

# I. Supervised Learning and Optimization

## ① Linear Regression

a hypothesis class:  $h_{\theta}(x) \rightarrow x^{(i)} \rightarrow y^{(i)}$

$$\text{Linear function: } h_{\theta}(x) = \theta^T x$$

$$\text{cost function: } J(\theta) = \frac{1}{2} \sum_i (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

(also called loss / penalty / objective function)

$\frac{1}{2} \sum_i (\theta^T x^{(i)} - y^{(i)})^2$  (可用 loop 算)

$= \frac{1}{2} (\theta^T X - Y)^T \cdot (\theta^T X - Y)$  (也可用向量形式)

problem formulation

$$\begin{array}{c} \theta^T \xrightarrow{\quad} \left( \begin{array}{c} x^{(i)} \\ \vdots \\ x^{(i)} \end{array} \right) - y^{(i)} \\ \left( \begin{array}{c} \theta^T \\ \vdots \\ \theta^T \end{array} \right) \left( \begin{array}{c} x^{(1)} \\ \vdots \\ x^{(n)} \end{array} \right) \\ \text{Train.X} \in \mathbb{R}^{14 \times 400} \end{array}$$

Function Minimization:

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \frac{\partial J(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix} \quad \text{且} \quad \frac{\partial J(\theta)}{\partial \theta_j} = \sum_i (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$g(j) = \frac{(\theta^T \cdot X - Y)}{14 \times 400} \cdot \begin{pmatrix} x_1^{(j)} \\ x_2^{(j)} \\ \vdots \\ x_{400}^{(j)} \end{pmatrix} \quad j=1, \dots, n \quad (n=14)$$

Exercise 1A: Linear Regression

$$\text{house\_data: } \begin{bmatrix} 1 & 1 & \cdots & 1 & 1 \\ 0.0449 & 6.4941 & \cdots & 1.0025 & 1.1517 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 22.900 & 16.100 & \cdots & 7.1 & 13.1000 \end{bmatrix} \quad 15 \times 506$$

- 列 1 表示 1 例，每行代表 1 feature

`size(data, 1)` ⇒ 返回 data 针对的行数

`size(data, 2)` ⇒ 返回 data 针对的列数

$$\text{train: } \begin{cases} \text{train.X} & 14 \times 400 \\ \text{train.y} & 1 \times 400 \end{cases} \quad \text{test: } \begin{cases} \text{test.X} & 14 \times 106 \\ \text{test.y} & 1 \times 106 \end{cases}$$

$$n: \boxed{\text{train.X}} \quad m: 400$$

$$\theta^T \cdot \text{test.X} \Rightarrow \text{test.y}$$

$$\theta \in \mathbb{R}^{14 \times 1}$$

## ② Logistic Regression

a hypothesis class  $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} = \sigma(\theta^T x) \rightarrow \text{sigmoid function}$

$$p(y=1|x) = h_{\theta}(x)$$

$$p(y=0|x) = 1 - h_{\theta}(x)$$

we label example  $\rightarrow 1, \text{ if } h_{\theta}(x) > 0.5$   
 $\rightarrow 0, \text{ otherwise}$

loss function:

$$J(\theta) = - \sum_{i=1}^m \left[ y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right]$$

若用梯度下降法(减少  $J(\theta)$ )

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_i (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \xrightarrow{\text{向量化写法}} \nabla_{\theta} J = \sum_i (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

5 linear regression 中梯度下降法  
梯度，只不過這裡的  $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$

$$[y^{(1)} \ y^{(2)} \dots \ y^{(m)}] \cdot \log \begin{bmatrix} h_\theta(x^{(1)}) \\ h_\theta(x^{(2)}) \\ \vdots \\ h_\theta(x^{(m)}) \end{bmatrix}$$

$$f = y * \log(h(x \cdot \theta)) + (1-y) * \log(1 - h(x \cdot \theta))$$

$$g \begin{bmatrix} g^{(1)} \\ \vdots \\ g^{(m)} \end{bmatrix} = \left[ h(\theta \cdot x) - y \right] \quad \text{temp} = \theta \cdot x$$

$\times m$  用向量化的写法就快多了。

## Exercise 1B

Mnist : D、1 级别

从  $X = \text{load MNIST Images ('...|train-images-idx3-ubyte')}$   $\rightarrow$  [ ] 什的标准化  
0均值, 1方差  
值都在 0~1 之间  
 $784 \times 60000$

$y = \text{load MNIST Labels ('...|train-labels-idx1-ubyte')}$   $\rightarrow$  [---]  $0-9$  的标签  
 $1 \times 60,000$

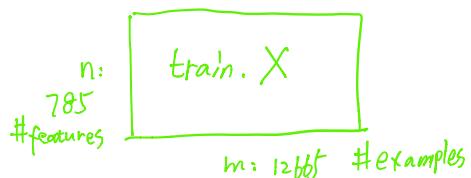
60,000 个中是 0 和 1 的共有 12665 个 (5923 个 0, 6742 个 1) / 60,000 个中有 0 和 1 共有 2115 个 ( $\frac{980}{1135}$ )

train.X  $784 \times 12665$  test.X  $784 \times 2115$  (Matlab 中也有 ( ).( ) 这种点操作, 个)

train.y  $1 \times 12665$  test.y  $1 \times 2115$  有点像 C 中的结构体)

增加一个 2 后

train.X  $785 \times 12665$  test.X  $785 \times 2115$



普通的方法复杂度：(自己特意写成最慢的方法)

$$\left( 785 \times 12665 + 12665 + 785 \times 12665 \times 785 + 785 \times 12665 \right) \times 100 \quad 7\text{小时} \rightarrow 20\text{次} \quad 12665$$

$$= 7.8 \times 10^{11} = 780 \text{ 亿元} \text{ (过于复杂)}$$

### ③ Debugging : Gradient Checking

$$g_i(\theta) = \frac{\partial}{\partial \theta_i} J(\theta) \quad \text{where } \vec{e}_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \leftarrow \text{with}$$

$$\theta^{(i+1)} = \theta + \varepsilon \times \vec{e}_i$$

$$\theta^{(i-1)} = \theta - \varepsilon \times \vec{e}_i$$

$$g_i(\theta) = \frac{J(\theta^{(i+1)}) - J(\theta^{(i-1)})}{2\varepsilon}$$

### ④ softmax regression (这个代码写的不是很正确)

是 logistic regression 的更一般化形式，针对多类别

$$h_{\theta}(x) = \begin{bmatrix} p(y=1|x; \theta) \\ p(y=2|x; \theta) \\ \vdots \\ p(y=k|x; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta^{(j)\top} x}} \begin{bmatrix} e^{\theta^{(1)\top} x} \\ e^{\theta^{(2)\top} x} \\ \vdots \\ e^{\theta^{(k)\top} x} \end{bmatrix}$$

$$\text{参数: } \theta = \begin{bmatrix} & & & \\ \theta^{(1)} & \theta^{(2)} & \theta^{(3)} & \cdots & \theta^{(K)} \\ & & & & \end{bmatrix}_{n \times k}$$

cost function:

在之前的 Logistic regression 中，我们的损失函数是

$$J(\theta) = - \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)})) \right]$$

$X_o: 7025 \times 1$	$785 \times 9$
$d: 7850 \times 1$	
$X: 785 \times 60000$	

$$J(\theta) = - \left[ \sum_{i=1}^m \sum_{k=0}^K I\{y^{(i)}=k\} \log P(y^{(i)}=k | X^{(i)}; \theta) \right]$$

now we have  $K$  classes

$$\text{so } J(\theta) = - \left[ \sum_{i=1}^m \sum_{k=1}^K I\{y^{(i)}=k\} \log P(y^{(i)}=k | X^{(i)}; \theta) \right]$$

$$= - \left[ \sum_{i=1}^m \sum_{k=1}^K I\{y^{(i)}=k\} \log \frac{e^{\theta^{(k)\top} X^{(i)}}}{\sum_{j=1}^K e^{\theta^{(j)\top} X^{(i)}}} \right]$$

$$I\{y^{(i)}=1\} \log \frac{e^{\theta^{(1)\top} X^{(i)}}}{\sum e^{\theta^{(j)\top} X^{(i)}}} + I\{y^{(i)}=2\} \log \frac{e^{\theta^{(2)\top} X^{(i)}}}{\sum e^{\theta^{(j)\top} X^{(i)}}}$$

