





覃雄派



# 提纲

HERS/77 OR CHINA

- 矩阵求导入门
- 矩阵求导实例



• 矩阵求导参考资料

#### The Matrix Cookbook

[ http://matrixcookbook.com ]

Kaare Brandt Petersen Michael Syskind Pedersen

Version: November 15, 2012

这里先进行简单入门,然后直接使用该cook book的一些结论

https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf







#### • 矩阵求导入门

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$oldsymbol{ heta} = egin{bmatrix} heta_0 \ heta_1 \ heta_2 \ dots \ heta_d \end{bmatrix}$$

按照"列"组织向量

$$\mathbf{x}^T \boldsymbol{\theta} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_d \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{bmatrix} = \theta_0 + x_1 \theta_1 + x_2 \theta_2 + \dots + x_d \theta_d$$



#### • 矩阵求导入门

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$oldsymbol{ heta} = egin{bmatrix} heta_0 \ heta_1 \ heta_2 \ dots \ heta_d \end{bmatrix}$$

$$\mathbf{x}^T \boldsymbol{\theta} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_d \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{bmatrix} = \theta_0 + x_1 \theta_1 + x_2 \theta_2 + \dots + x_d \theta_d$$

$$1.\theta_0+x_1\theta_1+x_2\theta_2+...+x_d\theta_d$$
 对 $\theta_0$ 、 $\theta_1$ 、 $\theta_2$ 、...、 $\theta_d$  求导

2.得到

$$1 \quad x_1 \quad x_2 \quad \dots \quad x_d$$

 $x_2$  3.按照列向量来组织为  $x_2$  :

• 矩阵求导入门

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$oldsymbol{ heta} = egin{bmatrix} heta_0 \ heta_1 \ heta_2 \ dots \ heta_d \end{bmatrix}$$

$$\mathbf{x}^T \boldsymbol{\theta} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_d \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{bmatrix} = \theta_0 + x_1 \theta_1 + x_2 \theta_2 + \dots + x_d \theta_d$$

$$1.\theta_0+x_1\theta_1+x_2\theta_2+...+x_d\theta_d$$
 对 $\theta_0$ 、 $\theta_1$ 、 $\theta_2$ 、...、 $\theta_d$  求导

2.得到

$$1 \quad x_1 \quad x_2 \quad \dots \quad x_d$$

3.按照列向量来组织为



$$\frac{\partial \mathbf{x}^T \boldsymbol{\theta}}{\partial \boldsymbol{\theta}} = \frac{\partial \boldsymbol{\theta}^T \mathbf{x}}{\partial \boldsymbol{\theta}} = \mathbf{x}$$

第一个重要公式

#### • 矩阵求导入门

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$oldsymbol{ heta} = egin{bmatrix} heta_0 \ heta_1 \ heta_2 \ dots \ heta_d \end{bmatrix}$$

$$\mathbf{x}^T \boldsymbol{\theta} = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_d \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{bmatrix} = \theta_0 + x_1 \theta_1 + x_2 \theta_2 + \dots + x_d \theta_d$$

$$\underline{\partial \boldsymbol{a}^T \mathbf{x}}$$

$$\theta_0 + x_1 \theta_1 + x_2 \theta_2 + ... + x_d \theta_d$$
  
对 $\theta_0 \ \theta_1 \ \theta_2 \ ... \ \theta_d$   
求导得到  
 $1 \ x_1 \ x_2 \ ... \ x_d$ 

按照列向量来组织为  $x_1$   $x_2$   $\vdots$   $x_{d}$ 

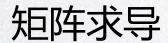
$$oldsymbol{x}$$
 列向量

$$\boldsymbol{a}$$
 列向量

$$rac{\partial oldsymbol{a}^T \mathbf{x}}{\partial \mathbf{x}} = rac{\partial \mathbf{x}^T oldsymbol{a}}{\partial \mathbf{x}} = oldsymbol{a}$$









• 矩阵求导入门

 $oldsymbol{x}$ 列向量

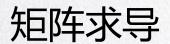
 $\boldsymbol{a}$  列向量

$$rac{\partial oldsymbol{a}^T \mathbf{x}}{\partial \mathbf{x}} = rac{\partial \mathbf{x}^T oldsymbol{a}}{\partial \mathbf{x}} = oldsymbol{a}$$



