# C2W3- Batch normalization

Normalization makes learning faster as it modified the contour line of cost function as a more round shape

$$\mu = \frac{1}{m} \sum_i x^{(i)}$$
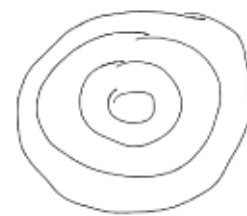
$$x = X - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i (x^{(i)})^2$$

$$X = X / \sigma^2$$

not normalized

normalized

<u>Batch normalization</u>:- previously we only normalized the input feature. But in batch normalization we will also normalize the activation unit (Z) of each neural node.

<u>Implementation</u>:-

For each hidden unit we will calculate $\tilde{z}$ as follows:-

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

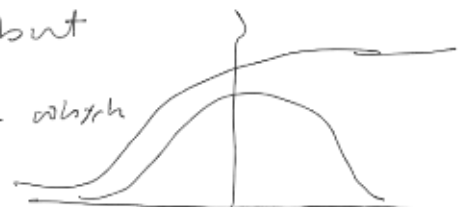$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

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

$$\tilde{z}^{(i)} = \gamma \, z_{norm} + \beta \quad \Big| \quad \gamma, \beta = \text{learnable parameters}$$

for next computation we will use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

the reason of using $\gamma$ & $\beta$ :-

we don't want to force the mean of $\tilde{z}^{(i)}$ to be zero, we might loose the nonlinearity of the sigmoid function $\gamma$ & $\beta$ doesn't set the mean=0 but it standarize the mean & variance which we can control now

# Adding batch norm to a network:-



$$x \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\substack{Batch \\ normalization}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \to a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{w^{[2]}, b^{[2]}} z^{[2]}$$

$$\beta^{[l]} = \beta^{[l]} - \alpha \, d\beta^{[l]}$$

# Working with minibatch:-

$$x^{\{1\}} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \to g(\tilde{z}^{[1]}) \to \ldots$$

$$x^{\{2\}} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\boxed{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \to g(\tilde{z}^{[1]}) \to \ldots$$

<span style="color:red">→ this BN takes the mean & variance from this $z^{[1]}$ only.</span>

$$\overline{z^{[l]}} = w^{[l]} a^{[l-1]} + \cancel{b^{[l]}}$$

<span style="color:red">→ During normalization of minibatch all constant will be canceled out. so this b will have no effect.</span>

<span style="color:red">if we implement batch norm we can eleminate the parameter b.</span>

implementing Gradient descend:-

for t = 1 ......... num of minibatches

   compute forward prop on $x^{\{1\}}$.

   In each hidden layer use BN to replace $z^{[1]}$ with $\tilde{z}^{[1]}$

   use backprop to compute $dw^{[1]}$, ~~$db^{[1]}$~~, $d\beta^{[1]}$, $d\gamma^{[1]}$

   Update parameter:-

$$w^{[1]} := w^{[1]} - \alpha \, dw^{[1]}.$$
$$\beta^{[1]} := \beta^{[1]} - \alpha \, d\beta^{[1]}.$$
$$\gamma^{[1]} := \gamma^{[1]} - \alpha \, d\gamma^{[1]}$$

This also works with momentum, adam/ RMSprop.

## Why batch norm works:-

① Makes the mean zero and provide a better shape for the contour line of cost function
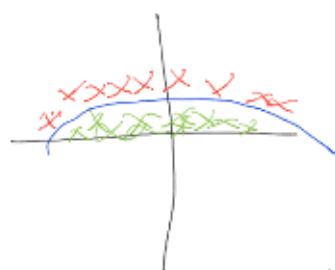
② Reduce the effect of covariate shift.

③ Get slight regularization effect as a side-effect

if we trained our network to identify only black cat then if we try to identify the non black cat then it will not perform better.
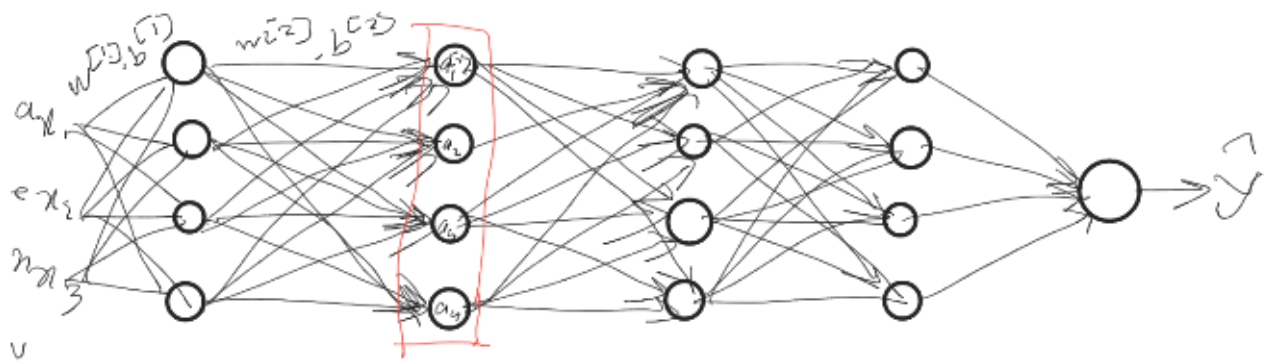


black cat                  non-black cat

:Even if black cat & non black cat use almost same function for decision boundary. We cannot use black cat model to non black cat model because of covariate shift. The idea of co-variate shift is if if we use mapping $(x \rightarrow y)$ & if the distribution of $x$ changes then we might need to retrain algorithm

# why co-variate shift is a problem of NN:-



if $w_{1/2}$ & $b_{1/2}$ changes then value of $a_i$ also changes that's co-variate shift. But if we use batch normalization then effect of covariate shift doesn't effect that much. Doesn't matter ow $a_i$ changes batch norm makes sure it's mean & variances stays the same

# Batch norm as regularization effects:-

• Each mini batch is scaled by mean/variance computed on just that minibatch.

• This adds some noise to the value of $z^{[l]}$ within that minibatch. So similar to dropout it adds some noise to each hidden layer activations.

• This has a slight regularization effects.

Regularization is just a side effect of batch norm.

its main job is the normalizing the data of hidden unit for faster learing.

## Batch norm at test time :-

in train time we have multiple train data so we can calculate mean, standard deviation thus normalize the data. But in test time we have only one example. so to calculate, $u, \sigma^2$ we do :-

Use the exponentially weighted average of $u, \sigma^2$ from minibatches $u^{[l]}$ = exponentially weighted average of $\{ u^{\{1\}[l]}, u^{\{2\}[l]} ... \}$
$\sigma^{2[l]}$ = exponentially weighted average of $\{ \sigma^{2\{1\}[l]}, \sigma^{2\{2\}[l]} ... \}$