

Week-2

Logistic Regression is a type of binary classification

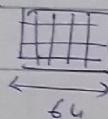
Binary Classification

example: Cat or NonCat

→ The image vectors^{pixel matrices} are combined for input vector

Consider 64p image

R



\downarrow^{64} similarly G & B

input vector $X = \begin{bmatrix} \vdots & \vdots & \vdots \\ R & G & B \\ \vdots & \vdots & \vdots \end{bmatrix}$ here, $n = n_x = 12288$ ($64 \times 64 \times 3$)

Single training example is represented as $(x, y) \quad x \in \mathbb{R}^n, y \in \{0, 1\}$
 m # training examples
 $(x^{(i)}, y^{(i)}) \rightarrow$ i^{th} training example ($i < m$)
 $M_{\text{train}}, M_{\text{test}}$

input matrix $X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \\ | & | & | & | & | \end{bmatrix}$ $x \in \mathbb{R}^{n_x \times m}$
 Training set \downarrow $y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$ $y \in \mathbb{R}^{1 \times m}$
 \downarrow $x \cdot \text{shape} = (n_x, m)$ \downarrow $y \cdot \text{shape} = (1, m)$

for this problem, given x we want \hat{y} .

where, $\hat{y} = P(y=1|x)$ [probability of a cat image for input x]

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

Output $\hat{y} = w^T x + b$ what used to be in linear Regress?

But, the output may be +ve or -ve.

So, instead, we use

$$\hat{y} = \sigma(w^T x + b) \text{ for controlling the output}$$

$$\sigma(z) = \begin{cases} 1 & \text{if } z \text{ is large} \approx \sigma(z) \approx 1 \\ \frac{1}{1+e^{-z}} & \text{if } z \text{ is small} \quad \sigma(z) \approx 0 \end{cases}$$

My Note:

[we take $w^T x$ to get a scalar product value where in we can define a line (x, y) . Though, we will use individual values of each vector to generate individual vectors. Next we combine these vectors in a matrix. In order to get \vec{x} & \vec{y} of similar size]

Loss (error) function:

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y})) \quad [\text{Why?}]$$

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

so, in order to minimize error, we need to have large $\log \hat{y}$.

In order to get

~~graph of log~~ $\log \hat{y} \neq 0$, we need to have $\hat{y} \neq 1$

$$\text{if } y=0 \quad L(\hat{y}, y) = -(\log(1-\hat{y}))$$

here, $\log(1-\hat{y})$ needs to be large

so, $(1-\hat{y})$ needs to be small

Thus, \hat{y} needs to be really small (very -ve)

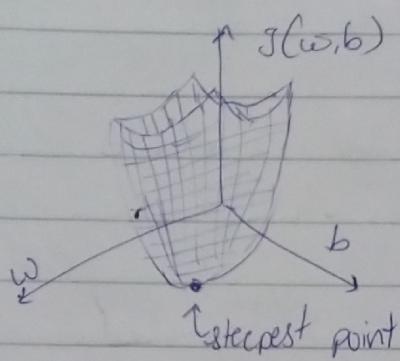
Cost function:

It indicates how good your training is doing.

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

So, we try to adjust the parameters in a way that cost function minimize

Gradient Descent



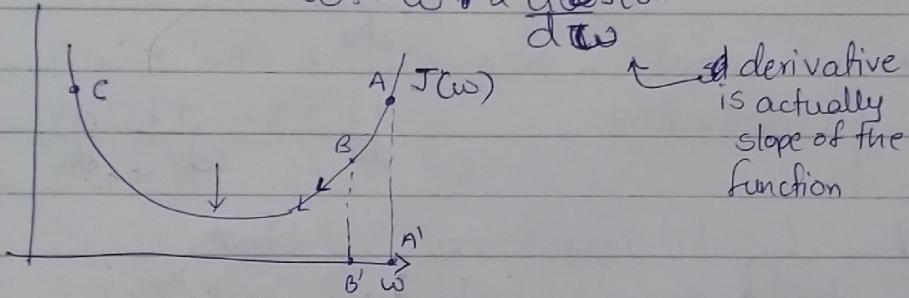
Cost function is a convex function
The gradient descent tries to minimize the cost function.

