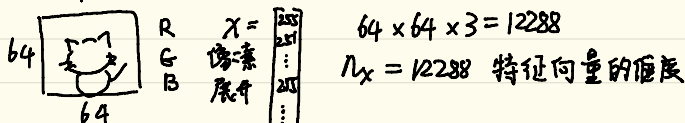


1.2 Course 1 Neural Networks And Deep Learning Week 2 Logistic Regression as a Neural Network

1.2 Video 1:

Binary Classification

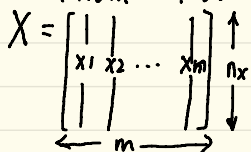
- Example cat: 1 not cat: 0



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

m training example: $\{(x^{(1)}, y^{(1)})\}, (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$m_{\text{train}} = \# \text{ training examples}$; $m_{\text{test}} = \# \text{ test examples}$



$X \in \mathbb{R}^{n_x \times m}$ 向量维度 \times 向量训练个数

$x.\text{shape} (n_x, m)$ python

$Y = [y^{(1)}, y^{(2)}, y^{(3)} \dots y^{(m)}]$

$y.\text{shape} (1, m)$ python

1.2 video 2:

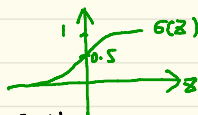
Logistic Regression

- Given x , want $\hat{y} = P(y=1|x)$ $0 \leq \hat{y} \leq 1$

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

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

\downarrow
sigmoid(z)



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

if $z \rightarrow \infty$

$\sigma(z) \rightarrow 1$

if $z \rightarrow -\infty$

$\sigma(z) \rightarrow 0$

另一种方式: $x_0 = 1$ $x \in \mathbb{R}^{n_x+1}$

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

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \rightarrow w$

1.2 Video 3:

Logistic Regression Cost Function

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

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

Loss (error) function: $\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ > 尽量小, 衡量 \hat{y} 与真实 y 的误差

error a single training example

$$\mathcal{L}(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log (1-\hat{y})] \text{ small}$$

$$\text{if } y=1 \Rightarrow \mathcal{L}(\hat{y}, y) = -\log \hat{y} \rightarrow \log \hat{y} \text{ large} \rightarrow \hat{y} \text{ large}$$

$$\rightarrow \hat{y} \in \{0, 1\} \rightarrow \hat{y} \rightarrow 1$$

$$\text{if } y=0 \Rightarrow \mathcal{L}(\hat{y}, y) = -\log (1-\hat{y}) \rightarrow \log (1-\hat{y}) \text{ large} \rightarrow \hat{y} \text{ small}$$

average of loss error

$$\rightarrow \hat{y} \in \{0, 1\} \rightarrow \hat{y} \rightarrow 0$$

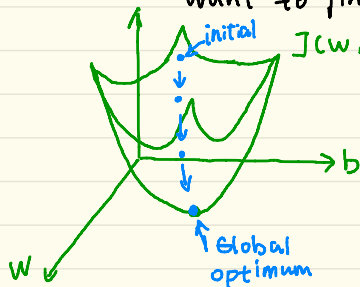
$$\text{Cost function} = J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

1.2 Video 4:

Gradient Descent

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

want to find w, b to minimize $J(w, b)$



凸函数: 定义在某个向量空间的凸子集 C 上的实值函数 f 对任意 $x, y \in \text{dom } f$, $0 \leq \theta \leq 1$ 有 $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$



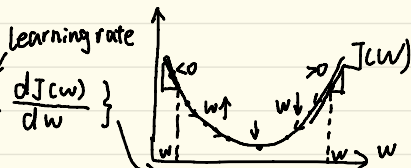
凸, 有最优解



非凸, 无最优解

忽略 b :
Repeat

$$\left\{ \begin{array}{l} w := w - \alpha \frac{dJ(w)}{dw} \\ \text{迭代更新} \end{array} \right\}$$



code use: dw'' 代表求导 $\Rightarrow w = w - \alpha dw$

考虑 b :

