

# **Deep Learning Notes**

**By Divyajeet Pala**

## Logistic Regression

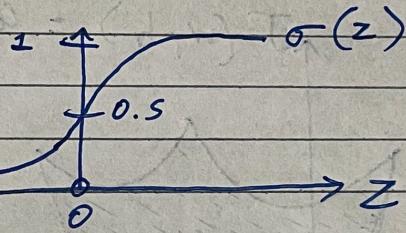
→ Input :  $x$ ,  $x \in \mathbb{R}^{n_x}$  i.e.,  $x$  is an  $n_x$  dimensional vector

Parameters :  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$   
 $n_x$  dim. vector      a real number

Output :  $\hat{y} = P(y=1 | x)$  i.e., probability  $y=1$  given  $x$

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

→ Sigmoid function



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

→ Training examples :  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ ,  
 want  $\hat{y}^{(i)} \approx y^{(i)}$ .

Note / superscript is in rounded bracket (i) denotes the  $i^{th}$  training example out of the total  $m$  examples.

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

convex ↑  
func.

→ Measures how well you're doing on a single training example

→ Cost function (measures performance w.r.t. whole dataset):

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

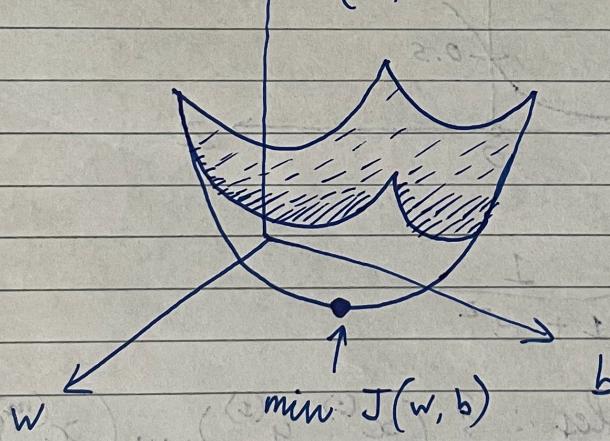
Note the -ve sign →

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

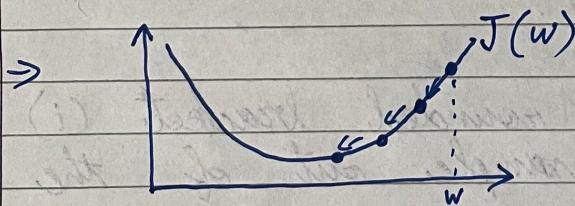
→ Gradient descent

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

→  $J(w, b)$  = convex function



3-D representation



Repeat {  
 $w := w - \alpha \frac{\partial J(w)}{\partial w}$   
 learning rate  
 slope ↑

[ := denotes updating the value]

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

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

→ In coding,  $\frac{\partial J}{\partial w}$  is assigned to 'dw' variable  
 &  $\frac{\partial J}{\partial b}$  will be 'db' variable.

→ Imp formulas:

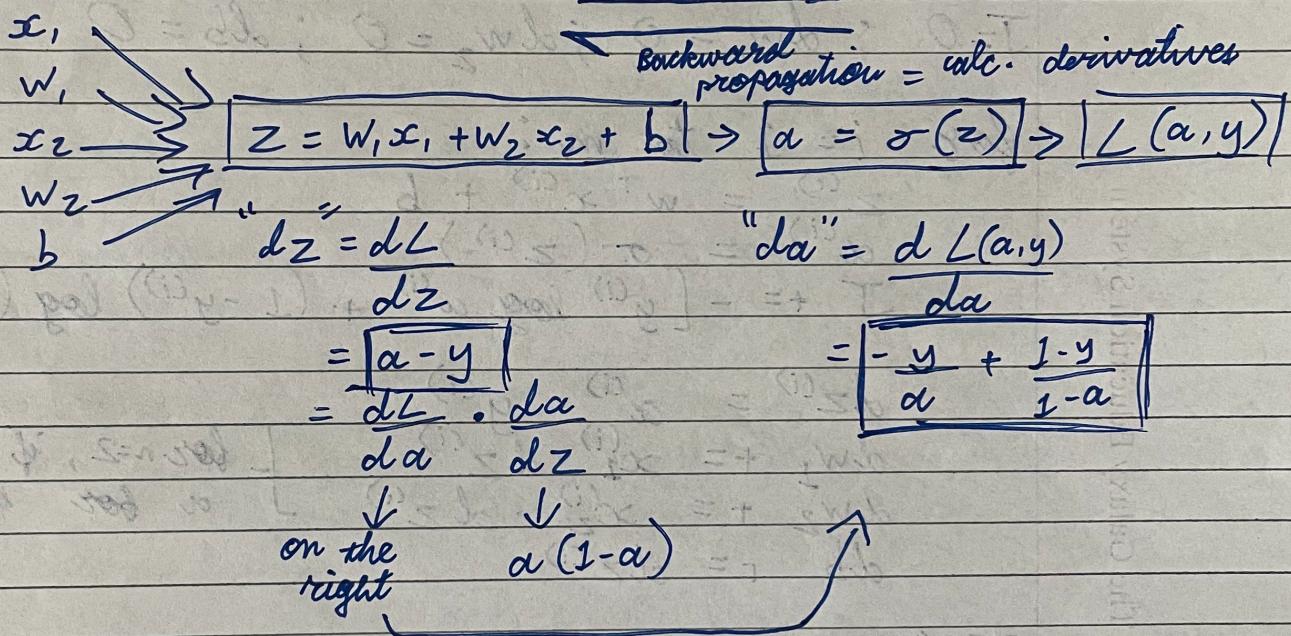
$$\rightarrow z = w^T x + b$$

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

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

→ Computational graph

Forward propagation = calculating values



$$\rightarrow \frac{\partial L}{\partial w_1} = "dw_1" = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_1} \quad \text{similarly,}$$

↓      ↓      ↓  
 coding       $\frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_1}$        $\frac{\partial L}{\partial w_2} = x_2 \cdot dz$   
 var. name       $= dz \cdot x_1$        $= "dw_2"$   
 ↑      ↑  
 coding var.

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

→ steps

① Calc.  $z$ , then  $a$ , then  $L(a, y)$ .

② Calc.  ~~$dz$~~  then  $dz^{(i)}$ , then  $dw_1, dw_2$ , and  $db$

③ Update  $w_1 := w_1 - \alpha dw_1$ ;  $w_2 := w_2 - \alpha dw_2$   
 $b := b - \alpha db$

④ Repeat

→ Implementation using for-loops

→ Just for understanding. In actual practice, vectorized form is implemented ~~as~~ cause for-loops slow af.