← Assume w as a real number here, for understanding

Q What does gradient descent do?

Repeat {

$$w := w - \alpha \frac{d}{dw} J(w)$$



Here, only 2-D graph is used for simplicity.

Initially value of w is A' and $J(w)$ is A . The derivative at A is positive, so, weights are updated, and we reach B .

At point C , the slope would be -ve and hence weight will be added.

Finally global minimum is achieved.

Similarly, in real scenario, gradient descent will change w and b as:

$J(w, b)$

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

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

since J is a function of w and b , instead of d , we use ∂ to denote partial derivation

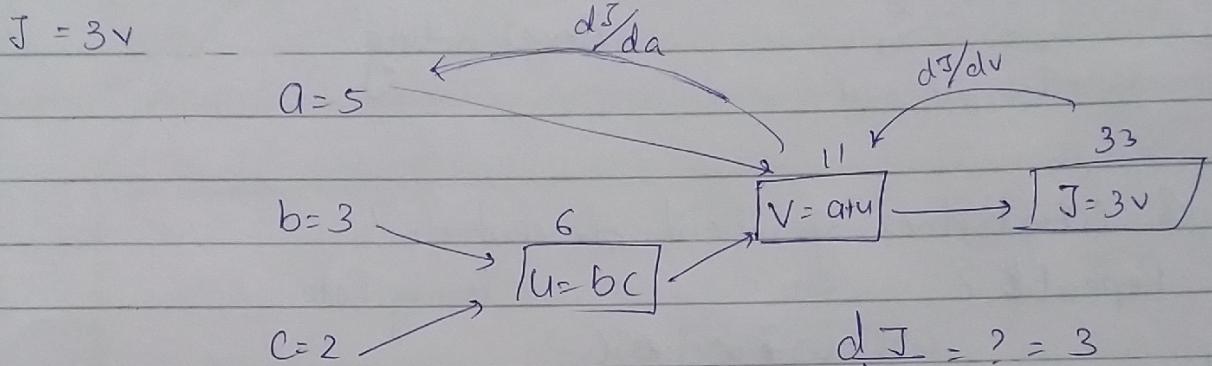
Computat'n Graph

$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



forward → get output

backward → get derivatives

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

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

chain rule

Now, we take logistic regression derivatives.

Assuming we have two variables x_1 and x_2

$$\begin{aligned}
 & x_1 \rightarrow w_1 \rightarrow Z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = a = \sigma(Z) \rightarrow L(a, y) \\
 & \frac{dL}{da} = \frac{dL}{dz} \times \frac{da}{dz} \quad (1) \\
 & \frac{dL}{da} = -\left[\frac{y}{a} - \frac{(1-y)}{1-a}\right] \quad (1)
 \end{aligned}$$

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

$$\frac{da}{dz} = \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right) = \frac{d}{dz} \left(1+e^{-z} \right)^{-1} = -\frac{d}{dz} \frac{(1+e^{-z})}{(1+e^{-z})^2}$$

$$= -\frac{(0+e^{-z})}{(1+e^{-z})^2} = \frac{e^{-z}}{(1+e^{-z})^2}$$

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

$$\begin{aligned}
 & = \frac{1}{1+e^{-z}} \times \left[\frac{(1+e^{-z})}{(1+e^{-z})} - \frac{1}{(1+e^{-z})} \right] = \sigma(z) \left[1 - \frac{1}{\sigma(z)} \right] \\
 & = a \left[1 - \frac{1}{a} \right] = (1)
 \end{aligned}$$

