

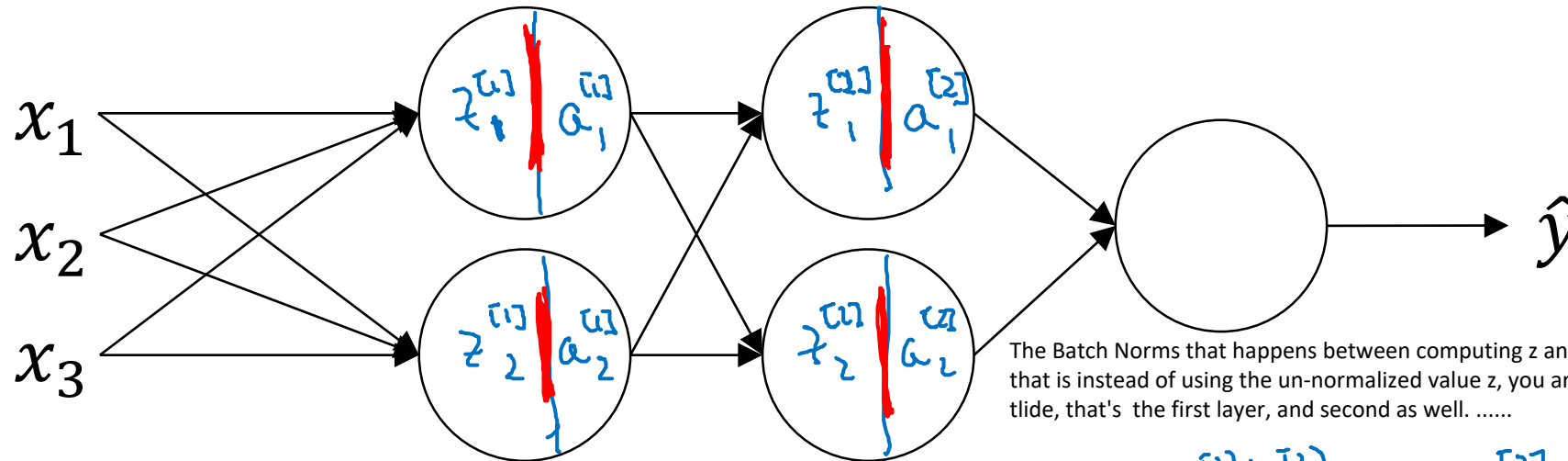


deeplearning.ai

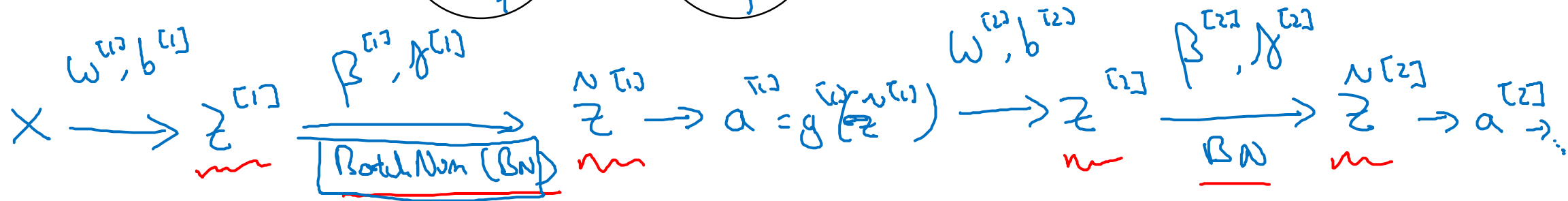
Batch Normalization

Fitting Batch Norm
into a neural network

Adding Batch Norm to a network



The Batch Norms that happens between computing z and computing a . And the intuition of that is instead of using the un-normalized value z , you are gonna use the normalized value z slide, that's the first layer, and second as well.



Parameters: $\left\{ \begin{array}{l} w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}, \\ \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]} \end{array} \right\}$

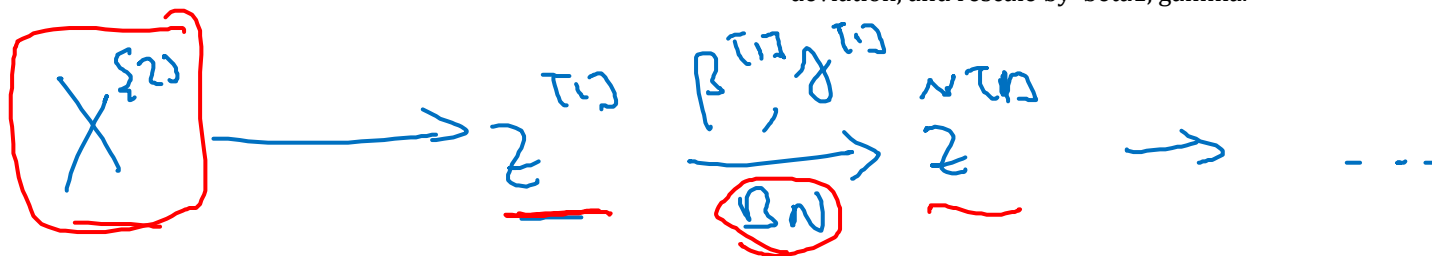
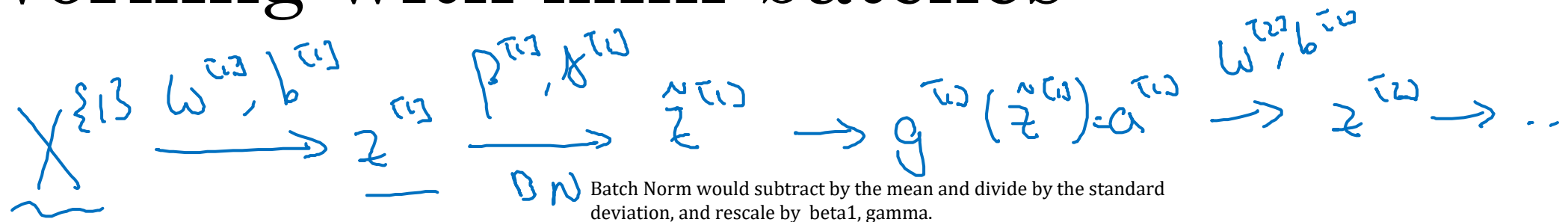
$$d\beta^{[2]} \quad \beta^{[2]} = \beta^{[1]} - \alpha d\beta^{[1]}$$