→ Code

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

for  $i = 1$  to  $m$ :

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

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

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

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

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

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

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

} for  $n=2$ , if more, apply a for loop.

$$J /= m$$

$$dw_1 /= m$$

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

$$dw_2 /= m$$

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

$$db /= m$$

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

→  $dw_1$ ,  $dw_2$ ,  $db$  are accumulators. At the end we take their average. Hence they don't have superscript  $(i)$ .

But  $dz^{(i)}$  is with w.r.t. just one e.g. hence the superscript  $(i)$ .

## Logistic Regression

→ Vectorizing Logistic Regression.

→ We need to calculate:

$$\begin{aligned} z^{(1)} &= w^T x^{(1)} + b & z^{(2)} &= w^T x^{(2)} + b & z^{(3)} &= w^T x^{(3)} + b \\ a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) & a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

i.e.,  $z$  &  $a$  for each training example

→ What we have is  $X$ ,

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

i.e., all the training inputs stacked together in different columns.

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

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

→ means  $X$  is a  $n_x \times m$  dimensional matrix

→ To get the  $z$  for all  $m$ -training inputs, a row vector is created of dim.  $1 \times m$ .

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \underbrace{w^T X}_{\substack{\text{also a} \\ \text{1} \times m \\ \text{vector}}} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_2 & \dots & w_m \end{bmatrix} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

which will give us,

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

Then add  $\begin{bmatrix} b & b & \dots & b \end{bmatrix}$  we get

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b & \dots & w^T x^{(m)} + b \end{bmatrix}$$

→ Python code for the above,

$$z = \text{np.dot}(w.T, x) + b \quad \text{--- [I]}$$

→ In the code,  $b = \mathbb{R}$ , But python will automatically expand it to fit the matrix obtained by the `np.dot()` operation i.e.,  $1 \times m$  matrix. This is called broadcasting.

→ To get 'a' of all the 'z',

Code

$$A = \text{sigmoid}(z) \quad \text{--- [II]}$$

where  $z$  is the vector calculated in [I] and `sigmoid()` is a user-defined function and  $A$  will be a matrix of dimension similar to that of  $z$ .

→ Calculating ~~derivative~~ derivatives with vectorization

$$\rightarrow dz^{(i)} = a^{(i)} - y^{(i)} ; \quad dz^{(2)} = a^{(2)} - y^{(2)} \\ \text{we want } dz = [dz^{(1)} \quad dz^{(2)} \dots \quad dz^{(m)}] \rightarrow \mathbb{R}^{1 \times m}$$

$$\rightarrow \text{We got } A = [a^{(1)} \dots a^{(m)}] \quad \& \quad Y = [y^{(1)} \dots y^{(m)}]$$

$$\rightarrow \text{so, } dz = A - Y \quad \Leftarrow \text{python code}$$

$$\rightarrow \text{As previously seen, } db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\ \text{so python code for that,}$$

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

→ Calculating derivatives with vectorization

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

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

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

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

Code

$$\rightarrow dw = 1/m * np.\dot{dot}(X, dz.T)$$

→ Complete code for vectorized logistic regression

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

$$A = sigmoid(z)$$

$$dz = A - Y$$

$$dw = 1/m * np.\dot{dot}(X, dz.T)$$

$$db = 1/m * np.sum(dz)$$

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

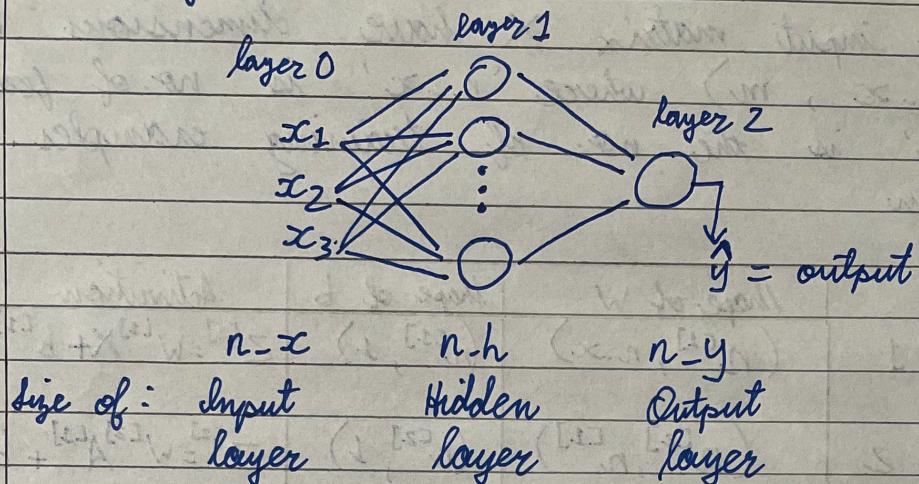
$$b := b - \alpha db$$

Note

The above code denotes one-iteration of gradient descent. For multiple iterations, a for-loop is necessary.

## Deep Neural Network

- ① Initialization of parameters  
→ 2-layer neural network



Parameters :  $w_1, b_1, w_2, b_2$

Dimensions of parameters :

$$\Rightarrow w_1 = (n-h, n-x)$$

$$\Rightarrow b_1 = (n-h, 1)$$

$$\Rightarrow w_2 = (n-y, n-h)$$

$$\Rightarrow b_2 = (n-y, 1)$$

where  $W_i$  corresponds to weight matrix of layer  $i$  and  $b_i$  corresponds to vector of bias for layer  $i$ .

$\Rightarrow$  Code (Done within function 'initialize\_parameters(n\_x, n\_h, n\_y)')

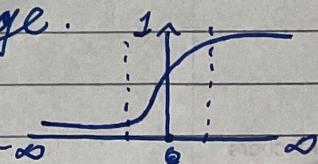
$w1 = np.random.rand(n-h, n-x) * 0.01$

$$b_1 = np \cdot \text{zeros} ((n-h, 1))$$

$$WZ = np.random.random(n-y, n-h) * 0.01$$

$$bz = np \cdot \text{zero}((n-y, 1))$$

Note '0.01' is multiplied to  $w_1$  so that the gradient remains large.   
 E.g. in sigmoid function,



we want to have no. between the dotted lines to have larger gradient and hence have better learning.

