

## Binary Classification



output for binary classification only has 2 classes, e.g. 0 or 1 (non-cat or cat)  
i.e.  $y = 0$  (non-cat) or  $1$  (cat)

Blue	Green	Red	how an image is represented in a computer (note this is way simplified)
255	134	93	22
255	231	42	22
123	94	83	2
34	44	187	92
34	76	232	124
67	83	194	202

an image is represented as 3 matrices with the size equal to the pixels of the image. For example, if the image is 64 pixels by 64 pixels, then it's represented as 3 64 by 64 matrices corresponding to the red, green and blue pixel intensity value for the image. So to unroll these pixel intensity values into feature vector, we will define a vector  $\mathbf{x}$  to store all these values in an orderly fashion.

$$\text{i.e. } \mathbf{x} = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix} \quad \text{then the dimension of } \mathbf{x} \text{ is } 64 \times 64 \times 3 = 12288 \\ \Rightarrow n_x \text{ or } n = 12288$$

Notation:

- a single training example  $(\mathbf{x}, y)$ , where  $\mathbf{x} \in \mathbb{R}^n$ ,  $y \in \{0, 1\}$
- $m$  training examples:  $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$ , where  $m = \# \text{ training examples. or } m = n_{\text{train}}$
- $n_{\text{test}} = \# \text{ test examples}$

input matrix:  $\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \end{bmatrix}_{n \times m} \quad \mathbf{x} \in \mathbb{R}^{n \times m} \quad \mathbf{x} \text{ shape} = (n, m)$

output matrix:  $\mathbf{Y} = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]_{1 \times m} \quad \mathbf{Y} \in \mathbb{R}^{1 \times m} \quad \mathbf{Y} \text{ shape} = (1, m)$

## Algorithm used for binary classification: logistic Regression

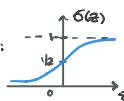
For the  $i$ -th example,

- Given  $\mathbf{x}^{(i)}$ , estimate of output  $y^{(i)}$  is  $\hat{y}^{(i)} = \Pr(y^{(i)}=1|\mathbf{x}^{(i)})$ , i.e. what is the probability of the image being a cat given  $\mathbf{x}^{(i)}$

- Suppose we have parameters  $\mathbf{w} \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ , and we model the estimated output as:

$$\hat{y}^{(i)} = \sigma(\mathbf{w}^T \mathbf{x}^{(i)} + b), \text{ where } \sigma \text{ is sigmoid function:}$$

which is a single layer NN, later we will show output of more layers NN



where does it come from:  $\Pr(y=1|\mathbf{x}) = \hat{y}^y (1-\hat{y})^{1-y}$   
 $\Pr(y=1|\mathbf{x}) = \hat{y} \Rightarrow \begin{cases} \text{if } y=1, \Pr(y=1|\mathbf{x}) = \hat{y} \approx 1 \\ \text{if } y=0, \Pr(y=0|\mathbf{x}) = 1-\hat{y} \approx 0 \end{cases}$

Sometimes there might be other notations:  
 $\mathbf{x}_0 = 1, \mathbf{x} \in \mathbb{R}^{n+1}$   
 $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x}), \text{ where } \mathbf{w} = \begin{bmatrix} \mathbf{w}_0 \\ \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_n \end{bmatrix} \in \mathbb{R}^{n+1}$

- Using  $\mathbf{z}^{(i)}$  to denote  $\mathbf{w}^T \mathbf{x}^{(i)} + b$ ,  $\hat{y}^{(i)} = \sigma(\mathbf{z}^{(i)}) = \frac{1}{1+e^{-\mathbf{z}^{(i)}}}$ , where  $\mathbf{z}^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b \quad \begin{cases} \text{if } \mathbf{z} \text{ is large, } \sigma(\mathbf{z}) \rightarrow 1 \\ \text{if } \mathbf{z} \text{ is small, } \sigma(\mathbf{z}) \rightarrow 0 \end{cases}$

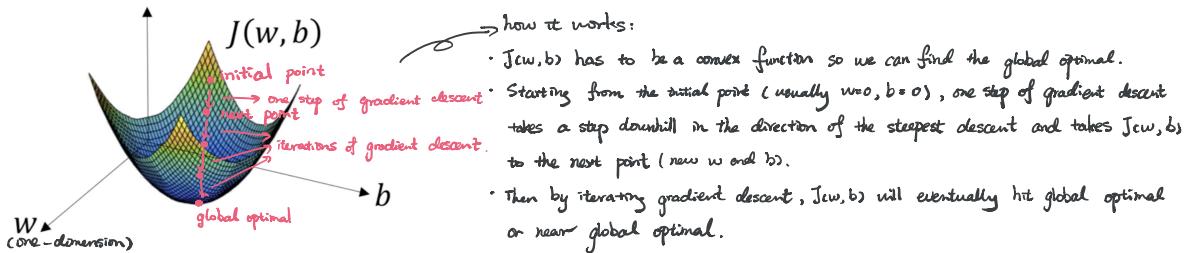
- To measure the performance of a single prediction, we define loss function:

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y})) \quad \begin{cases} \text{if } y=1, L(\hat{y}, y) = -\log \hat{y} \rightarrow \text{want } \hat{y} \approx 1 \Rightarrow \text{want } \log \hat{y} \text{ small} \\ \text{if } y=0, L(\hat{y}, y) = -\log(1-\hat{y}) \rightarrow \text{want } \hat{y} \approx 0 \Rightarrow \text{want } \log(1-\hat{y}) \text{ large} \Rightarrow \text{want } -\log(1-\hat{y}) \text{ small} \end{cases} \\ \Rightarrow \text{want } L(\hat{y}, y) \text{ small}$$

- Then for the entire dataset, we can measure the performance by cost function:

$$J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})] \quad \text{it's actually the log-likelihood function of } y|\mathbf{x}, \\ \text{i.e. } \sum_{i=1}^m \log \Pr(y^{(i)}|\mathbf{x}^{(i)}) = \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

Now we use an algorithm called **Gradient Descent** to find  $w$  and  $b$  that minimise cost function.



