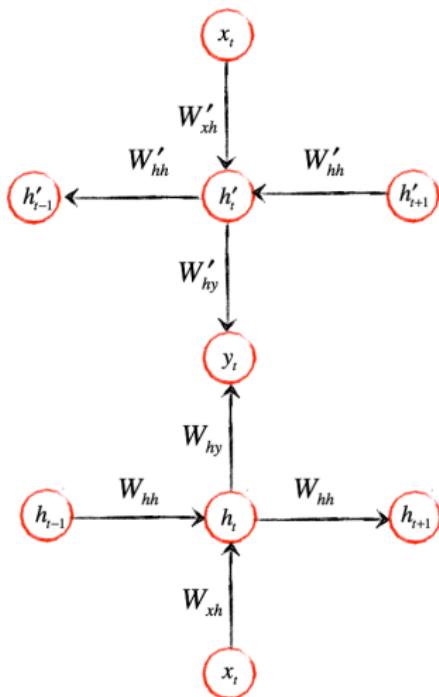
$$\delta_{n_t}^y = (\nabla_{y_{n_t}} L_{n_t}) \odot g'(v_{n_t})$$

$$\delta_{n_t}^h = (W_{hy})^T \delta_{n_t}^y \odot f'(u_{n_t})$$

网络架构



前向计算



$$u_t = W_{xh}x_t + W_{hh}h_{t-1}$$

$$u'_t = W'_{xh}x_t + W'_{hh}h'_{t-1}$$

$$h_t = f(u_t)$$

$$h'_t = f(u'_t)$$

$$v_t = W_{hy}h_t$$

$$v'_t = W'_{hy}h'_t$$

$$s_t = v_t + v'_t$$

$$y_t = g(s_t)$$

$$\nabla_{W_{hy}} L_t = \frac{\partial L_t}{\partial W_{hy}} = \frac{\partial L_t}{\partial v_t} \frac{\partial v_t}{\partial W_{hy}} = \delta_t^y \otimes h_t$$

$$\nabla_{W'_{hy}} L_t = \frac{\partial L_t}{\partial W'_{hy}} = \frac{\partial L_t}{\partial v'_t} \frac{\partial v'_t}{\partial W'_{hy}} = \delta_t'^y \otimes h'_t$$

$$\nabla_{W_{hh}} L_t = \frac{\partial L_t}{\partial W_{hh}} = \frac{\partial L_t}{\partial u_t} \frac{\partial u_t}{\partial W_{hh}} = \delta_t^h \otimes h_{t-1}$$

$$\nabla_{W'_{hh}} L_t = \frac{\partial L_t}{\partial W'_{hh}} = \frac{\partial L_t}{\partial u'_t} \frac{\partial u'_t}{\partial W'_{hh}} = \delta_t'^h \otimes h'_{t+1}$$

$$\nabla_{W_{xh}} L_t = \frac{\partial L_t}{\partial W_{xh}} = \frac{\partial L_t}{\partial u_t} \frac{\partial u_t}{\partial W_{xh}} = \delta_t^h \otimes x_t$$

$$\nabla_{W'_{xh}} L_t = \frac{\partial L_t}{\partial W'_{xh}} = \frac{\partial L_t}{\partial u'_t} \frac{\partial u'_t}{\partial W'_{xh}} = \delta_t'^h \otimes x_t$$

双向 RNN

反向传播：BPTT

 $t = 1$

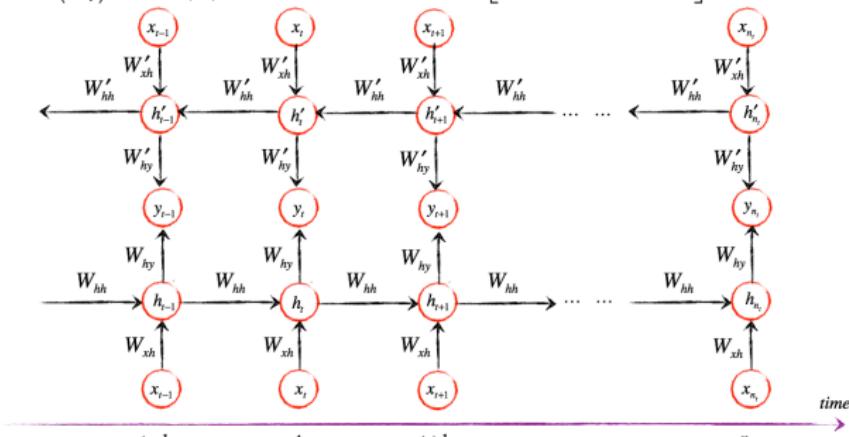
$$\delta_1^y = (\nabla_{y_1} L_1) \odot g'(s_1)$$

$$\delta_1'^h = (W'_{hy})^T \delta_1^y \odot f'(u'_1)$$

 $t = 2, 3, \dots, n_t$

$$\delta_t^y = (\nabla_{y_t} L_t) \odot g'(s_t)$$

$$\delta_t'^h = [(W'_{hy})^T \delta_t^y + (W'_{hh})^T \delta_{t-1}^h] \odot f'(u'_t)$$

 $t = n_t - 1, n_t - 2, \dots, 1$

$$\delta_t^y = (\nabla_{y_t} L_t) \odot g'(s_t)$$

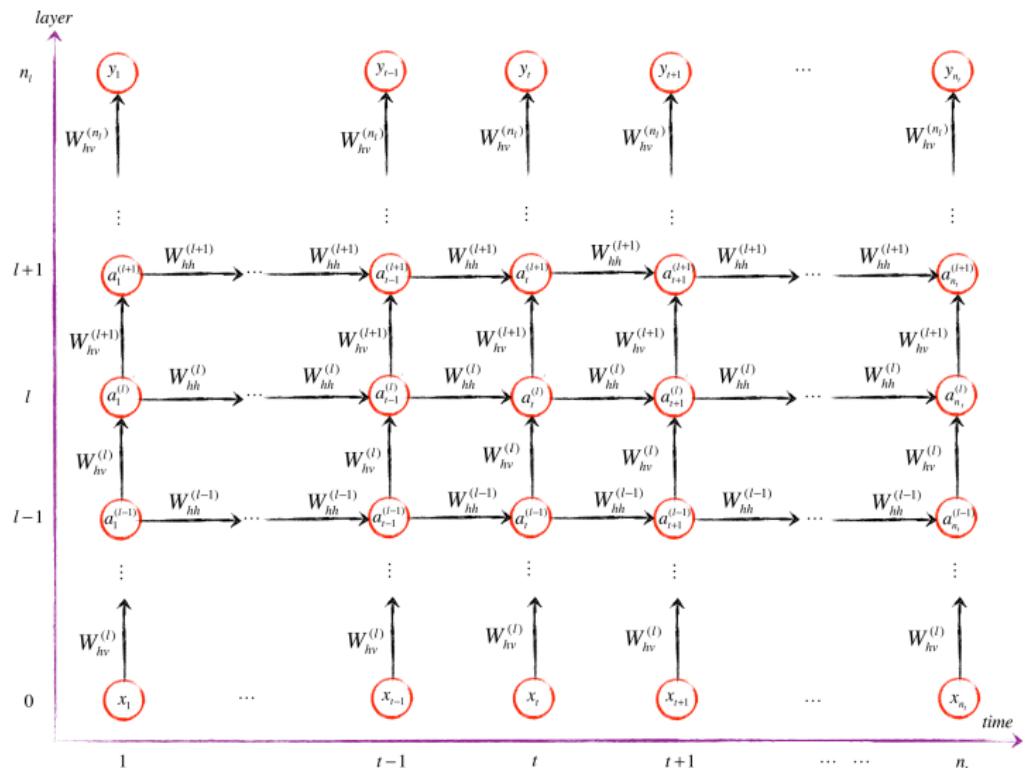
$$\delta_t^h = [(W_{hy})^T \delta_t^y + (W_{hh})^T \delta_{t+1}^h] \odot f'(u_t)$$

