

Week #2: Logistic Regression as a Neural Network

→ Logistic regression is an algorithm for binary classification.

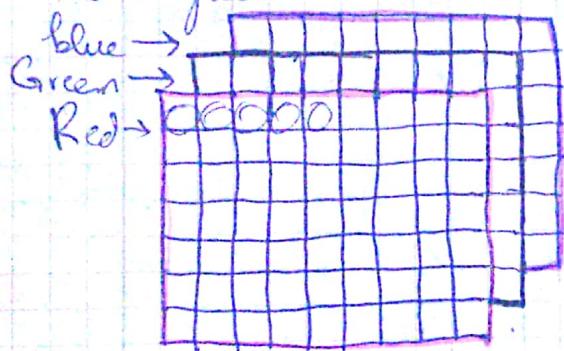
Let's set up an example:-

Suppose we have an image and we want our output to be 1 or 0.

where 1 denotes the label we intended is true and 0 denotes otherwise.

* An image is comprised of 3 matrices in the computer

red, green and blue.



Take the matrix values and append all of them into a vector.

long vector:-

$$X = \begin{bmatrix} \text{Red} \\ \text{Green} \\ \text{Blue} \end{bmatrix}$$

Suppose we have picture of 64×64 size, vector will be

$$64 \times 64 \times 3 \times 1$$

$$12288 \times 1$$

Goal:-

We want to train a model which can take this long feature vector X and tell us whether our image represents certain intended object or not.

Notations:-

One training example :- $(x, y) \in \mathbb{R}^n, y \in \{0, 1\}$

Training set comprises of m training examples:-

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

MF Mat

$N_{\text{test}} = \# \text{ test examples}$

- Finally to put all of training examples into a more compact notation, we are going to define a matrix X , as defined by taking the inputs and stacking them in columns.

$$X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(n)} \\ | & | & | \end{bmatrix} \quad \begin{matrix} \uparrow n_x \\ \downarrow \\ \leftarrow m \rightarrow \end{matrix}$$
$$X \in \mathbb{R}^{n_x \times m} \quad X.\text{shape} = n_x \times m$$

- To make implementation of neural networks easier we are gonna stack j 's too

$$y = [j^{(1)} \ j^{(2)} \ \dots \ j^{(m)}] \quad j \in \mathbb{R}^{1 \times m} \quad \Rightarrow y.\text{shape} = (1, m)$$

Logistic Regression:

- This algorithm is a supervised learning algorithm where in the output vector either all labels are zero or 1.

→ Given an image and then stacking all the values in a vector:

5. Give an input vector \vec{x} to an image,

$$\boxed{y = P(y=1 | \vec{x})} \quad \vec{x} \in \mathbb{R}^m$$

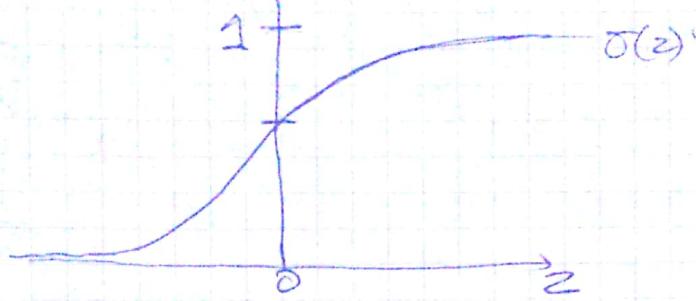
- Parameters of logistic regression:

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

So given a vector \vec{x} and parameters w and b
we want to know the output \hat{y} .

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

$$y \in \{0 \leq j \leq 1\}$$



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

If z is very large then e^z will be close to zero so $\sigma(z) \approx \frac{1}{1+0} = 1$.

If z is very small or it is very large negative number then e^z becomes very large and $\sigma(z)$ becomes very less.

When we implement the logistic regression, our job is to learn parameters w and b such that \hat{y} becomes a good estimate.