→ L-layer neural network

→ Let input matrix  $X$  have dimensions  $n^{[0]}$ ,  $(n-x, m)$  where ' $n-x$ ' is no. of features & ' $m$ ' is the no. of training examples.  
Then

	Shape of $W$	Shape of $b$	Activation
Layer 1	$(n^{[1]}, n-x)$	$(n^{[1]}, 1)$	$Z^{[1]} = W^{[1]}X + b^{[1]}$
Layer 2	$(n^{[2]}, n^{[1]})$	$(n^{[2]}, 1)$	$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$
⋮			
Layer $L$	$(n^{[L]}, n^{[L-1]})$	$(n^{[L]}, 1)$	$Z^{[L]} = W^{[L]}A^{[L-1]} + b^{[L]}$

$$\text{Shape of } Z^{[L]} = (n^{[L]}, m)$$

IMP → The Galaxy Education System  
 IMP → where  $n^{[i]}$  refers to size of  $i^{\text{th}}$  layer  
 → In general,  $W_i = (n^{[i]}, n^{[i-1]})$  &  $b_i = (n^{[i]}, 1)$   
 → The no. of units in different layers are stored in a variable 'layer-dims'.  
 For e.g. 'layer-dims' for the 2-layer NN would be  $[n-x, n-h, n-y]$ .

### Code

```

def initialize_parameters_deep(layer_dims):
    parameters = {}
    L = len(layer_dims) # no. of layers in NN

    for l in range(1, L):
        parameters["w" + str(l)] = np.random.randn(
            layer_dims[l], layer_dims[l-1]) * 0.01
        parameters["b" + str(l)] = np.zeros((layer_dims[l],
                                             1))

    return parameters
  
```

## (2) Functions for Forward Propagation

→ → Linear forward

→ The linear forward module (vectorized over all the examples) will compute :

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

where  $A^{[0]} = X$

→ The values of matrices / vectors  $W, A, b$  will be stored in a python tuple for computing the backwards pass efficiently in future.

### Code

```
def linear_forward(A, W, b):
    Z = np.dot(W, A) + b
    cache = (A, W, b)
    return Z, cache
```

→ Linear Activation Forward

→ Applying the activation function sigmoid or relu to the  $Z$ .

→ The function to be coded will return  $\sigma(z)$  or  $\text{relu}(z)$  depending on choice, and a activation cache that stores  $Z$ .

E.g. of activation function code

### Code

```
def sigmoid(Z):
```

$$A = 1 / (1 + \text{np.exp}(-z))$$

return A, Z

→ Hence the sigmoid () will return two items : the activation value "a" and a "cache" containing "z" (which will be fed to the backward function).

Can be called by ' A, activation-cache = sigmoid(z)'

→ Mathematical equation :  $A^{[l]} = g(z^{[l]}) = g(W^{[l]}A^{[l-1]} + b^{[l]})$   
where  $g$  can be sigmoid or relu

### Code

def linear\_activation\_forward(A\_prev, W, b, activation):

if activation == 'sigmoid':

    z, linear\_cache = linear\_forward(A\_prev, W, b)  
    A, activation\_cache = sigmoid(z)

elif activation == 'relu':

    z, linear\_cache = linear\_forward(A\_prev, W, b)  
    A, activation\_cache = relu(z)

cache = (linear\_cache, activation\_cache)

return A, cache

→ L-layer model

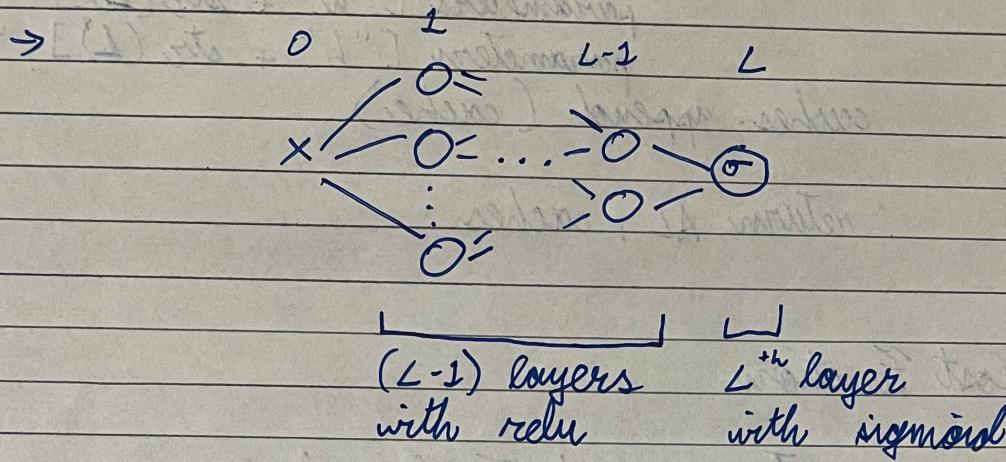


Fig.  $[\text{LINEAR} \rightarrow \text{RELU}] * (L-1) \rightarrow \text{LINEAR} \rightarrow \text{SIGMOID}$

→ Implementing forward propagation for the  $[\text{Linear} \rightarrow \text{Relu}] * (L-1) \rightarrow \text{Linear} \rightarrow \text{sigmoid}$  computation

Code

def L-model-forward(x, parameters):

caches = []

A = x

L = len(parameters) // 2 # no. of layers

# (Linear → Relu) \* (L-1)

# Loop starts from 1 as layer 0 is the input  
for l in range(1, L):

A-prev = A

A, cache = linear\_activation\_forward(

A-prev, parameters["W" + str(l)],  
parameters["b", str(l)], "relu")

caches.append(cache)

# (Linear → sigmoid) for final layer

$AL, cache = \text{linear\_activation\_forward}(A, \text{parameters}["W" + str(L)], \text{parameters}["b" + str(L)], "sigmoid")$   
 $caches.append(cache)$

return  $AL, caches$

## ⑤ Cost Function

→ Cross-entropy cost  $J$  formula:

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

### Code

def compute\_cost(AL, Y):

$m = Y.shape[1] \# \text{no. of training examples}$

$$\begin{aligned} \text{cost} = & -\frac{1}{m} * ((\text{np.dot}(Y, \text{np.log}(AL.T))) \\ & + (\text{np.dot}(1-Y, \text{np.log}(1-AL.T)))) \end{aligned}$$

$\text{cost} = \text{np.squeeze(cost)} \# \text{For To make sure } [1] \text{ turns to } 1.$

return cost

## ⑥ Backward propagation

→ np.sum (axis = 0/1) →

→ Linear Backward

→ For layer  $l$ , the linear part is:  $Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$   
 (followed by an activation)

Suppose we already have  $dZ^{[l]} = \frac{\partial L}{\partial Z^{[l]}}$ .  
 We want  $(dW^{[l]}, db^{[l]}, dA^{[l-1]})$ .

