# 神经网络和深度学习--神经网络基础

## 2-1 二分分类



对于二分类问题，大牛给出了一个小的Notation。

- 样本：(x,y)，训练样本包含m个；
- 其中x∈Rnx，表示样本x 包含nx个特征；
- y∈0,1，目标值属于0、1分类；
- 训练数据：{(x(1),y(1)),(x(2),y(2)),…,(x(m),y(m))}

# Notation

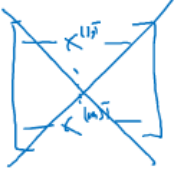$(x, y)$      $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$

$m$ training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})\}$

$M = M_{train}$      $M_{test}$ = #test examples.

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix} \Big\updownarrow n_x$$

$$\xcancel{\begin{matrix} x^{(1)} \\ x^{(m)} \end{matrix}}$$

$$Y = [y^{(1)} \ y^{(2)} \ \cdots , y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.shape = (1, m)$$

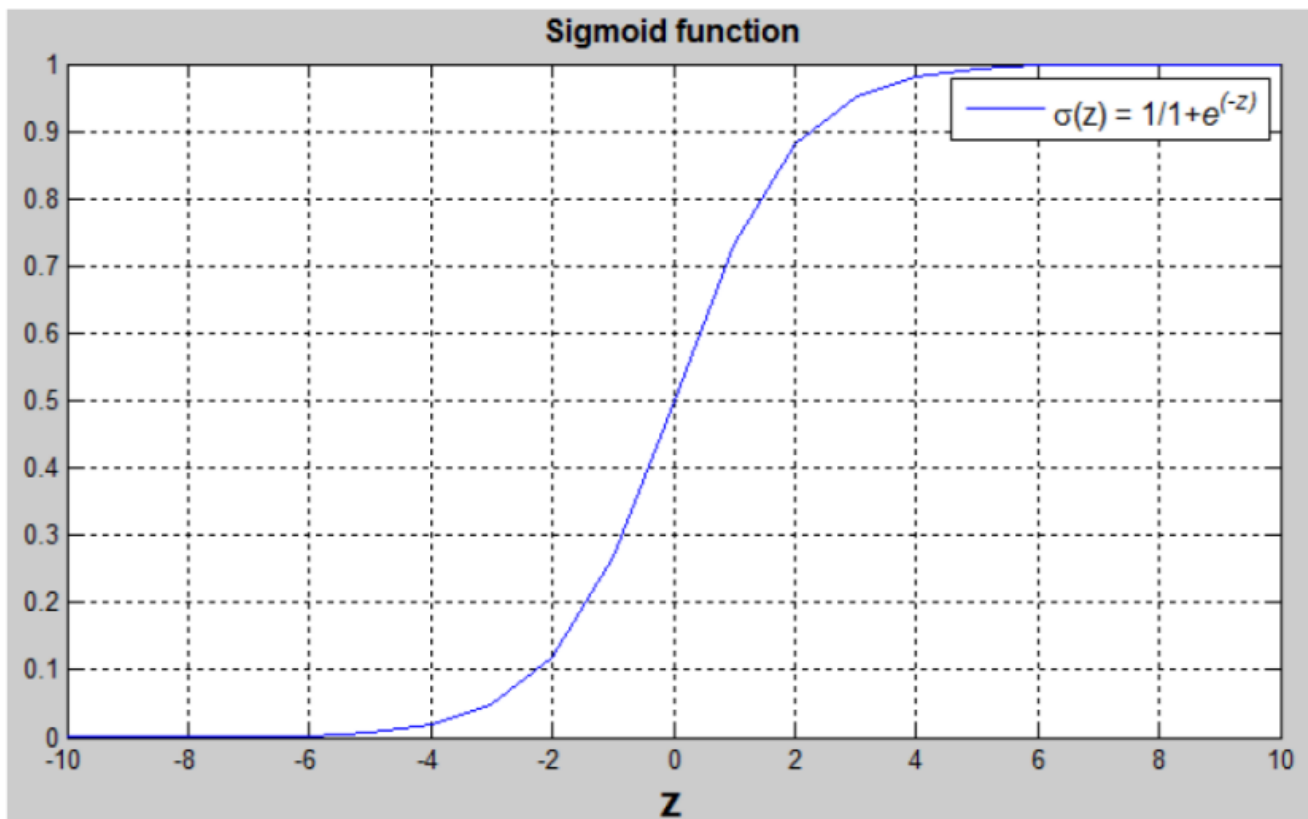$X \in \mathbb{R}^{n_x \times m}$      $X.shape = (n_x, m)$

