

Week 2

* Binary Classification.

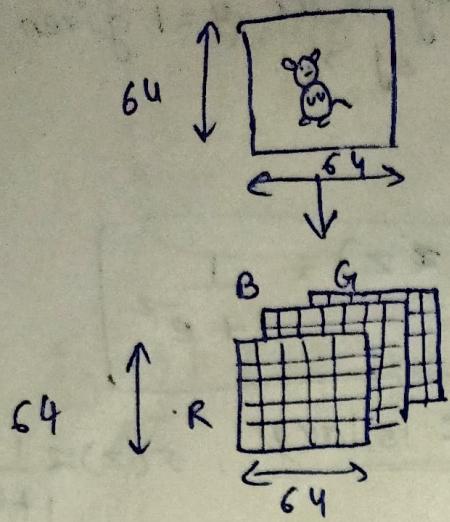
↳ logistic Regression algorithm.

Eg.

An image of a cat

output

1 (cat) vs 0 (non cat).



\rightarrow input feature vector x

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

data input

m train examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$m_{\text{train}}, m_{\text{test}} = \text{no. of test examples}$

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \\ | & | & | & & | \end{bmatrix} \quad \begin{matrix} n_x \\ \downarrow \\ m \end{matrix} \quad \underline{\underline{X \in \mathbb{R}^{n_x \times m}}}$$

$$Y = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(m)}] \Rightarrow \underline{\underline{Y \in \mathbb{R}^{1 \times m}}}$$

$$\left[\begin{array}{c} R_1 \\ R_2 \\ \vdots \\ R_{4096} \\ B_1 \\ B_2 \\ \vdots \\ B_{4096} \\ G_1 \\ G_2 \\ \vdots \\ G_{4096} \end{array} \right] \quad \left. \right\} 12288$$

$= n_x$

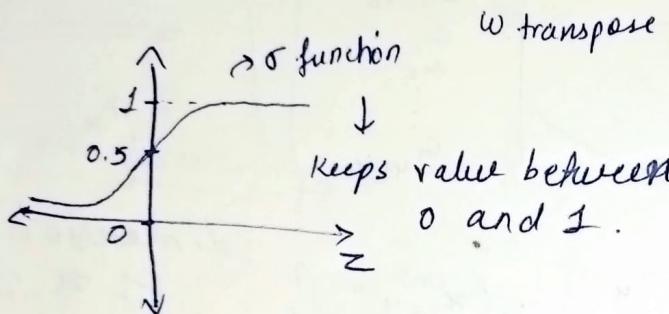
dimension of x .

• Logistic Regression.

Given x , $\hat{y} = P(y=1|x) \rightarrow$ Probability that $y=1$ given the value of x .
 ↓
 output $0 \leq \hat{y} \leq 1$
 $x \in \mathbb{R}^{n_x}$

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

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



w = weights (how much importance do we give to a particular feature)

b = bias (decides the threshold of the neuron)

W = no. of features \times no. of examples

$#_b$ = no. of neurons

• Logistic Regression cost function.

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

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

If z is large negative no., $\sigma(z) = 1 \approx 0$

Basically we need to find the values of parameters w and b such that \hat{y} will give an accurate prediction.

probability of y being equal to 1 -

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

Given $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$, we want

$$\hat{y}^{(i)} \approx y^{(i)}$$

Loss (error) function :- $L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$

If $y=1$: $L(\hat{y}, y) = -\log \hat{y} \leftarrow$ we want $\log \hat{y}$ large

\downarrow

should be small

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

$(0 \leq \hat{y} \leq 1) \Rightarrow \hat{y} \text{ should be closer to } 1$

If $y=0$: $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ $\log(1-\hat{y})$ should be large.

