

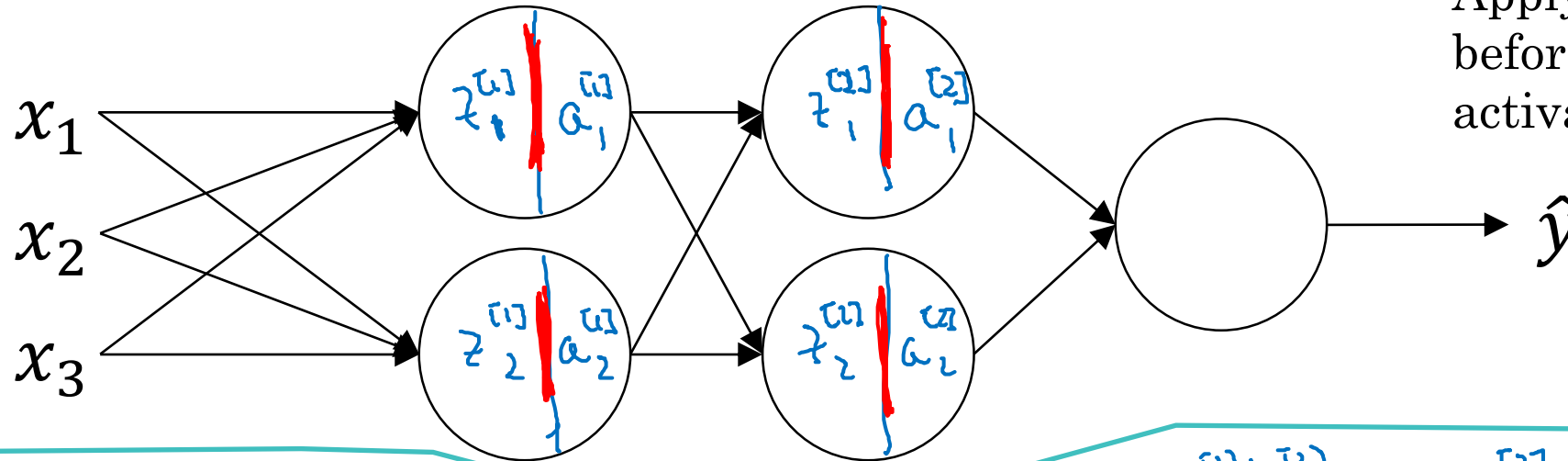


deeplearning.ai

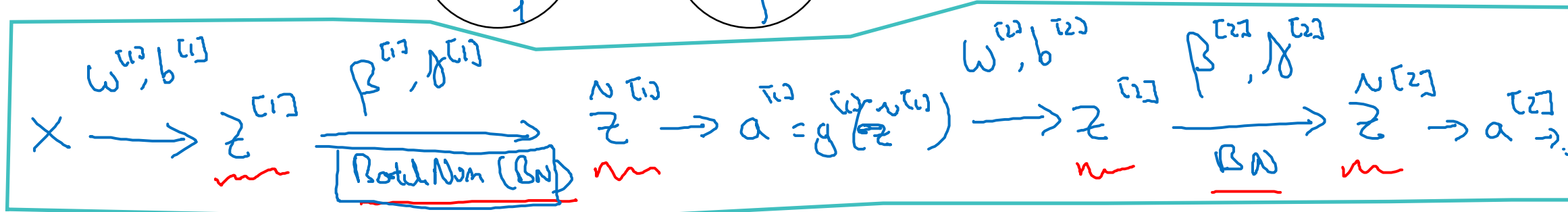
Batch Normalization

Fitting Batch Norm
into a neural network

Adding Batch Norm to a network



Apply Batch-norm before it goes into the activation function



Parameters: $\left\{ \begin{array}{l} w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)} \\ \rightarrow \underline{\beta}^{(1)}, \gamma^{(1)}, \underline{\beta}^{(2)}, \gamma^{(2)}, \dots, \underline{\beta}^{(L)}, \gamma^{(L)} \\ \rightarrow \underline{\beta} \end{array} \right\}$

$$d\beta^{(L)} \quad \beta = \beta - \alpha d\beta^{(L)}$$

tf.nn.batch-normalization ←

Working with mini-batches

$$\underline{X^{[1]}} \xrightarrow{W^{[1]}, b^{[1]}} \underline{z^{[1]}} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \underline{\tilde{z}^{[1]}} \rightarrow g^{[1]}(\tilde{z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} \underline{z^{[2]}} \rightarrow \dots$$

$$\boxed{\underline{X^{[2]}}} \rightarrow \underline{z^{[2]}} \xrightarrow[\boxed{\text{BN}}]{\beta^{[2]}, \gamma^{[2]}} \underline{\tilde{z}^{[2]}} \rightarrow \dots$$

$$\underline{X^{[2]}} \rightarrow \dots$$

parameter b is useless in this case since it will be eliminated by batch-norm anyway.

Parameters: $W^{[2]}, \cancel{b^{[2]}}, \beta^{[2]}, \gamma^{[2]}$

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ (n^{[2]}, 1) & (n^{[2]}, 1) & (n^{[2]}, 1) \end{matrix}$

$\begin{matrix} \uparrow \\ (n^{[2]}, 1) \end{matrix}$

$$\rightarrow \underline{z^{[2]}} = W^{[2]} a^{[1]} + \cancel{b^{[2]}}$$

$$\underline{z^{[2]}} = W^{[2]} a^{[1]}$$

$$\begin{matrix} \underline{z}_{\text{norm}}^{[2]} \\ \rightarrow \underline{\tilde{z}^{[2]}} = \gamma^{[2]} \underline{z}_{\text{norm}}^{[2]} + \boxed{\beta^{[2]}} \end{matrix}$$

Implementing gradient descent

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

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

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

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

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

Works w/ momentum, RMSprop, Adam.