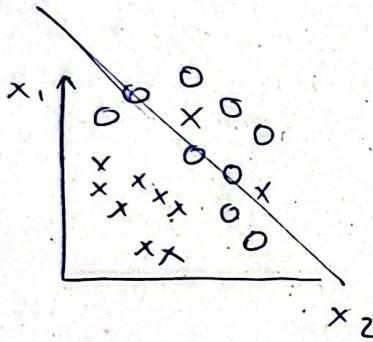


## Course 2

### Improving Deep Neural Networks:- Week 1

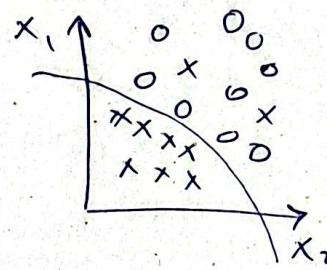
Hyper parameter :- it's the parameter that affects the resulting other parameters  $\rightarrow \alpha$ ,

#### Bias and Variance



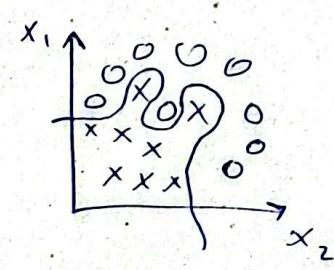
high bias

in this case there is high bias because the model is underfitting the data because the solution is biased thanks to our assumption that a linear decision boundary would fit.



Just right

the decision boundary is not underfitting nor overfitting the data so it actually has good bias and good variance.



high variance

here the decision boundary is overfitting the data resulting in a high variance

Train set  $\rightarrow 1\%$

15 %

0.05 %

Dev set  $\rightarrow 1.1\%$

16 %

1 %

in high dimensional plots you can't visualize the  
over fitting

D.B

high bias  
high variance  
over fitting

high bias  
underfitting

just right

Dev set Train set

80 %

15 %

$\rightarrow$  high variance.  $\rightarrow$  high bias.

## Basic Recipe for ML

if you have

High bias?

(training data overfitting)  
under

try

→ Bigger network

(more hidden units)

(more hidden layers)

Train the network longer

\* (change the Arch)

These options can fix the high bias problem.

High variance?

(Dev set performance?)

→ More Data

Regularization

\* (change the Arch).

if the train error is low but  
the Dev error is high then  
you're overfitting the training data

High variance

- \* in Deep Learning there is not much talk about the Bias Variance trade off because you don't have to affect one by changing the other like ML

## Regularization

$$[\omega \in \mathbb{R}^{n \times k}, b \in \mathbb{R}]$$

Logistic regression

$$\min_{\omega, b} J(\omega, b) = L(\hat{y}, y) + \underbrace{\frac{\lambda}{2m} \|\omega\|_2^2}_{\text{regularization parameter}}$$

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

$$\|\omega\|_2^2 = \sum_{j=1}^{n \times k} \omega_j^2 = \omega^\top \omega$$

$\rightarrow$  L2 regularization  
norm of  $\omega^2$

$$\text{L1 regularization} \rightarrow \frac{\lambda}{2m} \sum_{j=1}^{n \times k} |\omega_j| = \frac{\lambda}{2m} \|\omega\|_1$$

For a neural network

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

$$\|\omega^{[l]}\|^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (\omega_{i,j}^{[l]})^2 + \frac{\lambda}{2m} \sum_{l=1}^L \|\omega^{[l]}\|^2$$

$\hookrightarrow$  This is called the "Frobenius norm"

$$d\omega^{[l]} = (\text{from backprob}) + \frac{\lambda}{m} \omega^{[l]}$$

This is also called weight decay  $\hookrightarrow$  regularization term

$$\omega^{[l]} := \omega^{(l)} - \alpha \frac{d\omega^{[l]}}{\text{regularized}}$$

$$= \omega^{[l]} - \frac{\alpha \lambda}{m} \omega^{[l]} - \alpha (d\omega^{[l]})$$

## \*Drop out Regularization

→ In dropout regularization we drop some of the hidden units at random, dropout has multiple implementation methods.

### \* Inverted dropout :-

For Layer = 3. →  $l = 3$ ,  $\boxed{\text{keep\_prob} = 0.8}$

$d_3$  → dropout vector

$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep\_prob}$

→ each activation like ( $a_3$ ) has a value ranging from 0 to 1 so if  $a_3 < \text{keep\_prob}$  and  $\text{keep\_prob}$  is 0.8 that means that if  $a$  is ~~bigger~~ less than 0.8 set the  $d_3$  which is the dropout vector to 1 else set to 0

$a_3 = [0.5 \ 0.9 \ 0.4 \ 0.3]$ ,  $\text{keep\_prob} = 0.7$

$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep\_prob}$

`print(d3)` → we used `random.rand` because we don't actually have  $a_3$  but in implementation

\* to use the dropout vector :-  $\boxed{d_3 = a_3 < \text{keep\_prob}}$

$a_3 = \text{np.multiply}(a_3, d_3)$

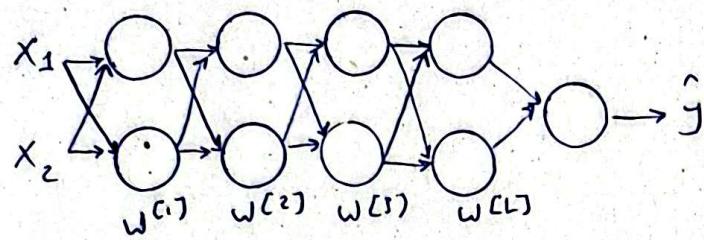
`print(a3)`

$\hookrightarrow [0.5 \ 0 \ 0.4 \ 0.3]$

$a_3 = a_3 / \text{keep\_prob}$

$\hookrightarrow$  to maintain the value of  $a$  which were reduced by eliminating one of the hidden units values.