 $\mathbf{b}^T \mathbf{A}$  看作一个整体

$$rac{\partial \mathbf{b}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}^T \mathbf{b}$$





• 矩阵求导入门

 $oldsymbol{x}$  列向量

a 列向量

$$rac{\partial oldsymbol{a}^T \mathbf{x}}{\partial \mathbf{x}} = rac{\partial \mathbf{x}^T oldsymbol{a}}{\partial \mathbf{x}} = oldsymbol{a}$$

有两个x,分别 求导,累加



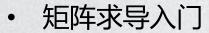


$$\mathbf{x}^T \mathbf{A} \mathbf{x}$$

$$rac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = oldsymbol{A} oldsymbol{x} + oldsymbol{A}^T oldsymbol{x}$$



AND CHINA



- $oldsymbol{x}$  列向量
- a 列向量

 $rac{\partial oldsymbol{a}^T \mathbf{x}}{\partial \mathbf{x}} = rac{\partial \mathbf{x}^T oldsymbol{a}}{\partial \mathbf{x}} = oldsymbol{a}$ 

有两个x,分别 求导,累加







 $\partial \mathbf{x}^T \mathbf{A} \mathbf{x}$ 

 $\partial \mathbf{x}$ 

 $= oldsymbol{A} oldsymbol{x} + oldsymbol{A}^T oldsymbol{x}$ 

 $2\boldsymbol{A}\boldsymbol{x}$ 

如果A是对称矩阵

 $\pi$ 





#### 在SVM的梯度下降法求解 过程中用到

#### Hinge loss → SVM

$$- \max(0,1-y_iw^Tx_i) + \frac{\lambda}{2}||w||^2$$

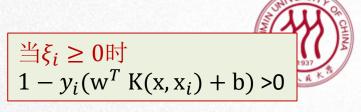
- · 对w求导
- 注意yi是标量, xi是一个列向量, w是一个列向量 (权重)
- $||w||^2$ 为w和自己点乘,即 $w^T w$
- 分两个部分

$$-\frac{\partial}{\partial w}\left[\max(0,1-y_{i}w^{T}x_{i})\right] = \begin{cases} -y_{i}x_{i}, & \text{if } 1-y_{i}w^{T}x_{i} \geq 0\\ 0, & \text{if } 1-y_{i}w^{T}x_{i} < 0 \end{cases} \qquad \underbrace{\frac{\partial \mathbf{a}^{T}\mathbf{x}}{\partial \mathbf{x}}} = \underbrace{\frac{\partial \mathbf{x}^{T}\mathbf{a}}{\partial \mathbf{x}}} = \mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{a}^{T}\mathbf{$$

分别对 $w^T$  w的两个w 求导,再加起来







- Hinge loss → kernelized SVM
  - $\frac{1}{2} \mathbf{w}^{T} K(x, x^{T}) \mathbf{w} + C \sum_{i=1}^{n} \max(0, 1 y_{i}(\mathbf{w}^{T} K(x, x_{i}) + \mathbf{b}))$ 
    - 对w求导, 对b求导

$$- \frac{\partial L}{\partial w} = K(x, x^T) w - C \sum_{i=1, \xi_i \ge 0}^n y_i K(x, x_i)$$



$$\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x}$$
$$K(\mathbf{x}, \mathbf{x}^T)$$
为一个对称矩阵

$$- \frac{\partial \mathbf{L}}{\partial \mathbf{b}} = -C \sum_{i=1, \xi_i \ge 0}^{n} y_i$$

• 注意b是标量, y<sub>i</sub>是标量

$$egin{aligned} rac{\partial \mathbf{x}^T oldsymbol{a}}{\partial \mathbf{x}} = oldsymbol{a} \end{aligned}$$

把 $K(x,x_i)$ 看作一个列向量







#### • 一元/多元线性回归的解析解

$$-\frac{1}{2n}(\mathbb{Y} - \mathbb{X}\theta)^T(\mathbb{Y} - \mathbb{X}\theta) = \frac{1}{2n}(\mathbb{Y}^T - \theta^T \mathbb{X}^T)(\mathbb{Y} - \mathbb{X}\theta)$$