$1-\hat{y} \rightarrow \text{large}$

$\hat{y} \rightarrow \text{small}$

$\hat{y} \approx 0$

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

average loss over

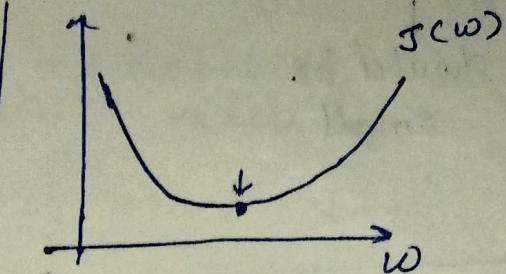
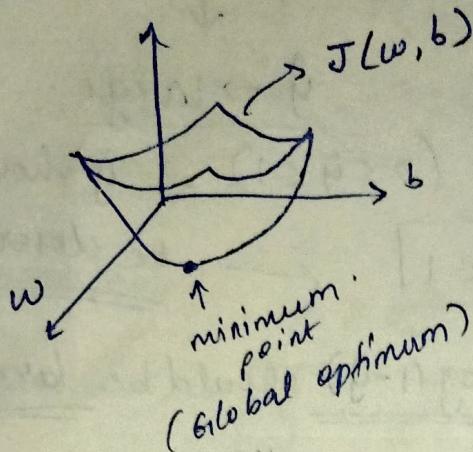
all the data examples

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

(whereas, Loss function is defined over a single example.)

task → We want to find parameters w, b for which the cost function will be minimum.

* Gradient Descent .



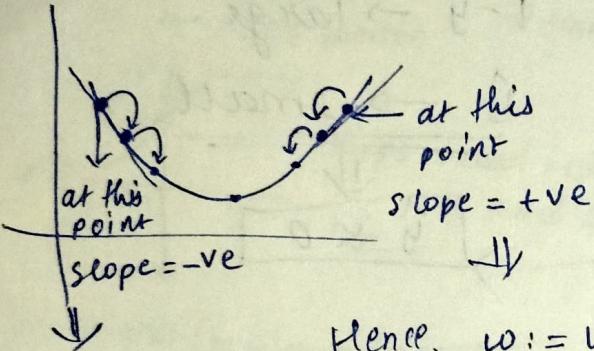
Repeat {

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

}

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

↑
slope



Hence, $w := w - \alpha dw \rightarrow$ we are subtracting from w .

eventually reaching the minimum point

Hence,

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

$$w := w + \alpha dw$$

we are adding to w . eventually reaching the minimum point.

for $J(w, b)$

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

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

Gradient descent tells us what changes to which weights matters the most or results in fastest decrease in cost function.

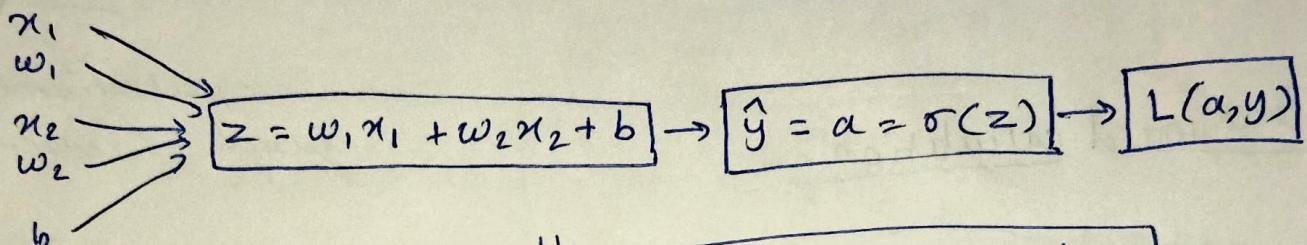
* Logistic Regression Gradient Descent:

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

$$\hat{y} = \sigma(z) = a$$

$$L(a, y) = - (y \log a + (1-y) \log(1-a))$$

e.g.:



$$\begin{aligned} w_1 &:= w_1 - \alpha \frac{\partial L}{\partial w_1} \\ w_2 &:= w_2 - \alpha \frac{\partial L}{\partial w_2} \\ b &:= b - \alpha \frac{\partial L}{\partial b} \end{aligned}$$

$$\frac{\partial L}{\partial w_1} = dw_1 = x_1 \cdot dz$$

$$\frac{\partial L}{\partial w_2} = dw_2 = x_2 \cdot dz$$

$$\frac{\partial L}{\partial b} = db = dz$$

$$dz = a - y$$

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

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

$$\frac{da}{dz} = a(1-a)$$

• Logistic regression on 'm' examples.

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

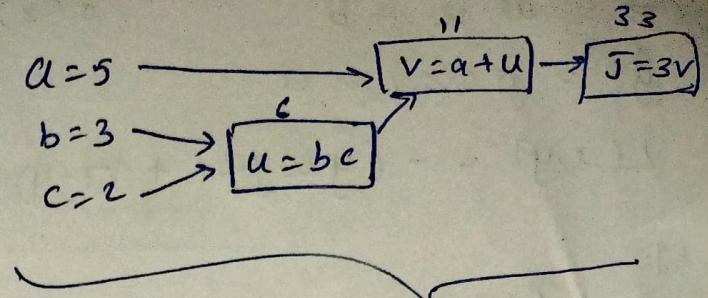
$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

* Computation graph.

