# Computational Intelligence

Samaneh Hosseini
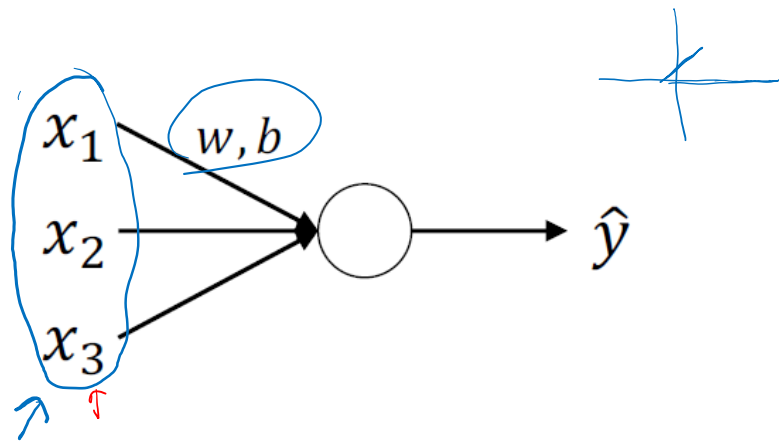
Isfahan University of Technology

# Outline

- Batch Normalization

  - Normalizing Activations in a Network

  - Fitting Batch Norm into a Neural Network

  - Why does Batch Norm work?

  - Batch Norm at Test Time

# Batch Normalization:
# Normalizing activations in a network

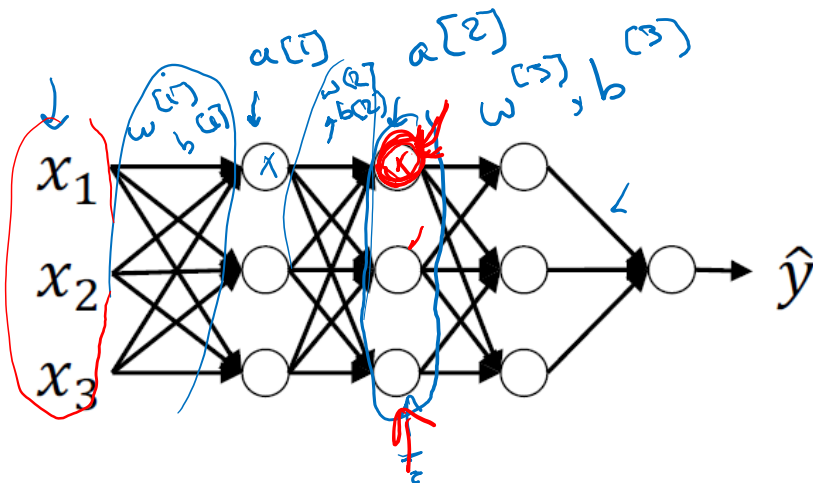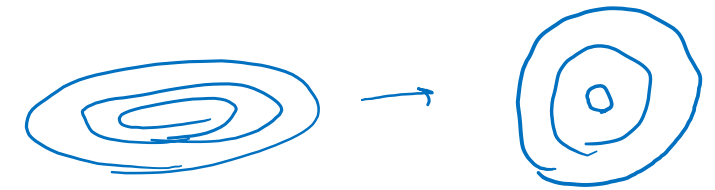# Normalizing inputs to speed up learning



$$\mu = \frac{1}{m} \sum_{i=1}^{m} X^{(i)}$$

$$\rightarrow X = X - \mu$$

element-wise

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} X^{(i)2}$$

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

Con we Normalize $a^{[2]}$ So

as to train $W^{[3]}, b^{[3]}$ faster

Normalize $Z^{[2]}$

64, 128, 512

# Implementing Batch Norm

$$\begin{bmatrix} Z^{(1)} & \cdots & Z^{(m)} \end{bmatrix}$$

$$\to Z^{[\ell](i)}$$

Given some intermediate value in NN

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

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

$$Z^{(i)}_{norm} = \frac{Z^{(i)} - \mu)}{\sqrt{\sigma^2 + \varepsilon}}$$

If

$$\gamma = \sqrt{\sigma^2 + \varepsilon}$$

$$\beta = \mu$$

then $\breve{Z}^{(i)} = Z^{(i)}$

$$\to \breve{Z}^{(i)} = \gamma Z^{(i)}_{norm} + \beta \qquad \text{learnable parameter} \atop \text{of Made}$$