对θ求导

$$- \frac{1}{2n} (\mathbb{Y}^T \mathbb{Y} - \mathbb{Y}^T \mathbb{X} \boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbb{X}^T \mathbb{Y} + \boldsymbol{\theta}^T \mathbb{X}^T \mathbb{X} \boldsymbol{\theta})$$

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}$$

$$-\frac{1}{2n}(0-X^TY-X^TY+X^TX\boldsymbol{\theta}+X^TX\boldsymbol{\theta})$$

分别对两个求导

$$\theta^{\mathrm{T}} \mathbb{X}^T \mathbb{X} \theta$$

$$\theta^{\mathrm{T}} \mathbb{X}^T \mathbb{X} \theta$$

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$$







 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>.<mark>></mark>A<sup>[1]</sup>

 $W^{[2]}, b^{[2]}$ 

输出层 1个神经元 Z<sup>[2]</sup>→A<sup>[2]</sup>



#### • 神经网络

- 
$$(1)[Z^{[1]} = W^{[1]}X + b^{[1]} = (64*10)(10*1) + (64*1) = (64*1)]$$

有1个样本,每个样本为10×1的列向量

- (2)  $A^{[1]} = \sigma(Z^{[1]}) = (64*1)$
- (3)  $Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$
- (4)  $\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$
- 处理过程 A<sup>[0]</sup> → Z<sup>[1]</sup>→A<sup>[1]</sup>→Z<sup>[2]</sup>→A<sup>[2]</sup>
- X即A<sup>[0]</sup>

#### 从输入层计算隐藏层

X 10元素 A<sup>[0]</sup>

 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 W<sup>[2]</sup>、b<sup>[2]</sup> Z<sup>[1]</sup> A<sup>[1]</sup> 输出层 1个神经元 Z<sup>[2]</sup>→A<sup>[2]</sup>



#### • 神经网络

- (1) Z<sup>[1]</sup> = W<sup>[1]</sup>X + b<sup>[1]</sup> = (64\*10)(10\*1) + (64\*1) = (64\*1)
- (2)  $A^{[1]} = \sigma(Z^{[1]}) = (64*1)$
- (3)  $Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$
- (4)  $\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$
- 处理过程 A<sup>[0]</sup> → Z<sup>[1]</sup>→A<sup>[1]</sup>→Z<sup>[2]</sup>→A<sup>[2]</sup>
- X即A<sup>[0]</sup>

有1个样本,每个样本为10×1的列向量

隐藏层的非线性传导



 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>→A<sup>[1]</sup> 输出层 W<sup>[2]</sup>、b<sup>[2]</sup> 1个神经元 Z<sup>[2]</sup>-▶A<sup>[2]</sup>



#### • 神经网络

- (1) 
$$Z^{[1]} = W^{[1]}X + b^{[1]} = (64*10)(10*1) + (64*1) = (64*1)$$

- (2)  $A^{[1]} = \sigma(Z^{[1]}) = (64*1)$ 

- (3) 
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$$

- (4) 
$$\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$$

- 处理过程 A<sup>[0]</sup> → Z<sup>[1]</sup>→A<sup>[1]</sup>→Z<sup>[2]</sup>→A<sup>[2]</sup>
- X即A<sup>[0]</sup>

有1个样本,每个样本为10×1的列向量

#### 从隐藏层计算输出层



 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>→A<sup>[1]</sup>

 $W^{[2]}, b^{[2]}$ 

输出层 1个神经元 Z<sup>[2]</sup>-▶A<sup>[2]</sup>



#### • 神经网络

- (1) 
$$Z^{[1]} = W^{[1]}X + b^{[1]} = (64*10)(10*1) + (64*1) = (64*1)$$

有1个样本,每个样本为10×1的列向量

- (2) 
$$A^{[1]} = \sigma(Z^{[1]}) = (64*1)$$

$$- (3) 7^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$$

- (4) 
$$\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$$

- 处理过程 A<sup>[0]</sup> → Z<sup>[1]</sup>→A<sup>[1]</sup>→Z<sup>[2]</sup>→A<sup>[2]</sup>
- X即A<sup>[0]</sup>