Forward :-

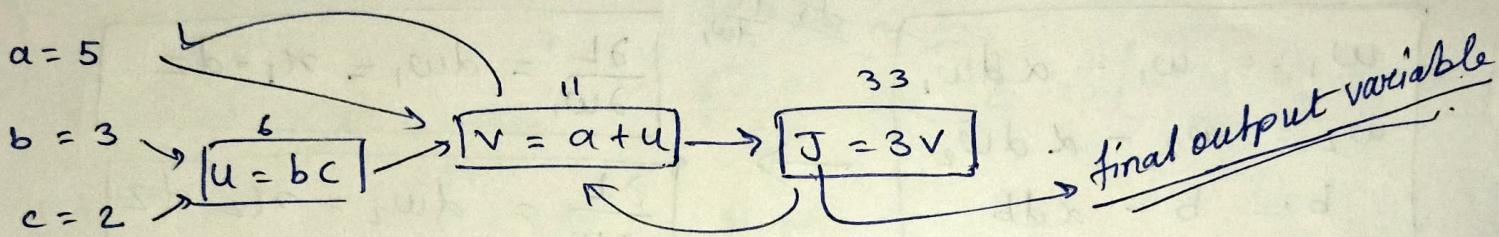
ex:- $J(a, b, c) = 3(a + \frac{bc}{u}) = 33$

$$\begin{array}{l} u = bc \\ v = a + u \\ J = 3v \end{array}$$



Computational graph (forward)

Backward calculation :-



$$\frac{dJ}{dv} = 3 \rightarrow \begin{aligned} J &= 3v \\ v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned} \quad \text{(increase in } J = 3 \text{ times increase in } v)$$

$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{da} \rightarrow \begin{aligned} a &= 5 \rightarrow 5.001 \\ v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned}$$

{ a → v → J }

Usually in code, we are interested in calculating the derivative of final output variable.

$$\frac{d \text{ Final Output Var}}{d \text{ var}}$$

→ this is too long.

so, we could replace it with just "dvar", which is derivative of final output variable with respect to 'var'.

* Vectorization

$$z = w^T x + b.$$

$$w = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad x = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad w, x \in \mathbb{R}^{n_x}$$

non-vectorized:

$$z = 0$$

for i in range (n_x):

$$z += w[i] * x[i]$$

$$z += b$$

vectorized:

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

$$w^T x.$$

This version is almost 3000 times faster! ~~in some cases~~.

e.g.: we want to apply exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

non-vectorized:

$$u = \text{np. zeros}((n, 1))$$

\rightarrow zero matrix

for i in range (n):

$$u[i] = \text{math.exp}(v[i])$$

vectorized:

import numpy as np

$$u = \text{np. exp}(v)$$

(Easy!)

Hence, in vectorized version,

$$\boxed{dw_1 = 0, dw_2 = 0} \times \rightarrow dw = \text{np.zeros}((n_m, 1))$$

and

$$\boxed{dw_1 = dw/m, dw_2 = dw/m} \times \rightarrow dw/m$$

for $(x^{(i)}, y^{(i)}) \rightarrow dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$

$$\Rightarrow \boxed{\frac{\partial J(\omega, b)}{\partial \omega_i} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L(a^{(i)}, y^{(i)})}{\partial \omega_i}}$$

$\underbrace{dw_i^{(i)}}$

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

for $i=1$ to m :

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

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

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

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

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

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

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

has many loops.
some large datasets
might have a larger
Then this method gets
inefficient.

$$J / m$$

$$dw_1 / m$$

$$dw_2 / m$$

$$db / m$$

Hence,

$$\omega_1 := \omega_1 - \alpha \frac{dw_1}{m}$$

$$\omega_2 := \omega_2 - \alpha \frac{dw_2}{m}$$

$$b := b - \alpha \frac{db}{m}$$

} 1 step of
gradient
descent.

To get rid of inefficient loops, we use Vectorization.