Cost Function:-

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

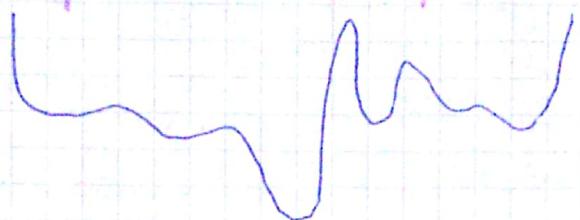
Given $\{(x^1, y^1), \dots, (x^m, y^m)\}$, and $y^{(i)} \approx \hat{y}^{(i)}$

Now let's see how can we define a loss(error) function to measure, how well our algorithm is doing.

One thing we can do is define loss by our predicted label \hat{y} and true label y as half the square error

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

- we could use this, but in logistic regression people don't normally use this, because in logistic regression our optimization problem becomes non-convex and we end up with multiple local optima.



- and in this case, gradient descent may not find the global optimum.

So in logistic regression we define another loss function that plays a similar role as squared function but gives us an optimization problem that is convex.

following is the loss function we use.

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

If $y=1$: $L(\hat{y}, y) = -\log \hat{y}$ \leftarrow want $\log \hat{y}$ large, want \hat{y} large
minimum loss shows that \hat{y} should be as ~~large~~
large as possible and clearly it can't be larger than 1 so
in other words we want \hat{y} as close as to 1.

If $y=0$ $L(\hat{y}, y) = -\log(1-\hat{y})$ \leftarrow want $\log(1-\hat{y})$ to be large
and thus want \hat{y} as small as possible.

→ The loss function was defined w.r.t one training example, it measures, how well you are doing on a single training example.

Now let's define Cost function which measures how well you are doing on the entire training examples.

Cost functions

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

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

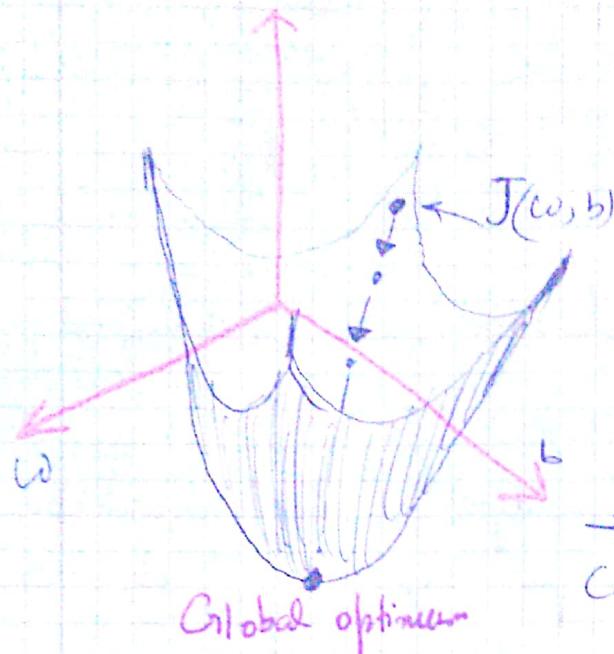
Loss function is applied to single training example and the cost function is the cost of your parameters, so we want to find the parameters which minimizes the cost function J .

Gradient Descent: Let's talk about, how can we use gradient descent on our algorithm to learn the parameters ω and b

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

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

We want to find w, b that minimizes $J(w, b)$.

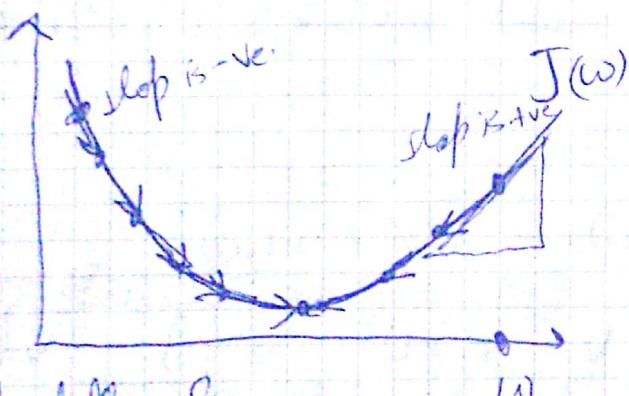


⇒ In practice w can be of higher dimension but for the sake of this example let's take them as one scalar.