输出层的非线性传导

X 10元素 A<sup>[0]</sup>

 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>→A<sup>[1]</sup>

 $W^{[2]}, b^{[2]}$ 

输出层 1个神经元 Z<sup>[2]</sup>→A<sup>[2]</sup> 损失函数J 二元交叉 熵

#### • 神经网络

- 处理过程 A<sup>[0]</sup> → Z<sup>[1]</sup>→A<sup>[1]</sup>→Z<sup>[2]</sup>→A<sup>[2]</sup>
- X即A<sup>[0]</sup>
- 最后用A<sup>[2]</sup>构造损失函数,,注意A<sup>[2]</sup>即预测值ŷ
  - 二值分类器(0/1)的交叉熵损失函数的形式为

$$J = -\frac{1}{n}((Y\log(A^{[2]}) + (1 - Y)\log(1 - A^{[2]}))$$

- Y=(1\*, 1)
- n=1

根据损失函数计算损失值

1个样本,每个 样本为10×1的 列向量

X 10元素 A<sup>[0]</sup>

 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>→A<sup>[1]</sup>

 $W^{[2]}, b^{[2]}$ 

输出层 1个神经元 Z<sup>[2]</sup>→A<sup>[2]</sup> 损失函数J 二元交叉 熵

#### • 神经网络

- (1) 
$$Z^{[1]} = W^{[1]}X + b^{[1]} = (64*10)(10*1) + (64*1) = (64*1)$$

- (2) 
$$A^{[1]} = \sigma(Z^{[1]}) = (64*1)$$

- (3) 
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$$

- (4) 
$$\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$$

- 处理过程A<sup>[0]</sup> → Z<sup>[1]</sup>→A<sup>[1]</sup>→Z<sup>[2]</sup>→A<sup>[2]</sup>
- 最后用 $A^{[2]}$ 构造损失函数,,注意 $A^{[2]}$ 即预测值 $\hat{y}$ 
  - 二值分类器(0/1)的交叉熵损失函数的形式为

• 
$$J = -\frac{1}{n}((Ylog(A^{[2]}) + (1 - Y)log(1 - A^{[2]}))$$

- 计算<mark>损失函数</mark>对W<sup>[1]</sup>、b<sup>[1]</sup>、W<sup>[2]</sup>、b<sup>[2]</sup>的导数

$$dZ^{[2]} = dA^{[2]}g'(Z^{[2]}) (1*1) (1*1)$$

$$dW^{[2]} = \frac{dJ}{dW^{[2]}} = \frac{dJ}{dA^{[2]}} \frac{dA^{[2]}}{dZ^{[2]}} \frac{dZ^{[2]}}{dW^{[2]}} = \frac{dJ}{dZ^{[2]}} \frac{dZ^{[2]}}{dW^{[2]}} = dZ^{[2]} (A^{[1]})^{T}$$

$$(1*1)(1*64)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T \\
\frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^T \\
\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} = \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^T$$

$$\frac{d\mathbf{Z}^{[2]}}{d\mathbf{W}^{[2]}} = \frac{d(\mathbf{W}^{[2]}\mathbf{A}^{[1]} + \mathbf{b}^{[2]})}{d\mathbf{W}^{[2]}} = \frac{d(\mathbf{E}\mathbf{W}^{[2]}\mathbf{A}^{[1]} + \mathbf{b}^{[2]})}{d\mathbf{W}^{[2]}} = E^T (A^{[1]})^T$$

E为单位矩阵

X 10元素 A<sup>[0]</sup>

 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>→A<sup>[1]</sup>

 $W^{[2]}, b^{[2]}$ 

输出层 1个神经元 Z<sup>[2]</sup>→A<sup>[2]</sup> 损失函数J 二元交叉 熵

#### • 神经网络

- (1) 
$$Z^{[1]} = W^{[1]}X + b^{[1]} = (64*10)(10*1) + (64*1) = (64*1)$$

- (2) 
$$A^{[1]} = \sigma(Z^{[1]}) = (64*1)$$

- (3) 
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$$

- (4) 
$$\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$$

- (5)最后用 $A^{[2]}$ 构造损失函数,,注意 $A^{[2]}$ 即预测值 $\hat{y}$ 
  - 二值分类器(0/1)的交叉熵损失函数的形式为