 $t = n_t$

$$\delta_{n_t}^y = (\nabla_{y_{n_t}} L_{n_t}) \odot g'(s_{n_t})$$

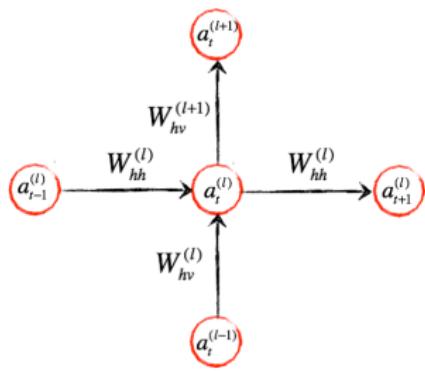
$$\delta_{n_t}^h = (W_{hy})^T \delta_{n_t}^y \odot f'(u_{n_t})$$

网络架构



深层 RNN

前向计算



$$u_t^{(l)} = W_{hv}^{(l)} a_t^{(l-1)} + W_{hh}^{(l)} a_{t-1}^{(l)}$$

$$a_t^{(l)} = f(u_t^{(l)})$$

$$a_0^{(l)} = \vec{0}$$

$$l = 1, 2, \dots, n_t$$

$$a_t^{(0)} = x_t$$

$$t = 1, 2, \dots, n_t$$

$$y_t = a_t^{(n_l)}$$

$$t = 1, 2, \dots, n_t$$

$$\nabla_{W_{hh}^{(l)}} L_t = \frac{\partial L_t}{W_{hh}^{(l)}} = \frac{\partial L_t}{u_t^{(l)}} \frac{u_t^{(l)}}{W_{hh}^{(l)}} = \delta_t^{(l)} \otimes a_{t-1}^{(l)}$$

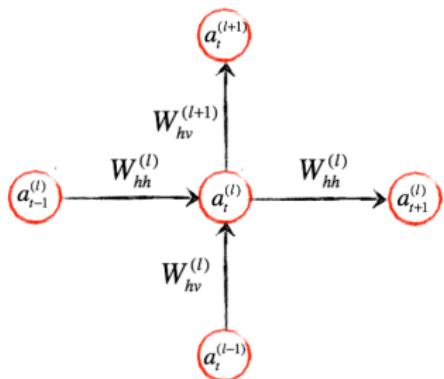
$$\nabla_{W_{hv}^{(l)}} L_t = \frac{\partial L_t}{W_{hv}^{(l)}} = \frac{\partial L_t}{u_t^{(l)}} \frac{u_t^{(l)}}{W_{hv}^{(l)}} = \delta_t^{(l)} \otimes a_t^{(l-1)}$$

$$\nabla_{W_{hh}^{(l)}} L = \sum_{t=1}^{n_l} \nabla_{W_{hh}^{(l)}} L_t$$

$$\nabla_{W_{hv}^{(l)}} L = \sum_{t=1}^{n_l} \nabla_{W_{hv}^{(l)}} L_t$$

深层 RNN

反向传播：BPTT



$$\delta_t^{(l)} = \vec{0} \quad \delta_t^{(n_l)} = \nabla_{y_t} L_t \odot f'(u_t^{(n_l)})$$

$$t > n_t \quad l = n_l$$

$$\delta_t^{(l)} = \left[\left(W_{hv}^{(l+1)} \right)^T \delta_t^{(l+1)} + \left(W_{hh}^{(l)} \right)^T \delta_{t+1}^{(l)} \right] \odot f'(u_t^{(l)})$$

