

Basics of neural network programming

Week 2 : Course 1

Binary classification

eg: we input an image from which we need to identify if it contains a cat (1) or not cat (0)

Any image is stored in a computer as three separate matrices corresponding to R G B channels of that particular image. eg: if the input image is 64×64 pixels then there will be three, 64×64 pixel matrices which will correspond to the pixel intensity values.

Then these are fed into an input feature vector x .

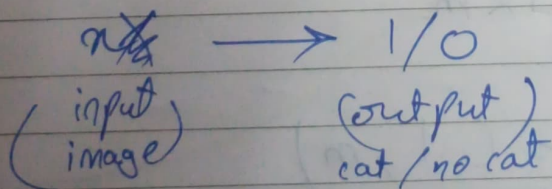
$$\text{i.e. } x = \begin{bmatrix} 255 \\ 302 \\ \vdots \\ 101 \\ 105 \\ \vdots \\ 209 \\ 212 \\ \vdots \end{bmatrix}$$

$$64 \times 64 \times 3 = 12288$$

$$\therefore n = n_x = 12288$$

↳ dimension of input feature x
(for brevity use n)

In short,



Notation:

data (n, y) $n \in \mathbb{R}^n$, $y \in \{0, 1\}$

A training example is represented by a pair (n, y) where, n is an n -dimensional feature vector and y is either 1 or 0.

The training set consists of ~~is~~ m total examples.
 \therefore The training sets will be written as ~~$(n_1, y_1), (n_2, y_2)$~~
 ~~$\dots, (n_m, y_m)$~~

$$\{(n^{(1)}, y^{(1)}), (n^{(2)}, y^{(2)}), \dots, (n^{(m)}, y^{(m)})\}$$

$$m = M_{\text{train}}, \quad M_{\text{test}} = \# \text{ test examples}$$

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

$\xleftarrow{\quad m \quad} \xrightarrow{\quad}$

$$\text{i.e. } X \in \mathbb{R}^{n_x \times m}$$

$$X.\text{shape} = (n_x, m)$$

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

$$\text{i.e. } Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

Logistic Regression

→ Used when output is either 1 or 0 i.e. in case of binary classification.

Given input feature vector x , we want an algorithm that outputs a prediction (\hat{y}) i.e. probability of $y=1$ given the input feature x .

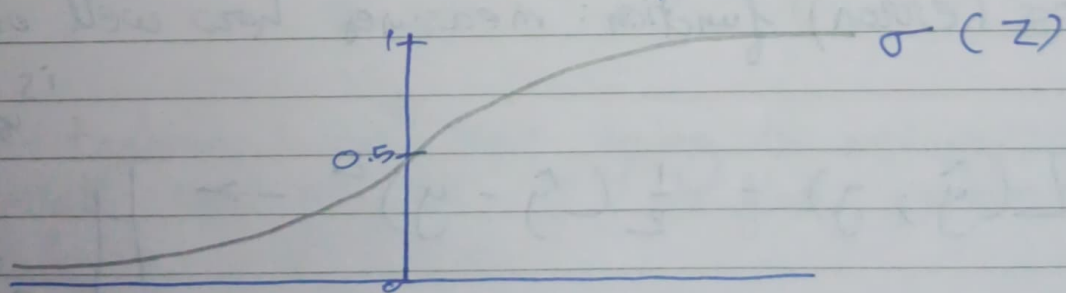
$$\hat{y} = P(y=1|x) \quad 0 \leq \hat{y} \leq 1$$

$$x \in \mathbb{R}^{n_x}$$

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

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

Only considering " $w^T x + b$ " would not make sense to probability as it would result in a very big number or negative number. Hence, we pass it to sigmoid function which will take care that our output is between 0 & 1.



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

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

If z is large negative number,

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{Big no}} \approx 0$$

'b' corresponds to inter-spectrum.

To train $w \in b$, we need to define cost f.

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

$$\text{where, } z^{(i)} = w^T x^{(i)} + b$$

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

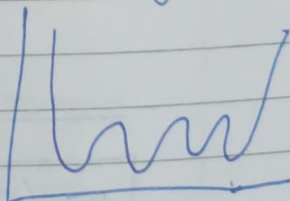
we want $\hat{y}^{(i)} \approx y^{(i)}$

* $x^{(i)}, y^{(i)}, z^{(i)}$ denote training data's i^{th} example

Loss (error) function: measures how well our algorithm is doing on single example

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

Local minima is not well defined.



convex function \rightarrow This would be preferred.

