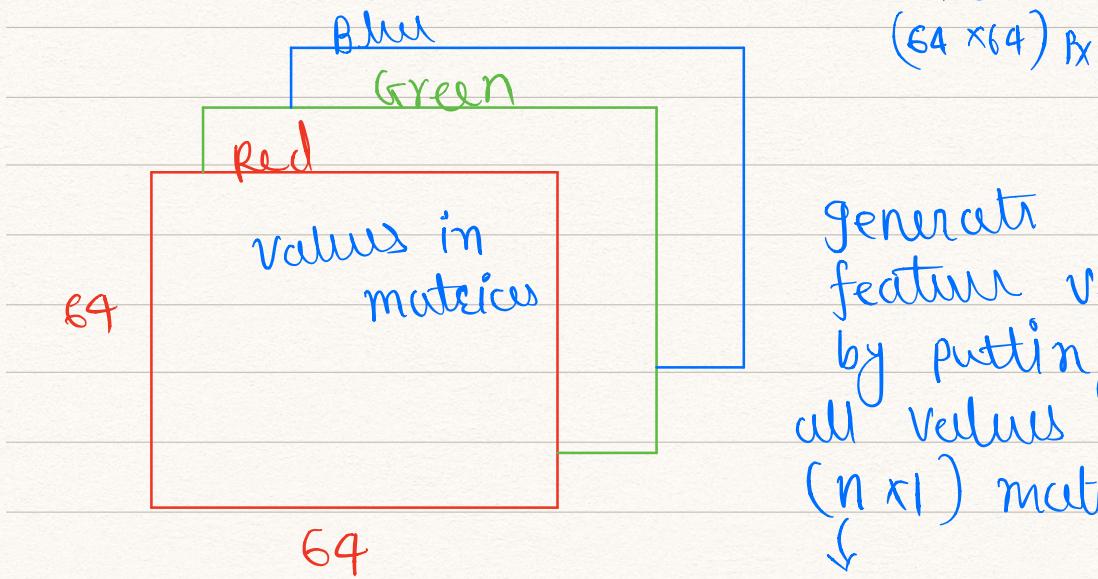


We can use for loop to go through entire training set but it will take lot of time complexity.

→ In Neural Networks we don't use for loops to go through training sets

Binary Classification :

(e.g.) Cat or No Cat in image



generate feature vector by putting all values in $(n \times 1)$ matrix
↓
 $64 \times 64 \times 3$

Here : (x, y)

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

m training Examples: $(x_1, y_1), (x_2, y_2), \dots$

$$X = \begin{bmatrix} & & & & \\ & & & & \\ x_1 & x_2 & x_3 & \dots & x_m \\ & & & & \\ & & & & \end{bmatrix} \begin{matrix} \\ \\ n_x \\ \\ \end{matrix}$$

m columns

$$Y = [y_1 \ y_2 \ \dots \ y_m]_{1 \times m}$$

Logistic Regression:

Given x , $\hat{y} = P(y=1|x)$
 $x \in \mathbb{R}^{n_x}$

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

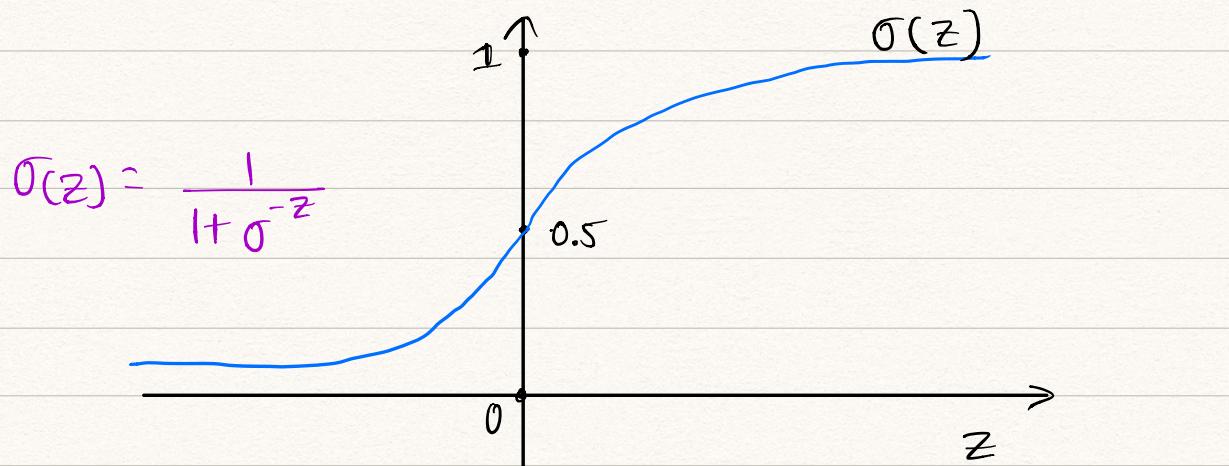
Output $\hat{y} = \frac{w^T x + b}{\text{not good}}$