$$l = n_l - 1, n_l - 2, \dots, 2, 1$$

$$t = n_l, n_l - 1, \dots, 2, 1$$

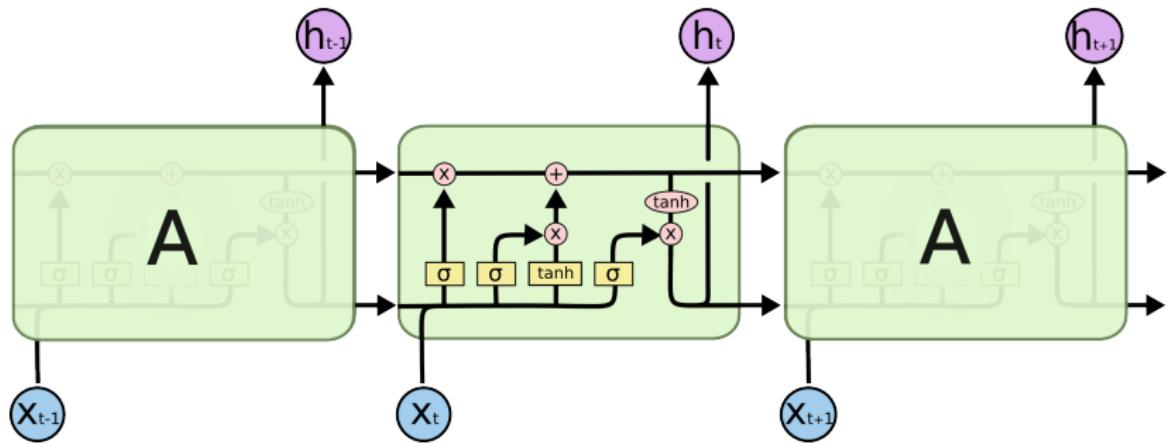
长短时记忆网络

- ① 网络架构
- ② 网络变种



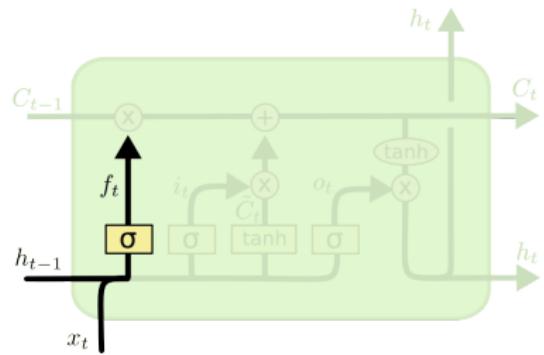
架构

网络架构



架构

遗忘门

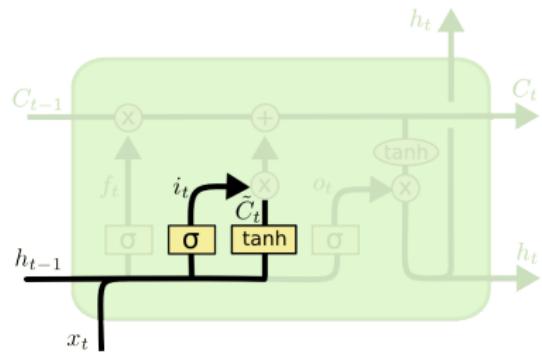


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$



架构

输入入门

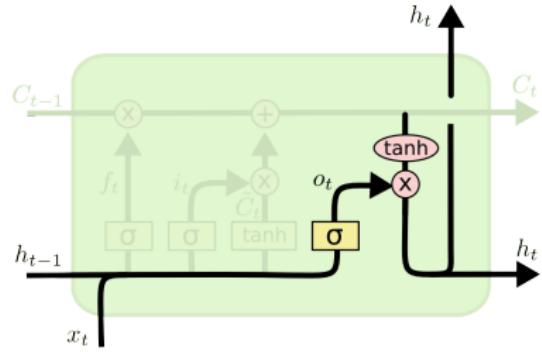


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

架构

输出门



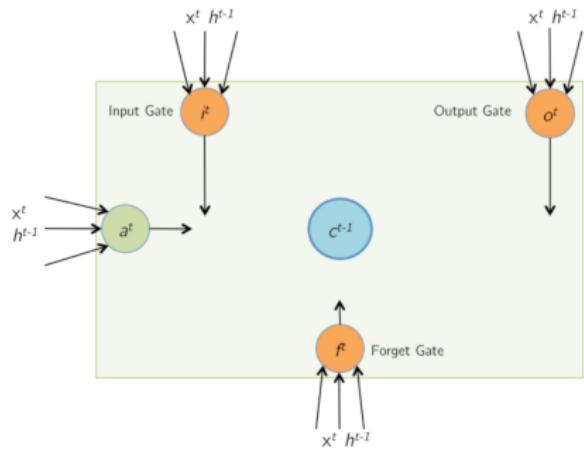
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

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

计算

初始状态: $t - 1$ 时刻 $x^t h^{t-1}$ $x^t h^{t-1}$ x^t
 h^{t-1} c^{t-1} o^t  $x^t h^{t-1}$

前向传播：输入与门计算



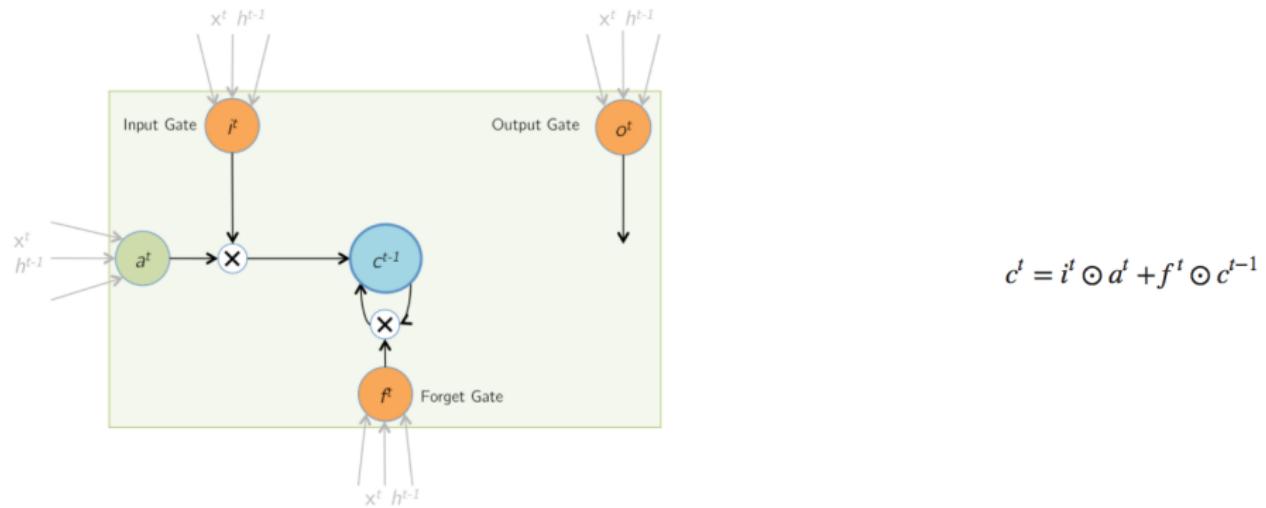
$$\begin{aligned}
 a^t &= \tanh(W_c x^t + U_c h^{t-1}) = \tanh(\hat{a}^t) \\
 i^t &= \sigma(W_i x^t + U_i h^{t-1}) = \sigma(\hat{i}^t) \\
 f^t &= \sigma(W_f x^t + U_f h^{t-1}) = \sigma(\hat{f}^t) \\
 o^t &= \sigma(W_o x^t + U_o h^{t-1}) = \sigma(\hat{o}^t)
 \end{aligned}$$

Ignoring the non-linearities,

$$\begin{aligned}
 z^t &= \begin{bmatrix} \hat{a}^t \\ \hat{i}^t \\ \hat{f}^t \\ \hat{o}^t \end{bmatrix} = \begin{bmatrix} W^c & U^c \\ W^i & U^i \\ W^f & U^f \\ W^o & U^o \end{bmatrix} \times \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix} \\
 &= W \times I^t
 \end{aligned}$$