$\stackrel{?}{=} \text{dL}$

$$\frac{dL}{dz} = \frac{dL}{da} \times \frac{da}{dz}$$

$$= \left[\frac{-y + (1-y)}{a(1-a)} \right] a[1-a] \quad \text{from i, ii, iii}$$

$$= -y(1-a) + (1-y)a$$

$$= -y + ay + \cancel{a} - \cancel{ay} a - ay$$

$$= \cancel{ay} - a - y \quad -(iv)$$

$$\frac{dL}{dw_1} = \frac{dL}{dz} \times \frac{dz}{dw_1} \quad -(v) \quad \& \text{ similar for } w_2$$

$$\frac{dz}{dw_1} = \frac{d}{dw_1} [\omega_1 x_1 + \omega_2 x_2 + b]$$

$$\frac{db}{dz} = 1 - v^{ii}$$

$$= x_1 \quad -(vi) \quad \text{similar for } w_2$$

(vii)

$$\therefore \frac{dL}{dw_1} = x_1(a-y) \quad \& \frac{dL}{dw_2} = x_2(a-y) \quad \text{from iv, v, vi}$$

$$\frac{dL}{db} = (a-y) \quad = dz \quad \text{from v}$$

Logistic Regression on m examples

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{L(z_i)}_{\text{Eqn 2}} \underbrace{L(a^{(i)}, y^{(i)})}_{(viii)}$$

$$\Rightarrow \frac{1}{m} \sum_{i=1}^m a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^T x^{(i)} + b)$$

Here, derivative of cost function will be,

$$\frac{d}{d\omega} J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \frac{dL}{d\omega} (a^{(i)}, y^{(i)})$$

$d\omega^{(i)}$ is a function of $(x^{(i)}, y \neq \hat{y})$ (vii) from

Vectorization

It is used to reduce the for loops in a code.

$$z = \text{np.dot}(w, x) + b$$

eliminating 1st for loop

$$J=0, dw_1=0, dw_2=0, db=0 \rightarrow dw=\text{np.zeros}(n_x, 1)$$

for weights

for i=1 to m:

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

$$o^{(i)} = g(z^{(i)})$$

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

$$dz^{(i)} = o^{(i)}(1-o^{(i)})$$

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

$$dw_2 += x_2^{(i)} dz^{(i)}$$

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

$$J = J/m \quad [dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m]$$

$$dw/m$$

eliminating
2nd for loop

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

$$x.\text{shape} = n_x \times m$$

$$w^T = \begin{bmatrix} | & | & | \\ w^T x^{(1)} & w^T x^{(2)} & \dots & w^T x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

$$w^T.\text{shape} = n_w \times m$$

$$Z = [z^{(1)} z^{(2)} \dots z^{(m)}] = w^T X + [b \ b \ \dots \ b] = \left[\begin{array}{c|c|c} w^T x^{(1)} + b & w^T x^{(2)} + b & \dots & w^T x^{(m)} + b \end{array} \right] \quad | \times m$$

$$a^{(1)} = g(z^{(1)}) \quad a^{(2)} = g(z^{(2)})$$

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = g(Z)$$

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

$$\begin{aligned} dw &= \frac{1}{m} X dz^T = \frac{1}{m} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix} \\ &= \frac{1}{m} \text{np.dot}(X, dz) \end{aligned}$$

$Z = \text{np.zeros}(1, m)$,

$Z = \text{np.dot}(w.T, x) + b$

$A = \sigma(Z)$

$J/\gamma Y/\text{fnl}$

$dZ = A - Y$

$dW = \frac{1}{m} \text{np.sum}(x * dZ.T)$

$dB = \frac{1}{m} \text{np.sum}(dZ)$

$w := w - \alpha \cdot \text{np.dot}$

$b := b - \alpha \cdot dB$

Broadcasting

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

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

Tip: use $\text{np.random.randn}(5, 1)$ instead of (5) rank 1 array should not be used
use assertions $\text{assert ca.shape == (5, 1)}$

Loss function Need
if $y=1$, we need

Programming Assignment

Propagate function

$w (2, 1)$

$x (2, 3)$

$y (1, 3)$

$(1, 3) \circ (1, 3)$

$$w = \begin{bmatrix} 1 \end{bmatrix}_{2 \times 1}$$

$$x = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}_{2 \times 3}$$

$$y = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}_{1 \times 3}$$

$$dB = \text{sum}(dZ, \text{axis}=1) \quad dW = x * dZ^T$$

$$A = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}^T \quad dZ = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}^T - \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}^T$$

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

$$= 1 \times 3 + b$$