can't be binary (0 or 1)

so apply **sigmoid function**

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

$\downarrow z$



In logistic Regression:

We try to find w & b so that \hat{y} become good estimate of chance of y being equal to 1.

Another Notation:

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

$$\Theta = \begin{bmatrix} \hat{y} \\ \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} = \sigma(\Theta^T x) \rightarrow b \quad w$$

Logistic Regression Cost Function:

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

optimization problem becomes non convex. might have local optimum so gradient decent might not find global optimum.

loss function

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

\downarrow
small as possible

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

\downarrow
large $\Rightarrow \hat{y} \rightarrow \text{large}$

if $y=0$

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

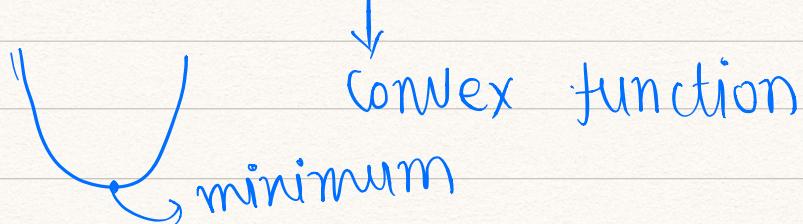
\downarrow
large $\rightarrow \hat{y} \rightarrow \underline{\text{small}}$

Cost Function:

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

Gradient Descent:

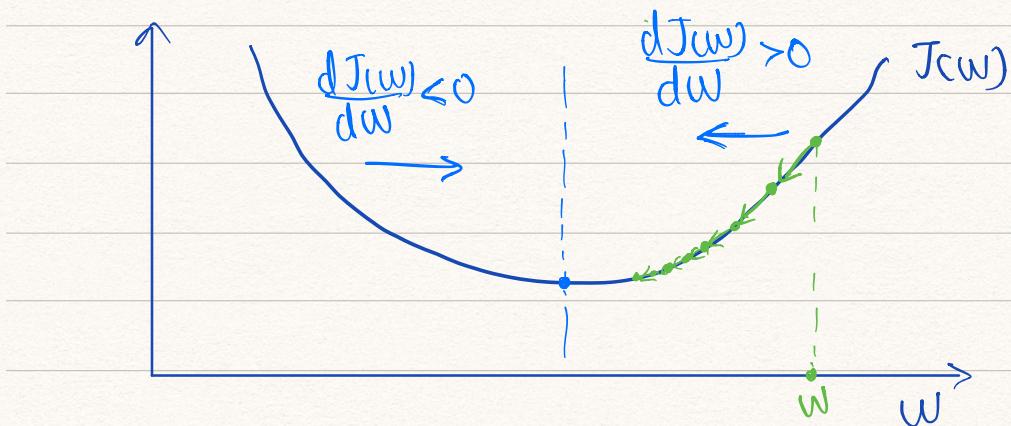
We want to find w, b for which $J(w, b)$ is minimum.



→ It will start at some initial point and then go down on the steepest descent in direction of $\nabla J(w, b)$.

→ It will converge to global minimum in case of convex function.

Gradient Descent for 1 D:



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

$\alpha \rightarrow$ learning rate

$$\frac{dJ(w)}{dw} \rightarrow \text{in code } \underline{dw}$$

Coding term: $w := w - \alpha dw$

Actual in 2D:

$$w := w - \alpha \frac{dJ(w, b)}{dw} \quad \begin{matrix} \text{partial} \\ \text{derivative} \end{matrix}$$

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

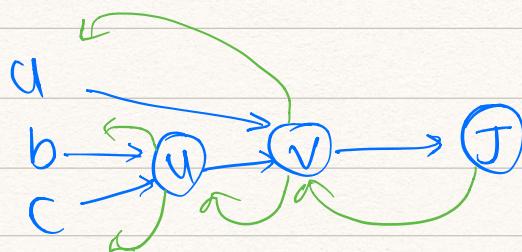
Computation Graph:

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

$$u = bc$$

$$J = 3v$$

$$v = a + u$$



$$u = bc$$

$$v = a + bc$$

$$J = 3v$$

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

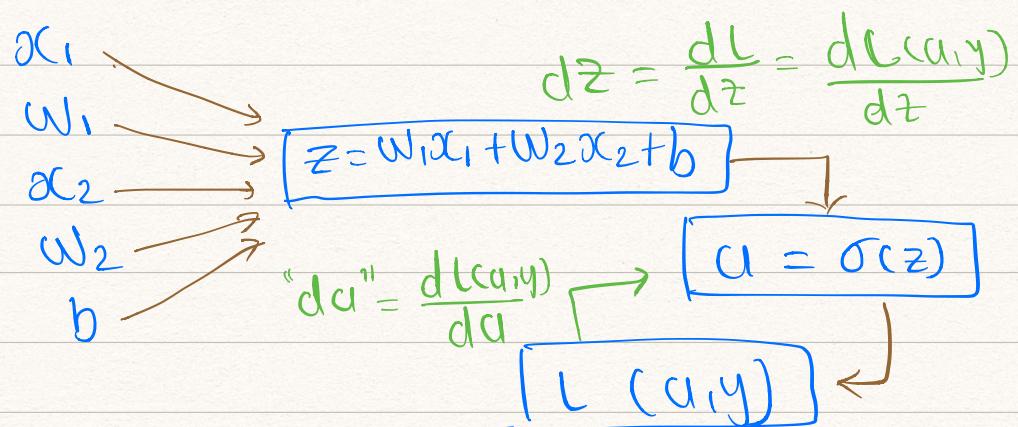
$$\frac{dJ}{da} = 3$$

Logistic Regression Gradient Descent:

$$z = \omega^T x + b$$

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

computational graph:



$$da = \left(\frac{-y}{a}\right) + \left(\frac{1-y}{1-a}\right)$$

$$dz = \frac{dL(a,y)}{dz} = \frac{dL}{da} \cdot \frac{da}{dz}$$

$$\frac{dL}{dw_1} = dw_1 = x_1 dz \quad dw_2 = x_2 dz$$

$$db = dz$$

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

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

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

Logistic Regression on m examples:

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

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

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

$$dw_i^{(i)} \rightarrow (x^{(i)}, y^{(i)})$$

Algorithm:

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

$$\left. \begin{array}{l} \text{for } i=1 \text{ to } m \\ \quad z^{(i)} = w^T x^{(i)} + b \\ \quad a^{(i)} = \sigma(z^{(i)}) \\ \quad J = -[y^{(i)} \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)})] \end{array} \right\}$$

for loop

for loop

$$\left\{ \begin{array}{l} dz^{(i)} = a^{(i)} - y^{(i)} \\ dw_1 + = x_1^{(i)} dz^{(i)} \\ dw_2 + = x_2^{(i)} dz^{(i)} \\ db + = dz^{(i)} \end{array} \right.$$

$\uparrow n=2$

$J | = m$

$dw_1 | = m$

$dw_2 | = m$

$db | = m$

accumulator

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

$$\begin{aligned} w_1 &:= w_1 - \alpha dw_1 \\ w_2 &:= w_2 - \alpha dw_2 \\ b &:= b - \alpha db \end{aligned}$$

Vectorization: get rid of for loop.

Logistic Regression Cost Function (Adv.)

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

$$y=0 \quad P(y|x) = \hat{y}$$

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

function with above boundary

conditions.

$$\log p(y|x) = y \underbrace{\log(\hat{y}) + (1-y)\log(1-\hat{y})}_{L(y, \hat{y})}$$

Cost Function:

$$P(\text{labeled in training set}) = \prod_{i=1}^n p(y^{(i)}|x^{(i)})$$

We have to maximize it.

Put log & maximize it.

$$\log P(\cdot) = \sum_{i=1}^n \underbrace{\log p(y^{(i)}|x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})}$$

Maximum likelihood Estimate:

choose parameter that
maximize $\log P(\cdot)$

$$= -\sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)})$$

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

\downarrow so get rid of negative sign.

$\frac{1}{m} \rightarrow$ for better scale

By minimizing cost function, we are carrying out max. likely hood estimation over logistic Reg. model under the assumption that our training examples are IID (Identical, Randomly distri.)