tf.nn.batch-normalization ←

→ **B** It has nothing to do with the hyperparameter beta that we had for momentum or for computing the variance exponentially weighted average.

Working with mini-batches

Batch Norm is usually applied with mini-batches of your training set. So the way you apply Batch Norm is you take your first mini-batch and compute z_1 ...

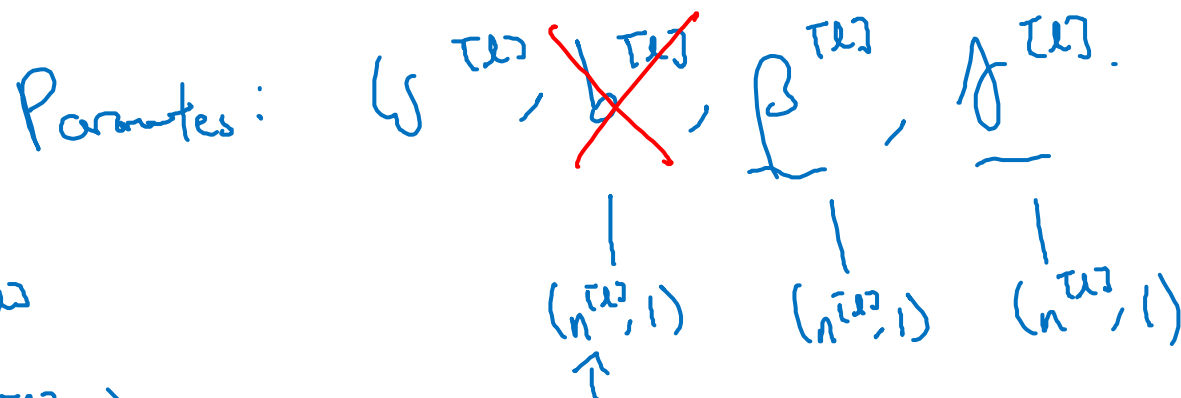


X^{l-1}

It's going to look at the mini-batch and normalize z_1 first to mean 0 and standard variance, and then rescale by β and γ .

any constant you add will get canceled out by mean subtraction step. So if you are using batch norm, you can actually eliminate that parameter. Or if you want, think of it as setting it permanently to 0.

kind of replaced with β_1 which is the parameter that controls that ends up affecting the shift or the bias terms.



$z^{l,l-1}$
 $(n^{l,l-1}, 1)$

because that's the number of hidden units you have and so β_1 and γ_1 are used to scale the mean and variance of each of the hidden units

$$\rightarrow \underline{z^{l,l}} = W^{l,l} a^{l,l-1} + \cancel{b^{l,l}}$$

$$z^{l,l} = W^{l,l} a^{l,l-1}$$

$$z_{\text{norm}}^{l,l} = \gamma^{l,l} z_{\text{norm}}^{l,l}$$

$$\rightarrow \tilde{z}^{l,l} = \gamma^{l,l} z_{\text{norm}}^{l,l} + \beta^{l,l}$$

You end up using this parameter β_1 in order to decide what's the mean of z slide 1.

Andrew Ng

Implementing gradient descent

for $t = 1 \dots \text{num Mini Batches}$

Compute forward pass on $X^{\{t\}}$.

In each hidden layer, use BN to repair $\underline{z}^{\{t\}}$ with $\underline{\hat{z}}^{\{t\}}$.

Use backprop to compute $\underline{dw}^{\{t\}}$, ~~$\underline{db}^{\{t\}}$~~ , $\underline{dp}^{\{t\}}$, $\underline{d\gamma}^{\{t\}}$

Update parameters

$$\left. \begin{aligned} w^{\{t\}} &:= w^{\{t-1\}} - \alpha dw^{\{t-1\}} \\ \beta^{\{t\}} &:= \beta^{\{t-1\}} - \alpha dp^{\{t-1\}} \\ \gamma^{\{t\}} &:= \dots \end{aligned} \right\}$$

← gradient descent

Works w/ momentum, RMSprop, Adam.

And so then it ensures that within that mini-batch, the value z end up with some normalized mean and variance of the version of z slide 1.

you can use the updates given by other optimization algorithms
Some of these other optimization algorithms as well can be used to update the parameters beta and gamma that Batch Norm added to algorithm