Hence we do not use the above function.

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

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

we want $\log \hat{y}$ large i.e. \hat{y} to be large (as close to 1)

If $y=0$: $L(\hat{y}, y) = -\log (1-\hat{y})$

we want $\log (1-\hat{y})$ large i.e. $1-\hat{y}$ to be large $\therefore \hat{y}$ close to 0.

Cost function: how well the algorithm is doing on entire training set.

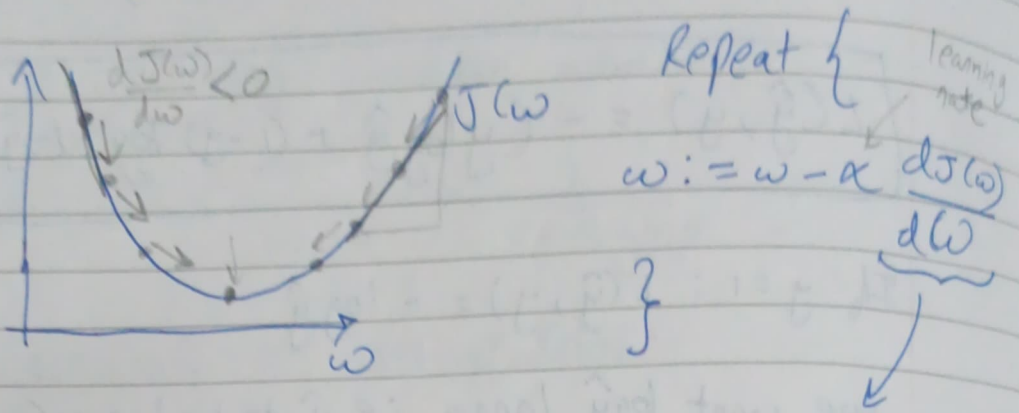
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

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

In training, we are going to minimize ~~to~~ overall cost function J to find parameters w & b .

Gradient descent algorithm is used for it.

Gradient Descent:



this derivate term in code is written as "dw."

$J(w, b)$

$$\begin{aligned}
 & \text{update } w := w - \alpha \frac{\partial J(w, b)}{\partial w} \\
 & b := b - \alpha \frac{\partial J(w, b)}{\partial b}
 \end{aligned}$$

actual formulae

partial derivative

written as "db" in code

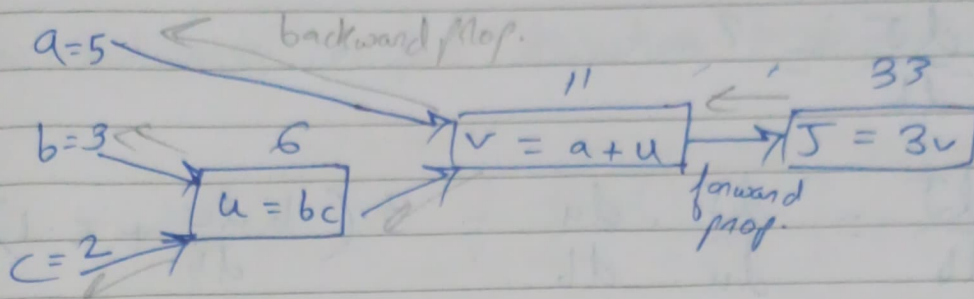
Computation graph:

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

Diagram showing the computation graph structure:

- $u = bc$
- $v = a + u$
- $J = 3v$

$$\begin{aligned}
 u &= bc \\
 v &= a + u \\
 J &= 3v
 \end{aligned}$$



Note: One step of backward propagation yields the derivative of final output

~~JS~~ Computing derivatives:

(prev example)

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

$$J = 3v$$

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

3 1

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

$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{da} \quad \dots \text{(chain rule)}$$

