Week 2

Binary classification

As the name suggests, this classification gives an output in binary form,i.e., 1 or 0.
This is achieved using feature vectors and classifier.

Feature Vector - It is an n-dimensional matrix containing all pixel values from the matrices of the image with red, green and blue pixels. Here, n is the number of pixel values in total. Each matrix (red, green or blue) has same size as the image.
Eg. Image of size 64*64 will give a feature vector with size (64*64*3, m), i.e., (n, m).
Here, n = no of total reshaped or un-rolled pixel values. [Total =r+g+b]
 m = Training data size

Training set :  matrix denoted by Capital X that contains total 'm' feature vectors;

X is the training set where x1, x2,... xm are feature vectors
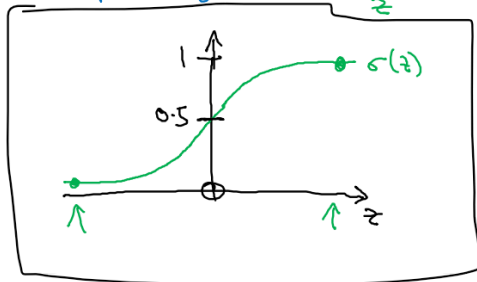Y is also introduced as stacking data in columns is conventional for implementing neural networks

## Logistic Regression

It is an algorithm that is used when we want output labells to be 0 or 1,i.e., binary

*Parameters* of the algorithm are W and b; W is n-dimensional vector and b is a real number.

Parameters : $\boxed{w} \in \mathbb{R}^{n_x}$ , $\boxed{b} \in \mathbb{R}$.

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

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

If $z$ large $\sigma(z) \approx \frac{1}{1+0} = 1$

If $z$ large negative number

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+ \text{Bignum}} \approx 0 .$$

y hat is output, and formula for is given in the image in terms of w and b ;Wt is w transpose
here, sigmoid function is used as it helps to get an output between 0 and 1. Note y is different for each training example

## Loss function

this function judges a models performance.

$\boxed{\text{Loss}}$ (error) function: $\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)$

$\boxed{\mathcal{L}(\hat{y}, y)} = - \left( \boxed{y \log \hat{y}} + (1-y) \log(1-\hat{y}) \right)$

If $\boxed{y=1}$ : $\mathcal{L}(\hat{y}, y) = - \log \hat{y}$ ← Want $\log \hat{y}$ large, want $\hat{y}$ large.

If $\boxed{y=0}$ : $\mathcal{L}(\hat{y}, y) = \frac{1}{\uparrow} \log(1-\hat{y})$ ← Want $\log 1-\hat{y}$ large .... want $\hat{y}$ small

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

The 1st formula besides the title in the image is actually a squared function [2 got erased]
This formula isnt feasible. The next formula below the title is used.

It computes error for single training eg., for entire training set, cost function is used.

*Cost function* is the avg of loss functions of the entire training set.

$$z = w^T x + b$$
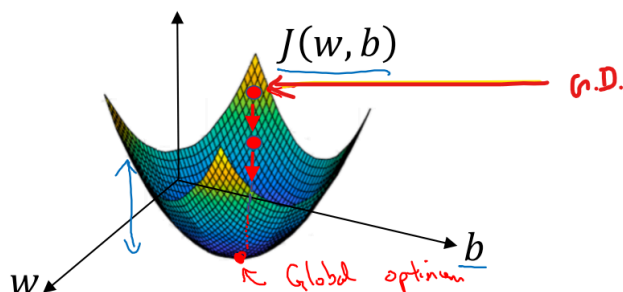$$\hat{y} = a = \sigma(z)$$
$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

## Cost function [ J(w,b) ]

## Gradient Descent

Recap: $\hat{y} = \sigma(w^T x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$

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

Want to find $w, b$ that minimize $J(w, b)$
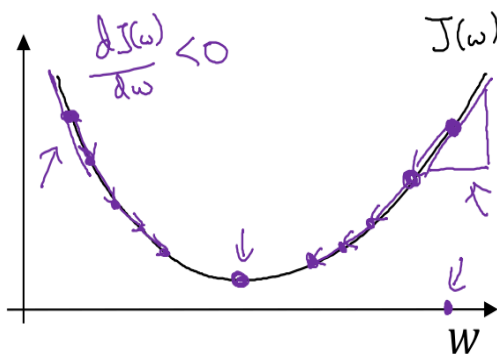


$J(w, b)$

G.D.