→ Cost function J is a convex function.

What we do, is initialize w and b to some random value like red blue dot on the graph after some iterations we converge to our global optimum.

Illustration:



Repeatedly {

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

let's say, we have a function that we want to minimize

Update on parameters.

until the algorithm converges

$\alpha \Rightarrow$ learning rate.

$$J(w, b)$$

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

partial derivative.

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

Computational graph:-

→ Computations of a neural networks are organized in terms of a forward pass or a forward propagation step, in which we compute the output of the Neural Network followed by the backward pass or backward propagation step, which we use to compute gradients.

Computational graph explains that, why it is organized this way.

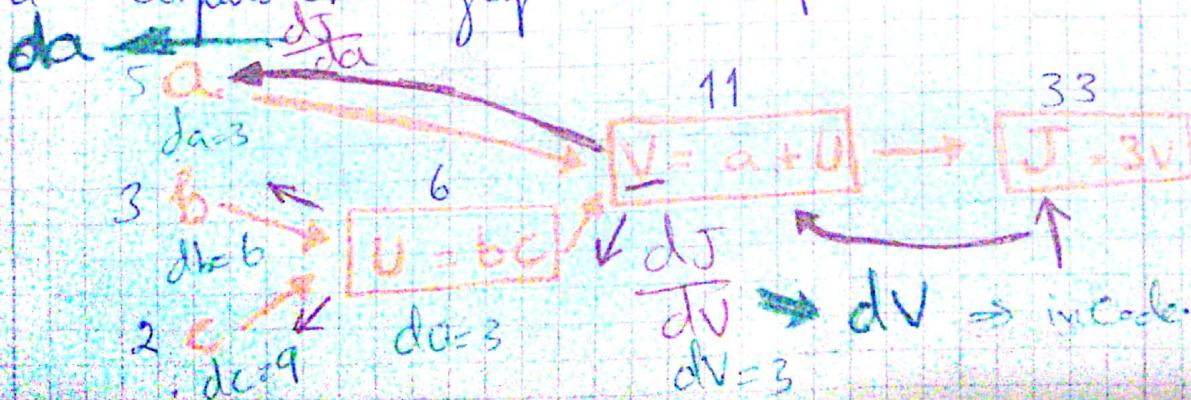
→ Let's use a simple example of logistic regression or a full blown neural network to illustrate Computational graph.

e.g

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

$$\textcircled{1} U = bc, \textcircled{2} V = a + U, \textcircled{3} J = 3V$$

we can take these 3 steps and draw them in a computational graph as follows:-



So, the computation graph comes in handy when there are some distinguished or some special output variables such as J in this case that needs to be optimized, and in case of logistic regression, J is of course the cost function that we are trying to minimize.

Derivatives with computation graphs:-

→ Keeping in consideration the last computation graph and let's say we want to compute derivative of J with respect to $V \frac{dJ}{dV}$, so that means if we were to change the value of V , how would it affect J .

$$J = 3V ; V = 11 \rightarrow 11.001$$

$$\frac{dJ}{dV} = 3 \quad J = 33 \rightarrow 33.003$$

Now what is $\frac{dJ}{da}$

$$a = 5 \rightarrow 5.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

→ by changing $a \frac{dV}{da}$, you end up changing V and change in the value of V will cause the value of J to also change. In calculus this is called the Chain Rule.