$$\begin{bmatrix} I\{y^{(1)}=1\} \\ I\{y^{(1)}=2\} \\ \vdots \\ I\{y^{(1)}=k\} \end{bmatrix}^T \begin{bmatrix} \log \frac{e^{\theta^{(1)\top} X^{(1)}}}{\sum e^{\theta^{(j)\top} X^{(1)}}} \\ \vdots \\ \log \frac{e^{\theta^{(1)\top} X^{(1)}}}{\sum e^{\theta^{(j)\top} X^{(1)}}} \end{bmatrix} + \begin{bmatrix} I\{y^{(2)}=1\} \\ I\{y^{(2)}=2\} \\ \vdots \\ I\{y^{(2)}=k\} \end{bmatrix}^T \begin{bmatrix} \log \frac{e^{\theta^{(2)\top} X^{(2)}}}{\sum e^{\theta^{(j)\top} X^{(2)}}} \\ \vdots \\ \log \frac{e^{\theta^{(2)\top} X^{(2)}}}{\sum e^{\theta^{(j)\top} X^{(2)}}} \end{bmatrix} + \dots + \begin{bmatrix} I\{y^{(m)}=1\} \\ I\{y^{(m)}=2\} \\ \vdots \\ I\{y^{(m)}=k\} \end{bmatrix}^T \begin{bmatrix} \log \frac{e^{\theta^{(m)\top} X^{(m)}}}{\sum e^{\theta^{(j)\top} X^{(m)}}} \\ \vdots \\ \log \frac{e^{\theta^{(m)\top} X^{(m)}}}{\sum e^{\theta^{(j)\top} X^{(m)}}} \end{bmatrix}$$

$$A = e^{(\theta^\top X)} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

sums  
↓ ↓ ↓ ↓  
↓ ↓ ↓ ↓  
↓ ↓ ↓ ↓

$B = 65 \times 65$

$$Y: \theta^\top X \quad 9 \times 60000$$

$$I = \text{sub2ind}(\text{size}(\theta^\top X), Y, 1: \text{size}(\theta^\top X, 2))$$

$$-\text{sum}(\theta^\top X(I)) = J$$

$$\theta = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \theta^{(1)} & \theta^{(2)} & \dots & \theta^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix}_{785 \times 9}$$

$$X: \begin{bmatrix} \dots \end{bmatrix}_{785 \times 60000}$$

$$Y: \begin{bmatrix} \dots \end{bmatrix}_{1 \times 60000}$$

梯度：

$$\nabla_{\theta^{(m)}} J(\theta) = - \sum_{i=1}^m \left[ X^{(i)} (I\{y^{(i)}=k\} - P(y^{(i)}=k | X^{(i)}; \theta)) \right]$$

(应该也可以从之前的  $J(\theta)$  的表达式推出来的)

暂且先不推导了。

$$\text{因为 } \frac{\partial J}{\partial \theta} = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix} \text{ 全部是一个向量}$$

$y = [x \ 60000]$   
debug of 现在的梯度匹配

向量形式表示：

$$g = \begin{bmatrix} 1 \\ g^{(1)} \\ | \\ g^{(2)} \\ | \\ g^{(3)} \\ | \\ \dots \\ g^{(L)} \\ | \end{bmatrix}$$

$$g(i) = \sum_{j=1}^L (x^{(j)}(1 \leq j \leq i) - p(y^{(j)}=1 | x^{(j)}; \theta))$$

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} 1 \{y^{(1)}=1\} - p(y^{(1)}=1 | x^{(1)}; \theta) \\ 1 \{y^{(2)}=1\} - p(y^{(2)}=1 | x^{(2)}; \theta) \\ \vdots \\ 1 \{y^{(m)}=1\} - p(y^{(m)}=1 | x^{(m)}; \theta) \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ x^{(m)} \\ | \end{bmatrix} \begin{bmatrix} 1 \{y^{(1)}=1\} - \frac{e^{\theta^{(1)\top} x^{(1)}}}{\sum_{j=1}^L e^{\theta^{(j)\top} x^{(j)}}} \\ 1 \{y^{(2)}=1\} - \frac{e^{\theta^{(2)\top} x^{(2)}}}{\sum_{j=1}^L e^{\theta^{(j)\top} x^{(j)}}} \\ \vdots \\ 1 \{y^{(m)}=1\} - \frac{e^{\theta^{(m)\top} x^{(m)}}}{\sum_{j=1}^L e^{\theta^{(j)\top} x^{(j)}}} \end{bmatrix}$$

$$g^{(2)} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \dots$$

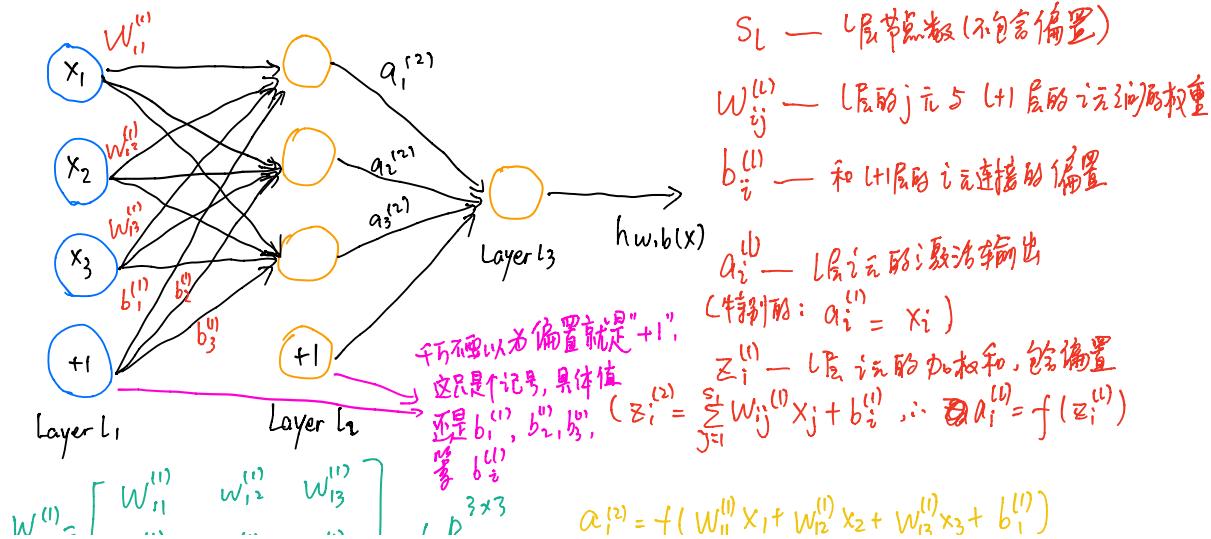
$$\begin{bmatrix} 1 \{y^{(1)}=2\} - \frac{e^{\theta^{(1)\top} x^{(1)}}}{\sum_{j=1}^L e^{\theta^{(j)\top} x^{(j)}}} \\ 1 \{y^{(2)}=2\} - \frac{e^{\theta^{(2)\top} x^{(2)}}}{\sum_{j=1}^L e^{\theta^{(j)\top} x^{(j)}}} \\ \vdots \\ 1 \{y^{(m)}=2\} - \frac{e^{\theta^{(m)\top} x^{(m)}}}{\sum_{j=1}^L e^{\theta^{(j)\top} x^{(j)}}} \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ x^{(m)} \\ | \end{bmatrix} k_1 + \begin{bmatrix} 1 \\ x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ x^{(m)} \\ | \end{bmatrix} k_2 + \begin{bmatrix} 1 \\ x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ x^{(m)} \\ | \end{bmatrix} k_3 + \dots + \begin{bmatrix} 1 \\ x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ x^{(m)} \\ | \end{bmatrix} k_m$$

$$= \begin{bmatrix} 1 \\ x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ x^{(m)} \\ | \end{bmatrix} \begin{pmatrix} k_1 \\ k_2 \\ k_3 \\ \vdots \\ k_m \end{pmatrix}$$

## I. Supervised Neural Networks

### ⑤ multi-layer neural networks



$$W^{(l)} = \begin{bmatrix} W_{11}^{(l)} & W_{12}^{(l)} & W_{13}^{(l)} \\ W_{21}^{(l)} & W_{22}^{(l)} & W_{23}^{(l)} \\ W_{31}^{(l)} & W_{32}^{(l)} & W_{33}^{(l)} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$b^{(l)} = [b_1^{(l)} \ b_2^{(l)} \ b_3^{(l)}] \in \mathbb{R}^{1 \times n}$$

element-wise: 向量之间对应元素的操作  
向量化的写法:

feedforward neural network  
前馈神经网络  
(没有 loop 或 cycle)

element-wise product

也叫作 hadamard product

$$\alpha_1^{(2)} = f(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)})$$

$$\alpha_2^{(2)} = f(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)})$$

$$\alpha_3^{(2)} = f(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)})$$

$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$\alpha^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}\alpha^{(2)} + b^{(2)}$$

$$h_{w,b}(x) = \alpha^{(3)} = f(z^{(3)})$$

总结来说

→ step forward propagation

$$J(W, b; x, y) = \frac{1}{2} \|h_{w,b}(x) - y\|^2 \text{ 对于单个样本 } (x, y) \text{ 来说, } J \text{ — cost function}$$

w<sub>b</sub>代表是 w<sub>b</sub>的函数 x<sub>j</sub>代表给 x(y) 样本时

given m examples:

$$J(W, b) = \left[ \frac{1}{m} \sum_{i=1}^m J(W, b; x^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (W_{j|i}^{(l)})^2$$