$b$

$w$

Global optimum

The formula given above the graph is of cost function, its avg of loss function

The graph represents a convex function whose height equals the cost function and its Global optimum gives values of w and b that help to minimise the J(w,b)
The red coloured point that is drawn down is due to iterations of gradient descent
G.d. is used to reduce the initial points of w and b to the global optimum

## Gradient Descent

$$\frac{dJ(w)}{dw} < 0$$

$J(w)$

$w$

update w

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

learning rate

"dw"

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

$$\frac{dJ(w)}{dw} = ?$$

$J(w,b)$

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

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

$$\frac{\partial J(w,b)}{\partial w}$$

$$\frac{\partial J(w,b)}{\partial b}$$
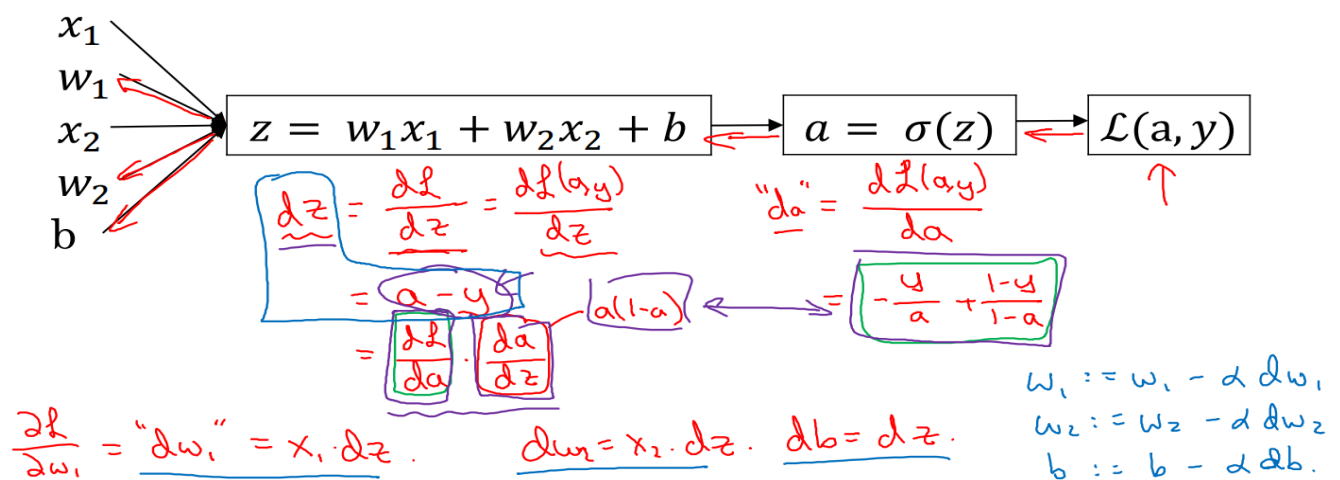
$\partial$

$\partial \leftarrow J$

"partial derivative"

$dw$

$db$

We also need to update w;  w updated w is denoted by  "w:"
the variable in code fore the derivative term
The Variable in code for the derivative present in the w formula is dw
dw  or dj/dw is the slope [ height/width = j/w ]  which if positive initially makes gradient descent to decrease.
J is a function of w and b as per the original formula, here, if J is supposed to be a function of b too  for that we use partial derivative of J

# Computational graphs

Helps to find the derivative of the final output variable at every step, with respect to some intermediate variable.

## Computing derivatives

$a = 5$  $\frac{dJ}{da}$  "da" $= 3$

$b = 3$

$c = 2$

$u = bc$ (6)

$v = a + u$ (11)

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

$J = 3v$ (33)

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

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

$3 \times 1$

$\frac{dv}{da} = 1$

$a \to v \to J$

$J = 3v$
$v = 11 \to 11.001$
$J = 33 \to 33.003$

$a = 5 \to 5.001$
$\to v = 11 \to 11.001$
$J = 33 \to 33.003$

$f(a) = 3a$
$\frac{df(a)}{da} = \frac{df}{da} = 3$
$J = 3v$
$\frac{dJ}{dv} = 3$

$\frac{d \text{ Final Output Var}}{d \text{ var}}$   $\frac{dJ d\text{var}}{}$  "dvar"

d var is used as variable in code to represent the derivatives of final output var with respect to various intermediate quantities
Suppose 2 parameters are present, that means 2 vals of x and w, and 1 val of b [ b is constt ]

Below is the computation graph for loss function

$x_1$
$w_1$
$x_2$
$w_2$
$b$

$z = w_1 x_1 + w_2 x_2 + b$   $a = \sigma(z)$   $\mathcal{L}(a, y)$

$dz = \frac{d\ell}{dz} = \frac{d\ell(a,y)}{dz}$

"da" $= \frac{d\mathcal{L}(a,y)}{da}$

$= a - y$

$= \frac{d\ell}{da} \cdot \frac{da}{dz}$

$a(1-a)$

$= -\frac{y}{a} + \frac{1-y}{1-a}$

$\frac{d\ell}{dw_1} = $ "dw₁" $= x_1 \cdot dz$.