To put this into mathematical representation:

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

↓

"learning rate", represented as " $\alpha w$ ", decides how big a step each iteration takes

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

↓

" $\alpha b$ "

where  $\frac{\partial w}{\partial w} < 0 \Rightarrow$  as  $w \uparrow$ ,  $J(w) \downarrow \Rightarrow$  want  $w \uparrow \Rightarrow -\alpha w > 0 \Rightarrow$  new  $w \uparrow$

$\frac{\partial w}{\partial w} > 0 \Rightarrow$  as  $w \uparrow$ ,  $J(w) \uparrow \Rightarrow$  want  $w \downarrow \Rightarrow -\alpha w < 0 \Rightarrow$  new  $w \downarrow$

To have an explicit expression for  $dw$  and  $db$ :

$$\text{let } a = \sigma(z) = \frac{1}{1+e^{-z}} \quad | \quad \frac{\partial a}{\partial z} = \frac{1}{(1+e^{-z})^2} \cdot e^{-z} = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} = a(1-a)$$

$$\frac{\partial l}{\partial a} = \frac{\partial l}{\partial a} \cdot (y \log a + (1-y) \log (1-a)) = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial l}{\partial z} = \frac{\partial l}{\partial a} \cdot \frac{\partial a}{\partial z} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) a(1-a) = -y(1-a) + (1-y)a = a - y$$

Then

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial z} \cdot \frac{\partial z}{\partial w_p} = (a - y) \cdot x_p, \quad p=1,2,\dots,n$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial z} \cdot \frac{\partial z}{\partial b} = (a - y)$$

Finally,

$$dw = \frac{\partial l}{\partial w} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m \frac{\partial l}{\partial w} (\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) \cdot x^{(i)}$$

$$db = \frac{\partial l}{\partial b} = \frac{1}{m} \sum_{i=1}^m \frac{\partial l}{\partial b} (\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m \frac{\partial l}{\partial b} (\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

Finally, we've come to the exciting **computing part**. Here we will show 2 methods of how to put logistic regression and Gradient Descent algorithm into coding:

Method 1 for-loop:

initialize with  $J=0$ ,  $dw_1=0, dw_2=0, \dots dw_n=0$ ,  $db=0$

For  $i=1$  to  $m$ ,

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

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

$t =$  the sum of  $m$  examples

$$J \leftarrow J - (y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)}))$$

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

$$dw_1 \leftarrow d^{(i)} \cdot x_1$$

$$dw_2 \leftarrow d^{(i)} \cdot x_2$$

$$\vdots$$

$$dw_n \leftarrow d^{(i)} \cdot x_n$$

$$db \leftarrow d^{(i)}$$

$$J = J/m \text{ and similarly for } dw_1, \dots, db$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$\vdots$$

$$w_n := w_n - \alpha dw_n$$

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

known as "feedforward"

known as "backpropagation"

\* Method 2 vectorization:

$$\text{Recall } X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m)} \end{bmatrix}_{n \times m} \in \mathbb{R}^{n \times m}, w = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}_{n \times 1} \in \mathbb{R}^n, Y = [y^{(1)} \dots y^{(m)}]_{1 \times m}$$

Then replacing the above codes, we have:

$$J = 0, dw = np.zeros((n, x, 1)), db = 0$$

$$Z = \underbrace{\np.dot(wT, x) + b}_{wT x + b} = [z^{(1)} \dots z^{(m)}]_{1 \times m} \quad (\text{note } b \in \mathbb{R} \text{ but this code still works thanks to Broadcasting in Python})$$

$$A = \sigma(Z) = [a^{(1)} \dots a^{(m)}]_{1 \times m}$$

$$dZ = A - Y = [a^{(1)} - y^{(1)} \dots a^{(m)} - y^{(m)}]_{1 \times m} = [dz^{(1)} \dots dz^{(m)}]_{1 \times m}$$

$$dW = \frac{1}{m} \times dZ.T = \frac{1}{m} \left[ \begin{bmatrix} z_1^{(1)} & \dots & z_1^{(m)} \end{bmatrix}_{n \times m} \begin{bmatrix} dz_1^{(1)} & \dots & dz_1^{(m)} \end{bmatrix}_{m \times 1} \right] = \frac{1}{m} \begin{bmatrix} \sum_i z_i^{(1)} dz_i^{(1)} \\ \vdots \\ \sum_i z_i^{(m)} dz_i^{(m)} \end{bmatrix}_{n \times 1} = \begin{bmatrix} dw_1 \\ \vdots \\ dw_n \end{bmatrix}_{n \times 1}$$

$$dw = w - dw$$

$$db = b - db$$

$$J = -\frac{1}{m} \cdot (\np.dot(Y, \log(A.T)) + \np.dot(1-Y, \log(1-A.T)).T)$$

Important: Always try to avoid coding with explicit expressions (i.e. for-loops) because it's too inefficient comparing to vectorization. For more details about vectorization, please view "Vectorization\_and\_Broadcasting.ipynb".