→ if a affects V which in turn affects J

then the amount that J changes is the product of how much V changes when you nudge a and how much J changes when you nudge V .

$$\frac{dJ}{da} = \frac{dJ}{dV} \cdot \frac{dV}{da}$$

$$\text{What is } \frac{dJ}{du} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{du}$$

$$U = 6 \rightarrow 6.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} = 6$$

$$b = 3 \rightarrow 3.001$$

$$U = b \cdot C = 6 \rightarrow 6.002 \quad C = 2$$

$$J = 33.006$$

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 9$$

$$3 \cdot 3$$

Logistic Regression Gradient Descent:-

Logistic Regression recap:-

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

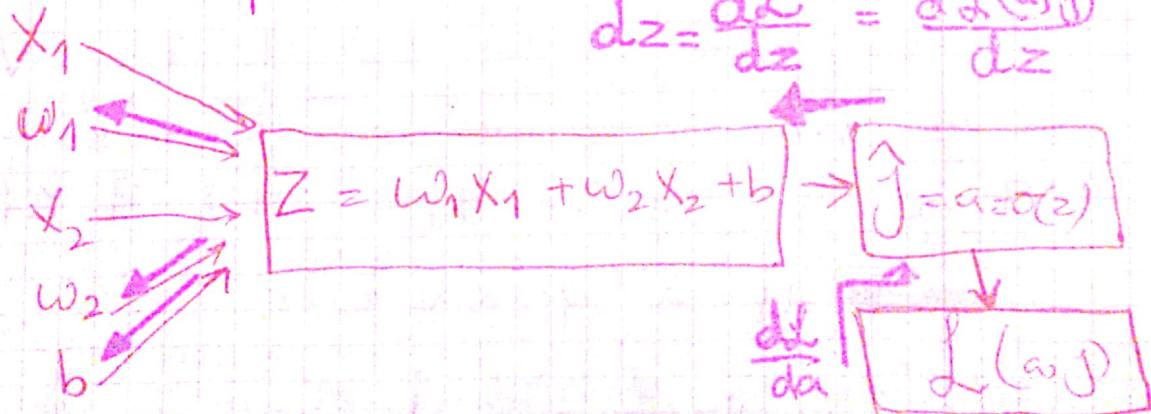
$$\hat{y} = a = \sigma(z) \cdot \text{output}$$

Loss function:

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

For the sake of this example let's say we have two features x_1 and x_2 .

So the output will be:



→ So, our job is to modify the parameters w_1 , w_2 and b in order to reduce the loss function $L(a, j)$.

→ First thing you do is compute the derivative of loss " $\frac{dL(a, j)}{da}$ " or da .

So let's do that :-

$$\frac{dL(a, j)}{da} = \frac{d}{da} [j \log(a) + (1-j) \log(1-a)].$$

$$da = -\frac{j}{a} + \frac{1-j}{1-a}$$

2nd step $\frac{dL}{dz} = \frac{dL}{da} \cdot \frac{da}{dz} \Rightarrow a(1-a)$

$$a = o(z) = \frac{1}{1+e^z} \Rightarrow \frac{da}{dz} = ? \frac{e^z}{(1+e^z)^2}$$

$$\frac{da}{dz} = \frac{1}{1+e^z} \cdot \left[1 - \frac{1}{1+e^z} \right] \Rightarrow a(1-a)$$

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

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

$$= -j + ja + a - a^2 j = a - j$$

$$\boxed{\frac{dL}{dz} = a - j}$$

Final Step :- go back to check how much do you need to change w_1, w_2, b .

$$\frac{\partial L_1}{\partial w_1} = dw_1 = x_1 \cdot dz, \quad dw_2 = x_2 \cdot dz, \quad db = dz$$

Now we perform update

$$w_1 = w_1 - \alpha d w_1$$

$$w_2 = w_2 - \alpha d w_2$$

$$b = b - \alpha d b$$

Logistic Regression with m examples

Cost function:

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

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

$$\frac{\partial}{\partial w_j} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_j} (L(a^{(i)}, y^{(i)})) \quad dw_1, dw_2, b$$

Algorithm

Let's make all of this compact into algorithm.

