



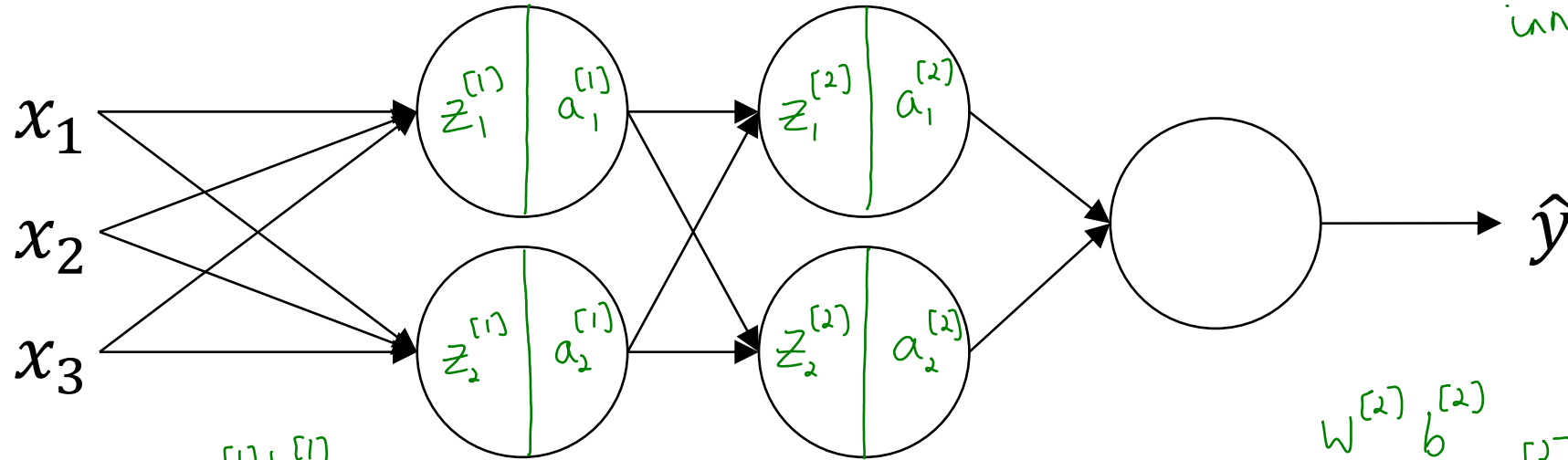
deeplearning.ai

Batch Normalization

Fitting Batch Norm
into a neural network

Adding Batch Norm to a network

performing batch normalization to inner hidden units



$$X \xrightarrow{W^{[1]} b^{[1]}} z^{[1]} \xrightarrow{\substack{b^{[1]}, y^{[1]} \\ \text{Batch norm}}} \tilde{z}^{[1]} \rightarrow a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{W^{[2]} b^{[2]}} z^{[2]} \xrightarrow{\substack{b^{[2]}, y^{[2]} \\ \text{Batch norm}}} \tilde{z}^{[2]} \rightarrow a^{[2]} = g^{[2]}(\tilde{z}^{[2]})$$

Basically, Batch norm happens b/w z & a 's computation

Params : $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]}$] For finding z
 $b^{[1]}, y^{[1]}, b^{[2]}, y^{[2]}, \dots, b^{[L]}, y^{[L]}$] For finding \tilde{z}

b here has nothing to do w/ b in gradient descent (Adam/RMS etc)