→ Mathematical formulas :

$$\rightarrow dW^{[l]} = \frac{\partial J}{\partial W^{[l]}} = \frac{1}{m} dZ^{[l]} (A^{[l-1] \cdot T})$$

Transpose of  $A^{[l-1]}$

$$\rightarrow db^{[l]} = \frac{\partial J}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m dZ^{[l](i)}$$

$$\rightarrow dA^{[l-1]} = \frac{\partial L}{\partial A^{[l-1]}} = W^{[l] \cdot T} \stackrel{\text{Transpose}}{\downarrow} dZ^{[l]}$$

Code

def linear\_backward(dZ, cache):

$A_{\text{prev}}, W, b = \text{cache}$

$m = A_{\text{prev}}. \text{shape}[1]$

$$dW = 1/m * \text{np.dot}(dZ, A_{\text{prev}}.T)$$

$$db = 1/m * \text{np.sum}(dZ, axis=1, keepdims=True)$$

$$dA_{\text{prev}} = \text{np.dot}(W.T, dZ)$$

return  $dA_{\text{prev}}, dW, db$

## → Linear Activation Backward

→ Two functions need to be user defined namely 'sigmoid-backward' & 'relu-backward'. The functions implement backward propagation for sigmoid & relu unit respectively.

→ Called by :

$$dZ = \text{sigmoid-backward}(dA, \text{activation\_cache})$$

OR

$$dZ = \text{relu-backward}(dA, \text{activation\_cache})$$

→ If  $g(\cdot)$  is the activation function, the sigmoid-backward & relu-backward compute :

$$dZ^{[i]} = dA^{[i]} * g'(Z^{[i]})$$

Code

```
def linear_activation_backward(dA, cache, activation):
    linear_cache, activation_cache = cache
```

if activation == "relu":

$dZ = \text{relu-backward}(dA, \text{activation\_cache})$   
 $dA\_prev, dW, db = \text{linear-backward}(dZ, \text{linear\_cache})$

elif activation == "sigmoid":

$dZ = \text{sigmoid-backward}(dA, \text{activation\_cache})$   
 $dA\_prev, dW, db = \text{linear-backward}(dZ, \text{linear\_cache})$

return  $dA\_prev, dW, db$

→ L - Model Backward

→ In L-model-forward(), a cache was stored containing ( $X, W, b \& Z$ ). Those will be used to compute gradients.

In L-model-backward(), we'll iterate through all the hidden layers backward.

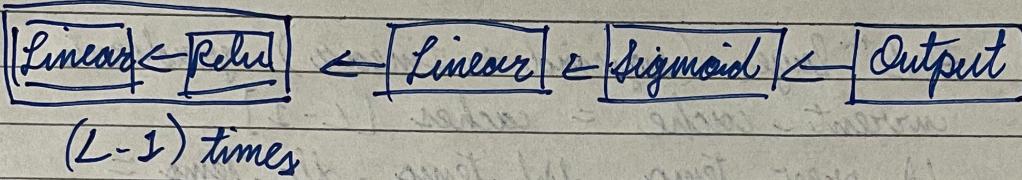


Fig. Backward Pass

→ We know  $A^{[L]} = \sigma(Z^{[L]})$ . We need  $dA_L = \frac{\partial L}{\partial A^{[L]}}$ .

$$\text{Loss} = - [y \log a + (1-y) \log (1-a)] \text{ hence, } \frac{\partial L}{\partial a}$$

$$\text{would be } \frac{\partial L}{\partial a} = \left[ - \frac{y}{a} + \frac{1-y}{1-a} \right]$$

$$\text{Therefore, } dA_L = - (\text{np.divide}(Y, AL) - \text{np.divide}(1-Y, 1-AL))$$

→ Each value of  $LA, dW$ , and  $db$  will be stored in the 'grads' dictionary.

Code

def L\_model\_backward(AL, Y, caches):

... ...

'caches' is list of caches containing:

- every cache of linear-activation-forward() with "relu" (it's caches[l], for l in range(L-1))
- the cache of linear-activation-forward() with "sigmoid" (it's caches[L-1])

grads = {}

$l = \text{len}(\text{caches})$

$m = \text{AL. shape}[1]$

$Y = Y \cdot \text{reshape}(\text{AL. shape})$

$$d\text{AL} = -(\text{np.divide}(Y, \text{AL}) - \text{np.divide}(1-Y, 1-\text{AL}))$$

#  $L^{th}$  layer (sigmoid  $\rightarrow$  Linear) gradients

current-cache = caches [ $L-1$ ]

dA-prev-temp, dW-temp, db-temp =

linear-activation-backward(dAL, current-cache, "sigmoid")

grads ["dA" + str( $L-1$ )] = dA-prev-temp

grads ["dW" + str( $L$ )] = dW-temp

grads ["db" + str( $L$ )] = db-temp

# Loop for  $l = L-2$  to  $l=0$

for  $l$  in reversed(range( $L-1$ )):

current-cache = caches [ $l$ ]

dA-prev-temp, dW-temp, db-temp =

linear-activation-backward(~~obt~~)

grads ["dA" + str( $l+1$ )], current-cache, "relu")

grads ["dA" + str( $l$ )] = dA-prev-temp

grads ["dW" + str( $l+1$ )] = dW-temp

grads ["db" + str( $l+1$ )] = db-temp

return grads

## Deep Learning

→ Updating parameters

→ Formula :

$$\begin{aligned} \Rightarrow w^{[l]} &= w^{[l]} - \alpha d_w^{[l]} \\ b^{[l]} &= b^{[l]} - \alpha d_b^{[l]} \end{aligned}$$

→ After updating parameters, stored in the 'parameters' dictionary.

Code

```
def update_parameters(params, grads, lr):  
    parameters = params.copy()  
    L = len(parameters) // 2
```

for l in range(L):

parameters["w" + str(l+1)] = parameters["w" + str(l+1)] - lr \* grads["dW" + str(l+1)]  
parameters["b" + str(l+1)] = parameters["b" + str(l+1)] - lr \* grads["db" + str(l+1)]

return parameters

# **Deep Learning Regularisation Notes**

**By Divyajeet Pala**

## Deep Learning Regularization

- Ways to solve a neural network overfitting on the data i.e., high variance problem :
  - Regularization
  - Get more training data
- Understanding Regularization based on Logistic regression.
- Cost function  $J$  for Logistic regression :

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

where  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$ ,  $m = \# \text{ no. of training e.g.}$

- When we want to regularize, we modify the equation as :

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

where  $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$   
   &  $\lambda$  = regularization parameter  $\uparrow$  vectorized form

- $\|w\|_2$  refers to the ~~L2~~  $L_2$  norm or Euclidean norm, where :

$$\|a\|_2 = \sqrt{\sum_{i=1}^n a_i^2}$$

Hence, the regularization used above is called  $L_2$  regularization.