参数要随机初始化

不能全设为 0

symmetry breaking

pseudo-code: 伪代码

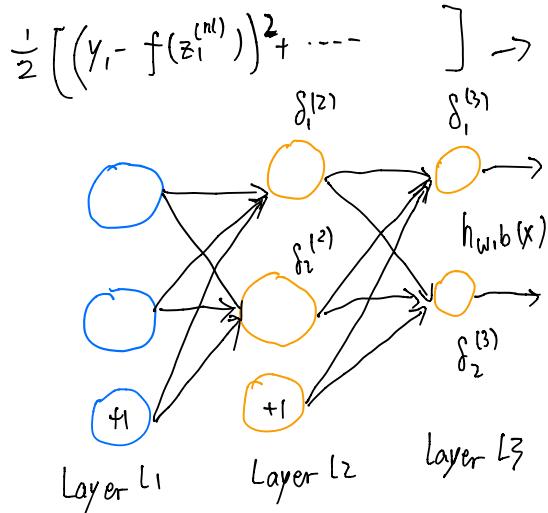
$$= \left[ \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{2} \|h_{w,b}(x^{(i)}) - y^{(i)}\|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (W_{j|i}^{(l)})^2$$

regularization term  
(weight decay term)  $\rightarrow$  通常对  $b^{(l)}$   
正则化相差不大

梯度下降迭代:

$$W_{j|i}^{(l)} = W_{j|i}^{(l)} - \alpha \frac{\partial}{\partial W_{j|i}^{(l)}} J(W, b)$$

$$b_i^{(l)} = b_i^{(l)} - \alpha \frac{\partial}{\partial b_i^{(l)}} J(w, b)$$



誤差  $J = \frac{1}{2} [(y_1 - f(z_1^{(3)}))^2 + (y_2 - f(z_2^{(3)}))^2]$

$$\begin{aligned} z_1^{(3)} &= W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + b_1^{(2)} \\ &= W_{11}^{(2)} f(z_1^{(2)}) + W_{12}^{(2)} f(z_2^{(2)}) + b_1^{(2)} \end{aligned}$$

$$\begin{aligned} z_2^{(3)} &= W_{21}^{(2)} a_1^{(2)} + W_{22}^{(2)} a_2^{(2)} + b_2^{(2)} \\ &= W_{21}^{(2)} f(z_1^{(2)}) + W_{22}^{(2)} f(z_2^{(2)}) + b_2^{(2)} \end{aligned}$$

$$\delta_1^{(2)} = \frac{\partial J}{\partial z_1^{(2)}} = \frac{\partial J}{\partial \delta_1^{(3)}} \cdot \frac{\partial \delta_1^{(3)}}{\partial z_1^{(2)}} = \delta_1^{(3)} \cdot (W_{11}^{(2)} f'(z_1^{(2)})) = \delta_2^{(3)} \cdot (W_{21}^{(2)} f'(z_1^{(2)}))$$

$$\frac{\partial J(w, b; x, y)}{\partial W_{11}^{(2)}} = a_1^{(2)} \delta_1^{(3)} = a_1^{(2)} (-1(y_1 - a_1^{(3)}) f'(z_1^{(3)})) \quad \checkmark$$

$$\frac{\partial J(w, b; x, y)}{\partial W_{11}^{(1)}} = a_1^{(1)} \delta_1^{(2)} = a_1^{(1)} (W_{11}^{(2)} \delta_1^{(3)} + W_{21}^{(2)} \delta_2^{(3)}) f'(z_1^{(2)}) \quad \checkmark$$

$$\frac{\partial J(w, b; x, y)}{\partial b_1^{(2)}} = \delta_1^{(3)} \quad \checkmark$$

$$\frac{\partial J(w, b; x, y)}{\partial b_2^{(2)}} = \delta_2^{(3)} \quad \checkmark$$

$$\begin{aligned} z^{(l+1)} &= W^{(l)} \otimes a^{(l)} + b^{(l)} \\ a^{(l+1)} &= f(z^{(l+1)}) \\ \delta_i^{(n_l)} & \end{aligned}$$

251784	10x756
786	10x1
784	10
δ <sup>(3)</sup>	10x1

$$\alpha_1^{(2)} = f(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + b_1^{(1)})$$

$$\alpha_2^{(2)} = f(W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + b_2^{(1)})$$

$$\alpha_1^{(3)} = h_{w,b}(x) = f(W_{11}^{(2)} \alpha_1^{(2)} + W_{12}^{(2)} \alpha_2^{(2)} + b_1^{(2)})$$

$$\alpha_2^{(3)} = h_{w,b}(x) = f(W_{21}^{(2)} \alpha_1^{(2)} + W_{22}^{(2)} \alpha_2^{(2)} + b_2^{(2)})$$

$$\delta_1^{(3)} = \frac{\partial J}{\partial z_1^{(3)}} = \frac{1}{2} \cdot 2(y_1 - f(z_1^{(3)})) \cdot -f'(z_1^{(3)})$$

$$= -(y_1 - \alpha_1^{(3)}) f'(z_1^{(3)})$$

$$\delta_2^{(3)} = -(y_2 - \alpha_2^{(3)}) f'(z_2^{(3)})$$

set:  $\rightarrow \delta_1^{(2)} = (W_{11}^{(2)} \delta_1^{(3)} + W_{21}^{(2)} \delta_2^{(3)}) f'(z_1^{(2)})$

$$\text{公式①: for } l = n-1, n-2, \dots, 1: \delta_i^{(l)} = \left( \sum_{j=1}^{S_{l+1}} w_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$$

$$\text{公式②: } \frac{\partial J(w, b; x, y)}{\partial w_{ij}^{(l)}} = a_j^{(l)} \delta_i^{(l+1)}$$

$$\frac{\partial J(w, b; x, y)}{\partial b_i^{(l)}} = \delta_i^{(l+1)}$$

$$\begin{aligned}\text{证明公式①: } \delta_i^{(l)} &= \frac{\partial J}{\partial z_i^{(l)}} = \sum_{j=1}^{S_{l+1}} \frac{\partial J}{\partial z_j^{(l+1)}} \frac{\partial z_j^{(l+1)}}{\partial z_i^{(l)}} \\ &= \sum_{j=1}^{S_{l+1}} \delta_j^{(l+1)} w_{ji}^{(l)} \underbrace{f'(z_i^{(l)})}_{\text{都有这一项, 提出来}} \\ &= \left( \sum_{j=1}^{S_{l+1}} w_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})\end{aligned}$$

$$\text{证明公式②: } \frac{\partial J}{\partial w_{ij}^{(l)}} = \frac{\partial J}{\partial z_i^{(l+1)}} \cdot \frac{\partial z_i^{(l+1)}}{\partial w_{ij}^{(l)}} = \delta_i^{(l+1)} \cdot a_j^{(l)}$$

$$\frac{\partial J}{\partial b_i^{(l)}} = \frac{\partial J}{\partial z_i^{(l+1)}} \frac{\partial z_i^{(l+1)}}{\partial b_i^{(l)}} = \delta_i^{(l+1)} \quad \begin{array}{l} \text{这一项是 1} \\ \text{下面是向量化} \end{array}$$

向量写法:

$$W^{(l)} : S_{l+1} \times S_l$$

$$\delta^{(l+1)} : S_{l+1} \times m$$

$$\delta^{(l)} : S_l \times m$$

两个  
维度相  
同  
 $z^{(l)} : S_l \times m$

同  
 $\delta^{(l)} :$

$m$  为 example 的数量。

本例中:  
Stack (cell)

$$W^{(1)}_{256 \times 784} \quad W^{(2)}_{10 \times 256}$$

$$b^{(1)}_{256 \times 1} \quad b^{(2)}_{10 \times 1}$$

$$[256 \times 784] \cdot [784 \times 60000]$$

$$= [256 \times 60000]$$

$$[10 \times 256] \cdot [256 \times 60000]$$

$$= [10 \times 60000]$$

Backpropagation	Algorithm 向量化(用“ $\bullet$ ”表示 hadamard product)
非向量化	
1. 前向传播, 算出 $L_2 \cup L_{nl} \cup L_0$ 的 $a_i, z_i$	1. 前向传播, 算出 $L_2 \cup L_{nl} \cup L_0$ 的 $a_i, z_i$
2. 对于层 $l_n$ (输出层), 设 $\delta_i^{(nl)} = -(y_i - a_i^{(nl)}) \cdot f'(z_i^{(nl)})$	2. 对于层 $l_n$ (输出层), 设 $\delta_i^{(nl)} = -(y_i - a_i^{(nl)}) \bullet f'(z_i^{(nl)})$
3. 对于层 $l_{n-1}, l_{n-2}, \dots, l_1$ , 设 $\delta_i^{(l)} = (\sum_{j=1}^{ S^{(l+1)} } W_{ji}^{(l)} \delta_j^{(l+1)}) f'(z_i^{(l)})$	3. 对于层 $l_{n-1}, l_{n-2}, \dots, l_1$ , 设 $\delta_i^{(l)} = ((W^{(l)})^\top \delta^{(l+1)}) \bullet f'(z_i^{(l)})$
4. $\frac{\partial}{\partial W_{ij}^{(l)}} J(w, b; x, y) = a_j^{(l)} \delta_i^{(l+1)}$	4. $\nabla_w J(w, b; x, y) = \delta^{(l+1)} (a^{(l)})^\top$
$\frac{\partial}{\partial b_i^{(l)}} J(w, b; x, y) = \delta_i^{(l+1)}$	$\nabla_b J(w, b; x, y) = \delta^{(l+1)}$
	用 BP 算梯度.