Repeat

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

1.2 Video 5.6

Derivatives 微积分基础

1.2 Video 7

computation graph

$$J(a, b, c) = 3(a + bc) \quad u = bc; v = a + u; J = 3v$$

$$a = 5$$

$$b = 3$$

$$c = 2$$

$$u = bc \rightarrow v = a + u \rightarrow J = 3v$$

正传播 Forward propagation

1.2 Video 8.

Derivatives with a Computation graph

$$a = 5$$

$$b = 3$$

$$c = 2$$

$$u = bc \rightarrow v = a + u \rightarrow J = 3v$$

backward propagation 反向传播: 求导

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

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

代码中可用 'dvar' 代表 $\frac{d(\text{Final output})}{d(\text{var})}$

1.2 video 9

Logistic Regression Gradient Descent

$$\begin{aligned} & x_1 \rightarrow \\ & \frac{dw_2}{dx_2} = x_2 \cdot \frac{dw_2}{dx_2} \rightarrow \\ & \frac{db}{db} = \frac{dw_2}{db} \rightarrow \\ & \frac{dw_1}{dw_1} = \frac{dz}{dz} \cdot \frac{dz}{dw_1} = dz \cdot x_1 \end{aligned}$$
$$z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = \sigma(z) \rightarrow \mathcal{L} = (\hat{y}, y)$$
$$\frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}}{d\hat{y}} \cdot \frac{d\hat{y}}{dz}$$
$$\frac{d\mathcal{L}}{d\hat{y}} = -\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})} = da$$
$$= \left[-\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right] \times \hat{y}(1-\hat{y}) \leftarrow \left(\frac{1}{1+e^{-z}} \right)'$$
$$= \hat{y} - y = a - y = dz$$

$$w_1 = w_1 - \alpha dw_1 \quad \text{update ...}$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

1.2 video 10

Gradient Descent on m Examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y) \Rightarrow a^{(i)} = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \ell(a^{(i)}, y^{(i)})$$

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

For $i=1$ to m 第1循环所有 training example m

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

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

$$\Rightarrow J_t = -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log (1 - a^{(i)})]$$

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

For $\left. \begin{aligned} dw_1 &+= dz^{(i)} \times x_1^{(i)} \\ dw_2 &+= dz^{(i)} \times x_2^{(i)} \end{aligned} \right\} n=2$ 第2循环所有特征 n

$$db += dz^{(i)}$$

$$\Rightarrow J \neq 0 \quad dw_1 \neq 0 \quad dw_2 \neq 0 \quad db \neq 0$$

$$\Rightarrow w_1 := w_1 - \alpha dw_1 \quad w_2 := w_2 - \alpha dw_2 \quad b := b - \alpha db$$

但数据集过大, 用For效率太低. 所以用 Vectorization
加速运算

For iter in range(1000): 迭代1000次梯度下降
np.dot(a, b) \Rightarrow a, b 为2个多维向量矩阵

去掉1st For: $z = np.dot(W.T, X) + b$ ★ "broadcasting"

$$dz = [dz^{(1)}, dz^{(2)}, \dots, dz^{(m)}]_{1 \times m}$$

$$A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = \sigma(z)$$

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

$$dz = A - Y$$

去掉2nd For: $db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} = \frac{1}{m} np.sum(dz)$

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

$$= \frac{1}{m} \begin{bmatrix} x_1 & \dots & x_m \\ | & & | \\ | & & | \end{bmatrix}_{n \times m} \times \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}_{m \times 1}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]_{n \times 1}$$

Python - broadcasting: $A_{3 \times 4} / B_{1 \times 4}$

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

$$\therefore (m, n) \overset{\pm}{\times} (1, n) \leftrightarrow (m, n)$$

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

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

bug: ~~X~~ $a = \text{np.random.randn}(5)$ Don't use
 $a.\text{shape} = (5,)$ \rightarrow 秩为1的数组而非矩阵

$\Rightarrow a = \text{np.random.randn}(5, 1)$
 $a.\text{shape} = (5, 1)$

~~或~~ $\text{assert}(a.\text{shape} == (5, 1))$

~~或~~ $a = a.\text{reshape}(5, 1)$

Shift + Enter. / Cell Run

1.2 video 18

Explanation of logic regression cost function

$$\text{If } y=1 \quad P(y|x) = \hat{y}$$

$$\text{If } y=0 \quad P(y|x) = 1 - \hat{y}$$

$$P(y|x) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

$\log P(y|x) \rightarrow$ 取 \log 令它单调

$$= y \log \hat{y} + (1-y) \log (1 - \hat{y})$$

想使它最小, 因此前加负号

\Rightarrow 最大似然估计