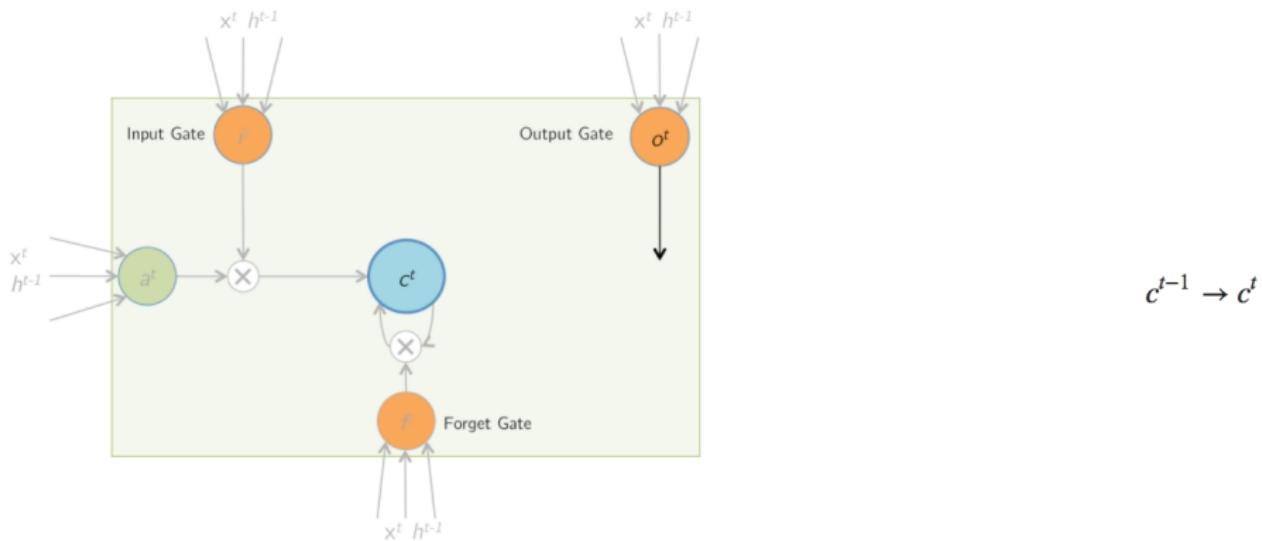
计算

前向传播：状态更新



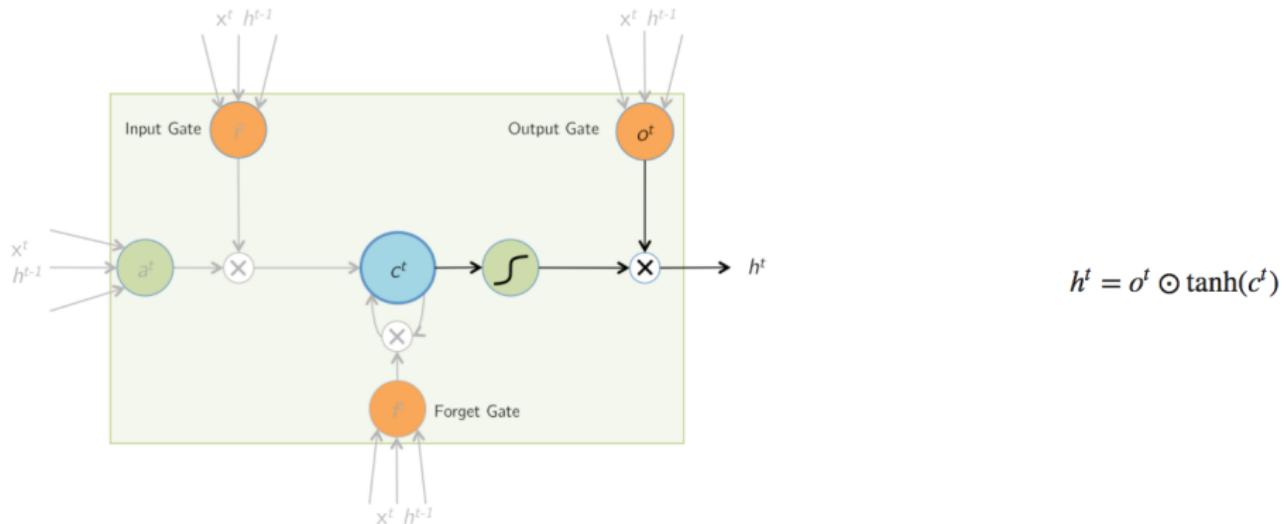
计算

前向传播：状态更新



计算

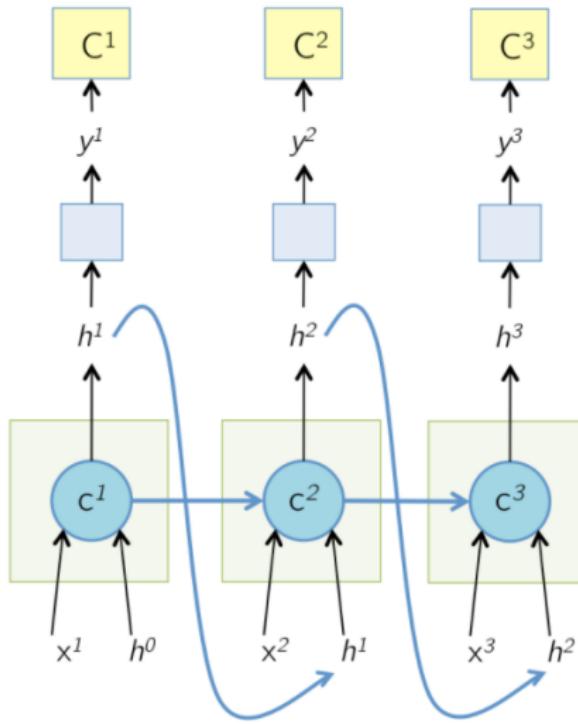
前向传播：输出





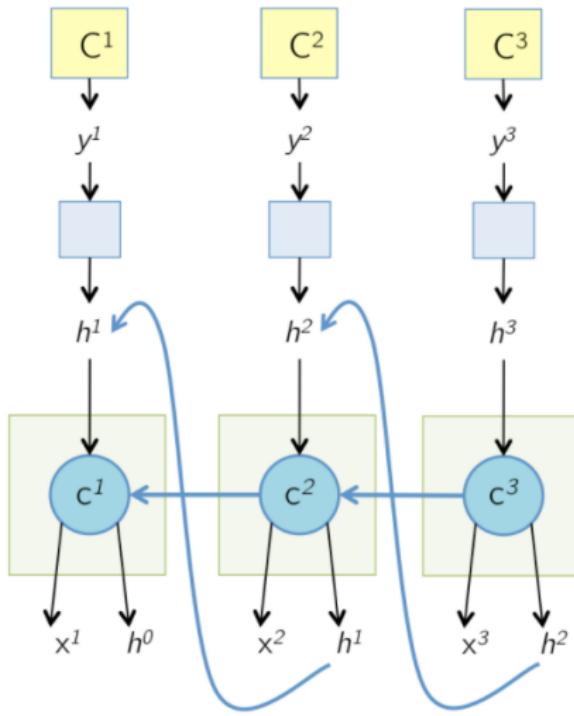
计算

前向传播



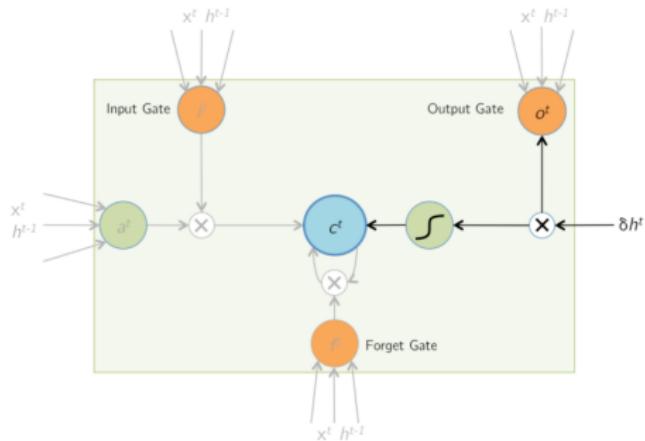
计算

反向传播



计算

反向传播



Forward Pass: $h^t = o^t \odot \tanh(c^t)$

Given $\delta h^t = \frac{\partial E}{\partial h^t}$, find $\delta o^t, \delta c^t$

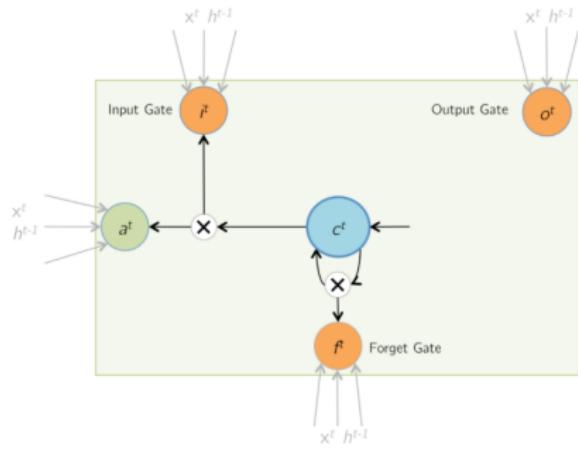
$$\begin{aligned}\frac{\partial E}{\partial o_i^t} &= \frac{\partial E}{\partial h_i^t} \cdot \frac{\partial h_i^t}{\partial o_i^t} \\ &= \delta h_i^t \cdot \tanh(c_i^t) \\ \therefore \delta o_i^t &\equiv \delta h_i^t \odot \tanh(c_i^t)\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial c_i^t} &= \frac{\partial E}{\partial h_i^t} \cdot \frac{\partial h_i^t}{\partial c_i^t} \\ &= \delta h_i^t \cdot o_i^t \cdot (1 - \tanh^2(c_i^t)) \\ \therefore \delta c^t &= \delta h^t \odot o^t \odot (1 - \tanh^2(c^t))\end{aligned}$$

Note that the $+$ above is so that this gradient is added to gradient from time step $(t + 1)$ (calculated on next slide, refer to the gradient accumulation mentioned in the previous slide)