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

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

→ "d var" in code

Similarly, $\frac{dJ}{dv} \rightarrow \text{"dv"} = 3$ in code

$$\text{"da"} = 3$$

$$\frac{dJ}{db} = ?$$

$$\begin{aligned} \frac{dJ}{db} &= \frac{dJ}{du} \cdot \frac{du}{db} = 3 \left(b \frac{dc}{db} + c \cdot 1 \right) \\ &= 3 (0 + (2)(1)) \\ &= 3 \times 2 \\ &= 6 \end{aligned}$$

$$\therefore \text{"} db \text{"} = 6$$

$$\text{|| by } \frac{dJ}{dc} = dc = 9 \quad \left\{ \begin{array}{l} \text{since } b=3 \\ \therefore 3 \times 3 = 9 \end{array} \right.$$

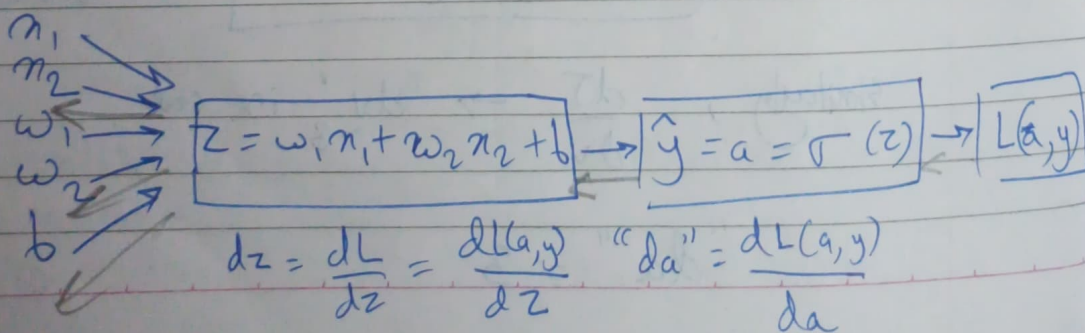
* In ~~linear~~ ^{logistic} regression, to compute final output go from left to right and to find derivatives go from right to left.

Gradient descent for logistic regression:

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

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

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



Go backwards i.e. find derivatives to decrease loss ℓ .

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

$$dz = \frac{dL}{dz} = \frac{dL}{da} \frac{da}{dz} = \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \cdot \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right)$$

$$= \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \cdot \left(\frac{+e^{-z}}{(1+e^{-z})^2} \right)$$

$$= \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \left(\frac{\frac{1}{a} - 1}{\frac{1}{a^2}} \right)$$

$$= \frac{-y(1-a) + a(1-y)}{a^2(1-a)} \left(\frac{1-a}{a} \right)^2$$

$$= -y + ay + 1 - 2y$$

$$= a - y$$

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

$$\frac{\partial L}{\partial \omega_1} = \text{"}\partial \omega_1\text{"} = \eta_1 \cdot dz$$

$$\frac{\partial L}{\partial \omega_2} = \text{"}\partial \omega_2\text{"} = \eta_2 \cdot dz, \quad db = dz$$

Then perform these updates:

$$\omega_1 := \omega_1 - \alpha d\omega_1$$

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

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

Logistic regression for m examples:

Vectorization: the art of getting rid of explicit for loops in code

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

$$\omega = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$\omega \in \mathbb{R}^{n \times 1} \quad x \in \mathbb{R}^{n \times 1}$$

Non-vectorized:

```

z = 0
for i in range(n):
    z += w[i] * x[i]
z += b
z += b
    
```

vectorized:

```

z = np.dot(w, x) + b
            $\omega^T x$ 
import numpy as np then
execute above code.
    
```

* Vectorized code also runs faster hence takes less time to execute.

examples of vectorization:

* Whenever possible, avoid explicit for-loops in logistic regression.

eg: compute $u = Av$

Non-vectorized:

$$u_i = \sum_j A_{ij} v_j$$

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

for i ...
for j ...

$$u[i] += A[i][j] * v[j]$$

Vectorized:

$$u = \text{np.dot}(A, v)$$

eg: $v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$, compute $u = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$

Non-vectorized:

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

for i in range(n):