### Batch Gradient Descent

1. set  $\Delta W^{(l)} := 0, \Delta b^{(l)} := 0$  for all layer  $l$
2. For  $i = 1$  to  $m$ 
  1. use backpropagation algorithm to compute  $\nabla_w J(w, b; x, y)$  and  $\nabla_b J(w, b; x, y)$
  2. set  $\Delta W^{(l)} := \Delta W^{(l)} + \nabla_w J(w, b; x, y)$   
 $\Delta b^{(l)} := \Delta b^{(l)} + \nabla_b J(w, b; x, y)$
3. update the parameters:  
 $w^{(l)} = w^{(l)} - \alpha [(\frac{1}{m} \Delta W^{(l)}) + \lambda w^{(l)}]$   
 $b^{(l)} = b^{(l)} - \alpha [\frac{1}{m} \Delta b^{(l)}]$

### Exercise: Supervised Neural Network

ei:

- input\_dim = 784
- output\_dim = 10
- layer\_sizes = [256, 10]
- lambda = 0

这一个  
Struct  
P

Stack (cell)  
((1), ..., (2), ...)

而 cell 里面的单元也可以为一个 struct, 可以用 {} 来索引

$$\text{activation\_fun} = \text{'logistic'}$$

$$J(\theta) = - \left[ \sum_{i=1}^m \sum_{k=1}^K \mathbb{1}\{y^{(i)} = k\} \log \frac{e^{\theta^{(k),T} h_{w,b}(x^{(i)})}}{\sum_{j=1}^K e^{\theta^{(j),T} h_{w,b}(x^{(i)})}} \right]$$

$$W^{(1)}_{256 \times 784} \quad v^{(1)}_{10 \times 256}$$

$$b^{(1)}_{256 \times 1} \quad b^{(2)}_{10 \times 1}$$

$$[256 \times 784] \cdot [784 \times 60000] \\ = [256 \times 60000]$$

$$[10 \times 256] \cdot [256 \times 60000] \\ = [10 \times 60000]$$

$$g^{(n)} = - \sum_{i=1}^m \left[ \mathbb{1}\{y^{(i)} = k\} - P(y^{(i)} = k | x^{(i)}; \theta) \right]$$

### III. Supervised Convolutional Neural Network

#### ⑥ Feature Extraction Using Convolution & Pooling

matlab中二维卷积的定义：convolve：

$$C(j, k) = \sum_{p=1}^P \sum_{q=1}^Q A(p, q) B(j-p+1, k-q+1)$$

P和q会遍历所有满足  $A(p, q)$  和  $B(j-p+1, k-q+1)$  的合法下标的值

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad \begin{bmatrix} 5 & 6 & 7 \\ 8 & 9 & 10 \end{bmatrix} \xrightarrow{\text{conv2}} \begin{bmatrix} 5 & 16 & 19 & 14 \\ 23 & 63 & 73 & 48 \\ 24 & 59 & 66 & 40 \end{bmatrix}$$

A                    B                    C

e.g.

$$C(1, 1) = \sum_{p=1}^P \sum_{q=1}^Q A(p, q) B(2-p, 2-q) = A(1, 1)B(1, 1) + \cancel{A(1, 2)B(1, 0)} \\ + \cancel{A(2, 1)B(0, 1)} + \cancel{A(2, 2)B(0, 0)} = 5$$

$$C(2, 2) = \sum_{p=1}^P \sum_{q=1}^Q A(p, q) B(3-p, 3-q) = A(1, 1)B(2, 2) + A(1, 2)B(2, 1) + A(2, 1)B(1, 2) \\ + A(2, 2)B(1, 1) = 9 + 16 + 18 + 20 = 63$$

:

$$\begin{matrix} A_{3 \times 3} \\ B_{2 \times 4} \end{matrix} \xrightarrow{\text{conv2}} C_{6 \times 6}$$

另一种直观的解释

1437

1427

1437

1437

$$\begin{bmatrix} 2 & [5] & 6 & 7 \\ & 8 & 9 & 10 \end{bmatrix} = 5 \quad \begin{bmatrix} 5 & [6] & 7 \\ 8 & 9 & 10 \end{bmatrix} = 16 \quad B: \begin{bmatrix} 5 & [6, 7] \\ 8 & 9 & 10 \end{bmatrix} = 19 \quad \begin{bmatrix} 5 & 6 & [7] \\ 8 & 9 & 10 \end{bmatrix} = 14$$

$$\begin{bmatrix} 4 & [3, 5] & 6 & 7 \\ 2 & 8 & 9 & 10 \end{bmatrix} = 23 \quad \begin{bmatrix} 5 & [6, 3] & 7 \\ 2 & 8 & 9 & 10 \end{bmatrix} = 63 \quad B: \begin{bmatrix} 5 & [6, 7] \\ 8 & [2, 9, 10] \end{bmatrix} = 73 \quad \begin{bmatrix} 5 & 6 & [7, 3] \\ 8 & 9 & [2, 10] \end{bmatrix} = 48$$

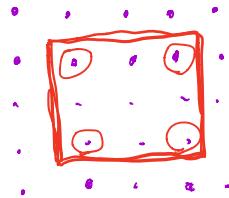
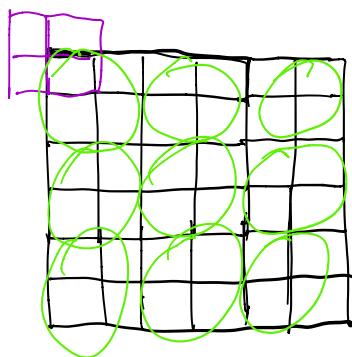
$$\begin{bmatrix} 5 & 6 & 7 \\ 4 & [3, 8] & 9 & 10 \\ 2 & 1 \end{bmatrix} = 24 \quad \begin{bmatrix} 5 & 6 & 7 \\ [8, 4, 3, 7] & 10 \\ 2 & 1 \end{bmatrix} = 57 \quad B: \begin{bmatrix} 5 & 6 & 7 \\ 8 & [7, 4, 3, 10] \\ 2 & 1 \end{bmatrix} = 66 \quad \begin{bmatrix} 5 & 6 & 7 \\ 8 & 9 & [10, 4, 3] \\ 2 & 1 \end{bmatrix} = 60$$

(numImages, numFilters, imageRow, imageCol) 等于  
 60000 100 28 28

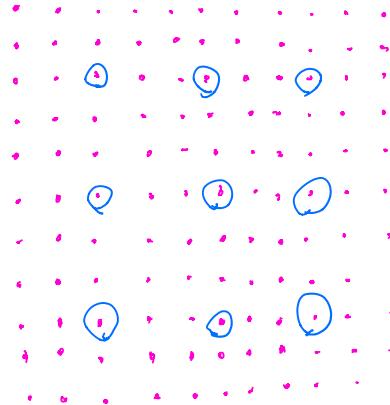
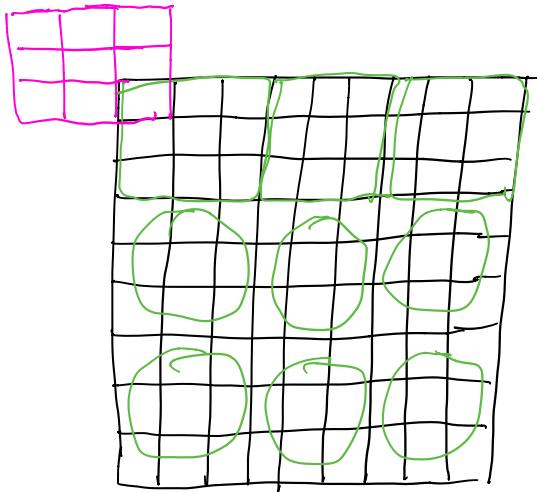
convolved Features 的 shape: (imageRow, imageCol, featureNum, imageNum)

pool:

21, 21, 100, 8



for i in range(poolDim):  
 $\left(ik * poolDim + 1 : i * poolDim, (j-1) * poolDim + 1 : j * poolDim\right)$



用 conv2 来 pool:  $(\text{poolDim}^{(i-1)} + 1, \text{poolDim}^{(j-1)} + 1)$

## ⑦ Optimization: Stochastic Gradient Descent L7 slide

good off the shelf 现成的

Batch Optimization Method: 能直接收敛至局部最优，实现容易，需要调整的超参数少。

Stochastic Gradient Descent: 不一次用整个数据集，只用一个或一些，解决反向传播耗时大的问题。  
通常用的是 mini-batch : 256，其中的学习率又要比 SGD 小

$$\theta = \theta - \alpha \nabla_{\theta} E[J(\theta)] \longrightarrow \theta = \theta - \alpha \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$$

学习率的正确选择是挺困难的

- 标准的做法：每两轮 epoch 后，减半  $\alpha$ （固定的值）
- 一个更好的：每 epoch 后，如果目标差小于阈值，则对  $\alpha$  进行退火。
- 另一个常用的：随着迭代次数  $t$  退火： $\alpha = \frac{\alpha}{t+1}$

SGD 输入数据 最好是打乱顺序的，不是刻意排列的  
randomly shuffle

Momentum (动量；势能；冲力)

→ 用 SGD 在 local optimum 附近会振荡，收敛较慢



steep wall 可以用 momentum :

$$v = \gamma v + \alpha \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$$

$\theta = \theta - v$   
↓ 连起来就是有  $\theta = \theta - \text{梯度的感覺}$

$\eta \in [0, 1]$ , generally,  $\eta$  is set to 0.5 until the initial learning stabilizes and then is increased to 0.9 or higher.

$D, V$  维数形状向量的向量,  $V$  是当前速度向量,  $D$  是参数向量

## ⑧ Convolutional Neural Network

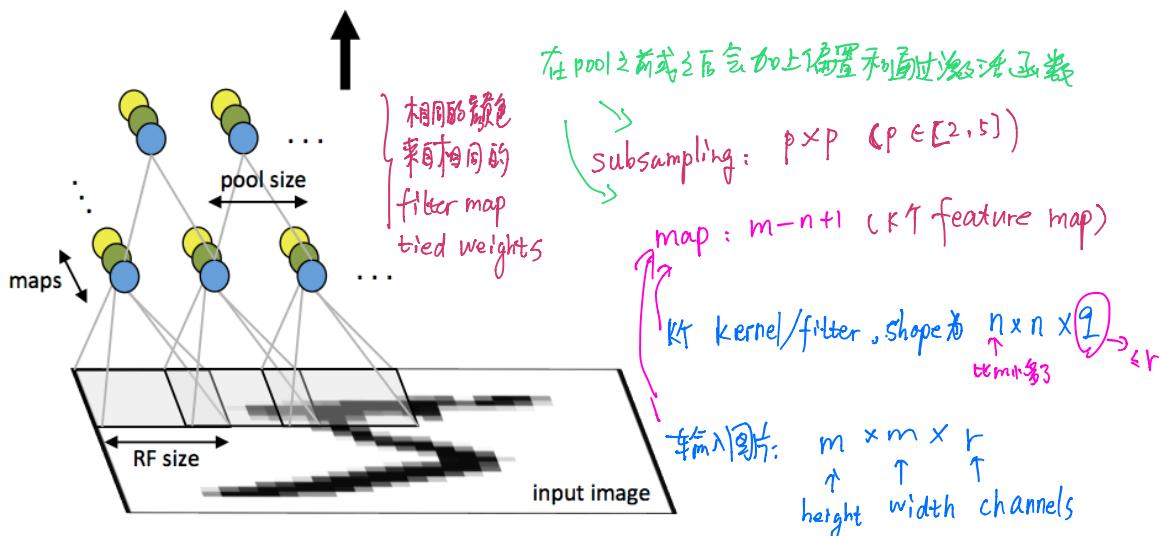
利用了图像的二维结构

CNN 在有相同隐藏神经元时, 参数为连接网络;

architecture

back propagation algorithm

Architecture:



## Back Propagation

w.r.t. —— with respect to 关于;涉及

记  $\delta^{(l+1)}$  为  $(l+1)$ -st layer 的 error term

如果  $l$  层与  $(l+1)$  层全连接的话, 则  $\delta^{(l)} = ((W^{(l)})^\top \delta^{(l+1)}) \cdot f'(z^{(l)})$

$$\begin{aligned} \text{梯度: } \nabla_{W^{(l)}} J(W, b; x, y) &= \delta^{(l+1)} (a^{(l)})^\top \\ \nabla_{b^{(l)}} J(W, b; x, y) &= \delta^{(l+1)} \end{aligned}$$

如果  $l$  层为卷积层与下采样层, 则  $\delta_K^{(l)} = \text{unsample} ((W_K^{(l)})^\top \delta_K^{(l+1)}) \cdot f'(z_K^{(l)})$

其中  $K$  表示 filter number

$$\text{梯度: } \nabla_{w_k^{(l)}} J(w, b; x, y) = \sum_{i=1}^m (a_i^{(l)}) * \text{tanh}(\delta_k^{(l+1)}, 2)$$

$$\nabla_{b_k^{(l)}} J(w, b; x, y) = \sum_{a, b} (\delta_k^{(l+1)})_{a, b}$$

↓  
Valid convolution

$a^{(l)}$  — the input image  
 $a^{(l)}$  — the input to the  $l$ -th layer  $\Rightarrow ??$

## ⑨ Exercise: Convolutional Neural Network

filterDim = 9

numFilters = 20 (网上说的是10)

poolDim = 2

cnnTrain.m

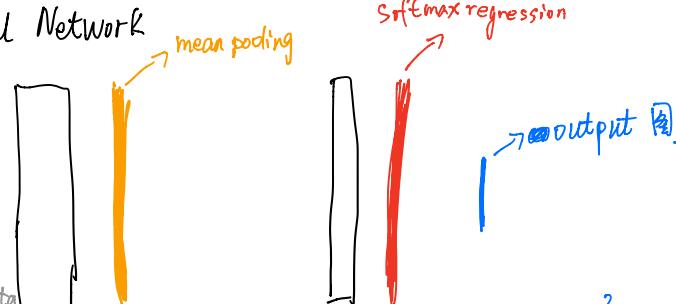
STEP0: Initialize Parameters and Load Data

STEP1: Implement convNet Objective convolutional layer

STEP2: Gradient Check  $W_c: (\text{filterDim}, \text{filterDim}, \text{numFilters})$  layer  $W_d: (\text{numClasses}, \text{hiddenSize})$

STEP3: Learn Parameters

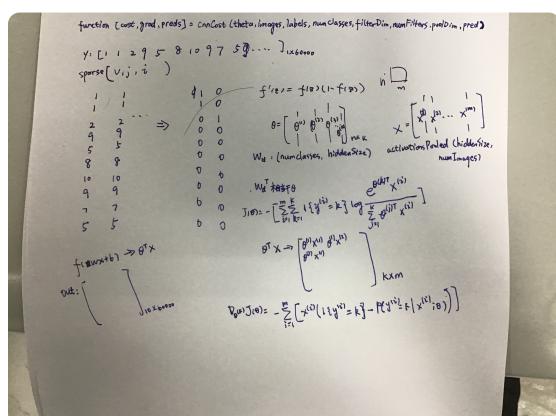
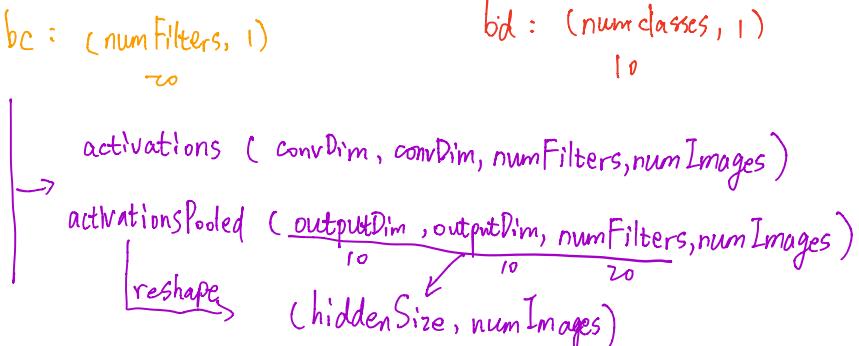
STEP4: Test



Step0a: Forward Propagation

convolution  $\rightarrow$  activation  $\rightarrow$  pooling

$\rightarrow$  activationPooled matrix  
 (every column  $\rightarrow$  an image)