计算

反向传播



Forward Pass: $c^t \equiv i^t \odot a^t + f^t \odot c^{t-1}$

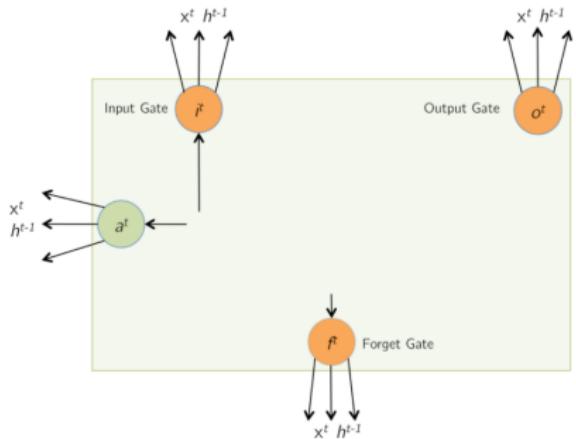
Given $\delta c^t = \frac{\partial E}{\partial c^t}$, find $\delta i^t, \delta a^t, \delta f^t, \delta c^{t-1}$

$$\begin{aligned} \frac{\partial E}{\partial t_i^t} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial c_i^t}{\partial i_i^t} \\ &= \delta c_i^t \cdot a_i^t \\ \therefore \delta t_i^t &= \delta c_i^t \odot a_i^t \end{aligned} \quad \begin{aligned} \frac{\partial E}{\partial f_i^t} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial}{\partial j} \\ &= \delta c_i^t \cdot c_i^{t-1} \\ \therefore \delta f_i^t &= \delta c_i^t \odot c_i^{t-1} \end{aligned}$$

$$\begin{aligned} \frac{\partial E}{\partial a_i^t} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial c_i^t}{\partial a_i^t} & \frac{\partial E}{\partial c_i^{t-1}} &= \frac{\partial E}{\partial c_i^t} \cdot \frac{\partial c_i^t}{\partial c_i^{t-1}} \\ &= \delta c_i^t \cdot i_i^t & &= \delta c_i^t \cdot f_i^t \\ \therefore \delta a^t &= \delta c^t \odot i^t & \therefore \delta c^{t-1} &= \delta c^t \odot f^t \end{aligned}$$

计算

反向传播



Forward Pass: $z^t = W \times I^t$
Given δz^t , find $\delta W^t, \delta h^{t-1}$

$$\delta I^t = W^T \times \delta z^t$$

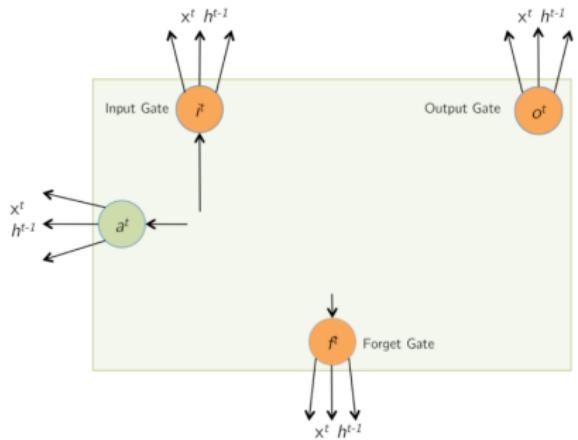
$$\text{As } I^t = \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix},$$

δh^{t-1} can be retrieved from δI^t

$$\delta W^t = \delta z^t \times (I^t)^T$$

计算

反向传播



Given δz^t , find $\delta W^t, \delta h^{t-1}$

$$\delta I^t = W^T \times \delta z^t$$

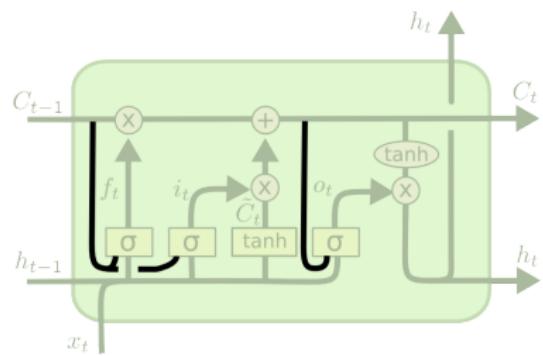
$$\text{As } I^t = \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix},$$

δh^{t-1} can be retrieved from δI^t

$$\delta W^t = \delta z^t \times (I^t)^T$$

变种

变种 1



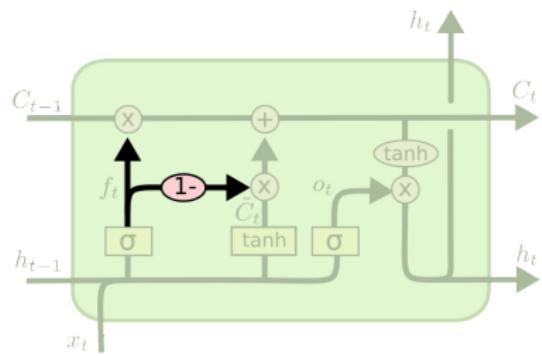
$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

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

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

变种

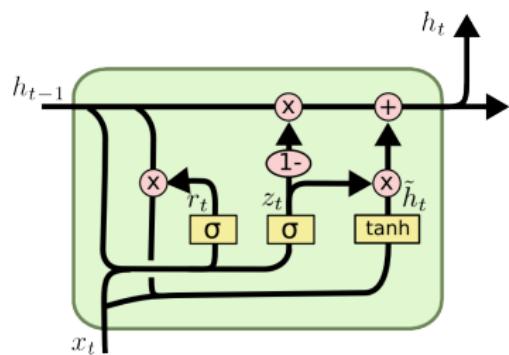
变种 2



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

变种

变种 3：GRU



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

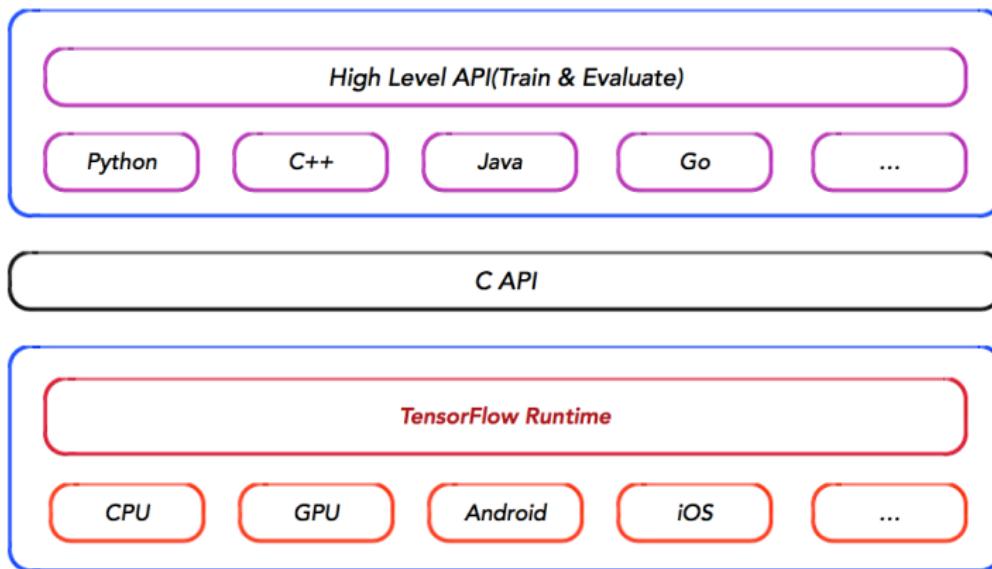
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