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

// initializing zero vector of n-dimensions

Vectorized:

import numpy as np

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

Vectorizing Logistic Regression:

$$z^{(1)} = \omega^T n^{(1)} + b$$

$$a^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = \omega^T n^{(2)} + b$$

$$a^{(2)} = \sigma(z^{(2)})$$

$$z^{(3)} = \omega^T n^{(3)} + b$$

$$a^{(3)} = \sigma(z^{(3)})$$

$$X = \begin{bmatrix} | & | & & | \\ n^{(1)} & n^{(2)} & \dots & n^{(m)} \\ | & | & & | \end{bmatrix} \quad (n_m, m)$$

$$R^{n_m \times m}$$

$\omega^T \rightarrow$ Row vector
according to
rules of matrix
multiplication
 $\rightarrow (1 \times m)$

$$Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = \omega^T X + [b \ b \ \dots \ b]$$

$$= [\omega^T n^{(1)} + b \ \omega^T n^{(2)} + b \ \dots \ \omega^T n^{(m)} + b]$$

python code: $zZ = \text{np.dot}(\omega.T, X) + b$

When we import the lib. numpy as np, we can use the np.dot to multiply matrices and use T to find transpose of matrix. Also, note that b is actually a $1 \times m$ vector (m -dimensional row vector) which is written as a real number. But python is smart enough to convert it to a vector when added to $\omega^T X$. This is called "Broadcasting".

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

Vectorizing LR's gradient computation:

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

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

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

$$dZ = A - Y = [a^{(1)} - y^{(1)} \quad a^{(2)} - y^{(2)} \quad \dots]$$

$$dw = 0$$

$$dw_+ = n^{(1)} dz^{(1)}$$

$$dw_+ = n^{(2)} dz^{(2)}$$

$$\vdots$$

$$dw_- = m$$

$$db = 0$$

$$db_+ = dz^{(1)}$$

$$db_+ = dz^{(2)}$$

$$\vdots$$

$$db_+ = dz^{(m)}$$

$$db_- = m$$

Vectorized:

$$\text{formulae: } dw = \frac{1}{m} X dZ^T$$

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

$$= \frac{1}{m} np. \text{dot}(X, dZ.T) = \frac{1}{m} np. \text{sum}(dZ)$$

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

$$= \frac{1}{m} [n^{(1)} dz^{(1)} + \dots + n^{(m)} dz^{(m)}] \quad n \times 1$$

Implementing logistic Regression using vectorization

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

for $i = 1$ to m :

$$z^{(i)} = w^T x^{(i)} + b \quad \text{--- (1)}$$

$$a^{(i)} = \sigma(z^{(i)}) \quad \text{--- (2)}$$

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

$$dz^{(i)} = a^{(i)} - y^{(i)} \quad \text{--- (3)}$$

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

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

$$db += dz^{(i)} \quad \text{--- (5)}$$

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

$$db = db/m$$

Vectorization:

$$\textcircled{1} \quad Z = w^T X + b$$

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

$$\textcircled{2} \quad A = \sigma(Z)$$

$$\textcircled{3} \quad dz = A - Y$$

$$\textcircled{4} \quad dw = \frac{1}{m} X \cdot dz^T$$

$$\textcircled{5} \quad db = \frac{1}{m} \text{np.sum}(dz)$$

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

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

Broadcasting example:

$$\text{ex 1)} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \Rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

$$\text{ex 2)} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = ?$$

2×3 1×3

⇒

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

2×3 2×3

$$\text{ex 3)} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = ?$$

2×3 2×1

⇒ 2×3

$$\Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

General principle:

$$\begin{array}{ccc} (m, n) & + & (1, n) \rightsquigarrow (m, n) \\ \text{matrix} & * & \text{matrix} \\ & / & (m, 1) \rightsquigarrow (m, n) \\ & & \text{matrix} \end{array}$$

$$(m, 1) + R$$

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

Python/numpy vectors

$a = \text{np.random.randn}(5)$

→ rank 1 array

(neither a row vector nor column vector)

$a.\text{shape} = (5,)$

Don't use to
implement logistic
regression

$a = \text{np.random.randn}(5, 1) \rightarrow a.\text{shape} = (5, 1)$ ✓
∴ column vector

$a = \text{np.random.randn}(1, 5) \rightarrow a.\text{shape} = (1, 5)$ ✓
∴ row vector