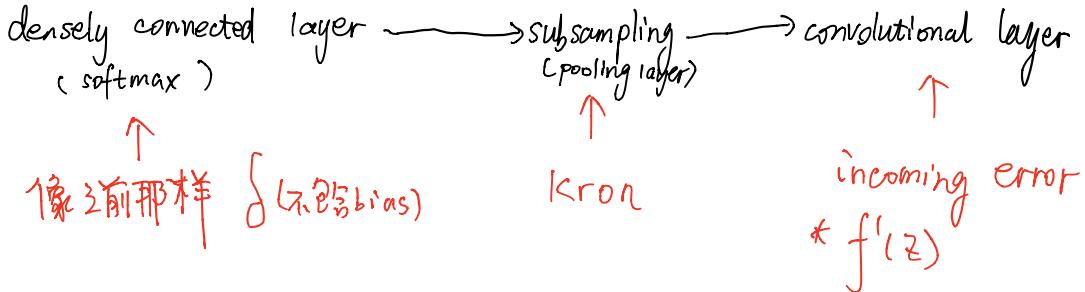
$W_d * \text{activationsPooled}$   
 在上偏置，经过激活函数

$\rightarrow \text{out} (\text{numClasses}, \text{numImages})$

纸上写的  
一点笔记

难点来了，卷积神经网络下的 back propagation (之前 multi-layer neural networks 的反向传播才慢慢懂。现在要来看 CNN 的 BP 了！坚持。)

跟 forward propagation 正好相反，这时的 error 要从



计算梯度的话 (Grad-Wc, Grad-bc, Grad-Wd, Grad-bd)

densely connected layer → subsampling → convolutional layer

↑  
还是和之前一样  
 $W-d, b-d$

因为卷积层有 error  
所以这里没有  
而且  
↑  
error 和原训练图像差  
 $W-c, b-c$   
↓

hadamard 积

梯度为 Ff 的 error 之和

如果当前层是全连接层:  $\delta^{i,L} = (W^{(L+1)})^T \delta^{i,L+1} \odot \sigma'(z^{i,L})$

如果当前层是卷积层:  $\delta^{i,L} = \delta^{i,L+1} * \text{rot 180 } (W^{(L+1)}) \odot \sigma'(z^{i,L})$

如果当前层是池化层:  $\delta^{i,L} = \text{upsample } (\delta^{i,L+1}) \odot \sigma'(z^{i,L})$

} for  $L=L-1 \rightarrow 2$

梯度下降更新权重

如果当前层是全连接层:  $W^L = W^L - \alpha \sum_{i=1}^m \delta^{i,L} (a^{i,L+1})^T, b^L = b^L - \alpha \sum_{i=1}^m \delta^{i,L}$

如果当前层是卷积层: 对每一个卷积核有:  $W^L = W^L - \alpha \sum_{i=1}^m \delta^{i,L} * a^{i,L+1}, b^L = b^L - \alpha \sum_{i=1}^m \sum_{u,v} (\delta^{i,L})_{u,v}$

} for  $L=2 \rightarrow L$

Step3: learn Parameters:

Batch method: L-BFGS  $\xrightarrow{\text{train}}$  a convolutional network  $\rightarrow$  一次迭代需花费几分钟  
或者更久

SGD + Momentum

iteration 和 epoch 的区别

velocity vector  $\rightarrow$   
parameter vector  $\rightarrow$

$$\text{SGD: } \theta = \theta - \alpha \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$$

$$\text{momentum: } v = \gamma v + \alpha \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$$

$$\theta = \theta - v$$

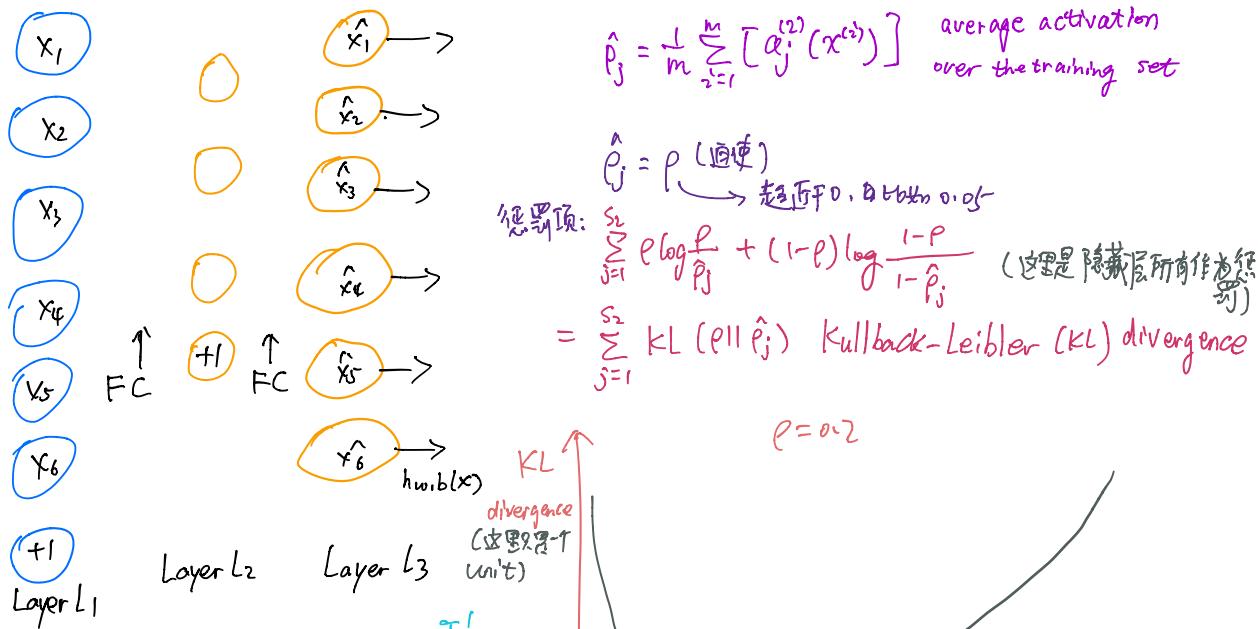
9745

| 257 513 ... 59649 | 59905

59745

①: Auto encoders:

An autoencoder neural network is an unsupervised learning algorithm that applies backpropagation, setting the target to be equal to the inputs. I.e., it uses  $y^{(1)} = x^{(1)}$



$$\delta_i^{(2)} = \left( \sum_{j=1}^{s_2} W_{ji}^{(2)} \delta_j^{(3)} \right) f'(z_i^{(2)}) \leftarrow \text{前}$$

$$\delta_i^{(2)} = \left( \left( \sum_{j=1}^{s_2} W_{ji}^{(2)} \delta_j^{(3)} \right) + \beta \left( -\frac{p}{\hat{p}_i} + \frac{1-p}{1-\hat{p}_i} \right) \right) f'(z_i^{(2)}) \leftarrow \text{现在}$$

$$\frac{\partial J'}{\partial z_i^{(2)}} = \beta \left( p \cdot \frac{\hat{p}_i}{p} \left( -\frac{p}{\hat{p}_i^2} \right) + (1-p) \cdot \frac{1-\hat{p}_i}{1-p} \cdot \frac{(1-p)}{(1-\hat{p}_i)^2} \right) f'(z_i^{(2)})$$

$$= \beta \left( -\frac{p}{\hat{p}_i} + \frac{1-p}{1-\hat{p}_i} \right)$$

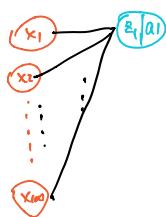
要想让 hidden unit  $i$  最大化地激活, 应该使输入  $x_j$  (for  $j=1 \dots 100$ ) 在

$$x_j = \frac{w_{ij}^{(1)}}{\sqrt{\sum_{i=1}^{100} (w_{ij}^{(1)})^2}} \quad (\text{试着自己证明一下})$$

visualization: one image per hidden unit.  $\xrightarrow{\text{learned}}$  hidden units learn to detect edges at different positions and orientations

$$a_i^{(2)} = f\left(\sum_{j=1}^{100} w_{ij}^{(1)} x_j + b_i^{(1)}\right)$$

$$\text{如果限制 } \|x\|^2 = \sum_{i=1}^{100} x_i^2 \leq 1$$



不考虑bias的情况下,  $a_i^{(2)} = f\left(\sum_{j=1}^{100} w_{ij}^{(1)} x_j\right)$  要让  $a_i^{(2)}$  尽可能趋近于1

$$(2) \sum_{j=1}^{100} w_{ij}^{(1)} x_j \text{ 要最大}$$

??