TensorFlow 实现

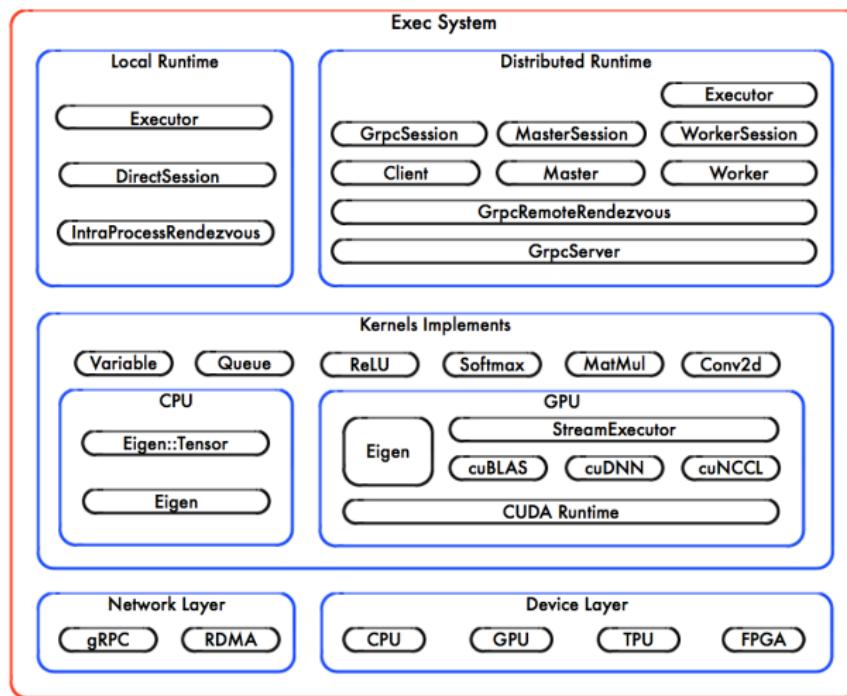
- ① 系统架构
- ② 自动微分

系统架构

架构



运行时

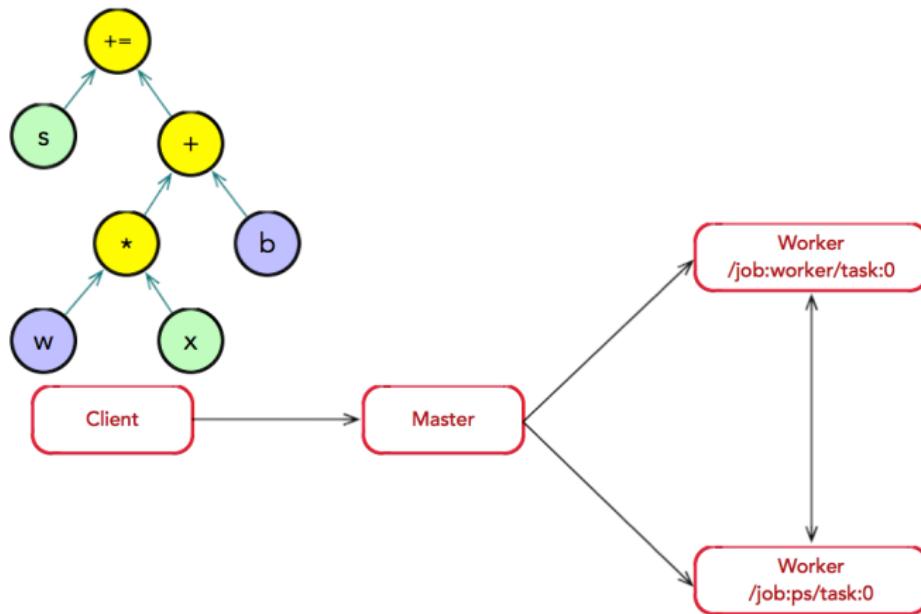


设计原则

- Deferred Execution
- Primitive OP
- Common Abstraction for Heterogeneous Accelerators