→ Start with $d w_1 = 0$; $d w_2 = 0$; $db = 0$

$$d w_i = \frac{\partial J}{\partial w_i} \quad \text{for } i=1 \text{ to } m \quad \star$$

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

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

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

$$d z^{(i)} = a^{(i)} - \hat{y}^{(i)}$$

$$d w_1 += x_1^{(i)} d z^{(i)} \quad \begin{matrix} \uparrow \\ \text{assuming there are} \\ \text{two features} \end{matrix}$$

$$d w_2 += x_2^{(i)} d z^{(i)} \quad \begin{matrix} \uparrow \\ \text{that's why 2 w's} \end{matrix}$$

$$d b += d z^{(i)}$$

$$J /= m \quad d w_1 /= m \quad d w_2 /= m \quad d b /= m$$

Now comes the update of gradient descent-

$$w_1 := w_1 - \alpha \frac{\partial J}{\partial w_1}$$

$$w_2 := w_2 - \alpha \frac{\partial J}{\partial w_2}$$

$$b := b - \alpha \frac{\partial J}{\partial b}$$

- Now problem with this algorithm is that we need two for loops to accomplish the task.
 - one is the main loop over all the training examples
 - other is for the update of all w 's i.e. parameters of logistic regression.

Vectorization is used to avoid all this expensive computation, let's talk about that.

→ What is Vectorization?

Vectorization is an art of getting rid of for loops from your code.

→ In Logistic regression we compute $Z = w^T x + b$

$$w = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}; x = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix} \quad w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x} \text{ Vectors.}$$

Non-Vectorized:

$z = 0$
for i in range(n_x):

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

$$z += b$$

Vectorized:

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

$$w^T x$$

→ More vectorization examples:

* Whenever possible avoid explicit for loops

$$\vec{U} = A \vec{V}$$

$$U_i = \sum_j A_{ij} V_j$$

Vectorized:

$$U = np.dot(A, V)$$

$$U = np.zeros(n, 1)$$

for $i=1 \dots n$

for $j=1 \dots m$

$$U[i] += A[i][j] * V[j]$$

another example:

Say you need to apply the exponential operation on every element of a matrix/vector.

$$\vec{V} = \begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$$

$$\rightarrow U = \begin{bmatrix} e^{V_1} \\ \vdots \\ e^{V_n} \end{bmatrix}$$

vectorized implementation
import numpy as np
 $U = np.exp(V)$

Let's see how can we get rid of one of the two for loops from our example. Or algorithm we developed.

→ Second loop can be avoided:

Instead of defining d_{w_1}, \dots, d_{w_n} as individual variables, introduce them as a vector \vec{dw} and instead individual updates, update them by using vector operations.

$$\vec{dw} := X^T \vec{dz}$$

$$dw / \vec{m}$$

Vectorizing logistic Regression

- Let's first examine the forward propagation step of logistic Regression.

Then from m training examples you need to do the following for 1st example -

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

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

for second then

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

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

⋮

$$\underline{z}^{(m)} = w^T x^{(m)} + b$$

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

X is a feature matrix for all the training examples we have.

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

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

$$= [w^T x^{(1)} + b \quad w^T x^{(2)} + b \quad \dots \quad w^T x^{(m)} + b]$$

how python sees it:

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

Now next thing is to compute $\hat{a}^{(1)}, \hat{a}^{(2)}, \dots, \hat{a}^{(m)}$

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

Now b is a real number, but remember in this condition python expands it and it becomes like $[b \ b \ \dots \ b]$

This is called Broadcasting -

Vectorizing Logistic Regression (Gradient Output):

We computed $d_2^{(1)} = \hat{a} - y^{(1)}$, $d_2^{(2)} = \hat{a} - y^{(2)}$...