$\Rightarrow L_1$  regularization

$$\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

where  $\|w\|_1$  is sum of the absolute value of weights.

$\Rightarrow$  On using  $L_1$  reg., w will end up being sparse i.e., w vector will have a lot of zeros in it.

Note In practice,  $L_2$  regularization is used more often.

$\Rightarrow$  Regularization in terms of Neural Network

$\Rightarrow$  For neural network,

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\text{where } \|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l-1]})^2 \because w : (n^{[l]}, n^{[l-1]})$$

i.e., sum of square of all weights

$\Rightarrow \|w\|_F^2$  is known as the "Frobenius norm" of the matrix w.

## Deep Learning Regularization

→ Implementing Gradient Descent with Regularization

→ Previously, we would calculate  $dW^{[l]}$ , which is  $\frac{\partial J}{\partial W^{[l]}}$  through the process of back-propagation.

→ Then update :

$$W^{[l]} := W^{[l]} - \alpha dW^{[l]}$$

→ Now after implementing regularization i.e., adding the  $\frac{\lambda}{2m} \sum_{l=1}^L \|W^{[l]}\|_F^2$  to the cost

function, we update modify the equation for gradient descent as follows :

$$\Rightarrow dW^{[l]} = (\text{value from backprop.}) + \frac{\lambda}{m} W^{[l]}$$

$$\Rightarrow W^{[l]} := W^{[l]} - \alpha dW^{[l]}$$

$$= W^{[l]} - \alpha (\text{backprop value} + \frac{\lambda}{m} W^{[l]})$$

$$= W^{[l]} - \frac{\alpha \lambda}{m} W^{[l]} - \alpha (\text{backprop value})$$

$$= \left(1 - \alpha \frac{\lambda}{m}\right) W^{[l]} - \alpha (\text{backprop value})$$

⇒ For this reason,  $L_2$  regularization is also known as "weight decay" since we are decreasing the original weights by multiplying with a number little less than 1.

# **Deep Learning Vanishing Gradient Notes**

**By Divyajeet Pala**

## Vanishing Gradient Problem

- Using more layers in a neural network causes more activation functions. And as we increase that, increases the gradients of loss function which leads to zero. What this causes is that the weights initial layers will not update much ~~e.g.~~ will update by infinitesimally small amount. So the network won't improve much.
- Partial solution
- One of the ways in which it can be partially solved is by initializing the weights in a different manner.
- $Z = w_1 x_1 + w_2 x_2 + \dots$   
 We want  $Z$  to not blow up and nor become too small.
- For that, larger the  $n$  (no. of layers),  $w_i$  should be smaller. Since  $Z = \sum_{i=1}^n w_i x_i$ , if we are adding more number of  $w_i x_i$  terms due to more no. of layers, we want each of the  $w_i x_i$  terms to be smaller.
- To solve the issue (partially), we can set the variance of  $w_i$  to  $\frac{1}{n}$  where  $n$  = no. of input features going into the neuron.

$$\text{Var}(w_i) = \frac{1}{n} \quad \text{where } n = \text{no. of input features going into the neuron}$$

To do this, initialize as follows :

$$w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$$

→ If the activation function used is Relu,

$\text{Var}(w_i) = \frac{2}{n}$  works better. It can be achieved by :

$$w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{[l-1]}}\right)$$

→ For activation function tanh :

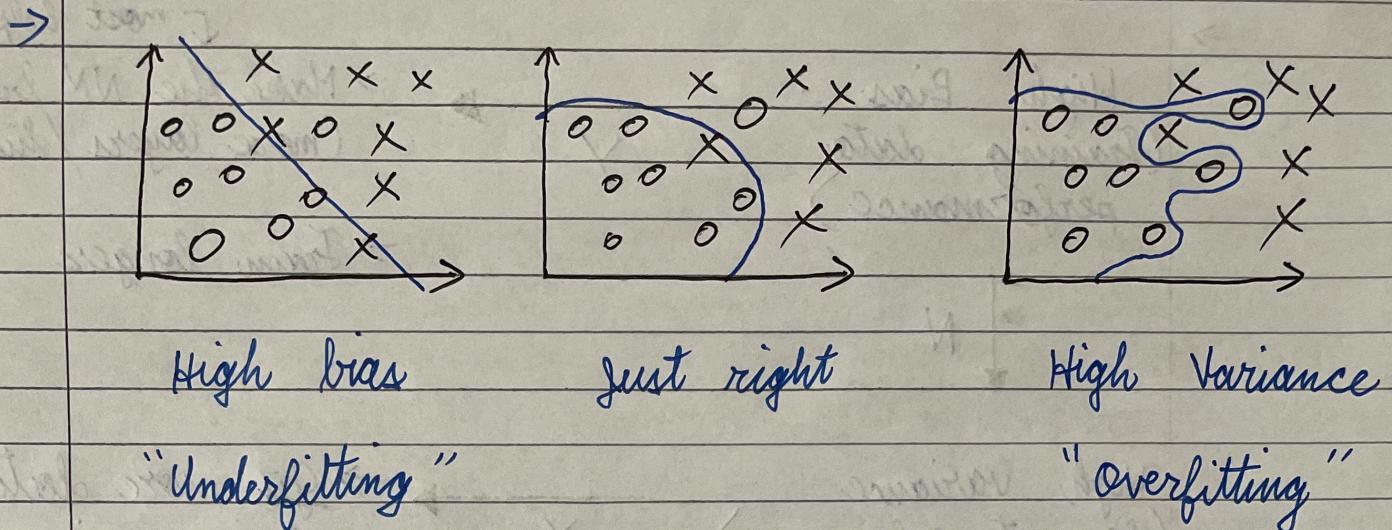
$$\text{Var}(w_i) = \frac{1}{n}$$

$$w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$$

Xavier ~~at~~ initialization  $\rightarrow$

# **Deep Learning Bias and Variance Notes**

**By Divyajeet Pala**



The Galaxy Education System

Train set error:	1%	15%	15%	0.5%
Dev set error:	11%	16%	30%	1%

Assuming optimal : High High High bias Low  
error to be variance Bias & high bias &  
0% variance variance

→ The train set error gives a sense of the bias problem and the dev set error gives a sense of variance problem.