* Vectorizing Logistic Regression

→ Let's say we have n data examples.

$$z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b \quad z^{(3)} = w^T x^{(3)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad a^{(3)} = \sigma(z^{(3)}) \quad \dots$$

$$X = \left[\begin{array}{c|c|c|c} 1 & | & x^{(1)} & | \\ \hline 1 & | & x^{(2)} & | \\ \hline 1 & | & \dots & | \\ \hline 1 & | & x^{(m)} & | \end{array} \right] \quad (n_x \times m) \quad \mid \quad w = \left[\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right] \quad \{ n_x \quad w^T = \left[\begin{array}{c} \dots \\ \hline \dots \\ \hline \dots \\ \hline n_x \end{array} \right]$$

$$[z^{(1)} \ z^{(2)} \ z^{(3)} \ \dots \ z^{(m)}] = w^T X + [b \ b \ b \ \dots \ b] \quad \downarrow \quad \underbrace{1 \times m}_{1 \times m}$$

$$w^T \quad 1 \times n_x \quad \left[\begin{array}{c|c|c|c} 1 & | & x^{(1)} & | \\ \hline 1 & | & x^{(2)} & | \\ \hline 1 & | & \dots & | \\ \hline 1 & | & x^{(m)} & | \end{array} \right]_{n_x \times m} = \left[\begin{array}{c} \dots \\ \hline \dots \\ \hline \dots \\ \hline 1 \times m \end{array} \right]$$

$$\therefore \underbrace{[z^{(1)} \ z^{(2)} \ z^{(3)} \ \dots \ z^{(m)}]}_Z = \underbrace{[w^T x^{(1)} + b]}_{z_1} \quad \underbrace{[w^T x^{(2)} + b]}_{z_2} \quad \dots \quad \underbrace{[w^T x^{(m)} + b]}_{z_m}$$

↳ This can be done using:- ~~np.zdot~~

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

↓
Transpose
↓
a real number

→ here, python automatically adds b to all numbers in the $w^T x^T$ matrix.

↑
This is called broadcasting

* Vectorizing Logistic Regression

$$dz^{(1)} = \hat{a}^{(1)} - y^{(1)}, dz^{(2)} = \hat{a}^{(2)} - y^{(2)} \dots$$

$$dz = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]$$

$$A = [\hat{a}^{(1)} \ \hat{a}^{(2)} \ \dots \ \hat{a}^{(m)}], Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$\rightarrow dz = A - Y = [\hat{a}^{(1)} - y^{(1)} \ \hat{a}^{(2)} - y^{(2)} \ \dots \ \hat{a}^{(m)} - y^{(m)}]$$

$$\frac{dw}{dz} = \frac{\sum_{i=1}^m x^{(i)} \cdot dz^{(i)}}{m}$$

$$\frac{db}{dz} = \frac{\sum_{i=1}^m dz^{(i)}}{m}$$

$$\rightarrow db = [np \cdot \text{sum}(dz)]/m$$

$$\rightarrow dw = \frac{1}{m} \cdot X \cdot dz^T$$

$$= \frac{1}{m} \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$dw = \frac{1}{m} \begin{bmatrix} x^1 dz^1 + x^2 dz^2 + \dots + x^m dz^m \\ \vdots \end{bmatrix}_{n_x \times 1}$$

~~By hand~~

Updated : Implementing Logistic Regression:

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

$$= np \cdot \text{dot}(w \cdot T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{m} X \cdot dZ^T$$

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

$$\text{Gradient Descent : - } \left. \begin{array}{l} w := w - \alpha dw \\ b := b - \alpha db \end{array} \right\} \begin{array}{l} \text{we might need a for loop} \\ \text{here tho.} \end{array}$$

Broadcasting in Python

7

	Apples	Beef	Eggs	Potatoes
carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
fat	1.8	135.0	99.0	0.9

We want to find % of
calories from carb,
protein, fat.

$$\text{cal} = A \cdot \text{sum}(\text{axis}=0)$$

percentage = $100 * \frac{A / (\text{cal.reshape}(1,4))}{\text{cal}}$

horizontal
vertical addition
(axis=1 → horizontal)
 Broadcasting

this is not really necessary
since $A \cdot \text{sum}(\text{axis}=0)$ already
made the cal → (1,4).

Broadcasting example:-

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(2,3)} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{(1,3) \rightarrow (2,3)} \rightarrow \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(2,3)} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}_{(2,1) \rightarrow (2,3)} \overline{\begin{bmatrix} 100 & 100 \\ 200 & 200 \end{bmatrix}} \rightarrow \begin{bmatrix} 1001 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$