$dw_2 = x_2 \cdot dz$.  $db = dz$.

$w_1 := w_1 - \alpha \, dw_1$
$w_2 := w_2 - \alpha \, dw_2$
$b := b - \alpha \, db$.

Remember derivative of loss wrt  z is a-y
The derivatives of loss remain constt, like formulae do.

In the below pic it is shown which variables vary for m no. of training
examples [ subscript i associated all variables ]

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

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

$$(x^{(i)}, y^{(i)})$$

$$dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$$

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

$$dw_1^{(i)} - (x^{(i)}, y^{(i)})$$

The above figure gives the formula of loss for m  training egs. and
its derivative for some parameter w1

Below we initialise J, w and b to 0 and the following figure shows
how to implement the formulas for the m training examples

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

$\rightarrow$ For $i = 1$ to $m$

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

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

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

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

$dw_3$
$\vdots$
$dw_n$

$\left[ \begin{array}{l} dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)} \\ dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \\ db \mathrel{+}= dz^{(i)} \end{array} \right]$  $n = 2$

$$J {/}{=} m \;\Leftarrow$$

$$dw_1 {/}{=} m \; ; \; dw_2 {/}{=} m \; ; \; db {/}{=} m. \;\Leftarrow$$

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

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

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

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

Vectorization

this is applicable only works for small datasets.

In the above fig, 2 for loops run -
1 for the Training egs
Another for the parameters w1 and w2

Even though in neural networks we use *vectorisation* instead of for loops

After W(s)are found, we update them
W(s) are accumulative variables which's why there are  no subscripts with them
derivative dz_i helps calculate the loss function for each training eg

**Vectorisation**

Colab nbk for vectorisation function, diff bw this function and normal for loop [execution time]

co **Untitled14.ipynb**

## Logistic regression derivatives

$J = 0$, $dw1 = 0$, $dw2 = 0$, $db = 0$     $dw = np.zeros((n_x, 1))$

→ for $i = 1$ to m:

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

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

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

for $j=1..n_x$
$\quad dz^{(i)} = a^{(i)} - y^{(i)}$

for $dw_j += ?$
$\quad dw_1 += x_1^{(i)} dz^{(i)}$   $n_x = 2$     $dw += x^{(i)} dz^{(i)}$

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

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

$J = J/m$, $dw_1 = dw_1/m$, $dw_2 = dw_2/m$, $db = db/m$

$\quad\quad\quad\quad\quad dw \mathrel{/}= m.$

**The above** fig is for calculating derivatives by vectorisation
The original version was by non-vec, we then cancelled the areas, where np functions and other vectorisation tactics can be used
Modifications-
Initialisations replaced with np.zeros
The for loop for $j=1..n_x$ and $w[j] += x\_1[i]*dz[i]$
**is replaced by** $dw += x[i]*dz[i]$
So, now only 1 for loop present [ was 2 earlier ]
At the end, $dw/=m$ replaces $dw\_1/=m$ and dw2 too

# Vectorizing Logistic Regression

$$z^{(1)} = w^T x^{(1)} + b$$
$$a^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = w^T x^{(2)} + b$$
$$a^{(2)} = \sigma(z^{(2)})$$

$$z^{(3)} = w^T x^{(3)} + b$$
$$a^{(3)} = \sigma(z^{(3)})$$

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

$(n_x, m)$

$\mathbb{R}^{n_x \times m}$

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

$$Z = \begin{bmatrix} z^{(1)} & z^{(1)} & \cdots & z^{(m)} \end{bmatrix} = w^T X + \begin{bmatrix} b & b & \cdots & b \end{bmatrix}$$

$1 \times m$

$$= \begin{bmatrix} w^T x^{(1)} + b \end{bmatrix} \begin{bmatrix} w^T x^{(1)} + b \end{bmatrix} \cdots w^T x^{(m)} + b$$

$1 \times m$

$$Z = np.dot(w.T, x) + b$$

$(1,1)$   $\mathbb{R}$   "Broadcasting"

$$A = \begin{bmatrix} a^{(1)} & a^{(1)} & \cdots & a^{(m)} \end{bmatrix} = \sigma(Z)$$

Instead of operating individually on z1, z2,z3 ,lly,for a(s)
We use the Z that equals dot product of wT and x and adds b

Refer the pdf directly for rest - codes and syntax