## 2.2 逻辑回归

逻辑回归是一个用于监督学习问题的学习算法，输出值不是0就是1.

逻辑回归中用到的参数：

The parameters used in Logistic regression are:

- The input features vector: $x \in \mathbb{R}^{n_x}$, where $n_x$ is the number of features
- The training label: $y \in 0, 1$
- The weights: $w \in \mathbb{R}^{n_x}$, where $n_x$ is the number of features
- The threshold: $b \in \mathbb{R}$
- The output: $\hat{y} = \sigma(w^T x + b)$
- Sigmoid function: s = $\sigma(w^T x + b) = \sigma(z) = \dfrac{1}{1 + e^{-z}}$
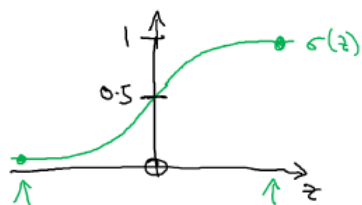
sigmold函数的图像

Sigmoid function

- 如果z趋向于比较大的正数，那么sigmold等于1。
- 如果z趋向于比较大的负数，那么sigmold等于0。
- 如果z等于零，那么sigmold等于0.5。

# Logistic Regression

Given $x$, want $\hat{y} = P(y=1|x)$

$x \in \mathbb{R}^{n_x}$

$0 \leq \hat{y} \leq 1$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$.

Output $\hat{y} = \sigma(\underbrace{w^T x + b}_{z})$



$x_0 = 1$, $x \in \mathbb{R}^{n_x+1}$

$\hat{y} = \sigma(\theta^T x)$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \begin{array}{l} \}b \leftarrow \\ \\ \}w \leftarrow \end{array}$$

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

If $z$ large $\sigma(z) \approx \frac{1}{1+0} = 1$

If $z$ large negativ numbe

$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+Bignum} \approx 0$

Andrew Ng

# Logistic Regression cost function

$\rightarrow \hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$, where $\sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$    $z^{(i)} = w^T x^{(i)} + b$

Given $\{(x^{(1)}, y^{(1)}),\ldots,(x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$.

$\left. \begin{array}{l} x^{(i)} \\ y^{(i)} \\ z^{(i)} \end{array} \right\}$ i-th example.

Loss (error) function:   $\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y}-y)^2$

$\boxed{\mathcal{L}(\hat{y}, y)} = -(y \log \hat{y} + (1-y) \log(1-\hat{y})) \leftarrow$

If $y=1$: $\mathcal{L}(\hat{y}, y) = -\log \hat{y}$   $\leftarrow$ Want $\log \hat{y}$ large, want $\hat{y}$ large.

If $y=0$: $\mathcal{L}(\hat{y}, y) = -\log(1-\hat{y})$ $\leftarrow$ Want $\log 1-\hat{y}$ large $\ldots$ want $\hat{y}$ small

$\boxed{\text{Cost}}$ function: $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right]$

Andrew Ng

## 2.3 逻辑回归损失函数

损失函数的公式

Loss (error) function:

The loss function measures the discrepancy between the prediction ($\hat{y}^{(i)}$) and the desired output ($y^{(i)}$). In other words, the loss function computes the error for a single training example.