→ Note that if the optimal error for the above e.g. was 15%, the second model would be best and would be considered as low bias & low variance.

→ To solve the issue of high bias and/or high variance, a few things can be done. Basic recipe of ML ~~summary~~ tells a few of those ways.

→

High Bias  
(Training data performance)

Y

↙ most helpful

→ Make the NN bigger  
(more layers / hidden units)



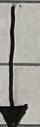
N

→ Train longer

High variance  
(Dev set performance)

Y

→ Get more data & train again



N

⇒ Regularization

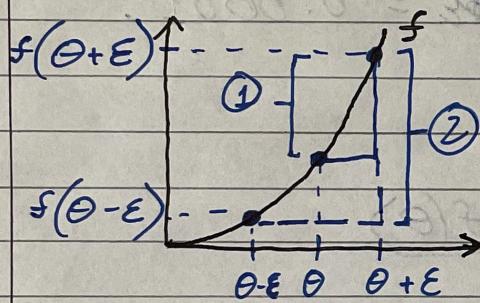
Done

# **Deep Learning Gradient Checking Notes**

**By Divyajeet Pala**

## Checking derivative computation

- Derivatives are calculated when implementing back propagation in a neural network to update the weight matrices and biases.
- There is a test called gradient checking that can help make sure that the implementation of back prop is correct.
- E.g.  $f(\theta) = \theta^3$



① → The method with which gradient is actually calculated.

② → Method proposed for checking the gradient

if  $\epsilon = 0.01$ ,

$$\theta = 1$$

$$\theta + \epsilon = 1.01$$

$$\theta - \epsilon = 0.99$$

$$\Rightarrow ② = f(\theta + \epsilon) - f(\theta - \epsilon)$$

$$\Rightarrow \text{width of the bigger } \Delta \text{ is } [\theta + \epsilon - (\theta - \epsilon)] = 2\epsilon.$$

$$\Rightarrow \text{Gradient} = \frac{\text{Height}}{\text{Width}} = \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx f'(\theta)$$

$$= \frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001$$

$$\Rightarrow g(\theta) = f'(\theta) = 3\theta^2 = 3(1)^2 = 3$$

$$\text{Approx. error} \Rightarrow 3.0001 - 3 = 0.0001$$

- As the error is very less, the two-sided difference gradient checking can be used to evaluate the derivatives.

→ Why two-sided difference  $(\theta + \epsilon, \theta - \epsilon)$  is better than  $(\theta, \theta + \epsilon)$  one-sided difference?

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \leftarrow \text{Formal definition of derivative}$$

The error for above is  $\overset{\text{big-O}}{\sim} (\epsilon^2)$  where  $\epsilon$  is a value tending to 0 (zero) i.e., very small value.

E.g. if  $\epsilon = 0.01$ , error = 0.0001.

For one-sided diff.,

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$$

the error tends to be  $\overset{\text{big-O}}{\sim} (\epsilon)$ .

E.g. For  $\epsilon = 0.01$ , error = 0.01

which is significantly more than 0.0001.

→ Therefore, one-sided difference is a much less accurate approximation when compared to two-sided difference.

For the same reason, two-sided difference is used when doing gradient checking.

→ By comparing the derivative from back-prop with the calculated two-sided difference, the accuracy can be checked of the back prop.

## Gradient checking

→ Implementation

→  $w^{[i]}$  refers to the weight matrix of layer  $i$ .  
 $b^{[i]}$  refers to bias of layer  $i$ .

→ Take  $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$  and reshape into a big vector  $\Theta$ .  $\Theta$  will contain all the values of  $w^{[1]} \dots b^{[L]}$ .

Therefore, the cost function  $J$  which earlier was a function of all the different weights & bias, will now be a function of  $\Theta$ .

Earlier :  $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]})$   
 Now :  $J(\Theta)$

→ In the same way as before, concatenate  $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$  and reshape into a big vector  $d\Theta$ .

shape of  $\Theta$  &  $d\Theta$  will be same.

→ Now checking whether  $d\Theta$  is the gradient of  $J(\Theta)$ .

→ Gradient checking

Code (Pseudo)

for each  $i$  in  $\Theta$  : # iterating over the vector  $\Theta$

$$d\Theta_{\text{approx.}}[i] = \frac{J(\underline{\Theta_1, \Theta_2, \dots, \Theta_i + \epsilon, \dots}) - J(\underline{\Theta_1, \dots, \Theta_i - \epsilon, \dots})}{2\epsilon}$$

→ Basically getting a vector  $d\Theta_{\text{approx.}}$  where the  $i^{\text{th}}$  element is two-sided difference of  $J(\Theta)$  in which only the  $i^{\text{th}}$  element is added or subtracted with  $\epsilon$  & then two-sided difference is calculated.

→ **Formulas :**

$$d\theta_{approx}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\text{so, } d\theta_{approx}[i] \approx \frac{\partial J}{\partial \theta_i}.$$

Dimension of  $d\theta_{approx}$ . will be same as that of  $d\theta$ .

Aim is to check whether  $d\theta_{approx}$  &  $d\theta$  are almost equal or not.

To do that euclidean distance between them is calculated.

$$\rightarrow \text{Check } \frac{\|d\theta_{approx.} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} \quad \text{--- (i)}$$

where  $\|d\theta_{app.} - d\theta\|_2$  is the L2 norm i.e., sum of squares of element-wise difference and then taking the root of answer.

→ The numerator is the euclid. distance and denominator is used to normalize the values.

→ In practice,  $\epsilon = 10^{-7}$  will be good.

so with  $\epsilon = 10^{-7}$ ,

→ if the formula (i) gives  $10^{-7}$  or smaller, the derivative approximation is correct ✓.

→ if (i)  $\approx 10^{-5}$ , okayish (can work but double-checking necessary).

→ if (i)  $\approx 10^{-3}$ , the implementation of back prop has a bug.

# **Deep Learning Optimisation Algorithms Notes**

**By Divyajeet Pala**

## Optimization Algorithms

### Exponentially moving averages EMA (EMA)

- A moving average shows an average of data-points for a certain number of time periods. It's called 'moving' because each data-point is calculated using data from previous  $\times$  number of time periods. By averaging prior data, MA smooth the price data to form a trend following indicator. They don't predict price direction, but rather define the current direction.
- E.g. Let  $\theta_1, \theta_2, \theta_3 \dots$  be temperatures of a place where  $\theta_i$  is the temperature on  $i$ -th day.