To vectorize, we define a new variable

$$\vec{d}_2 = [d_2^{(1)} \ d_2^{(2)} \ \dots \ d_2^{(m)}]_{1 \times m}$$

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

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

$$d_2 = A - Y$$

In the previous implementation:-

$$\left. \begin{array}{l} dw = 0 \\ dw_s += X^{(1)} d_2^{(1)} \\ dw_o += X^{(2)} d_2^{(2)} \\ dw_b / m \end{array} \right\} \quad \left. \begin{array}{l} db = 0 \\ db += d_2^{(1)} \\ db += d_2^{(2)} \\ db / m \end{array} \right\}$$

Vectorized Implementation:-

$$db = \frac{1}{m} \sum_{i=1}^m d_2^{(i)}$$
$$= \frac{1}{m} np \text{ sum}(d_2)$$

$$dw = \frac{1}{m} X d_2^T$$
$$= \frac{1}{m} \left[\begin{matrix} X^{(1)} & \dots & X^{(m)} \end{matrix} \right] \begin{bmatrix} d_2^{(1)} \\ \vdots \\ d_2^{(m)} \end{bmatrix}$$

$$dw = \frac{1}{m} [X^T d_2^{(1)} + \dots + X^T d_2^{(m)}]$$

Summing all this up and let's condense up the vectorized version of algorithm.

for i in range(100):

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

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

$$A = \sigma(Z)$$

$$dZ = A \cdot Y$$

$$dw = \frac{1}{m} X dZ$$

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

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

Single iteration of Gradient descent

Broadcasting in Python:-

Let's motivate broadcasting by an example:-

→ Calories from Carbs, Proteins, Fats, in 100g of different foods

	Apple	Beef	Eggs	Potatoes
Carbs	56.0	0.0	44.0	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

Goal is to calculate % of Calories from Carb, protein and Fat. Can you do this with explicit for loop.

⇒ Cal = A.sum(axis=0)

Print (Cal) [591. 239. 155.4, 79.1]

Per = 100 * A / Cal.reshape(1, 4)

Print (Per)

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m, n)} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{(1, n)} \rightarrow (m, n)$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}$$

if $(m, n) + (1, n) \rightarrow (m, n)$ then operation element wise.

If conversely:

$$(m, n) + (m, 1) \rightarrow (m, n)$$
 then elementwise operation

Note on Python / Numpy codes

$$a = np.random.rand(5, 1) \Rightarrow \text{Preferred}$$

$$a = np.random.rand(5) \left. \begin{array}{l} a.shape = (5,) \\ \text{rank 1 array} \end{array} \right\} \text{dots use}$$

Logistic Regression Cost function:-

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

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

$$\text{if } y=1 : P(y| x) = \hat{y}$$

$$\text{if } y=0 : P(y| x) = 1 - \hat{y}$$

Let's summarize the above 2 equations into a single one

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

we can write the probability as follows, do why?

Let's check out.

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

2nd Case.

if $y=0$ $P(y|x) = (\hat{y})^0 \cdot (1-\hat{y})^{(0-y)} = 1-\hat{y}$

because log function is strictly monotonically increasing function we are maximizing $P(y|x)$

$$\log P(y|x) = \log[\hat{y} \cdot (1-\hat{y})^{(1-y)}]$$

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

$$\boxed{\log[P(y|x)] = y \log(\hat{y}) + (1-y) \log(1-\hat{y})}$$

and it is equal to the negative of the loss function.

This is equal to the negative of the loss function because in the logistic regression we are trying to minimize the loss by increasing the probability of doing so.

$$\text{So } -L(\hat{y}, y) \downarrow = \log[P(y|x)] \uparrow$$

Cost on m examples:-

$$\log P(\text{labels in training set}) = \log \prod_{i=1}^m P(y^i|x^i)$$

$$= \sum_{i=1}^m \log(P(y^i|x^i))$$

Maximum likelihood estimate

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

$\log P(\text{labels in training set})$

$$= -\sum_{i=1}^m L(g_i^v, \hat{y}).$$

Maximum
likelihood
estimation.

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

minimize