$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2}(\hat{y}^{(i)} - y^{(i)})^2$

$L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)})\log(1 - \hat{y}^{(i)}))$

- If $y^{(i)} = 1$: $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$ where $\log(\hat{y}^{(i)})$ and $\hat{y}^{(i)}$ should be close to 1
- If $y^{(i)} = 0$: $L(\hat{y}^{(i)}, y^{(i)}) = -\log(1 - \hat{y}^{(i)})$ where $\log(1 - \hat{y}^{(i)})$ and $\hat{y}^{(i)}$ should be close to 0

但是，对于logistic regression 来说，一般不适用平方错误来作为Loss Function，这是因为上面的平方错误损失函数一般是非凸函数（non-convex），其在使用低度下降算法的时候，容易得到局部最优解，而不是全局最优解。因此要选择凸函数。

- 我们的目标是最小化样本点的损失Loss Function，损失函数是针对单个样本点的。

代价函数

Cost function

The cost function is the average of the loss function of the entire training set. We are going to find the parameters $w \ and \ b$ that minimize the overall cost function.

$$J(w,b) = \frac{1}{m}\sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m}\sum_{i=1}^{m}[(y^{(i)}\log(\hat{y}^{(i)}) + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})]$$

代价函数是损失函数的平均值。

## 2.4 梯度下降法

用梯度下降法（Gradient Descent）算法来最小化Cost function，以计算出合适的w和b的值。



Andrew Ng

每次迭代更新的修正表达式：

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

在程序代码中，我们通常使用

dw来表示$\frac{\partial J(w, b)}{\partial w}$，用db来表示$\frac{\partial J(w, b)}{\partial b}$。

# 2.5 基本的神经网络编程--导数

理解导数

导数的直观理解就是经过某点的直线的斜率

## Intuition about derivatives

$f(a) = 3a$



$a = 2 \quad f(a) = 6$

$a = 2.001 \quad f(a) = 6.003$

slope (derivation) of $f(a)$ at $a=2$ is 3

$\frac{0.003}{0.001} \quad \frac{height}{width}$

$a = 5 \quad f(a) = 15$

$a = 5.001 \quad f(a) = 15.003$

slope at $a=5$ is also 3

$\frac{d \, f(a)}{da} = 3 = \frac{d}{da} f(a)$

$0.001$
$0.0000001$
$0.00000001$

Andrew Ng

## 2.6 更过导数的例子

## Intuition about derivatives

$f(a) = a^2$



$\frac{d}{da} a^2 = 2a$

$0.001$

$(2a) \times 0.001$

$0.001$
$0.00000....01$

$a = 2 \quad f(a) = 4$

$a = 2.001 \quad f(a) \approx 4.004$

$(4.004001)$

slope (derivation) of $f(a)$ at $a=2$ is 4.

$\boxed{\frac{d}{da} f(a) = 4}$ when $\boxed{a=2}$

$a = 5 \quad f(a) = 25$

$a = 5.001 \quad f(a) \approx 25.010$

$\boxed{\frac{d}{da} f(a) = 10}$ when $\boxed{a=5}$

$\frac{d}{da} f(a) = \frac{d}{da} a^2 = \boxed{2a}$

Andrew Ng

**更多倒数的例子**

# More derivative examples

$f(a) = a^2$      $\dfrac{d}{da} f(a) = \underset{4}{2a}$      $a = 2$    $f(a) = 4$

                                             $a = 2.001$    $f(a) \approx 4.004$

$f(a) = a^3$      $\dfrac{d}{da} f(a) = \underset{3 \times 2^2 = 12}{3a^2}$      $a = 2$    $f(a) = 8$

                                             $a = 2.001$    $f(a) \approx 8.012$

$f(a) = \log_e(a)$      $\dfrac{d}{da} f(a) = \dfrac{1}{a}$      $a = 2$    $f(a) \approx 0.69315$

    $\ln(a)$                                    $a = 2.001$    $f(a) \approx 0.69265$

          $\dfrac{\ln(a)}{0.0005}$,   $0.001$     $\dfrac{d}{da} f(a) = \boxed{\dfrac{1}{2}}$      $0.0005$    $0.0005$