Calculating moving averages :

→ Initializing  $V_0 = 0$ .

$$\Rightarrow V_1 = 0.9 V_0 + 0.1 \theta_1 \leftarrow \begin{array}{l} \text{That day's temperature} \\ \text{previous value} \end{array}$$

$$\Rightarrow V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$\Rightarrow V_3 = 0.9 V_2 + 0.1 \theta_3$$

⋮

$$\Rightarrow V_t = 0.9 V_{t-1} + 0.1 \theta_t \leftarrow \text{general formula}$$

→ Formula

$$\Rightarrow V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

→ In the e.g.  $\beta = 0.9$ .

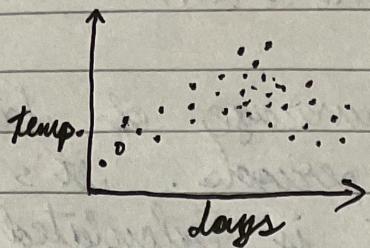
→ Here,  $V_t$  <sup>approx.</sup> averages over  $\frac{1}{1-\beta}$  number of days'

temperature.

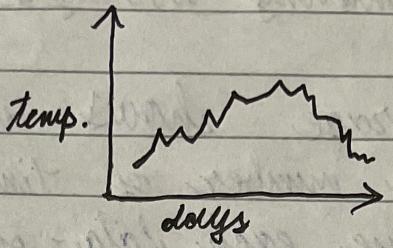
So, when  $\beta = 0.9$ ,  $V_t$  would approximately average over 10 days of temperatures.

→ If  $\beta = 0.98$ ,  $V_t \approx 50$  days.

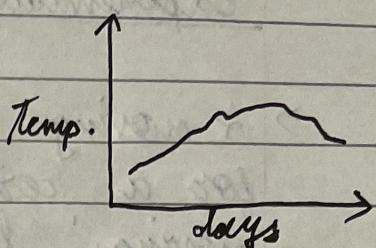
→



Scatter plot of  
temp v/s days



$V_t$  graph when  
 $\beta = 0.9 \approx 10$  days  
average



$V_t$  when  
 $\beta = 0.98 \approx 50$   
days avg.

- As seen, when  $\beta$  is high, the graph is smoother. Downside of having high  $\beta$  is that the EMA formula, adapts more slowly to the temperature changes.

The reason :

$$\Rightarrow V_t = \underbrace{\beta(V_{t-1})}_{\text{Previous value}} + (1-\beta)\theta_t \quad \text{Current value}$$

when  $\beta$  is large, it gives more weight to  $V_{t-1}$  i.e., the previous values, and a smaller weight to the  $\theta_t$  i.e., current temperature.

- Intuition behind ~~EAT~~ EMA in the next page.

Exponentially moving averages (EMA)

$$\Rightarrow V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$\Rightarrow$  Writing the equations for  $V_t$  in decreasing order of  $t$  with  $\beta = 0.9$ :

$$\Rightarrow V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

$$V_{99} = 0.9 V_{98} + 0.1 \theta_{99}$$

$$V_{98} = 0.9 V_{97} + 0.1 \theta_{98}$$

...

$\Rightarrow$  Substituting  $V_{99}$  in the eq. of  $V_{100}$  we get,

$$V_{100} = \cancel{0.9} \cdot 0.1 \theta_{100} + 0.9 (0.9 V_{98} + 0.1 \theta_{99}) \quad \text{--- } \textcircled{i}$$

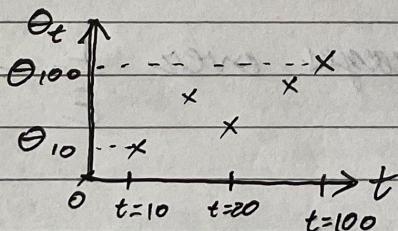
$\Rightarrow$  Subst.  $V_{98}$  in  $\textcircled{i}$ ,

$$V_{100} = 0.1 \theta_{100} + 0.9 [0.1 \theta_{99} + 0.9 (0.1 \theta_{98} + 0.9 V_{97})]$$

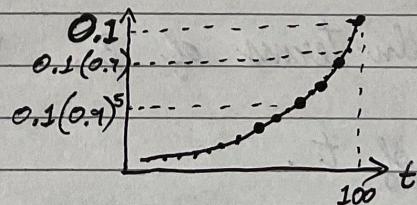
$\Rightarrow$  Expanding, we get :

$$V_{100} = 0.1 \theta_{100} + 0.1 (0.9) (\theta_{99}) + 0.1 (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97}$$

$\Rightarrow \theta_t$  v/s  $t$  scatter plot



Scatter plot of the datapoints



An exponentially decaying graph for the  $0.1 (0.9)^t$  part of the equation i.e., co-efficients

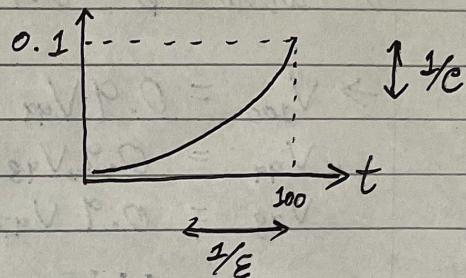
$\Rightarrow$  For calculating,  $\cancel{\theta_{100}}, V_{100}$ , sum up the element-wise product of

the two functions.

→ The sum of the co-efficients ( $0.1 \times 0.9^x$  part of the term), is approx. 1.

→ Let  $\beta = 1 - \epsilon$ ,

$$(1 - \epsilon)^{\frac{t}{\epsilon}} = \frac{1}{e}$$



Basically, the equation tells us that it takes  $\frac{1}{\epsilon}$  days for the height of the graph to decay to  $\frac{1}{e} \sim \frac{1}{3}$  of the peak.