## ① PCA Whitening

Principal components analysis (PCA)  $\xrightarrow{\text{降维}} \text{无监督特征学习}$

covariance matrix + 协方差矩阵

eigen vector / eigenvector / 特征向量, 本征向量

$\Sigma$  - covariance matrix standard notation

$$\Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} (x^{(i)})^T \xrightarrow{\substack{\text{他张特征向量} \\ \text{组成的矩阵U}}} U = \begin{bmatrix} | & | & | \\ u_1 & u_2 & \dots & u_n \\ | & | & | \end{bmatrix} \Rightarrow \text{正交矩阵}$$

principal  
↑  
 $u_1, u_2, \dots, u_n$   
对称  
↓  
特征值

$u_1^T x$  —  $x$  在  $u_1$  方向上的投影

$u_2^T x$  —  $x$  在  $u_2$  方向上的投影

$$X : \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m)} \\ | & | & \dots & | \end{bmatrix}_{n \times m}$$

$$\begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_n^T \end{bmatrix} X$$

Rotating

$$X_{\text{rot}} = U^T X = \begin{bmatrix} u_1^T X \\ u_2^T X \\ \vdots \\ u_n^T X \end{bmatrix}_{n \times m}$$

(在这里  $n=2$  罢了)

$$X = U X_{\text{rot}}_{n \times m}$$

$$U^T \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}_{n \times m}$$

Reducing

$$\tilde{x} = \begin{bmatrix} x_{\text{rot},1} \\ \vdots \\ x_{\text{rot},k} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \approx \begin{bmatrix} x_{\text{rot},1} \\ \vdots \\ x_{\text{rot},k} \\ x_{\text{rot},k+1} \\ \vdots \\ x_{\text{rot},n} \end{bmatrix} = x_{\text{rot}}$$

$x_{\text{rot},1}^{(i)} = u_i^T x^{(i)}$  去掉没用的特征.

也把  $\tilde{x} = \begin{bmatrix} x_{\text{rot},1} \\ \vdots \\ x_{\text{rot},k} \end{bmatrix}_{K \times n} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{bmatrix}$

Recovering

$$\hat{x} = U \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \vdots \\ \hat{x}_k \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \sum_{i=1}^k u_i \hat{x}_i$$

(实际上我们是让  $\hat{x} \in \mathbb{R}^k$  乘以  $U$  的前  $K$  列)

Number  $K$  of components to retain

用这个描述 percentage of variance retained : 入为特征值,  $\lambda_j \rightarrow$  入降序排列

$$\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^n \lambda_j}$$

图像有一种性质叫stationarity

通常用PCA的时候, 都要让各特征对应的值有0均值, 相同的变化范围

对于natural images, 一般不需要做第②步: variance normalization

只需做第①步:

for all  $j$ :

$$\mu^{(i)} := \frac{1}{n} \sum_{j=1}^n x_j^{(i)} \quad x_j^{(i)} := x_j^{(i)} - \mu^{(i)}$$

跟PCA算法联系相对紧密的一个预处理步骤叫：whitening / spherering

$$X_{\text{PCA white}, i} = \frac{\underline{x}_{\text{rot}, i}}{\sqrt{\lambda_i}}, \text{使协方差矩阵为单位矩阵 I;}$$

作用  
使输入冗余降低：  
①各特征相关性降低  
②各特征有相同的方差

如果 R 为正交矩阵（即  $R^T R = R R^T = I$ ）

但  $R X_{\text{PCA white}}$  也有单位协方差矩阵

$$\underline{x}_{\text{ZCA white}} = U \underline{x}_{\text{PCA white}} \quad (\text{且 } R=U), \text{这样使 } \underline{x}_{\text{ZCA white}} \text{ 最接近 } \underline{x} \quad (\text{原始图像})$$

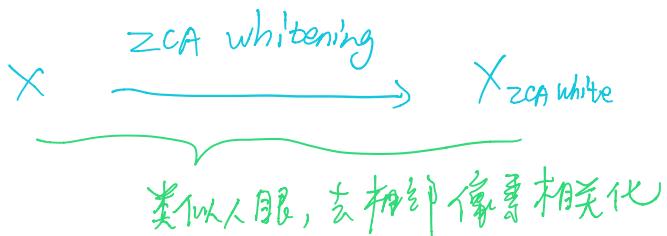
ZCA whitening 一般不单独使用

## Regularization

有些  $\lambda_i$  接近于零， $\frac{1}{\sqrt{\lambda_i}} \Rightarrow$  会很大，所以加入固定值  $\epsilon$

$$\text{使: } X_{\text{PCA white}, i} = \frac{\underline{x}_{\text{rot}, i}}{\sqrt{\lambda_i + \epsilon}} \quad \epsilon \in (-1, 1]$$

↓  
 $\epsilon$  大约是  $10^{-5}$  (也有平滑, 即通滤波的反梯)



## (12) sparse coding

目标：找到一系列的基底向量  $\phi_i$ ，使我们可以用其线性组合来表示一个输入向量

$$x = \sum_{i=1}^k a_i \phi_i \quad x \in \mathbb{R}^n$$

$\rightarrow (K > n)$

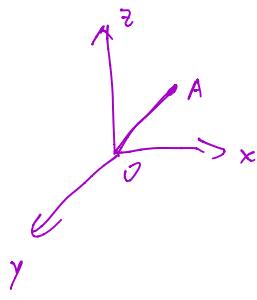
Learn an over-complete set of basis vectors to represent input vectors

$\downarrow$  sparsity  $\rightarrow$  要使系数尽可能多的为0  
 $\rightarrow$  防止对同一个  $x$  表示时， $a_i$  的不唯一

就好像在三维空间中  $\vec{OA} = (2, 3, 4)$ ， $i, j, k$  分别是  $x, y, z$  轴方向上的基底向量

$$\text{本来 } \vec{OA} = 2\vec{i} + 3\vec{j} + 4\vec{k} = \sum_{i=1}^3 a_i \phi_i$$

但是现在要用 4 个或更多向基底来表示  $\rightarrow$  over-complete  
 过完备的



定义 sparse coding cost function ( $m$  个输入向量)

$$\underset{a_i^{(j)}, \phi_i}{\text{minimize}} \underbrace{\sum_{j=1}^m \left\| x^{(j)} - \sum_{i=1}^k a_i^{(j)} \phi_i \right\|^2}_{\text{这块体现“无监督”}} + \lambda \sum_{i=1}^k S(a_i^{(j)})$$

$S(\cdot)$  当  $a_i$  远离0时就惩罚越厉害

惩罚项的几种选择：

①  $L_0$  norm  $S(a_i) = 1 (|a_i| > 0)$  不可导且难以优化

②  $L_1$  norm  $S(a_i) = |a_i| \quad \}$  常用

③ log  $S(a_i) = \log(1 + a_i^2)$

一个问题， $\lambda \sum_{i=1}^k S(a_i^{(j)})$  可以很小，当  $a_i$  减小时，但这会使  $\left\| x^{(j)} - \sum_{i=1}^k a_i^{(j)} \phi_i \right\|^2$  中的  $\phi_i$  增大，为了 solve this，加一个限制： $\|\phi_i\|^2 < \text{常数 } C$ 。所以完整的目标函数是：

$$\underset{a_i^{(j)}, \phi_i}{\text{minimize}} \sum_{j=1}^m \left\| x^{(j)} - \sum_{i=1}^k a_i^{(j)} \phi_i \right\|^2 + \lambda \sum_{i=1}^k S(a_i^{(j)})$$

$$\text{subject to} \quad \|\phi_i\|^2 \leq C, \forall i = 1, \dots, K$$

## probabilistic interpretation

sparse coding → a generative model (概率上的可解释性)

$L_1$  sparsity penalty  $\xrightarrow{\text{ai convex}} \text{共轭梯度法 (conjugate gradient method)}$

$L_2$  norm constraint  $\xrightarrow{\text{phi convex}} \text{拉格朗日对偶法 (Lagrange dual)}$

→ 来一个新数据得重新优化来得到 "coefficients". — computationally expensive

## ICA (independent component analysis) 独立成分分析

不是线性无关的基底，而是“正交基底” (orthonormal basis)

$$\hookrightarrow \text{basis}(\phi_1, \dots, \phi_n) \left\{ \begin{array}{l} \phi_i \cdot \phi_j = 0 \text{ if } i \neq j \\ \phi_i \cdot \phi_j = 1 \text{ if } i = j \end{array} \right.$$

W (列向量之间是一系列基底，正交且稀疏)

sparse coding : mapping features to raw data  $x = \sum_{j=1}^k \alpha_j \phi_j$

ICA : mapping raw data  $x$  to features

orthonormal ICA objective :

$$\text{minimize } \|Wx\|,$$

$$\text{s.t. } WW^\top = I$$

① 学到的基底数 < 输入维度,

即 under-complete basis

优化目标函数，使用梯度下降法:

② 数据必须是 ZCA whitened  
with no regularization

repeat until done:

1.  $W \leftarrow W - \alpha \nabla_W \|Wx\|$ , (只会用 line-search 算法来改变着进行梯度下降)

2.  $W \leftarrow \text{proj}_U W$   $U$  是满足  $UU^\top = I$  的矩阵空间 ( $W \in (WW^\top)^{-\frac{1}{2}} U$ )