## Vanishing and exploding gradients.



$g(z) = z \rightarrow$  Linear activation function.

$$z^L = w \cdot x + y$$

$$\hat{y} = w^{(1)} w^{(2)} \dots w^{(L)} \cdot x$$

$$z^{(1)} = w^{(1)} \cdot x$$

$$a^{(1)} = g(z^{(1)}) = z^{(1)}$$

$$a^{(2)} = g(z^{(2)}) = z^{(2)}$$

$$z^{(2)} = w^{(2)} \cdot z^{(1)}$$

$$z^{(L)} = \hat{y} = w^{(L)} w^{(L-1)} w^{(L-2)} \dots w^{(3)} w^{(2)} w^{(1)} X$$

$\underbrace{\qquad\qquad\qquad}_{z^{(1)}}$   
 $\underbrace{\qquad\qquad\qquad}_{z^{(2)}}$   
 $\underbrace{\qquad\qquad\qquad}_{z^{(3)}}$

Let's suppose that

$$w^{(L)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$\hat{y} = w^{(L)} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} \cdot X$$

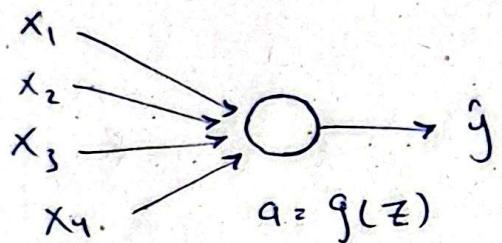
this weight

matrix is usually a bit larger than the rest.

since this equation will result in multiplying  $x$  with  $1.5^{(L-1)}$  times then the value of  $\hat{y}$  will likely explode and if we do the same setting the value of  $w^{(L)}$  to  $\begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$  then the value of  $\hat{y}$  will decrease until it vanishes.

\* a partial solution to the problem of vanishing and exploding gradients is

• Careful choice of the random initialization parameters.



\*  $z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$

→ to solve the problem of vanishing and exploding gradients for a neural network with  $l$  number of layers

\* if  $l$  is large you want to minimize  $w^{(l)}$

→ because the vanishing gradient can change as the number of layers get bigger.

$$z = \sum_{i=1}^l w_i x_i$$

→ a way to deal with this problem is setting the variance of  $(w_i)$  to be  $\frac{2}{n}$ , this helps with the gradients stabilizing and not growing or diminishing.

how to set  $\text{Var}(w_i) = \frac{2}{n}$  ?

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

setting the  $\text{Var}(w_i)$  to  $\frac{2}{n}$  is required when dealing with ReLU because ReLU only activates half the inputs.

\* this only helps to stabilize  $z$  and reduce the vanishing and exploding gradient's problem.

\* Why did we decide by  $n^{[l-1]}$  when setting the variance to be  $\frac{2}{n}$ ?

↳ because in large neural nets the layer  $\underline{L}$  has  $n-1$  inputs so the  $n^{[l-1]}$  is the correct input to layer  $L$ .

\* Why this solution helps ~~solve~~ with the vanishing and exploding gradients problem?

↳ we can control the variance of  $Z$  through the variance of  $w_i$  and  $x_i$ :

$$Z = \sum_{i=1}^n w_i x_i$$

$$\text{Var}(Z) = n \cdot \text{Var}(w_i) \cdot \text{Var}(x_i) \rightarrow \text{assuring } x \text{ has mean } 0 \text{ and } \sigma \text{ of } 1.$$

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

↳ we set the  $\text{Var}(w_i) = \frac{1}{n}$  to ensure that the variance of  $Z$  does not grow as the value of  $n$  increases and this stabilizes  $Z$ .

Why scaling  $w_i$  by  $\text{np.sqrt}(\frac{2}{n})$  works?

\* np.random.rand generates  $N(0, 1)$  with Gaussian distribution

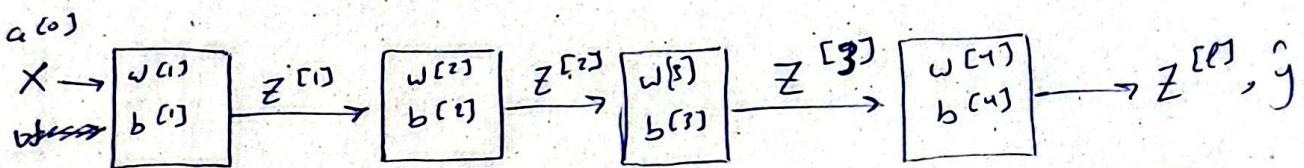
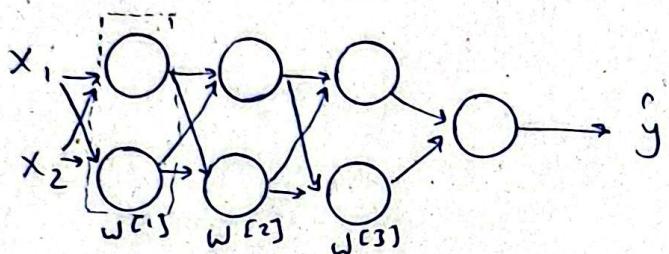
$$\text{Var}(w_i) = (\underbrace{\sqrt{\frac{2}{n}}}_\text{multiplies the samples})^2 \cdot \text{Var}(0, 1) = \frac{2}{n}$$

↳ this multiplies the samples

this is why we multiplied the samples with the square root of  $\frac{2}{n}$  because the variance of  $\sqrt{\frac{2}{n}}$  =  $(\sqrt{\frac{2}{n}})^2 = \frac{2}{n}$

\* why stabilizing  $Z$  is the solution that we have to help reduce this problem?

↳ because the value that is being fed to the neural network's layers is  $Z$



→ since  $Z$  is the value that is being fed because  $Z = \sum_{i=1}^n w_i x_i$  and the main component responsible for the problem of vanishing and exploding gradients is  $w$  (supposing  $x$  has mean 0 and  $\sigma$  of 1).

then ~~scales~~ scaling  $w$  scales the value of  $Z$  helping maintaining stable distribution and ensuring the signal does not diminish or explode through random initialization of  $w$  with a certain distribution that has variance  $\text{Var}(w) = \frac{1}{n}$

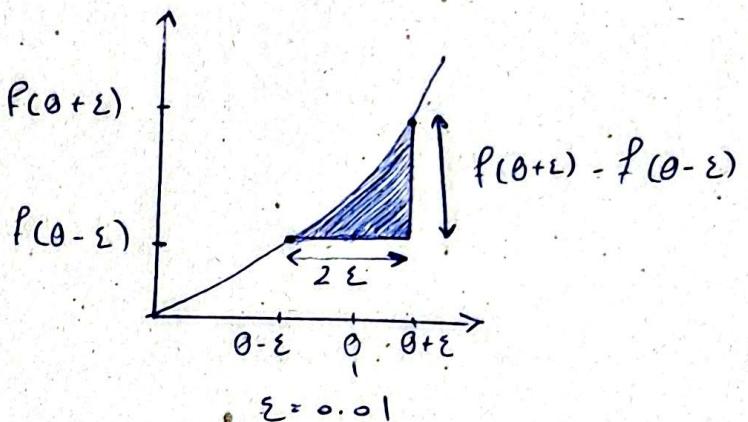
$$\boxed{\text{Var}(Z) = n \cdot \text{Var}(w) \cdot \text{Var}(x)}$$

## Numerical approximation of gradients

$$\text{slop} = \frac{f(\theta + \varepsilon) - f(\theta - \varepsilon)}{2\varepsilon} \approx g(\theta)$$

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

~~= 3.0001~~



$$f(\theta) = \theta^3 \rightarrow g(\theta) = f'(\theta) = 3\theta^2 = 3 \times 0.001^2 = 3$$

\* approximation error =  $3.0001 - 3 = 0.0001$

↳ From this we can be confident that  $g(\theta) = 3\theta^2$  is a correct implementation of the derivative of ~~f~~  $f(\theta) = \theta^3$ .

## Gradient checking

- First step of implementing gradient checking is taking all our parameters like  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$  and reshaping (concatenating) them into 1 vector  $\theta$ .

\*  $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$

- next we take the parameters  $dw^{[1]}, db^{[1]}, \dots, dw^{[L]}, db^{[L]}$  from backprob and concatenate them into 1 vector  $d\theta$

- The question now is  $\rightarrow$  Is  $d\theta$  the gradient  $J(\theta)$

\*  $J(\theta) = J(\theta_1, \theta_2, \theta_3, \theta_4, \dots)$

for each  $i$ :

$$d\theta \text{ approx } [i] = \frac{J(\theta_1, \theta_2, \dots, \cancel{\theta_i}, \theta_i + \varepsilon) - J(\theta_1, \theta_2, \dots, \cancel{\theta_i}, \theta_i - \varepsilon)}{2\varepsilon}$$

$$\therefore d\theta[i] = \frac{\partial J}{\partial \theta_i}$$

now we need to check ;  $f \boxed{d\theta_{approx} \approx d\theta}$

\* Check  $\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2}$

- the main thing that we did is calculate the numerical gradient for each  $\theta_i$ .

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

and doing that for all each  $\theta_i$ .

now we have  $d\theta_{approx} = \left[ \frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \frac{\partial J}{\partial \theta_3}, \dots, \frac{\partial J}{\partial \theta_n} \right]$

↳ after each time we calculated  $d\theta_{approx}[i]$   
we added the result to the vector  $d\theta_{approx}$

- now we can go check the difference between our numerical gradient and the gradient we got from backprob, we can measure them by calculating the euclidean distance between them.

check  $\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} \approx 10^{-7}$  for  $\epsilon = 10^{-7}$

- if we have used  $\epsilon = 10^{-7}$  then the diffrence is supposed to be around  $10^{-7}$ , if else then we need to check our individual components.

## Course 2 & week 2

### Mini batch GD :-

Mini batch gradient descent is just about processing mini or smaller batches or chunks of the training data at a time because processing the whole dataset would consume a lot of time and resources.

\* a mini-batch has  $>1$  and  $\leq m$  size

$$1 < \text{mini-batch-size} < m$$

if mini-batch = m then it's just batch GD

and if m=1 then it's called stochastic GD.

\* Mini batches typically have a size that is the exponential of 2. Like  $2^2$  or  $2^4$  or  $2^6$  so the sizes could be 4, 16, 32, 64, 128, 256, 512 etc with the most common mini batch size = 512

$$X_{(n \times m)} = [x^{(1)}, x^{(2)}, \dots, x^{(1000)} | x^{(1001)}, \dots, x^{(2000)}] \quad x^{\{1\}} \quad x^{\{2\}} \quad \dots \quad x^{\{3\}}$$

$$Y_{(1, m)} = [y^{(1)}, y^{(2)}, \dots, y^{(1000)} | y^{(1001)}, \dots, y^{(2000)}] \quad y^{\{1\}} \quad y^{\{2\}} \quad \dots \quad y^{\{3\}}$$

## Mini-batch GD implementation

→ for  $t = 1, \dots, \text{last mini-batch}$  num-mini-batch.

→ forward prop on  $X^{[t]}$

$$\rightarrow Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$\rightarrow A^{[1]} = g^{[1]}(Z^{[1]}) \quad \left. \begin{array}{l} \text{vectorized} \\ \text{implementation} \end{array} \right\}$$

$$\rightarrow A_L = g^{[L]}(Z^{[L]})$$

$$\rightarrow \text{compute cost } J^{\text{avg}} = \frac{1}{\text{size of minibatch}} \sum_{i=1}^{\text{size of minibatch}} L(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2000} \sum \|W^{(l)}\|_F^2$$

↳ This is the cost on only 1 mini-batch

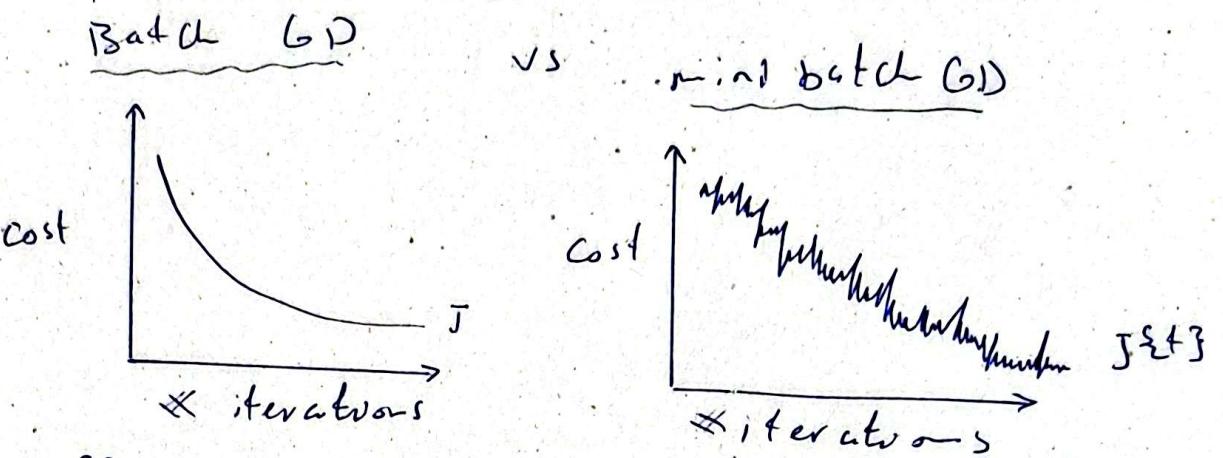
→ Go back prop to compute gradients with respect to  $J^{[t]}$

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

$$\rightarrow b^{[l]} = b^{[l]} - \alpha db^{[l]}$$

This is all called one epoch or a single pass through the training data or in this case a single pass through all the mini-batches.

\* Mini-batch GD allows you to make modifications to the weights and biases based on each mini-batch without needing to go through the whole training data to just make one update.



→ the difference between the two graphs is that when we're using batch gradient descent we have to go through all the training data just to make one weight update and that surely smoothes out the curve because the weight update is based on all the training data but with mini-batch GD the update is based on separate smaller size mini-batches which make the training and weight updates prone to more bias but surely after running through all the mini-batches it gets to the same cost as batch GD; if not smaller.

## exponentially weighted averages

$$\theta_1 = 40^{\circ}$$

$$\theta_2 = 49^{\circ}$$

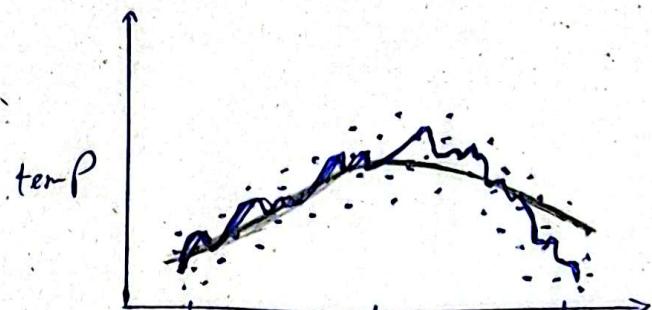
$$\theta_3 = 45^{\circ}$$

:

$$\theta_{180} = 60^{\circ}$$

$$\theta_{181} = 56^{\circ}$$

:



$$\rightarrow \beta = 0.98 \text{ days}$$

$$\rightarrow \beta = 0.9$$

$$v_0 = \theta$$

$$v_1 = 0.9 v_0 + 0.1 \theta_1$$

$$v_2 = 0.9 v_1 + 0.1 \theta_2$$

$$v_3 = 0.9 v_2 + 0.1 \theta_3$$

:

$$v_t = 0.9 v_{t-1} + 0.1 \theta_t$$

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t$$

→ This is the general form for the exponentially weighted averages algorithm that sometimes performs better than Gradient descent.

$v_t$  averages over approximately  $\frac{1}{1-\beta}$  training examples

so if  $\beta = 0.9$  then it's averaging over  $\frac{1}{1-0.9} = 10$  days

→ if you increase the value of  $\beta$  then you're averaging over more past training examples.

$$\beta = 0.98 \rightarrow \frac{1}{1-0.98} = \underline{50 \text{ days}}$$

## Understanding exponentially weighted averages -

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

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

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

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

$$V_{100} = 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 (0.1 \theta_{98} + 0.9 ( \dots \dots \dots$$

$$\begin{aligned} V_{100} &= 0.1 \theta_{100} + 0.1 * 0.9 \cancel{\theta_{99}} + 0.1 * 0.9 * 0.9 \theta_{98} \dots \\ &= 0.1 \theta_{100} + 0.1 (0.9) \theta_{99} + 0.1 (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} \\ &\quad + \dots \dots \dots \end{aligned}$$

So basically exponentially weighted moving averages is just computed the averages of two points with giving respect to the previous points, the number of previous points that you want to compute the e.w.m.a of is going to be set based on the value of  $\beta$  like we said before but because the first value does not have any training examples before it we need to correct its bias.

## Bias Correction

When we implement the E.M.W.A we compute  $v_t$  as follows

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1-\beta) \theta_1$$

$$\text{but } v_0 = 0$$

$$\text{so } v_1 = (1-\beta) \theta_1$$

$$v_2 = \beta v_1 + (1-\beta) \theta_2$$
  
very low estimate

So when we are computing the exponentially weighted moving averages and we start with  $v_0 = 0$  then we're estimating a very poor value of  $v_1, v_2, v_3$  till it finally converges at  $v_t(\frac{1}{1-\beta})$  or the  $v_{\text{number of days}}$  we're averaging over. So to correct this bias instead of computing  $v_t$  we'll compute  $\frac{v_t}{(1-\beta^t)}$  so:-

$$t=2 \rightarrow 1-\beta^t = 1-(0.9)^2 = 0.19$$

$$v_t = \frac{v_t}{0.19} = \frac{0.09 \theta_1 + 0.1 \theta_2}{0.19} = 0.473 \theta_1 + 0.526 \theta_2$$

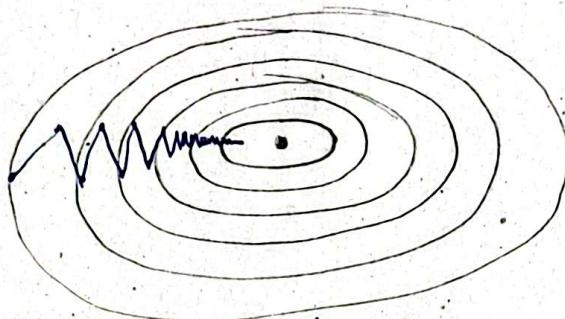
So this modification for the first few examples and later when  $t = 50$  for example the  $(1-\beta^t) = (1-\beta^{50}) = 0.999$  which is close to 1

so  $v_t \approx v_t$  so it corrects over the first few examples then finally goes to  $v_t \approx v_t$

$$t=1000 \rightarrow (1-\beta^{1000}) = 1 \quad \boxed{v_t = \frac{v_t}{1} = v_t}$$

## Gradient descent with momentum

→ a better idea than using GD alone is combining it with exponentially weighted moving averages which makes GD work faster almost every time.



This contour plot shows that when we're using GD it takes a lot of steps to get to the lowest cost and you wouldn't even be able to use a larger  $\alpha \rightarrow$  learning rate to prevent it from overshooting but if we plug in the exponentially weighted moving averages in the process of updating the parameters to take into consideration the parameters of the previous mini-batches parameters this would smooth out the oscillations and yield better results altogether.

on iteration  $t =$

compute  $d_w, d_b$

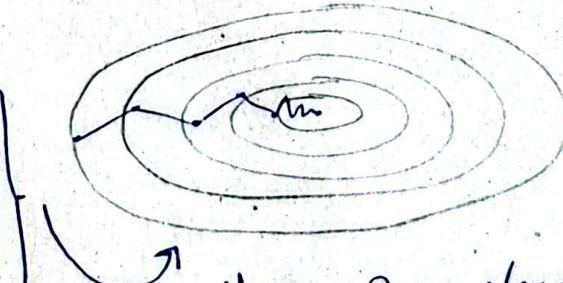
then  $v_{dw} = 0, v_{db} = 0$

$$v_{dw} = \beta v_{dw} + (1-\beta) dw$$

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

$$\omega = \omega - \alpha v_{dw}$$

$$b = b - \alpha v_{db}$$



Hyperparameters  
 $\alpha, \beta$

## RMS Prop $\rightarrow$ (Root mean squared error prob)

on iteration  $t$ :

Compute  $d_w, d_b$  on  $t$ :

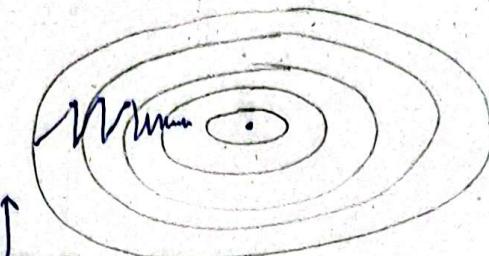
$$sdw = \beta_2 sdw + (1-\beta_2) dw^2$$

$$sdb = \beta_2 sdb + (1-\beta_2) db^2$$

$$w = w - \alpha \frac{dw}{\sqrt{sdw} + \epsilon}$$

$$b = b - \alpha \frac{db}{\sqrt{sdb} + \epsilon}$$

$$\epsilon = 10^{-8}$$



in this plot we wanted to slow the learning in direction  $b$  and maintain the same speed or increase it in direction  $w$  to get to the lowest Cost faster.

$\rightarrow$  so our original goal was to minimize the learning in the direction  $b$ , so we wanted to decrease the value of  $b$  which is why we divided it by  $\sqrt{sdb}$

because:  $\sqrt{sdb}$  here is a large number so:-

$b = b - \alpha \frac{db}{\sqrt{sdb}}$  is dividing  $db$  by a large number

so  $b - \alpha \left( \frac{db}{\text{large number}} \right) = \frac{db}{\text{large num}} = \text{small num}$

$\rightarrow b - \alpha (\text{small num})$

So overall we're taking smaller steps in the  $b$  direction hence the updates to  $b$  are small and the updates to  $w$  are big which smooths out the learning like in momentum.

$\rightarrow$  Finally we add  $\epsilon = 10^{-8}$  to  $b = \alpha \frac{db}{\sqrt{sdb} + \epsilon}$  just in case  $sdb$  or  $sdw$  are very small you don't end up exploding the gradients.

ADAM → adaptive moment estimation.

Combining the exponentially weighted moving averages with RMS prop make an algorithm which is even better than both of them separately.

$$\nabla dW = 0, SdW = 0, \nabla bd = 0, Sdb = 0$$

on iteration 0.

Compute  $dW, db$  on current t:-

$$\begin{aligned} \nabla dW &= \beta_1, \nabla dW + (1 - \beta_1) dW \\ \nabla bd &= \beta_1, \nabla db + (1 - \beta_1) db \end{aligned} \quad \left. \begin{array}{l} \text{momentum} \\ \text{RMS prob} \end{array} \right\}$$

$$\begin{aligned} SdW &= \beta_2 SdW + (1 - \beta_2) dW^2 \\ Sdb &= \beta_2 Sdb + (1 - \beta_2) db^2 \end{aligned} \quad \left. \begin{array}{l} \text{RMS prob} \\ \text{RMS prob} \end{array} \right\}$$

$$\nabla_{dW}^{\text{corrected}} = \frac{\nabla dW}{(1 - \beta_1)^t}, \nabla_{db}^{\text{corrected}} = \frac{\nabla db}{(1 - \beta_1)^t}$$

$$S_{dW}^{\text{corrected}} = \frac{SdW}{(1 - \beta_2)^t}, S_{db}^{\text{corrected}} = \frac{Sdb}{(1 - \beta_2)^t}$$

Finally,

$$w = w - \alpha \frac{\nabla_{dW}^{\text{corrected}}}{\sqrt{S_{dW}^{\text{corrected}}}} + \epsilon \quad \rightarrow 10^{-8}$$

and same for  $b$

→ This algorithm is the combination of GD, exponentially weighted moving average, and lastly RMS prop.

Hyperparameters  $\alpha, \beta_1, \beta_2, \epsilon$ ,

$$\begin{array}{ll} 0.9 & 0.999 \\ dW & dW^2 \end{array} \quad 10^{-8}$$

- \* So far the Hyper Parameters that we've come across include:-

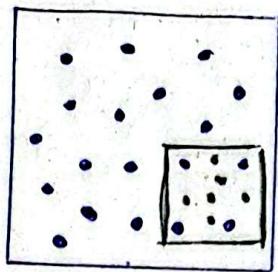
- $\alpha$
- $\beta_1, \beta_2, \epsilon$
- number of layers
- number of hidden units for the layers
- Learning rate decay
- mini-batch size.

- \* in the tuning process some hyperparameters are more important to tune than others
  - $\alpha$  is the most important hyperparameter to tune
  - next in importance :-
    - \* if we're using momentum then  $\beta$  is a good start hidden units.  
then maybe the mini-batch size and the number of layers
    - \* if we're using adam then the default values for  $\beta_1, \beta_2$  and  $\epsilon$  are pretty much ok and you don't have to tune them.
  - Lastly the least important hyperparameters are the Learning rate decay, the number of layers and the  $\beta, \beta_2, \epsilon$ ; if you're using ADAM.

- \* in the hyperparameter tuning process it's better to search or try out values at random instead of using Grid search because based on the importance of the change in value of the hyperparameters, Grid search could be wasting a lot of time and computational power.

## Coarse to fine sampling scheme :-

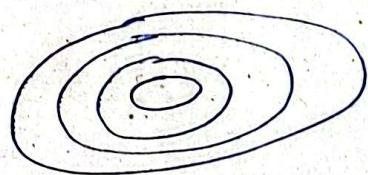
- \* In the Coarse to fine scheme we look at some randomly sampled values for hyperparameters and if we find some points that are near to each other and result in Lowering the cost function or reducing the time i.e. better hyper parameters then we might want to zone in and look for better values in that selected region.



## Normalizing activations :-

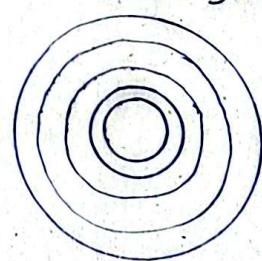
- it has been said before that normalizing the input  $X$  speeds up the training process because it helps GD to converge faster as it smooths and rounds the contours of the cost function plot hence making it more convex or so

From this shape :

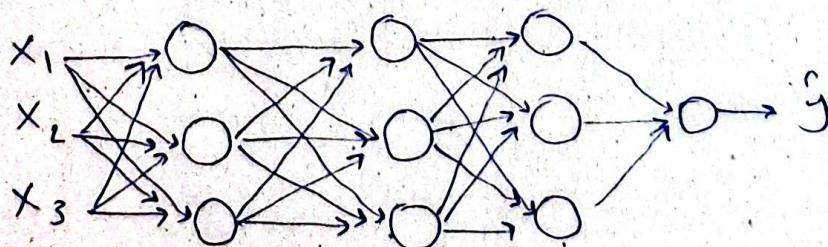


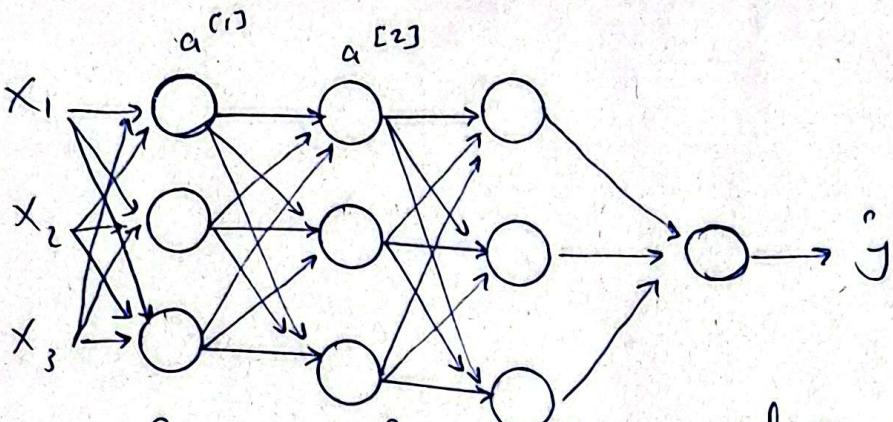
to

this shape



- Wouldn't also help to normalize the output of each and every layer as they are the input to the next layer so it also speeds up the training process :- each layer





- we found before that normalizing  $X$  would help Gradient descent find the best parameters for  $w^{[1]}$  and  $b^{[1]}$ , so wouldn't normalizing  $Z^{[2]}$  also help speeding up the search for parameters  $w^{[2]}$  and  $b^{[2]}$  that yield the lowest cost.
- This is called batch normalization

### Implementing batch norm.

\* Given some intermediate values in the NN  $Z^{(1)}, \dots, Z^{(n)}$

$$\mu = \frac{1}{m} \sum_{i=1}^m Z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (Z^{(i)} - \mu)^2$$

$$Z_{\text{norm}}^{(i)} = \frac{Z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

now all the  $Z$  values have  $\mu = 0$  and  $\sigma^2 = 1$

→ now we add some parameters to control the distribution of our inputs

$$\tilde{Z}^{(i)} = \gamma Z_{\text{norm}}^{(i)} + \beta$$

\* Gamma  $\gamma$  and beta  $\beta$  are learnable parameters.

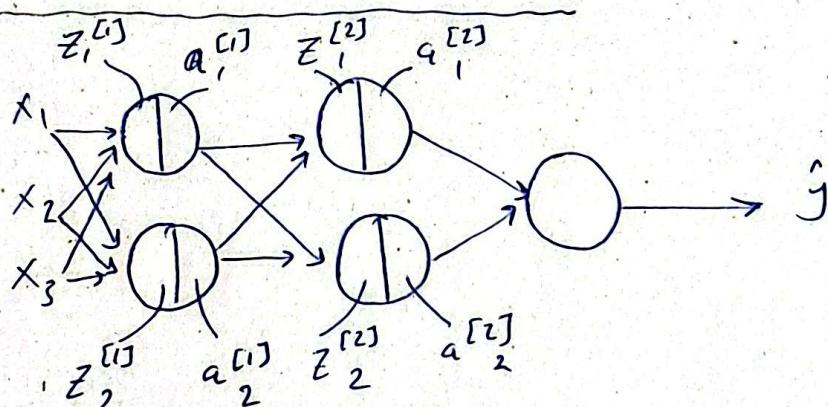
\* how does having the parameters  $\gamma$  and  $\beta$  control the distribution of our inputs?

↳ for instance if we set  $\gamma = \sqrt{\delta^2 + \epsilon}$  and  $\beta = \mu$   
then

$$\tilde{Z}_{\text{norm}} = \sqrt{\delta^2 + \epsilon} \frac{Z_{\text{norm}} - \mu}{\sqrt{\delta^2 + \epsilon}} + \mu = Z_{\text{norm}}$$

which is the original distribution so the parameters  $\gamma$  and  $\beta$  control the mean and the standard deviation i.e. the spread of the data points and their center.

Fitting batch norm in a NN



$$x \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{Batch norm}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \xrightarrow{g^{[1]}(\tilde{z}^{[1]})} a^{[1]}$$

$$a^{[2]} \xleftarrow{g^{[2]}(z^{[2]})} z^{[2]} \xleftarrow[\text{Batch norm}]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]} \xleftarrow{w^{[2]}, b^{[2]}} z^{[2]}$$

Parameters  $\rightarrow w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$   
 $\beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]}$

## Working with mini batches

- \* since we're using minibatches now, we have to redesign our batch normalization for it to work on our mini batches.

$$x^{[1]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[B.N.]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \xrightarrow{g^{[1]}(\tilde{z}^{[1]})} a^{[1]} \rightarrow \text{etc...}$$

First minibatch

$$x^{[2]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[2]} \xrightarrow[B.N.]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[2]} \xrightarrow{g^{[1]}(\tilde{z}^{[2]})} a^{[2]} \rightarrow \text{etc...}$$

⋮

$$x^{[m]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[m]} \xrightarrow[B.N.]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[m]} \xrightarrow{g^{[1]}(\tilde{z}^{[m]})} a^{[m]} \rightarrow \text{etc...}$$

- \* and that is called one epoch.

- since when we're rescaling  $z^{[l]}$  by  $\beta^{[l]}$  and  $\gamma^{[l]}$  and we're subtracting the  $\beta^{[l]}$  and setting  $z^{[l]}$  to  $\mu=0$  and  $\sigma^2=1$  or so, we're omitting the  $b$  value since we're subtracting a raw number.
- since  $z^{[l]}$  has a  $\mu=0$  then  $b^{[l]}$  does not matter anymore so we're gonna drop it from the parameters.

Parameters  $\rightarrow w^{[l]}, \beta^{[l]}, \gamma^{[l]}$

$$z^{[l]} = w^{[l]} * a^{[l-1]}$$

$$\boxed{\tilde{z}^{[l]}_{\text{norm}} = \gamma^{[l]} z^{[l]} + \beta^{[l]}}$$

use these values  $d\gamma^{[l]}, d\beta^{[l]}, d\mu^{[l]}$  to compute Adam or GD.

for t=1 ... run minibatches

→ compute forward prop

→ compute B.N to replace  $z^{[l]}$  to  $\tilde{z}^{[l]}$

→ compute back prop

$$* w^{[l]} = w^{[l]} - \alpha d w^{[l]} \rightarrow ①$$

$$* \beta^{[l]} = \beta^{[l]} - \alpha d \beta^{[l]} \rightarrow ②$$

$$* \gamma^{[l]} = \gamma^{[l]} - \alpha d \gamma^{[l]} \rightarrow ③$$

## Batch norm at test time

→ since mini batch normalization works with multiple data points to figure out  $\mu$  and  $\sigma^2$  so it can normalize the distribution of this specific minibatch and this not an option at test since we're dealing with individual data points so we have to compute with an estimate for the mean  $\mu$  and the S.D  $\sigma$  to normalize the test point so that the classifier i.e. NN can handle them correctly.

→ First we use the exponentially weighted averages across all minibatches for  $\mu$  and  $\sigma^2$ .

$$X^{[1]}, X^{[2]}, X^{[3]} \rightarrow L=1$$

$$\downarrow \mu^{[1]}, \mu^{[2]}, \mu^{[3]}$$

these are the means for the minibatches.  
now we calculate the exponentially weighted average for  $\mu$  and we're going to end up with some value for the last mean.

$$\tilde{\mu}_0 = 0, \tilde{\mu}_t = \beta \tilde{\mu}_0 + (1-\beta) \tilde{\mu}_0$$

$$\tilde{\mu}_t = \beta \tilde{\mu}_{t-1} + (1-\beta) \tilde{\mu}_t$$

until we reach:

$$\boxed{\tilde{\mu}_{\text{num-mini-batches}} = \tilde{\mu}_n = \beta \tilde{\mu}_{n-1} + (1-\beta) \tilde{\mu}_n}$$

→  $n$  referring to the number of minibatches

\* We do the same for  $\sigma^2$

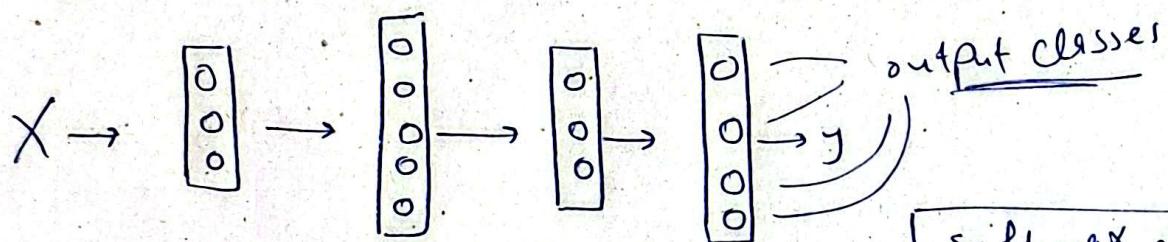
then we use the values that we got from this process  
\* meaning the last  $\mu$  and  $\sigma^2$  and compute.

$$Z_{\text{norm}}^{(i)(c)} = \frac{Z^{(i)(c)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \rightarrow \tilde{Z}^{(i)(c)} = \gamma \cdot Z_{\text{norm}}^{(i)} + \beta$$

↳ this  $\tilde{Z}$  is calculated before each layer in the Linear forward process per test sample, so it is normalized just like we did to the training data..

### Soft max function

→ soft max function is used when the classes of the model are more than 2  $\rightarrow > 2$  so we can't use the sigmoid any more or use the binary cross entropy in general



$$Z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

Activations  $\rightarrow$

$$t = e^{(Z^{[L]})}$$

$$a^{[L]} = \frac{e^{(Z^{[L]})}}{\sum_{j=1}^n t_j}, \quad a_i^{[L]} = \frac{t_i}{\sum_{j=1}^n t_j}$$

output layer units

$$Z^{[1]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.9 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$\sum_{i=1}^4 t_i = 176.3 \rightarrow a^{[L]} = \frac{t}{176.3} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

→ the model classifies that it's gonna be the first class with probability of being  $\frac{1}{4} = 84.2\%$

soft max algorithm takes the output of the final layer and maps it to some Probabilities.

## Loss function for soft max

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

Our trained algorithm outputs that the class is a cat meanwhile it actually isn't, so how would we calculate the cost in this case?  
 → in Logistic regression we used to calculate the cost as follows:

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

This used to work because we only had 2 classes but now we have number of classes = C.

→ the Loss function for the softmax algorithm:

$$L(\hat{y}, y) = - \sum_{j=1}^C y_j \log \hat{y}_j$$

\* In the previous case we had  $y = [0 \ 1 \ 0 \ 0]$

$$\begin{aligned} \text{so } L(\hat{y}, y) &= - (0 \log \hat{y}_1) - y_2 \log \hat{y}_2 - 0 \log \hat{y}_3 - 0 \log \hat{y}_4 \\ &= - y_2 \log \hat{y}_2 = \underline{-\log \hat{y}_2} \end{aligned}$$

So the only way to minimize  $L(\hat{y}, y)$  the cost function is to maximize the value of  $\hat{y}_2$  since  $-\log(1) = 0$

so this is exactly what happens when we're training the model; the only way to minimize the cost function is to maximize the probability of the true class.

$$J = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}) \quad \left| \begin{array}{l} Y^{(4, m)} \\ Y^{(4, m)} \end{array} \right.$$

(w, b, ...)