

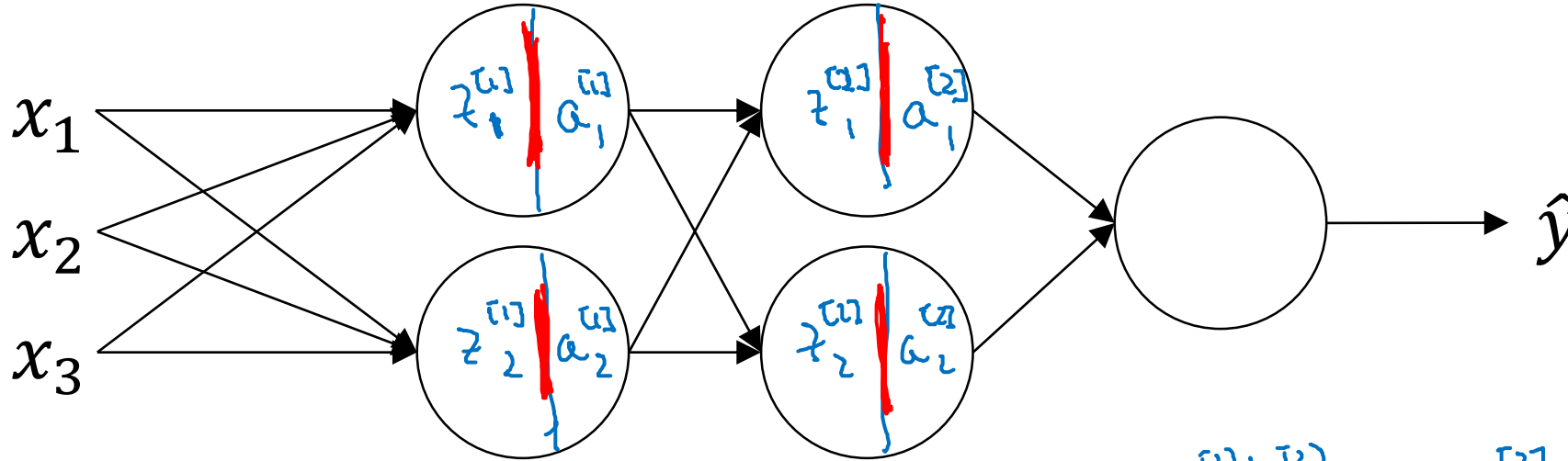


deeplearning.ai

Batch Normalization

Fitting Batch Norm
into a neural network

Adding Batch Norm to a network



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

$a = g(z)$

Parameters: $\left\{ W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots, W^{(L)}, b^{(L)} \right\}$
 $\rightarrow \underline{\beta^{(1)}}, \underline{\gamma^{(1)}}, \underline{\beta^{(2)}}, \underline{\gamma^{(2)}}, \dots, \underline{\beta^{(L)}}, \underline{\gamma^{(L)}} \}$
 $\rightarrow \underline{\beta}$

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

tf.nn.batch-normalization ←

Working with mini-batches

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

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

$$X^{\{2\}} \rightarrow \dots$$

Parameters: $W^{\{1\}}, \cancel{b^{\{1\}}}, \beta^{\{1\}}, \gamma^{\{1\}}$

\downarrow \downarrow \downarrow
 $(n^{\{1\}}, 1)$ $(n^{\{1\}}, 1)$ $(n^{\{1\}}, 1)$

\uparrow

$z^{\{1\}}_{(n^{\{1\}}, 1)}$

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

$$z^{\{2\}} = W^{\{2\}} a^{\{1\}}$$

$$z^{\{2\}}_{\text{norm}}$$

$$\rightarrow \underline{\tilde{z}^{\{2\}}} = \gamma^{\{2\}} z^{\{2\}}_{\text{norm}} + \boxed{\beta^{\{2\}}}$$

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

Works w/ momentum, RMSprop, Adam.