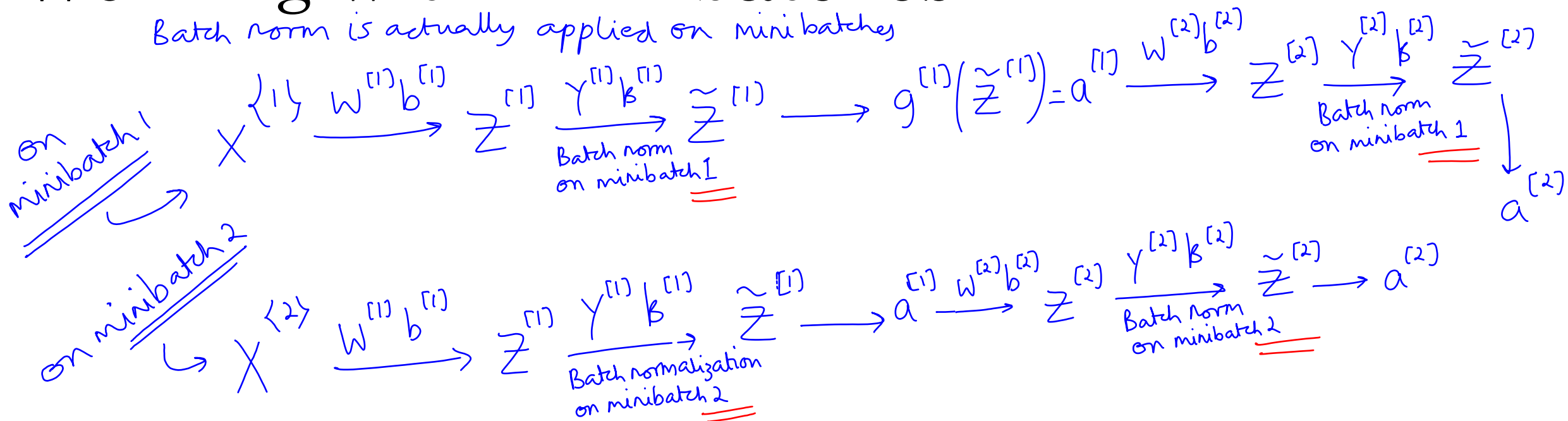
Now we will have to find db & dy & update the model

$$b^{[1]} = b^{[1]} - \alpha \cdot db^{[1]}$$

$$y^{[1]} = y^{[1]} - \alpha \cdot dy^{[1]}$$

Working with mini-batches

Batch norm is actually applied on minibatches



$\dim b^{(L)}$
 $\& Y^{(L)}$
 $= (n^{(L)}, 1)$

Why?

$\tilde{Z}^{(L)} = Y^{(L)} \cdot Z_{\text{norm}}^{(L)} + b^{(L)}$

$Z_{\text{norm}}^{(L)}$ is same as $Z^{(L)}$

$Z^{(L)} = W^{(L)} a^{(L-1)} + b^{(L)}$
 $\rightarrow \dim Z^{(L)} = (n^{(L)} \times 1)$

There is one catch
 We said Params were $W^{(L)}, b^{(L)}, \beta^{(L)}, \gamma^{(L)}$
 In reality, it is $W^{(L)}, \beta^{(L)}, \gamma^{(L)}$ (no $b^{(L)}$) - why?

$Z^{(L)} = W^{(L)} \cdot a^{(L-1)} + b^{(L)}$

Now when we find

$Z_{\text{norm}}^{(L)}, \text{ we do } Z_{\text{norm}}^{(L)} = \frac{Z^{(L)} - \mu}{\sigma^2}$

Where $\mu = \frac{1}{M} \sum_{i=1}^M Z^{(L)}_{(i)}$ for layer L

$\mu = \frac{1}{M} \sum_{i=1}^M (W^{(L)} a^{(L-1)} + b^{(L)})$

$\mu = \frac{1}{M} \sum_{i=1}^M (W^{(L)} a^{(L-1)}) + \frac{1}{M} \times M \cdot b^{(L)}$

∴ When we do $Z_{\text{norm}}^{(L)}$, each $Z^{(L)}$ has a $b^{(L)}$
 component which cancels out because we do $Z - \mu$
 ∴ No need of optimizing for $b^{(L)}$, which is
 never used, now from $Z_{\text{norm}}^{(L)}$, we find $\tilde{Z}^{(L)}$

Implementing gradient descent

for $t=1$ to Minibatches

compute fwd prop on $X^{[t]}$

In each hidden layer, use batch norm to replace $Z^{[l]}$ with $\tilde{Z}^{[l]}$

use backprop, to find $dw^{[l]}$, ~~$db^{[l]}$~~ , $dy^{[l]}$

update the params of Gradient descent

$$\left. \begin{aligned} W^{[l]} &= W^{[l]} - \alpha \cdot dw^{[l]} \\ b^{[l]} &= b^{[l]} - \alpha \cdot db^{[l]} \\ y^{[l]} &= y^{[l]} - \alpha \cdot dy^{[l]} \end{aligned} \right\} \begin{aligned} &\text{- No need to find } b^{[l]} \\ &\text{- Can use RMS Prop/Adam} \\ &\text{to update weights.} \end{aligned}$$