

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	4
123	94	83	2	192
34	44	187	92	34
34	76	232	124	94
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.

i.e. $\mathbf{x} = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}$ then the dimension of \mathbf{x} is $64 \times 64 \times 3 = 12288$
 $\Rightarrow n_x$ 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 $m = \#$ training examples. or $m = m_{train}$
- m_{test} = # test examples

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

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

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

- Using $z^{(i)}$ to denote $\mathbf{w}^T \mathbf{x}^{(i)} + b$, $\hat{y}^{(i)} = \frac{1}{1+e^{-z^{(i)}}}$, where $z^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b$

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

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

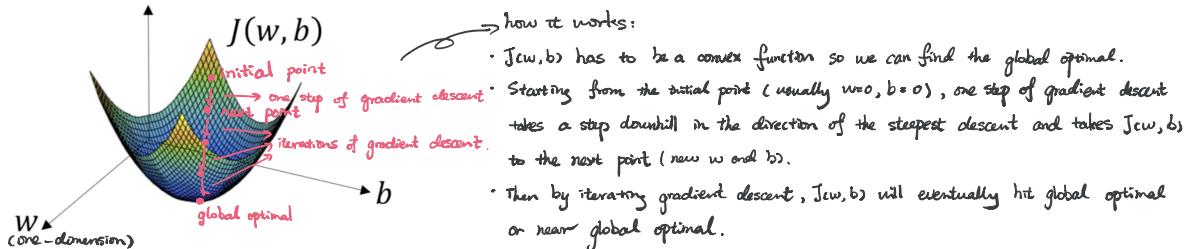
$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

$\left\{ \begin{array}{l} \text{if } y=1, L(\hat{y}, y) = -\log \hat{y} \rightarrow \text{want } \hat{y} \approx 1 \Rightarrow \text{want } \log \hat{y} \text{ small} \\ \text{if } y=0, L(\hat{y}, y) = -\log (1-\hat{y}) \rightarrow \text{want } \hat{y} \approx 0 \Rightarrow \text{want } \log (1-\hat{y}) \text{ large} \Rightarrow \text{want } -\log (1-\hat{y}) \text{ small} \end{array} \right.$
 $\Rightarrow \text{want } L(\hat{y}, y) \text{ small}$

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

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

- 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 " αw ", decides how big a step each iteration takes
 ↓ "dw" sets the direction of the steepest descent to reach global optimal

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

↓ "db"

To have an explicit expression for dw and db :

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

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

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

Then

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

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

Finally,

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

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

Then in a loop we need to define:

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

For $i=1$ to m ,

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

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

\leftarrow \sum of m examples

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

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

$$dw_1 += d\hat{z}^{(i)} \cdot x_1$$

$$dw_2 += d\hat{z}^{(i)} \cdot x_2$$

$$\vdots$$

$$dw_n += d\hat{z}^{(i)} \cdot x_n$$

$$db += d\hat{z}^{(i)}$$

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

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

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

$$\vdots$$

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

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

} known as "feedforward"

} known as "backpropagation"

But putting this explicit expressions in the code is too inefficient, so we need **Vectorization**:

for example, vectorization of z : $z = \underbrace{\text{np.dot}(w^T, x)}_{w^T x} + b$

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

Then replacing the above codes, we have:

$$J = 0, dw = \text{np.zeros}(c_n \cdot x, 1), db = 0$$

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

$$A = \phi(z) = [a^{(1)} \cdots a^{(m)}]_{1 \times m}$$

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

$$dw = \frac{1}{m} \times dz^T = \frac{1}{m} \left[x^{(1)} \cdots x^{(m)} \right]_{m \times n} \left[\frac{dz^{(1)}}{dw^{(1)}} \cdots \frac{dz^{(m)}}{dw^{(m)}} \right]_{m \times 1} = \frac{1}{m} \begin{bmatrix} z_1^{(1)} dz^{(1)} \\ \vdots \\ z_1^{(m)} dz^{(m)} \end{bmatrix}_{n \times 1} = \begin{bmatrix} dw_1 \\ \vdots \\ dw_n \end{bmatrix}_{n \times 1}$$

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

$$w := w - dw$$

$$b := b - db$$

$$J = -\text{np.dot}(Y, \log(A)^T) - \text{np.dot}(1 - Y, \log(1 - A))^T$$