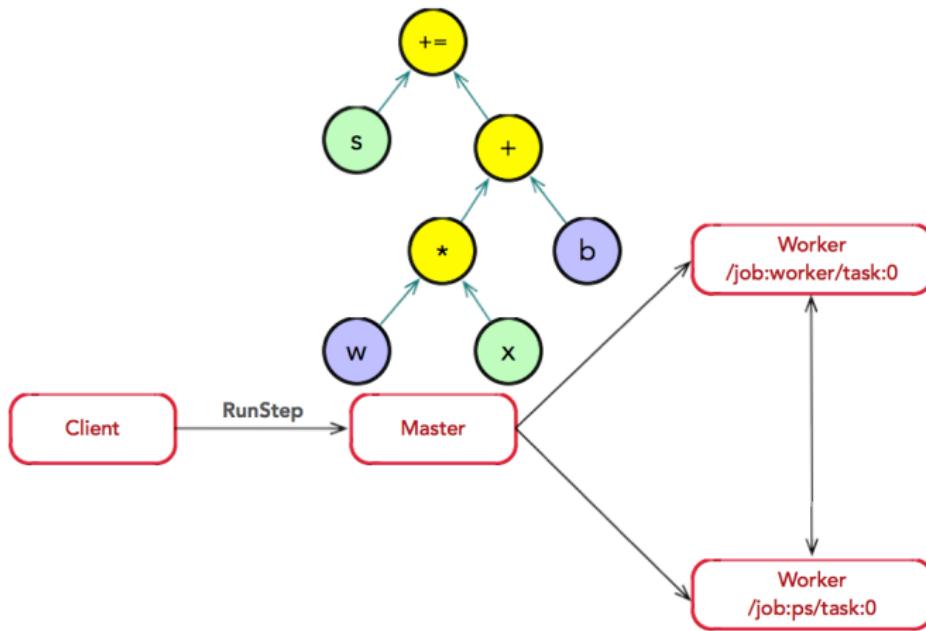
计算图

图构造



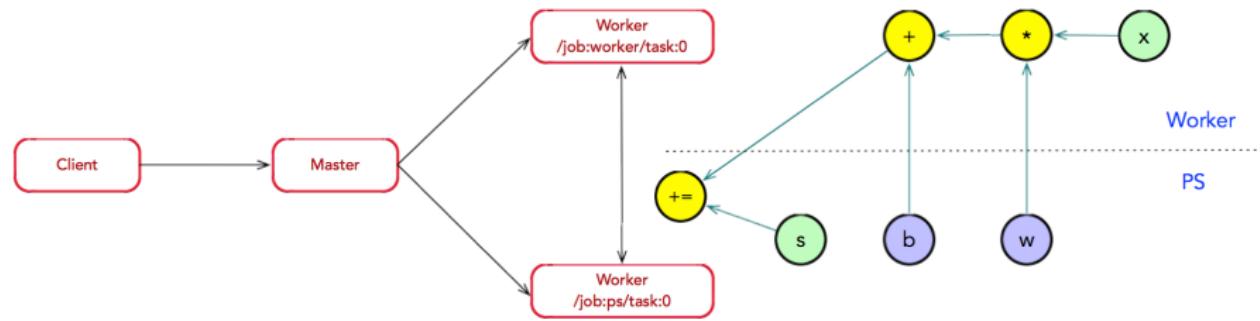
计算图

图执行



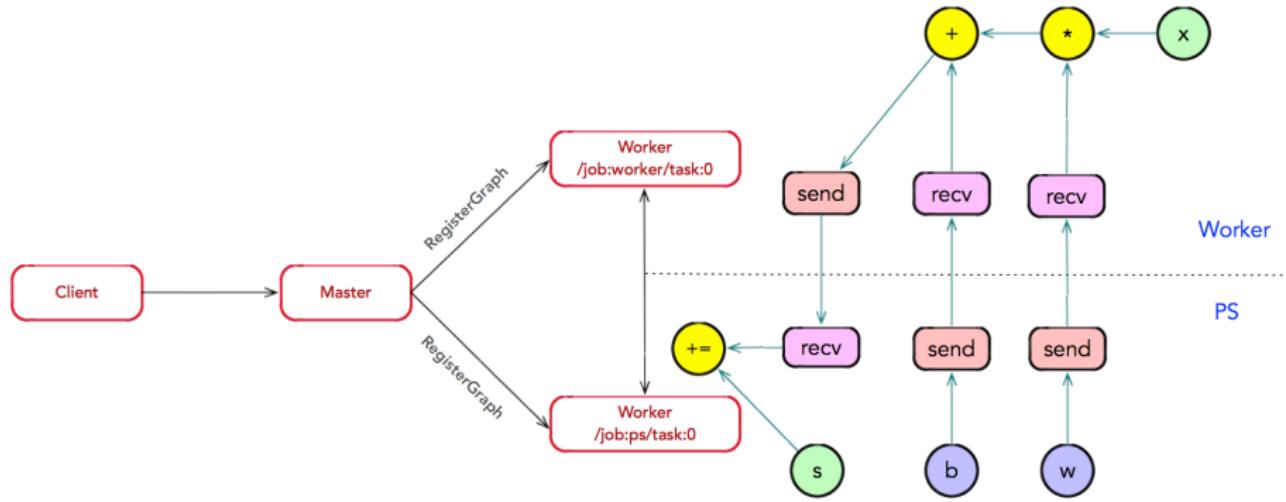
计算图

图分裂



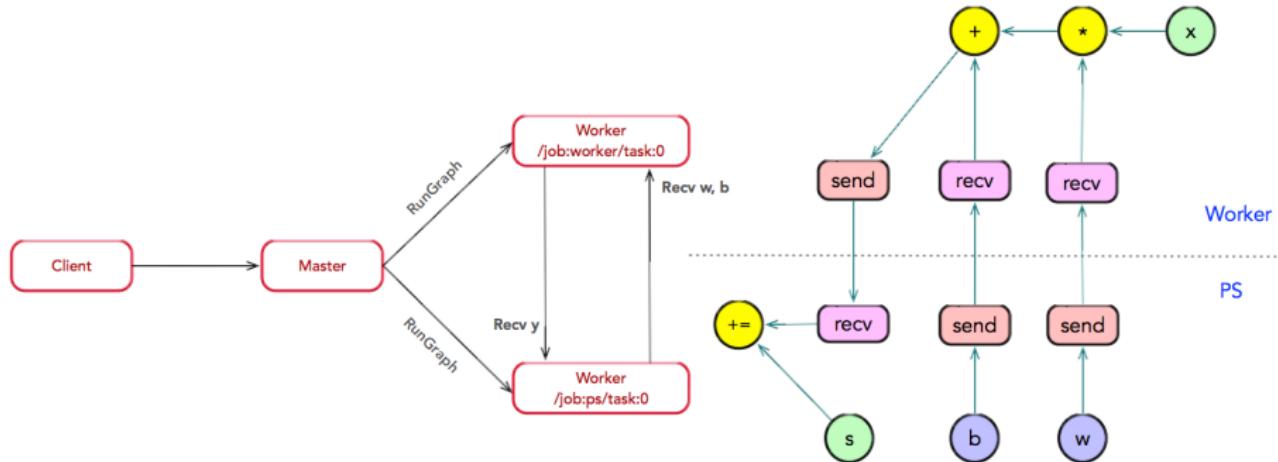
计算图

注册图



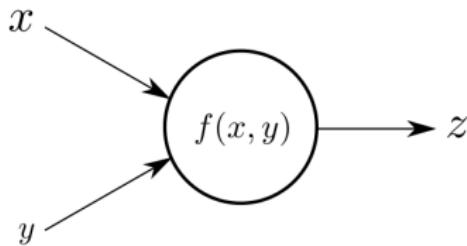
计算图

执行图



链式法则

Forwardpass



Backwardpass

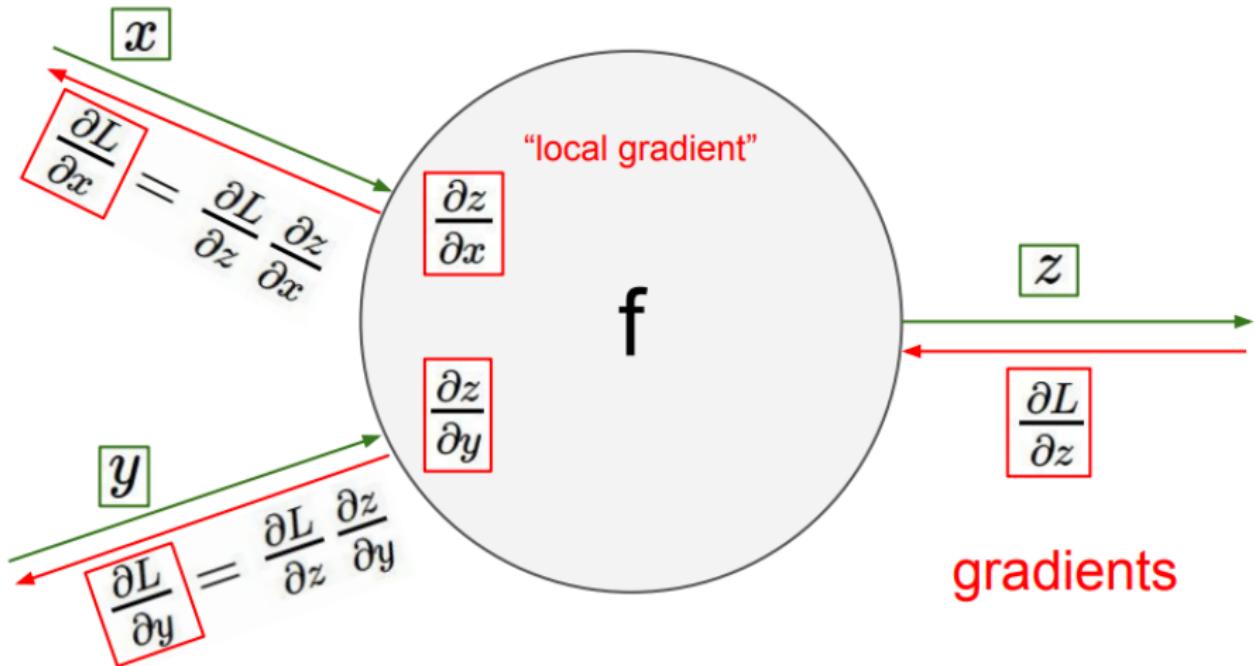
$$\frac{dL}{dx} = \frac{dL}{dz} \frac{dz}{dx}$$