ZCA :  $U \cdot \text{PCA}$

$$U \cdot (\quad) U^\top \cdot x$$

⇒ 这里过程中好像还有待补充说明：为什么像  
ZCA whitening

## RICA (Reconstruction ICA)

用 soft reconstruction penalty 来代替 ICA 的正交约束, 为了解决 ICA 的缺陷

$$\min_w \lambda \|w\|_1 + \frac{1}{2} \|W^T W x - x\|_2^2$$

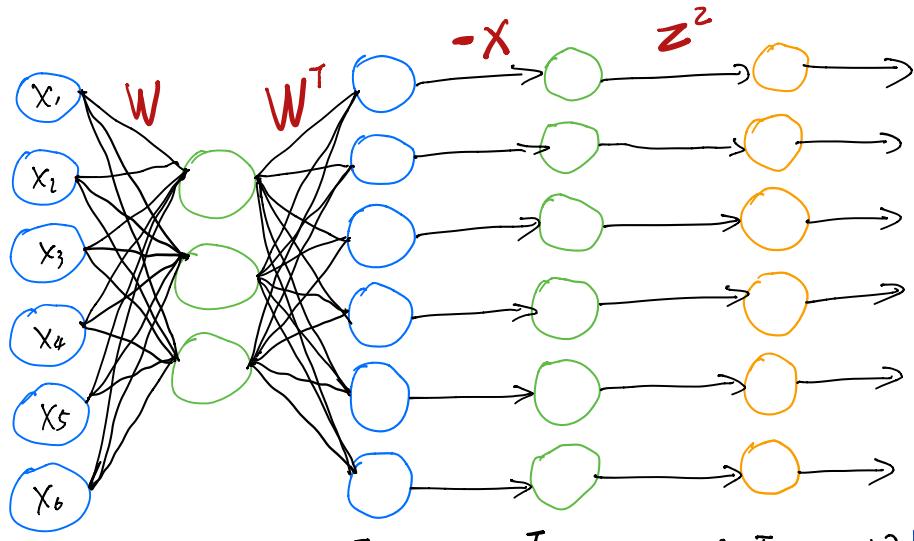
sparse autoencoders :

$$\min_w \lambda \|\sigma(wx)\|_1 + \frac{1}{2} \|\sigma(w^T \sigma(wx)) - x\|_2^2$$

退化控制

相差个 sigmoid  
非线性

RICA 和 over-complete basis  
和 non-white data  
要更鲁棒.



	$X$	$WX$	$W^T WX$	$W^T W x - x$	$(W^T W x - x)^2$
层	1	2	3	4	
权重	$w$	$w^T$	I	N/A	

激活函数

$f(z_i) = z_i$	$f(z_i) = z_i$	$f(z_i) = z_i - x_i$	$f(z_i) = z_i^2$
$f'(z_i) = 1$	$f'(z_i) = 1$	$f'(z_i) = 1$	$f'(z_i) = 2z_i$

在神经网络中多次出现  $W$ ,  
对其梯度就是对每一个  $w$  的简单加和.

$$\Delta = (W^T \delta^{(2)}) \cdot 1 \quad ((W^T)^T \delta^{(1)}) \cdot 1 \quad (I^T \delta^{(1)}) \cdot 1 \quad f'(z_i) = 2z_i$$

for  $\delta^{(1)}$

$$\text{Input } z \text{ to this layer} \quad x \quad Wx \quad W^T Wx \quad W^T Wx - x$$

$$\nabla_{w^T} F = \delta^{(3)}(a^{(2)})^T = 2(w^T w x - x) \cdot (w x)^T$$

$$\nabla_w F = \delta^{(2)}(a^{(1)})^T = (w)(2(w^T w x - x))(x)^T$$

$$\min_w \lambda \|w x\|_1 + \frac{1}{2} \|w^T w x - x\|_2^2$$

$\left. \begin{array}{l} \nabla_w F = \nabla_w F + (\nabla_{w^T} F)^T \\ = (w)(2(w^T w x - x))(x)^T + 2(w x)(w^T w x - x)^T \end{array} \right\}$   
↑ 3 个维度: 50 × 81  
 $w_{50 \times 81}$   
 $x_{81 \times 10000}$

$x$	$w x$	$(w x)^2$	$(w x)^2 + \varepsilon$	$\sqrt{(w x)^2 + \varepsilon}$
<small>输入</small>	1	2	3	4
<small>权重</small>	$w$	1	1	1
<small>激活函数</small>	$f(z_i) = z_i$	$f(z_i) = z_i^2$	$f(z_i) = z_i + \varepsilon_i$	$f(z_i) = \sqrt{z_i}$
<small>导数</small>	$f'(z_i) = 1$	$f'(z_i) = 2z_i$	$f'(z_i) = 1$	$f'(z_i) = \frac{1}{2\sqrt{z_i}}$

Delta  $(w^T \delta^{(2)}) \cdot 1 \quad (I^T \delta^{(3)}) \cdot 2 w x \quad (I^T \delta^{(4)}) \cdot 1 \quad f'(z_i) = \frac{1}{2\sqrt{z_i}}$

Input  $z$   
 to this layer  $x$   $w x$   $(w x)^2$   $(w x)^2 + \varepsilon$

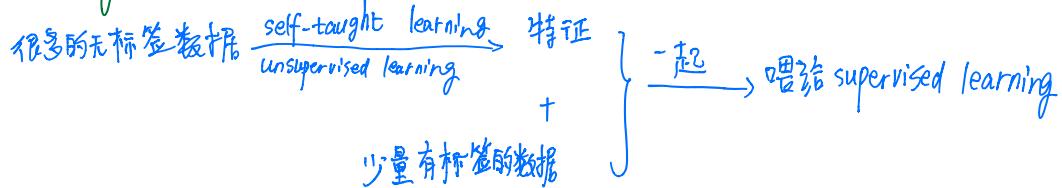
$$\nabla_w \|w x\|_1 = (\delta^{(2)}(a^{(1)})^T) = \frac{2 w x}{2\sqrt{(w x)^2 + \varepsilon}} x^T \Rightarrow \frac{w x x^T}{\sqrt{(w x)^2 + \varepsilon}}$$

$\Rightarrow$  在编程时还要注意,  
 一下下, 看源码写好了.

$\uparrow$   
 $\delta^{(2)}$        $(a^{(1)})^T$

### ⑬ self-taught learning

Something's it's not who has the best algorithm that wins; it's who has the most data

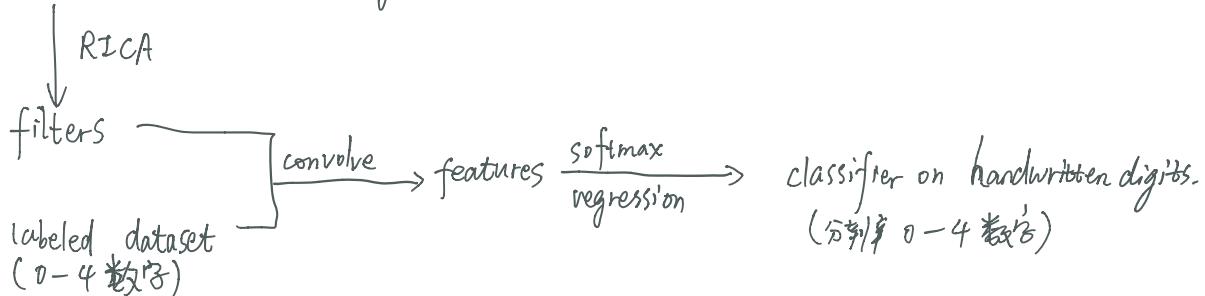


self-taught learning setting: 无标签的数据 和 有标签的数据不一定有相同的分布

semi-supervised learning: 无标签数据 依然要和有标签数据 同分布

exercise

patches (from unlabeled training set of handwritten digits) (0-9 数字)





questions:  $\leftarrow$  这是看CS229 note时产生的问题

- ① page. 13.  $\Sigma^2$
- ② page. 15  $\tau$  bandwidth parameter
- ③ page 17 why choose sigmoid function?
- ④ page 19 linear regression update rule is identical with logistic regression update rule? why?

$$p(x|\eta) = h(x)g(\eta) e^{\eta^T u(x)} \Rightarrow \text{指数形式} \quad p(y|\eta) = b(y) e^{\eta^T T(y) - a(y)}$$

$$p(x|\mu) = \mu^x (1-\mu)^{1-x}$$

$$\begin{aligned} & e^{\ln \mu^x (1-\mu)^{1-x}} \\ &= e^{\ln \mu^x (1-\mu)^{-x} \cdot (1-\mu)} \\ &= (1-\mu)^{-x} e^{x \ln \left( \frac{\mu}{1-\mu} \right)} \end{aligned}$$

$$\text{设 } \eta = \ln \frac{\mu}{1-\mu}, \text{ 则 } \mu = \sigma(\eta)$$

$$p(x|\mu) = [1 - \sigma(\eta)] e^{\eta \cdot x}$$