• 
$$J = -\frac{1}{n}((Ylog(A^{[2]}) + (1 - Y)log(1 - A^{[2]}))$$

- 计算损失函数对W<sup>[1]</sup>、b<sup>[1]</sup>、W<sup>[2]</sup>、b<sup>[2]</sup>的导数

$$\frac{dJ}{db^{[2]}} = \frac{dJ}{dA^{[2]}} \frac{dA^{[2]}}{dZ^{[2]}} \frac{dZ^{[2]}}{db^{[2]}} 
= [A^{[2]} - Y][1] = [A^{[2]} - Y] 
= dZ^{[2]} 
(1*1)$$

$$\frac{dJ}{dA^{[1]}} = \frac{dJ}{dA^{[2]}} \frac{dA^{[2]}}{dZ^{[2]}} \frac{dZ^{[2]}}{dA^{[1]}} = \frac{dJ}{dZ^{[2]}} \frac{dZ^{[2]}}{dA^{[1]}} 
= dZ^{[2]} (W^{[2]})^T \mathbb{R} (W^{[2]})^T dZ^{[2]} 
(64*1)(1*1)$$

$$\frac{d\mathbf{Z}^{[2]}}{dA^{[1]}} = \frac{d(\mathbf{W}^{[2]}\mathbf{A}^{[1]} + \mathbf{b}^{[2]})}{dA^{[1]}} = (\mathbf{W}^{[2]})^T$$

X 10元素 A<sup>[0]</sup>

 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>→A<sup>[1]</sup>

 $W^{[2]}, b^{[2]}$ 

输出层 1个神经元 Z<sup>[2]</sup>→A<sup>[2]</sup> 损失函数J 二元交叉 熵

#### • 神经网络

- (1) 
$$Z^{[1]} = W^{[1]}X + b^{[1]} = (64*10)(10*1) + (64*1) = (64*1)$$

- (2) 
$$A^{[1]} = \sigma(Z^{[1]}) = (64*1)$$

- (3) 
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$$

- (4) 
$$\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$$

- (5)最后用 $A^{[2]}$ 构造损失函数,,注意 $A^{[2]}$ 即预测值 $\hat{y}$ 
  - 二值分类器(0/1)的交叉熵损失函数的形式为

• 
$$J = -\frac{1}{n}((Ylog(A^{[2]}) + (1 - Y)log(1 - A^{[2]}))$$

- 计算损失函数对W<sup>[1]</sup>、b<sup>[1]</sup>、W<sup>[2]</sup>、b<sup>[2]</sup>的导数

$$dZ^{[1]} = dA^{[1]}g'(Z^{[1]})$$

$$(64*1) (64*1)$$

$$dW^{[1]} = \frac{dJ}{dW^{[1]}} = \frac{dJ}{dA^{[2]}} \frac{dA^{[2]}}{dZ^{[2]}} \frac{dZ^{[2]}}{dA^{[1]}} \frac{dZ^{[1]}}{dW^{[1]}}$$

$$= \frac{dJ}{dZ^{[1]}} \frac{dZ^{[1]}}{dW^{[1]}}$$

$$= dZ^{[1]} (A^{[0]})^{T}$$

$$(64*1)(1*10)$$

$$\begin{array}{lll} \frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} & = & \mathbf{a} \mathbf{b}^T \\ \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} & = & \mathbf{b} \mathbf{a}^T \\ \frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} & = & \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{a}}{\partial \mathbf{X}} & = & \mathbf{a} \mathbf{a}^T \end{array}$$

$$\frac{d\mathbf{Z}^{[1]}}{d\mathbf{W}^{[1]}} = \frac{d(\mathbf{W}^{[1]}\mathbf{A}^{[0]} + \mathbf{b}^{[1]})}{d\mathbf{W}^{[1]}} = \frac{d(\mathbf{E}\mathbf{W}^{[1]}\mathbf{A}^{[0]} + \mathbf{b}^{[1]})}{d\mathbf{W}^{[1]}} = E^T (A^{[0]})^T$$