$$\frac{dL}{dy} = \frac{dL}{dz} \frac{dz}{dy}$$

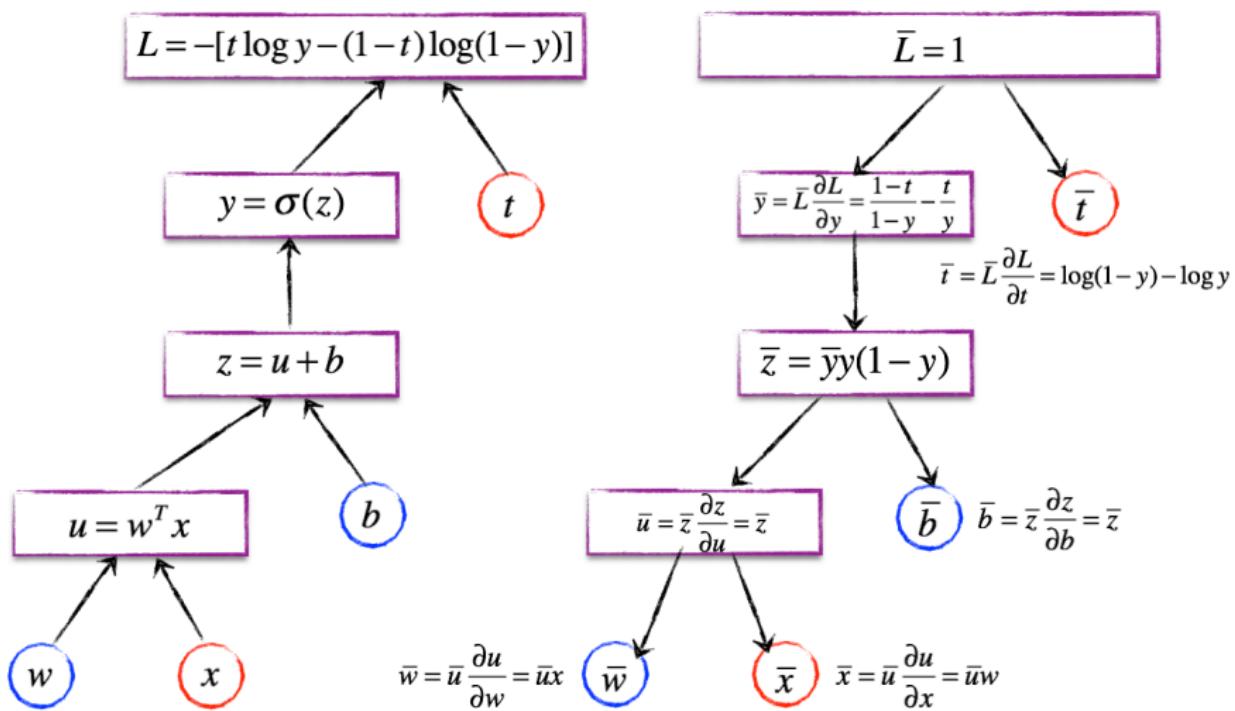
$$\frac{dL}{dz}$$

自动微分

局部梯度



例子：逻辑回归



优化器

```
class Optimizer(object):
    def minimize(self, loss, var_list=None, global_step=None):
        grads_and_vars = self.compute_gradients(
            loss, var_list=var_list)
        return self.apply_gradients(
            grads_and_vars,
            global_step=global_step)

    def compute_gradients(loss, var_list):
        grads = gradients(loss, var_list, grad)
        return list(zip(grads, var_list))

    def gradients(loss, var_list, grads=1):
        ops_and_grads = {}
        for op in reversed_graph(loss).topological_sort():
            grad = op.grad_fn(grad)
            ops_and_grads[op] = grad
        return [ops_and_grads.get(var) for var in var_list]
```

梯度函数

```
@ops.RegisterGradient("op_name")
def grad_func(op, grad):
    """construct gradient subgraph for an op type.
    Returns:
        A list of gradients, one per each input of op.
    """
    return cons_grad_subgraph(op, grad)
```

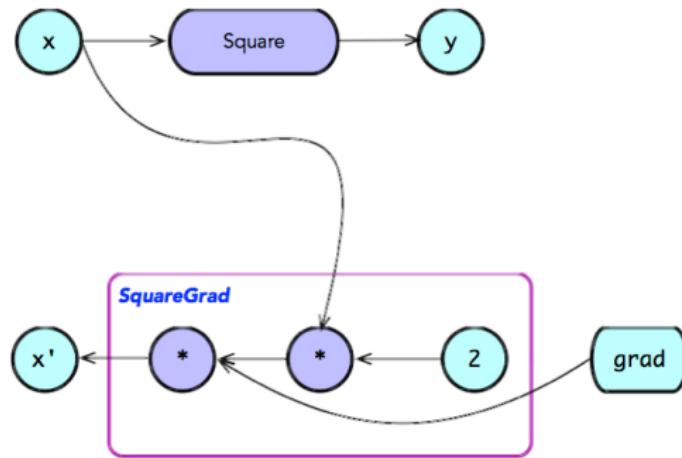
$$(y_1, y_2, \dots, y_m) = f(x_1, x_2, \dots, x_n)$$

$$\left(\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2}, \dots, \frac{\partial L}{\partial x_n}\right) = g\left(x_1, x_2, \dots, x_n; \frac{\partial L}{\partial y_1}, \frac{\partial L}{\partial y_2}, \dots, \frac{\partial L}{\partial y_n}\right)$$

自动微分

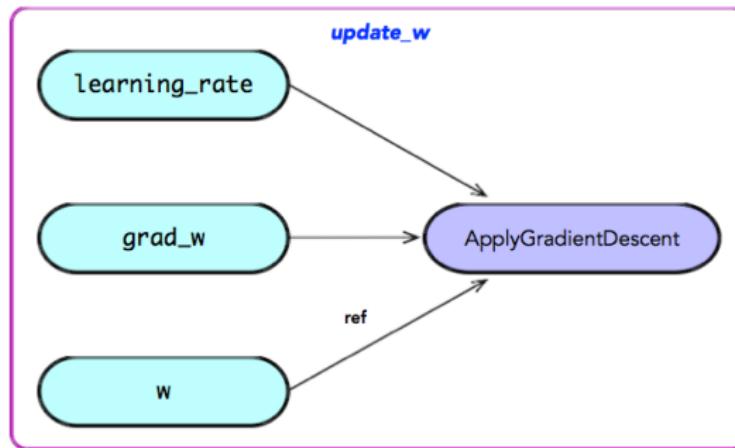
例子：square

```
@ops.RegisterGradient("Square")
def SquareGrad(op, grad):
    x = op.inputs[0]
    with ops.control_dependencies([grad.op]):
        return grad * (2.0 * x)
```



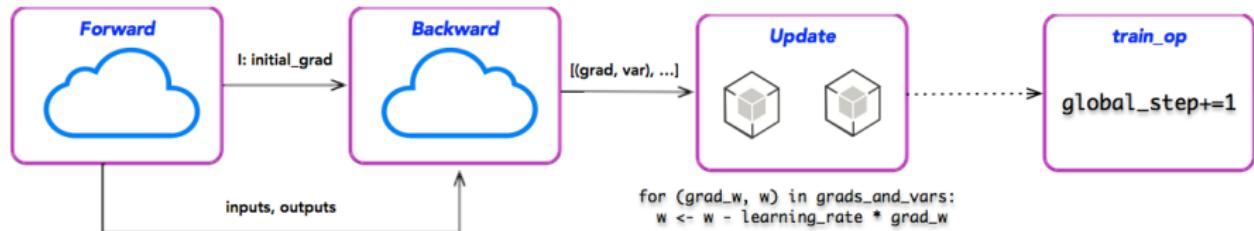
应用梯度

```
def apply_gradients(grads_and_vars, learning_rate):
    for (grad, var) in grads_and_vars:
        apply_gradient_descent(learning_rate, grad, var)
```



自动微分

关键路径: RunStep



致谢

答疑





致谢

谢谢