So, for  $\beta = 0.9$ ,  $\epsilon = 0.1$ . Hence it takes  $\frac{1}{\epsilon} = 10$  days for the peak of the graph to come down to  $\frac{1}{e}$  or approx.  $\frac{1}{3}$  of the peak.

$$(0.9)^{\frac{10}{0.1}} \Rightarrow 0.9^{10} \approx \frac{1}{e}$$

values before

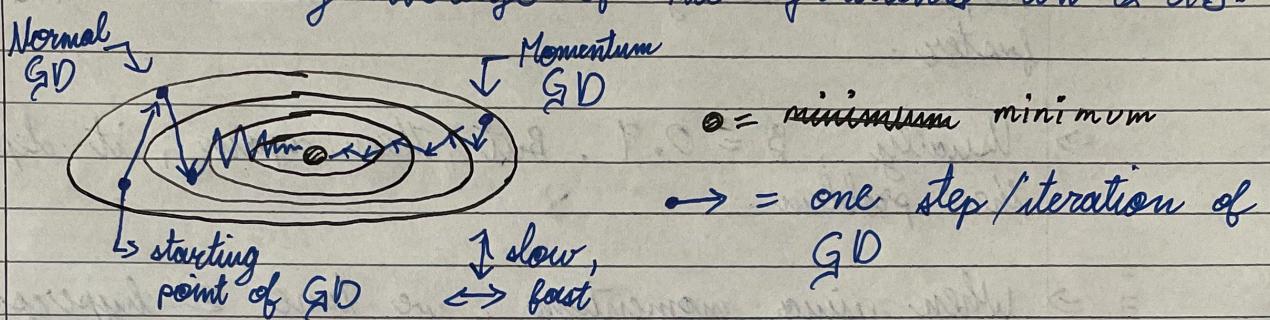
Hence we can say that after 10 days, contribute a very less weight  $\ll \frac{1}{e}$  (peak), to the current moving average value. Hence we say that  $\beta = 0.9$ , averages over 10 days of temperature.

In terms of  $\beta$ ,  $\sqrt{t}$  averages over  $\frac{1}{\epsilon} = \frac{1}{1-\beta}$  values of  $t$ .

→ We could just add the values of 10 days and divide with 10 to get average, but that is computationally and memory-wise really expensive though more accurate. EMA is more efficient when no. of variables is high.

### Gradient Descent with Momentum

→ Idea: The basic idea is that instead of updating weights ( $w$ ) & biases ( $b$ ) with  $dW$  &  $db$ , we update them with the exponentially moving average of the gradients  $dW$  &  $db$ .



we want learning to be slower vertically, & faster horizontally.

→ Applying momentum:

On iteration  $t$ :

Compute  $dW$ ,  $db$  on the current mini-batch/batch:

$$\cancel{V_{dw}} = \cancel{dW}$$

$$\rightarrow V_{dw} = \beta V_{dw} + (1 - \beta) dW$$

$$V_{db} = \beta V_{db} + (1 - \beta) db$$

$$w = w - \alpha V_{dw}, b = b - \alpha V_{db}$$

where  $\alpha$  = learning rate,  $\beta$  = EMA term

→ Intuition:

Think of the function as a bowl shaped function. We roll down a ball, the terms  $dW$ ,  $db$  act as acceleration, and  $V_{dw}$ ,  $V_{db}$  (being multiplied  $\beta$ ) act as velocity. So due to acceleration, our ball moves faster & faster towards the center instead of going here-and-there much.  $\beta$  acts as friction and

prevents the ball from speeding infinitely.

So due to momentum, gradient descent algorithm now doesn't take the steps independently of the previous steps. Instead previous steps momentum helps in making the ball reach the minima point faster.

- Usually,  $\beta = 0.9$ . But, then again, it depends on the problem.
- When using momentum, we tune 2 hyperparameters namely,  $\alpha$  &  $\beta$ . Where  $\alpha$  is the learning rate and  $\beta$  is the momentum term.

$$w_k(\beta - 1) + \gamma v_k = w_k$$

$$d_k(\beta - 1) + \gamma v_k = d_k$$

$$\gamma v_k - d_k = d_k - w_k$$

## RMS prop

→ Idea : Updating params  $w$  &  $b$  by  $dW$  &  $db$  divided by a numbers terms that would be small, if the gradient of the term ( $w$  or  $b$ ) would favour in reaching the minima and would be large if the direction of the gradient of the term ( $w$  or  $b$ ) would not favour reaching the minima.

→ On iteration  $t$  :

Compute  $dW$ ,  $db$  for the mini-batch/batch :

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2 \leftarrow \text{element-wise squaring}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \leftarrow$$

$$w := w - \alpha \frac{dW}{\sqrt{S_{dw}}}, b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

Note  $\beta_2$  is just  $\beta$ , the  $\beta_2$  just separates the  $\beta$  of RMS prop from that of momentum.

→ Suppose, the direction of gradients of  $w$  favour in reaching the minima. Then, while calculating  $S_{dw}$ ,  $(dW)^2$  will turn out to be small. Hence  $dW$  will be divided with a relatively small no.. so the updates in  $w$  direction will be larger as the update ( $dW$ ) is divided with a small number ( $\sqrt{S_{dw}}$ ).

If direction of gradient of  $b$  is not favourable,  $S_{db}$  will be a larger number as  $(db)^2$  will be large.

Hence, the update will be smaller as the denominator ( $\sqrt{S_{db}}$ ) will be large.

Adam Optimization

- ⇒ Idea : Combining momentum with RMSprop.
- ⇒ Note :  $\beta_1$  corresponds to  $\beta$  of momentum, and  $\beta_2$  to  $\beta$  of RMSprop.
- ⇒ Initialize  $V_{dw} = 0$ ,  $S_{dw} = 0$ ,  $V_{db} = 0$ ,  $S_{db} = 0$

On iteration  $t$ :

Compute  $dW$ ,  $db$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dW$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

# Bias correction for all terms

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1 - (\beta_1)^t}, \quad V_{db}^{\text{corrected}} = \frac{V_{db}}{1 - (\beta_1)^t}$$

$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{1 - (\beta_2)^t}, \quad S_{db}^{\text{corrected}} = \frac{S_{db}}{1 - (\beta_2)^t}$$

# Updating  $W$  and  $b$

$$W := W - \alpha \frac{V_{dw}^{\text{corr.}}}{\sqrt{S_{dw}^{\text{corr.}}} + \epsilon}$$

$$b := b - \alpha \frac{V_{db}^{\text{corr.}}}{\sqrt{S_{db}^{\text{corr.}}} + \epsilon}$$

Note  $\epsilon$  is a small value which prevents the term from exploding in case the  $\sqrt{S^{\text{corr.}}}$  term becomes very small.

→ Hyperparameter tuning :

→  $\alpha$  (learning rate) : Needs to be tuned according to the problem

→  $B_1$  (momentum parameter) : 0.9 (by default) but can be changed

→  $B_2$  (RMS prop parameter) : 0.999 (by default) is optimum but can be changed

→  $\epsilon$  :  $10^{-8}$  (Doesn't contribute much to model's performance)

→ Learning Rate Decay

Idea : As the GD gets closer to the minima, the learning rate decreases so as to provide a smoother slope and prevent over-shooting.

$$\text{Formula} : \alpha = \frac{1 \times \alpha_0}{1 + (\text{decay-rate} \times \text{epoch.num})}$$

$$\text{E.g. } \alpha_0 = 0.2, \text{ decay} = 1$$

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
:	: