

# DI.Ai\_Optimization algorithms

Optimization algorithms are what we use to speed up the process of learning though they may affect the accuracy of the model sometimes, they are not meant to increase the accuracy of the model usually

## a) Mini Batch Gradient Descent

$$X = \underbrace{[x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)}]}_{\substack{\{1\} \\ X}} \ \underbrace{\left[ \begin{array}{c|c|c} x^{(1001)} & \dots & x^{(2000)} \end{array} \right]}_{\substack{\{2\} \\ X}} \ \dots \ \underbrace{\left[ \begin{array}{c|c|c} \dots & \dots & x^{(m)} \end{array} \right]}_{\substack{\{m\} \\ X}} \quad (n_x, m)$$
$$Y = \underbrace{[y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)}]}_{\substack{\{1\} \\ Y}} \ \underbrace{\left[ \begin{array}{c|c|c} y^{(1001)} & \dots & y^{(2000)} \end{array} \right]}_{\substack{\{2\} \\ Y}} \ \dots \ \underbrace{\left[ \begin{array}{c|c|c} \dots & \dots & y^{(m)} \end{array} \right]}_{\substack{\{m\} \\ Y}} \quad (1, m)$$

- In batch gradient descent all the instances are vectorized to form a vector  $x$  with dimension  $n \times m$  where  $n$  is the no of features and  $m$  the no of instances
- each instance is represented as  $x^{(i)}$  we will create a new vector  $x^{\{i\}}$  comprising of  $x^{(i)}$  to  $x^{(1000)}$  which is called as our minibatch

for  $t = 1, \dots, 5000$

use  $X$   
(as if mini)

Forward prop on  $X^{t_{00}}$ .

$$z^{[l]} = W^{[l]} X^{t_{00}} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

Vectorized implementation  
(1000 example)

$$\text{Compute cost } J^{t_{00}} = \frac{1}{1000} \sum_{i=1}^L l(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$$

Backprop to compute gradients wrt  $J^{t_{00}}$  (using  $(X^{t_{00}}, Y^{t_{00}})$ )

$$W^{[l]} := W^{[l]} - \alpha \delta W^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha \delta b^{[l]}$$

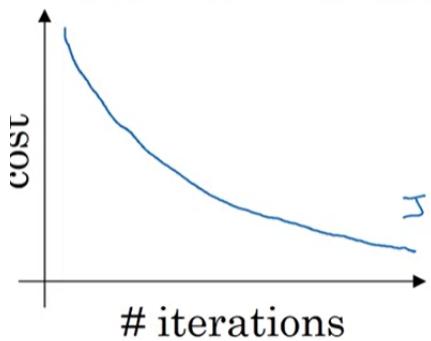
}

"1 epoch"

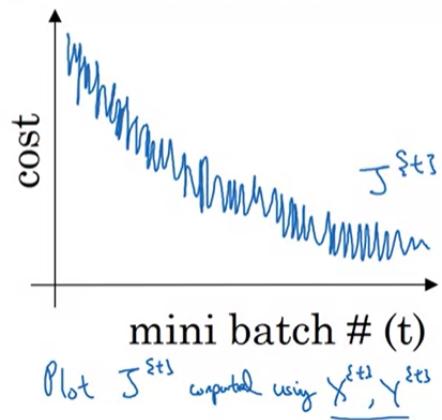
... much training set

- first we must run backward propagation and forward propagation for each mini batch individually
- so we use for  $t$  in range (1000):
- we run our forward propagation steps just for the mini batch we have  $z^{[1]}$  is  $t[2] \times 5000$
- We compute the loss function just for this mini batch, then use  $J^{t_{00}}$  to update weights of the layer for the neural network
- Note that  $z^{[l]}$  and  $a^{[l]}$  are separate for each mini batch but the weights and biases of each layer is accessed by all mini batches.
- once all mini batches has been accessed atleast once a epoch is completed
- The weight and bias update we do for each minibatch is called a gradient step. Each epoch of mini batch gradient descent consists of 5000 gradient steps compared to 1 step in batch gradient descent.
- Mini batch converges faster than a batch gradient descent.

Batch gradient descent

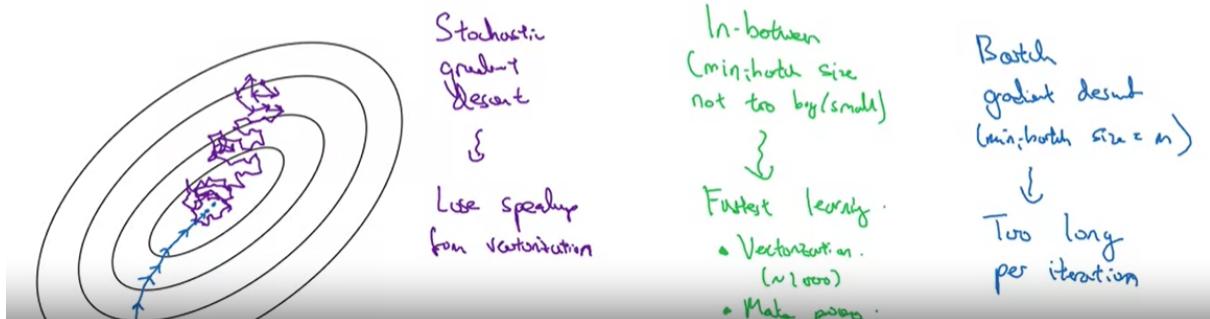


Mini-batch gradient descent



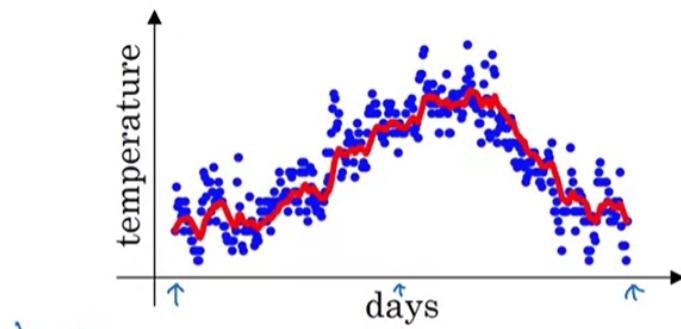
the cost vs mini batch is highly noisy because of obvious reasons

- Stochastic is not the fastest model because we loose the speedup we get by using vectorization.



## b) Exponentially weighted averages

- this is just a statistical concept for now



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

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

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

:

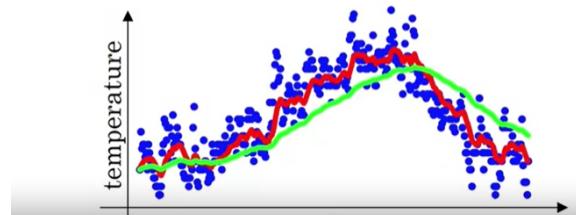
$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

- we take a moving average of all existing instances but giving more weightage to recent values.
- even though technically all pre existing values are considered in reality only the last

$\frac{1}{1-\beta}$  will affect the average.hence its said to be a moving average

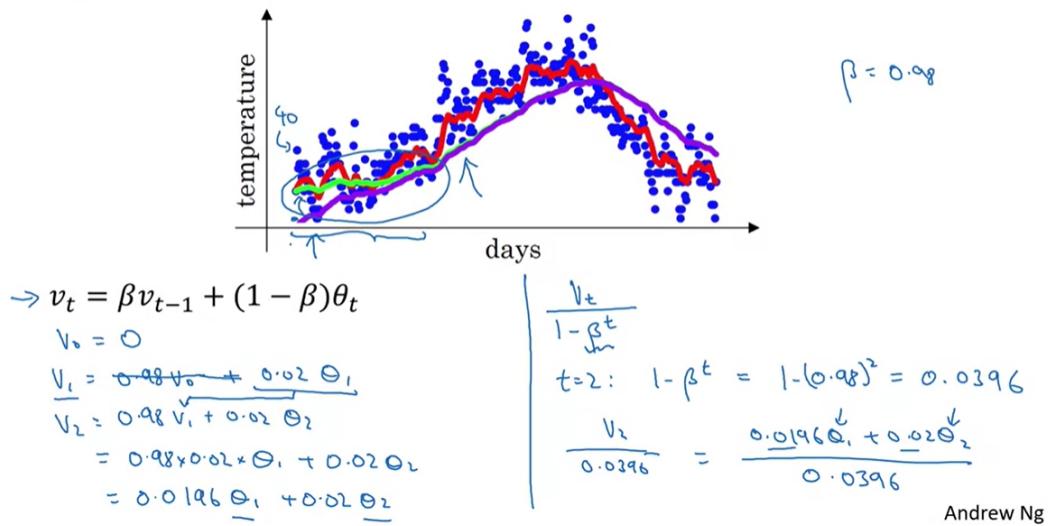
- Higher value of beta means that the average is taken over a larger window leading to more smoothening out

$$\begin{aligned} \beta = 0.9 & : \approx 10 \text{ day tapering} \\ \beta = 0.98 & : \approx 50 \text{ day} \end{aligned}$$



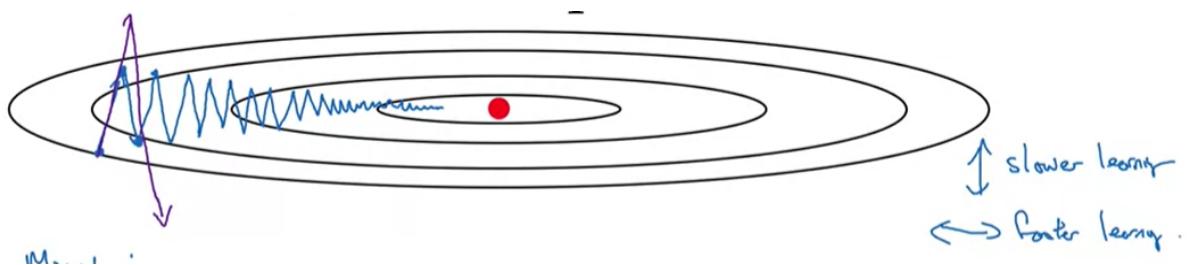
- It can be observed that for higher values such as b=0.98 the average takes place over an window of 50 instances the initial value of the average will be very low ( first 50 instances) hence we use bias correction

## Bias correction



- we use  $\frac{V^t}{1-\beta^T}$  instead of  $V_t$  it will converge to  $V_t$  for higher values of  $t$

## c) Momentum



In this case we need to have faster learning in the horizontal direction and slower learning in the vertical direction. We use Momentum which is a application of exponentially weighted average.

- If we can take the average of the individual steps we could neglect up and down oscillations.
- Hence we use exponentially weighted averages instead of the dw and db (i.e the gradients of loss function)

On iteration  $t$ :

Compute  $\Delta w, \Delta b$  on current mini-batch.

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

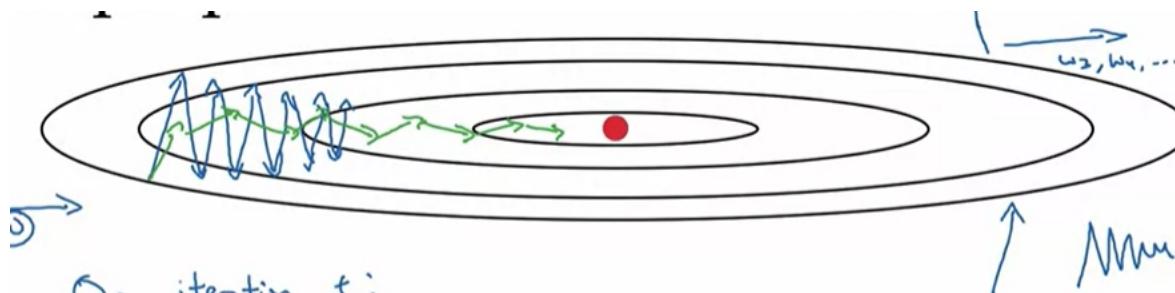
$$v_{db} = \beta v_{db} + (1-\beta) \underline{\Delta b}$$

$\nwarrow$  velocity       $\nearrow$  acceleration

$$w := w - \alpha v_{dw}, \quad b := b - \alpha v_{db}$$

- Note we average out  $\Delta w$  from different minibatches in this case but we can also do it from different iterations.
- we can consider that the  $v_{dw}$  represents the velocity of a body beta the friction that applies on it  $\Delta w$  as the new acceleration it gets.
- $\beta$  is a hyperparameter like  $\alpha$

## d) RMSprop



On iteration  $t$ :

Compute  $\Delta w, \Delta b$  on current mini-batch

$$\underline{s_{dw}} = \beta_2 \underline{s_{dw}} + (1-\beta_2) \underline{\Delta w^2} \quad \leftarrow \text{element-wise}$$

$$\rightarrow \underline{s_{db}} = \beta_2 \underline{s_{db}} + (1-\beta_2) \underline{\Delta b^2} \quad \leftarrow \text{large}$$

$$w := w - \alpha \frac{\underline{\Delta w}}{\sqrt{\underline{s_{dw}} + \epsilon}} \quad \leftarrow$$

$$b := b - \alpha \frac{\underline{\Delta b}}{\sqrt{\underline{s_{db}} + \epsilon}} \quad \leftarrow$$

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

- RMS prop or root mean square prop is different from momentum in the sense that it solves the issue of oscillation by adapting the learning rate according to the parameter, instead of changing the differential.

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

Rmsprop      Momentum

- we use  $dw^2$  and  $db^2$  since  $db$  is very high and  $dw$  is very low (change in  $b$  is more than change in  $w$ )  $S_{dw}$  is smaller hence  $w$  is learned more compared to  $b$ .
- we add the term  $E$  to ensure that we don't divide by zero

## c+d) Adam Optimizer

- When we use both RMS prop and Momentum together we get adam optimizer which stands for adaptive moment optimizer.

$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$

On iteration  $t$ :

Compute  $dw, db$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \leftarrow \text{"moment"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \leftarrow \text{"RMSprop"} \beta_2$$

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

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

$$w := w - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}$$

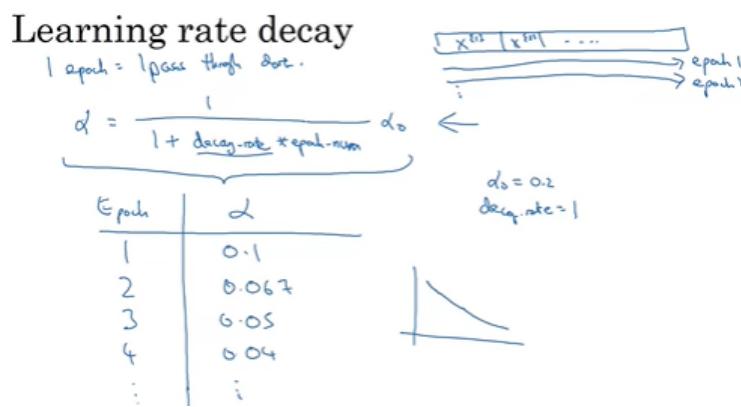
$$b := b - \alpha \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

- Hyperparameters for this model are  $\alpha, \beta_1, \beta_2, E$

## e) Learning Rate Decay



- we can decrease the learning rate over epochs to reach convergence faster

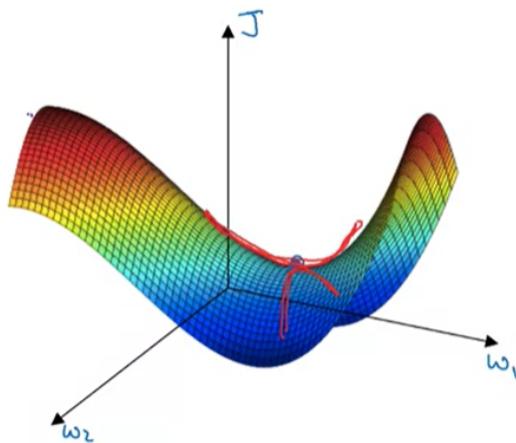


- decay is also an hyperparameter

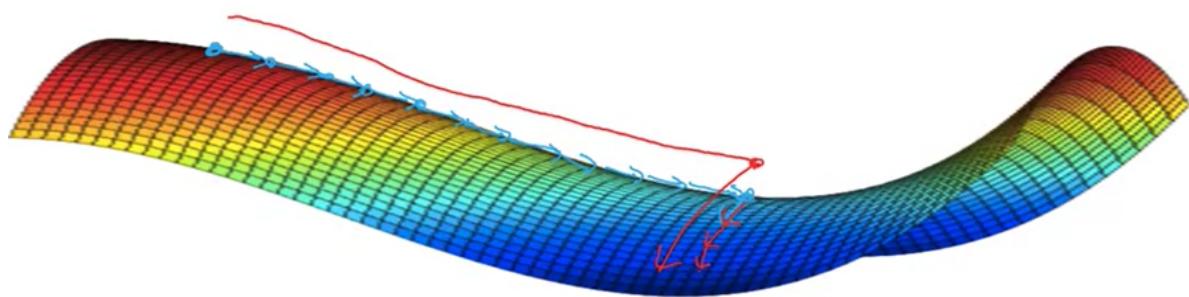
## f) Local optima

- unlike ml models are have low dimensional, dl models deal with high dimensional data
- Lets say our model deals with 20000 parameters a local minima can occur in our cost vs parameters hyper plot only if a minima occurs in the sectional plots of each dimension
- The chance of all these sectional plot having a minima is  $\frac{1}{2}^{20000}$  which is extremely low.
- One global minima will exist but the chances of our model running into a local minima is negligible.

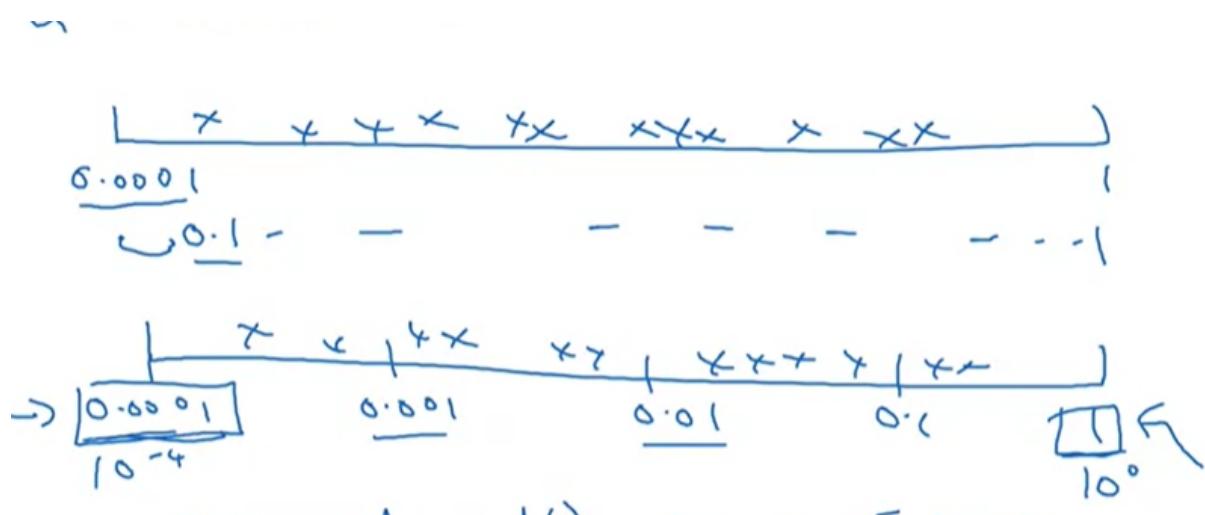
- Instead our model can run into multiple saddle points where some of the parameters reach a minima but all of the do not.



- plateaus instead slowdown our models



## g) Hyperparameter Tuning



$$\beta = 0.9 \dots 0.999$$

$\downarrow$        $\downarrow$

10      1000

---

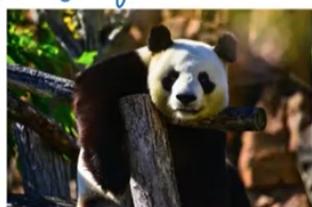
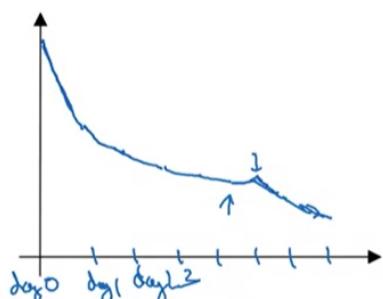

$$-\beta = 0.1 \dots 0.001$$

$\beta: 0.900 \rightarrow 0.9005$

$\beta: 0.999 \rightarrow 0.9995$

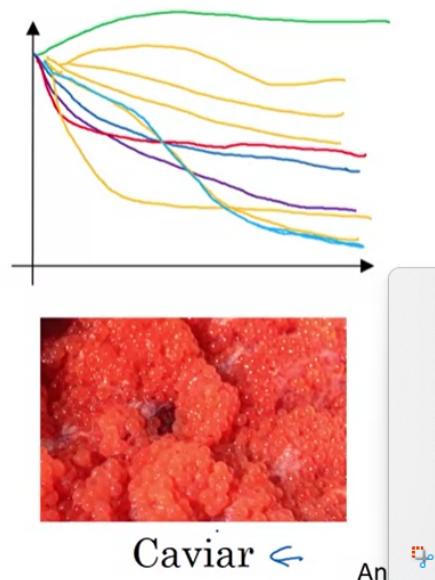
$$\begin{aligned} & \overbrace{\quad \quad \quad \quad \quad}^{\leftarrow} \\ & 0.9 \qquad \qquad \qquad 0.999 \\ & \overbrace{\quad \quad \quad}^{\leftarrow} \qquad \qquad \qquad \overbrace{\quad \quad \quad}^{\leftarrow} \\ & 0.9 \qquad 0.99 \qquad 0.999 \\ & \overbrace{\quad \quad \quad}^{\leftarrow} \qquad \qquad \qquad \overbrace{\quad \quad \quad}^{\leftarrow} \\ & 0.1 \qquad 0.01 \qquad 0.001 \\ & \overbrace{\quad \quad \quad}^{\leftarrow} \qquad \qquad \qquad \overbrace{\quad \quad \quad}^{\leftarrow} \\ & 10^{-1} \qquad \qquad \qquad 10^{-3} \\ & r \in [-3, -1] \\ & 1 - \beta = 10^r \\ & \beta = 1 - 10^r \end{aligned}$$

## Babysitting one model



Panda  $\leftarrow$

## Training many models in parallel



Caviar  $\leftarrow$

An

## f) Batch Normalization

- In batch normalization we normalize the values of  $z^i$  which come out of a layer.

Given some intermediate values in NN

$$\underbrace{z^{(1)}, \dots, z^{(n)}}_{z^{[l]}(:)}$$

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

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

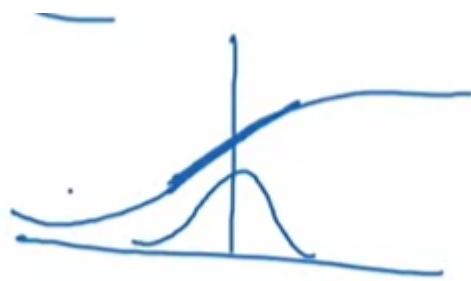
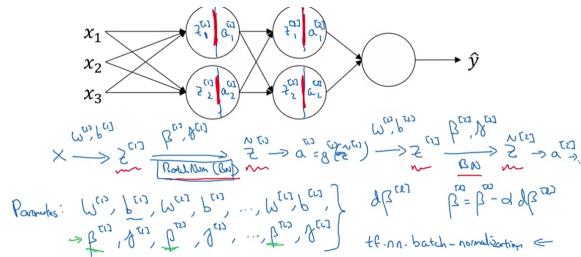
$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

- But we don't use z norm directly instead we use

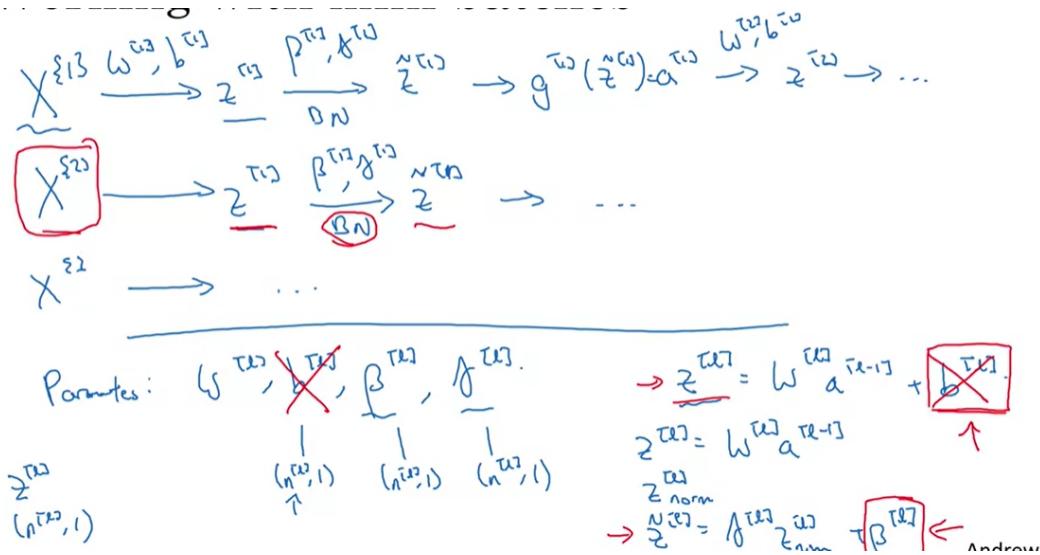
$$\hat{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

where gamma and beta are **Scale and Shift**-learnable parameters of the model

- we can't directly use , z norm is centered at z=0 with variance 1, hence only the linear part of the activation function will be used which makes the model linear.



- adding batch normalisation to a batch gradient descent
- when we use batch normalization bias becomes redundant as z norm won't change when bias is changed(if you add or subtract a constant mean remains same).
- Beta has same dimension as b,z,z norm which is  $[n(l), m]$  j must have  $[n(l), n[l(l)]]$  but its somehow different in coursera



- Covariance shift-occurs when the task is same but the model is not able to generalize to a different distribution of data.e.g)black cat vs dogs and white cat vs dogs
- if we learn a mapping  $X \rightarrow Y$  if distribution of  $X$  changes then we have to retrain  $y$ .
- occurs because the inputs received by the hidden layers in the middle vary a lot leading to them being highly dependent on the initial layers
- Batch normalizing ensures that the inputs that each layer receives always have a fixed mean and variance i.e more stable making the later layers more robust.
- Batch Normalization does different changes in different mini batches leading to slight regularization effect.
- When we Normalize the values in input layer it is done outside of the model as a transformation on the input data. Hence it can be done both in test and train data
- However batch normalization requires the data of all instances in a batch and if we don't normalize the values of  $z$  during test time the outputs may change a lot.

$$\begin{aligned} \mu &= \frac{1}{m} \sum_i z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \\ z_{\text{norm}}^{(i)} &= \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ \tilde{z}^{(i)} &= \gamma z_{\text{norm}}^{(i)} + \beta \end{aligned}$$

$\mu, \sigma^2$ : estimate using exponentially weighted average (across mini-batches).  
 $x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots$   
 $\downarrow$   
 $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}, \dots$   
 $\theta_1, \theta_2, \theta_3, \dots$   
 $\tilde{z}_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}}$   
 $\tilde{z} = \gamma z_{\text{norm}} + \beta$

Andrew Ng