## 2.7 计算图

# Computation Graph

$J(a, b, c) = 3(\underset{J}{\underset{V}{\underbrace{a + \underset{u}{\underbrace{bc}}}}}) = 3(5 + 3 \times 2) = 33$

$u = bc$
$V = a + u$
$J = 3v$

链式法则求导

# Computing derivatives

$$a = 5$$
$$\frac{dJ}{da} \quad \text{`da`} = 3$$

$$b = 3$$

$$u = bc$$

$$c = 2$$

11

$$v = a + u \longrightarrow J = 3v$$

33

$$\frac{dJ}{dv} \quad \text{`dv`} = 3$$

$$\frac{dJ}{dv} = ? = 3$$

$$a \to v \to J$$

$$J = 3v$$
$$v = 11 \to 11.001$$
$$J = 33 \to 33.003$$

$$\frac{dJ}{da} = 3 = \frac{dJ}{dv}\frac{dv}{da}$$

$$3 \times 1$$

$$\frac{dv}{da} = 1$$

$$a = 5 \to 5.001$$
$$\to v = 11 \to 11.001$$
$$J = 33 \to 33.003$$

$$f(a) = 3a$$

$$\frac{df(a)}{da} = \frac{df}{da} = 3$$

$$J = 3v$$
$$\frac{dJ}{dv} = 3$$

$$\frac{d \, Final \, Output \, Var}{d \, var}$$

$$\frac{dJ dvar}{} \to \text{`dvar`}$$

Andrew Ng

## 2.9 逻辑回归中的梯度下降法

逻辑回归损失函数回顾

# Logistic regression recap

$$\to z = w^T x + b$$
$$\to \hat{y} = a = \sigma(z)$$
$$\to \mathcal{L}(a, y) = -(y \log(a) + (1-y)\log(1-a))$$

$$x_1$$
$$w_1$$
$$x_2$$
$$w_2$$
$$b$$
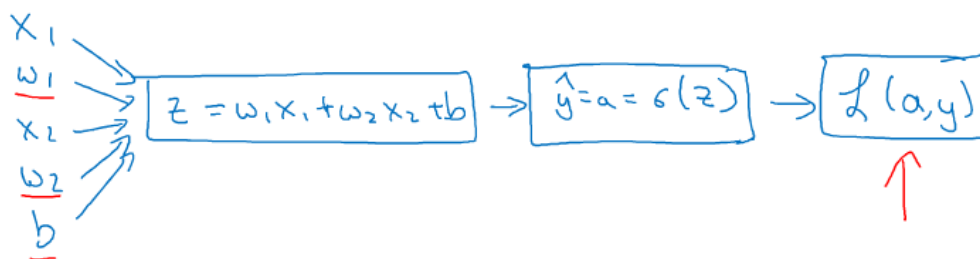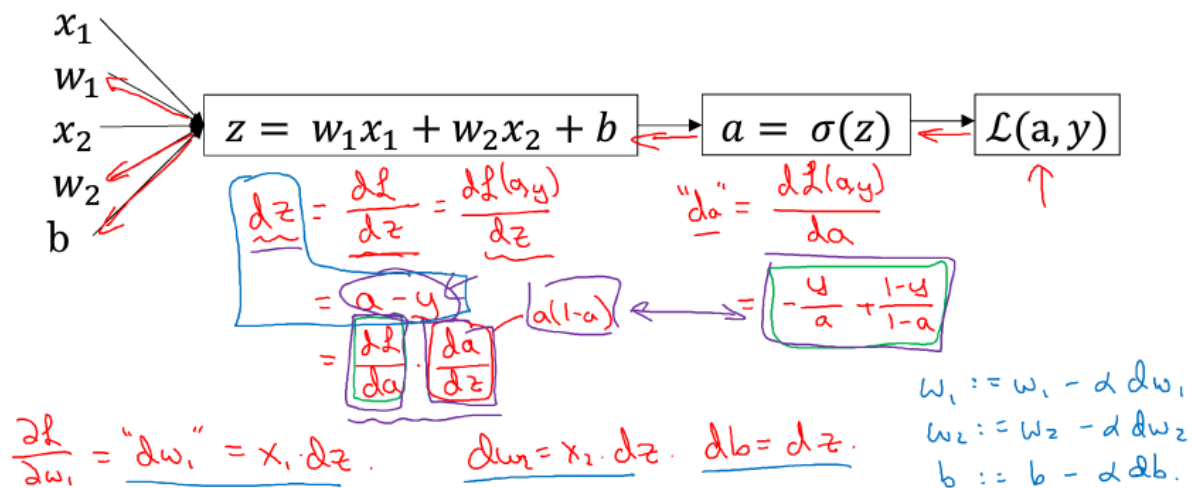
$$z = w_1 x_1 + w_2 x_2 + b \to \hat{y} = a = \sigma(z) \to \mathcal{L}(a, y)$$

Andrew Ng

逻辑回归求导

# Logistic regression derivatives

$x_1$
$w_1$
$x_2$
$w_2$
$b$

$$z = w_1 x_1 + w_2 x_2 + b \leftarrow a = \sigma(z) \leftarrow \mathcal{L}(a, y)$$

$$\underline{dz} = \frac{d\ell}{dz} = \frac{d\mathcal{L}(a,y)}{dz} \qquad "da" = \frac{d\mathcal{L}(a,y)}{da}$$

$$= a - y$$

$$= \frac{d\ell}{da} \cdot \frac{da}{dz} \qquad a(1-a) \leftarrow = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{d\ell}{dw_1} = "dw_1" = x_1 \cdot dz. \qquad dw_2 = x_2 \cdot dz. \qquad db = dz.$$

$$w_1 := w_1 - \alpha \, dw_1$$
$$w_2 := w_2 - \alpha \, dw_2$$
$$b := b - \alpha \, db.$$

Andrew Ng

## 2-10 m个样本的梯度下降法

# Logistic regression on $m$ examples

$$J = 0 ; \quad dw_1 = 0 ; \quad dw_2 = 0 ; \quad db = 0$$

$\rightarrow$ For $i = 1$ to $m$

$\quad z^{(i)} = w^T x^{(i)} + b$

$\quad a^{(i)} = \sigma(z^{(i)})$

$\quad J += -\left[ y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)}) \right]$

$\quad \underline{dz^{(i)}} = a^{(i)} - y^{(i)}$

$\quad dw_1 += x_1^{(i)} \, dz^{(i)}$

$\quad dw_2 += x_2^{(i)} \, dz^{(i)} \qquad \Big\} \, n = 2$

$dw_3 \atop \vdots \atop dw_n$ $\quad db += dz^{(i)}$

$J /= m$

$dw_1 /= m ; \quad dw_2 /= m ; \quad db /= m.$

$$dw_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha \, \underline{dw_1}$$
$$w_2 := w_2 - \alpha \, \underline{dw_2}$$
$$b := b - \alpha \, \underline{db}.$$

Vectorization

Andrew Ng

对m个样本来说，其Cost function表达式如下：

$$z^{(i)} = w^T x^{(i)} + b$$
$$\hat{y}^{(i)} = a^{(i)} = \sigma(z^{(i)})$$
$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$$

Cost function 关于w和b的偏导数可以写成所有样本点偏导数和的平均形式：

$$dw_1 = \frac{1}{m} \sum_{i=1}^{m} x_1^{(i)} (a^{(i)} - y^{(i)})$$

$$db = \frac{1}{m} \sum_{i=1}^{m} (a^{(i)} - y^{(i)})$$

## 2.14 向量化logistics回归的梯度输出

在深度学习的算法中，我们通常拥有大量的数据，在程序的编写过程中，应该尽最大可能的少使用loop循环语句，利用python可以实现矩阵运算，进而来提高程序的运行速度，避免for循环的使用。

什么是向量化？



### 逻辑回归向量化

输入矩阵X：$( n_x, m )$

权重矩阵W：$( n_x, 1 )$

偏置b：为一个常数

输出矩阵Y：$( 1, m )$

所有m个样本的线性输出Z可以用矩阵表示：

$$Z = w^T X + b$$

python 代码

```
Z = np.dot(w,T,X) + b
A = sigold(Z)
```

## 2.12 更多的向量化的例子--神经网络

# Neural network programming guideline

## Whenever possible, avoid explicit for-loops.

$u = Av$

$u_i = \sum_i \sum_j A_{ij} v_j$

$u = np.zeros((n,1))$
for i ...
    for j ...
       u[i] += A[i][j] * v[j]

$u = np.dot(A,v)$

# Vectors and matrix valued functions

Say you need to apply the exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_c} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

```
u = np.zeros((n,1))
for i in range(n):
    u[i]=math.exp(v[i])
```

import numpy as np
$u = np.exp(v)$

np.log(v)
np.abs(v)
np.maximum(v,0)
v**2
1/v

# Logistic regression derivatives

$J = 0,$ $\boxed{dw1 = 0,\ dw2 = 0},$ $db = 0$ $\quad dw = np.zeros((n\text{-}x, 1))$

$\rightarrow$ for i = 1 to n:

$\qquad z^{(i)} = w^T x^{(i)} + b$

$\qquad a^{(i)} = \sigma(z^{(i)})$

$\qquad J \mathrel{+}= -\left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})\right]$

for j=1...$n_x$  $\qquad dz^{(i)} = a^{(i)}(1 - a^{(i)})$

$dw_j \mathrel{+}=$ $\qquad \boxed{\begin{array}{l} dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)} \\ dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \end{array}}$ $\quad n_x = 2 \qquad dw \mathrel{+}= x^{(i)} dz^{(i)}$

$\qquad db \mathrel{+}= dz^{(i)}$

$J = J/m,$ $\boxed{dw_1 = dw_1/m,\ dw_2 = dw_2/m},$ $db = db/m$

$\qquad\qquad\qquad dw\ /= m.$

# 2.13 向量化逻辑回归的梯度计算

## Vectorizing Logistic Regression

$\underline{dz^{(1)} = a^{(1)} - y^{(1)}} \qquad \underline{dz^{(2)} = a^{(2)} - y^{(2)}} \quad \cdots$

$\boxed{\circled{dz} = \underbrace{[dz^{(1)}\ dz^{(2)} \cdots dz^{(m)}]}_{1 \times m}} \leftarrow$

$A = [a^{(1)} \cdots a^{(m)}]. \qquad Y = [y^{(1)} \cdots y^{(m)}]$

$\rightarrow dz = A - Y = [a^{(1)} - y^{(1)} \quad a^{(2)} - y^{(2)} \quad \cdots]$

$\rightarrow dw = 0$ $\qquad db = 0$

$\boxed{\begin{array}{l} dw \mathrel{+}= x^{(1)} dz^{(1)} \\ dw \mathrel{+}= x^{(2)} dz^{(2)} \\ \vdots \end{array}}$ $\quad \begin{array}{l} db \mathrel{+}= \circled{dz^{(1)}} \\ db \mathrel{+}= \circled{dz^{(2)}} \\ \vdots \\ db \mathrel{+}= \circled{dz^{(m)}} \end{array}$

$dw/= m$ $\qquad db/= m.$

$db = \frac{1}{m} \sum_{i=1}^{m} dz^{(i)}$

$\qquad = \frac{1}{m} np.sum(dz)$

$\boxed{dw = \frac{1}{m} X\, dz^T}$

$= \frac{1}{m} \begin{bmatrix} x^{(1)} \cdots x^{(m)} \\ | \quad\quad | \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$

$= \frac{1}{m} \underbrace{[x^{(1)} dz^{(1)} + \cdots + x^{(m)} dz^{(m)}]}_{n_x \times 1}$

# Implementing Logistic Regression

$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0$

for i = 1 to m:

$\quad z^{(i)} = w^T x^{(i)} + b \quad\leftarrow$

$\quad a^{(i)} = \sigma(z^{(i)}) \quad\leftarrow$

$\quad J \mathrel{+}= -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$

$\quad dz^{(i)} = a^{(i)} - y^{(i)} \quad\leftarrow$

$\quad \left.\begin{array}{l} dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)} \\ dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \end{array}\right\} \quad dw \mathrel{+}= x^{(i)} \ast dz^{(i)}$

$\quad db \mathrel{+}= dz^{(i)}$

$J = J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m$

$db = db/m$

for iter in range (1000): $\leftarrow$

$Z = w^T X + b$

$\quad = np.dot(w.T, X) + b$

$A = \sigma(Z)$

$dZ = A - Y$

$dw = \frac{1}{m} X \, dZ^T$

$db = \frac{1}{m} np.sum(dz)$

$w := w - \alpha \, dw$

$b := b - \alpha \, db$

Andrew Ng

```
Z = np.dot(w,T,X) + b
A = sigmold(Z)
dZ = A-Y
dw = 1/m*np.dot(X,dZ.T)
db = i/m*np.sum(dZ)

w = w - alpha*dw
b = b - alpha*db
```

# 2.15 python 中的广播

# Broadcasting example

Calories from Carbs, Proteins, Fats in 100g of different foods:

|        | Apples | Beef  | Eggs | Potatoes |
|--------|--------|-------|------|----------|
| Carb   | 56.0   | 0.0   | 4.4  | 68.0     |
| Protein| 1.2    | 104.0 | 52.0 | 8.0      |
| Fat    | 1.8    | 135.0 | 99.0 | 0.9      |

$= A$

$(3,4)$

59 cal

$\frac{56}{59} \approx 94.9\%$

Calculate % of calories from Carb, Protein, Fat. Can you do this without explicit for-loop?

```
cal = A.sum(axis = 0)
percentage = 100*A/(cal.reshape(1,4))
```

$\uparrow(3,4)$  /  $(1,4)$

# Broadcasting example

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} 100$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

$(m,n)$ $(2,3)$ $(1,n) \leadsto (m,n)$ $(2,3)$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} =$$

$(m,n)$

$(m,1)$
$\downarrow$
$(m,n)$

## General Principle

$(m, n)$  
matrix

$$+ \quad -\quad *\quad /$$

$(1, n) \quad \leadsto \quad (m, n)$

$(m, 1) \quad \leadsto \quad (m, n)$

$(m, 1) \quad + \quad \mathbb{R}$

$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad + \quad 100 \quad = \quad \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$

$[1 \ 2 \ 3] \quad + \quad 100 \quad = \quad [101 \quad 102 \quad 103]$

Matlab/Octave: bsxfun

## 2.16 python numpy 向量的说明

- 虽然在Python有广播的机制，但是在Python程序中，为了保证矩阵运算的正确性，可以使用reshape()函数来对矩阵设定所需要进行计算的维度，这是个好的习惯；
- 如果用下列语句来定义一个向量，则这条语句生成的a的维度为（5，），既不是行向量也不是列向量，称为秩（rank）为1的array，如果对a进行转置，则会得到a本身，这在计算中会给我们带来一些问题。

## Python / numpy vectors

```
import numpy as np

a = np.random.randn(5)

a = np.random.randn((5,1))

a = np.random.randn((1,5))

assert(a.shape = (5,1))
```

- 如果需要定义（5，1）或者（1，5）向量，要使用下面标准的语句：

```
a = np.random.randn(5,1)
b = np.random.randn(1,5)
```

- 可以使用assert语句对向量或数组的维度进行判断。assert会对内嵌语句进行判断，即判断a的维度是不是（5，1），如果不是，则程序在此处停止。使用assert语句也是一种很好的习惯，能够帮助我们及时检查、发现语句是否正确。

```
assert(a.shape == (5,1))
```

- 可以使用reshape函数对数组设定所需的维度

```
a.reshape((5,1))
```