E为单位矩阵

X 10元素 A<sup>[0]</sup>

 $W^{[1]}, b^{[1]}$ 

隐藏层 64神经元 Z<sup>[1]</sup>→A<sup>[1]</sup>

 $W^{[2]}, b^{[2]}$ 

输出层 1个神经元 Z<sup>[2]</sup>→A<sup>[2]</sup> 损失函数J 二元交叉 熵

#### • 神经网络

- (1)  $Z^{[1]} = W^{[1]}X + b^{[1]} = (64*10)(10*1) + (64*1) = (64*1)$
- (2)  $A^{[1]} = \sigma(Z^{[1]}) = (64*1)$
- (3)  $Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} = (1*64)(64*1) + (1*1) = (1*1)$
- (4)  $\hat{y} = A^{[2]} = \sigma(Z^{[2]}) = (1*1)$
- 处理过程A<sup>[0]</sup> → Z<sup>[1]</sup>→A<sup>[1]</sup>→Z<sup>[2]</sup>→A<sup>[2]</sup>
- (5)最后用 $A^{[2]}$ 构造损失函数,,注意 $A^{[2]}$ 即预测值 $\hat{y}$ 
  - 二值分类器(0/1)的交叉熵损失函数的形式为
  - $J = -\frac{1}{n}((Ylog(A^{[2]}) + (1 Y)log(1 A^{[2]}))$
- 计算损失函数对W<sup>[1]</sup>、b<sup>[1]</sup>、W<sup>[2]</sup>、b<sup>[2]</sup>的导数







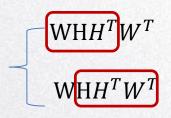
- NMF梯度 (矩阵求导)
- 假设我们把V分解为W和H
- $L = \min \frac{1}{2} ||V WH||^2$
- $\frac{\partial L}{\partial w_{ik}} = \left[ -\left[ (V WH)H^T \right]_{ik} \right]$
- $\frac{\partial L}{\partial \mathbf{h}_{k,i}} = -[W^T(V WH)]_{k,j}$

#### Frobenius norm

$$\frac{\partial}{\partial \mathbf{X}}||\mathbf{X}||_{\mathrm{F}}^{2} = \frac{\partial}{\partial \mathbf{X}}\mathrm{Tr}(\mathbf{X}\mathbf{X}^{H}) = 2\mathbf{X}$$

$$||\mathbf{A}||_{\mathrm{F}} = \sqrt{\sum_{ij} |A_{ij}|^2} = \sqrt{\mathrm{Tr}(\mathbf{A}\mathbf{A}^H)}$$
 (Frobenius)

$$(V - WH)(V - WH)^{T} =$$
  
 $(V - WH)(V^{T} - H^{T}W^{T}) =$   
 $VV^{T} - WHV^{T} - VH^{T}W^{T} + WHH^{T}W^{T}$   
对着W求导,得到  
 $0-VH^{T}-VH^{T}+WHH^{T}+WHH^{T}$   
=-2(V-WH)  $H^{T}$ 





- NMF梯度 (矩阵求导)
- 假设我们把V分解为W和H
- $L = \min \frac{1}{2} ||V WH||^2$
- $\frac{\partial L}{\partial w_{ik}} = -[(V WH)H^T]_{ik}$
- $\frac{\partial L}{\partial \mathbf{h}_{kj}} = \begin{bmatrix} -[W^T(V WH)]_{kj} \\ (V WH)^T(V WH) = \end{bmatrix}$

#### Frobenius norm

$$\frac{\partial}{\partial \mathbf{X}}||\mathbf{X}||_{\mathrm{F}}^{2} = \frac{\partial}{\partial \mathbf{X}}\mathrm{Tr}(\mathbf{X}\mathbf{X}^{H}) = 2\mathbf{X}$$

$$||\mathbf{A}||_{\mathrm{F}} = \sqrt{\sum_{ij} |A_{ij}|^2} = \sqrt{\mathrm{Tr}(\mathbf{A}\mathbf{A}^H)}$$
 (Frobenius)

$$(V - WH)^{T}(V - WH) =$$
 $(V^{T} - H^{T}W^{T})(V - WH) =$ 
 $V^{T}V - H^{T}W^{T}V - V^{T}WH + H^{T}W^{T}WH$ 
对着H求导,得到
 $0 - W^{T}V - W^{T}V + W^{T}WH + W^{T}WH$ 
=-  $2W^{T}(V - WH)$ 