Use $\breve{Z}^{[\ell](i)}$ insted of $Z^{(i)}$

# Batch Normalization
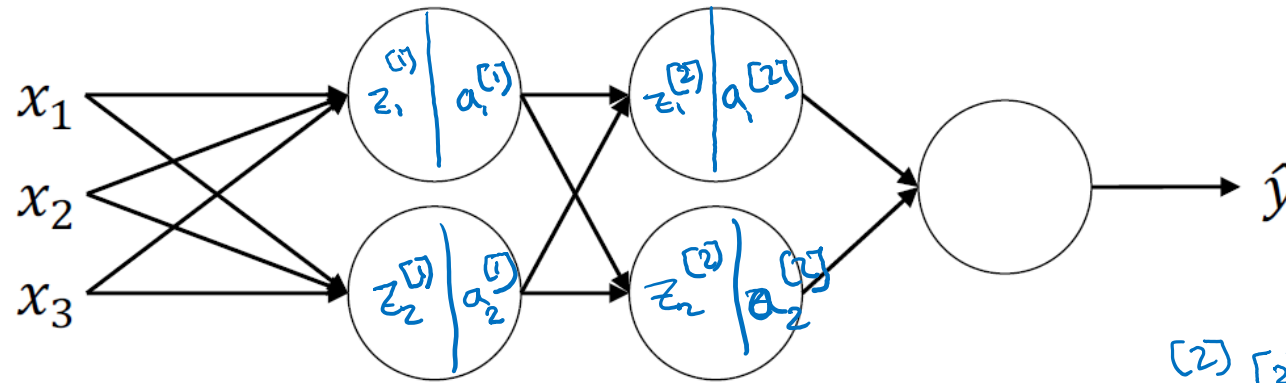# Fitting Batch Norm into a Neural Network

# Adding Batch Norm to a network



$$X \xrightarrow{\omega^{[1]}, b^{[1]}} Z^{[1]} \underbrace{\xrightarrow{\beta^{[1]}, \gamma^{[1]}}}_{\text{Batch Norm (BN)}} \tilde{Z}^{[1]} \longrightarrow a^{[1]} = g^{[1]}(\tilde{Z}^{[1]}) \xrightarrow{\omega^{[2]}, b^{[2]}} Z^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{Z}^{[2]} \rightarrow a^{[2]} \cdots$$

$$\left. \begin{array}{c} \omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]} \cdots \omega^{[L]}, b^{[L]} \\ \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]} \qquad \beta^{[L]}, \gamma^{[L]} \end{array} \right\} \quad d\beta^{[\ell]} \qquad \beta^{[\ell]} = \beta^{[\ell]} - \alpha \, d\beta^{[\ell]}$$

$$tf.nn.batch\_normalization(- \quad )$$

# Working with mini-batches

$$\text{batch} \rightarrow X^{\{1\}} \xrightarrow{\omega,b^{[1]}} z^{[1]} \xrightarrow{\beta^{[1]},\gamma^{[1]}} \tilde{z}^{[1]} \rightarrow g^{[1]}(\tilde{z}^{[1]}) = a^{[1]} \xrightarrow{\omega^{[2]},b^{[2]}} z^{[2]} \leftarrow - - -$$

$$\text{batch} \rightarrow X^{\{2\}} \rightarrow z^{[1]} \boxed{\beta^{[1]},\gamma^{[1]}} \tilde{z}^{[1]} \rightarrow - - - -$$

$$\underbrace{}_{BN}$$

$$\rightarrow X^{\{1\}} \rightarrow$$

---

$$\omega^{[l]} , \cancel{b^{[l]}} , \beta^{[l]} , \gamma^{[l]}$$

$$z$$

$$\downarrow \qquad \downarrow$$

$$(n^{[l]},1) \qquad (n^{[l]},1)$$

$$z - \mu$$

$$z^{[l]}$$

$$(n^{[l]},1)$$

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

$$z^{[l]} = \omega^{[l]} a^{[l-1]}$$

$$\rightarrow z^{[l]}_{norm}$$

$$\tilde{z} = \gamma^{[l]} z^{[l]}_{nom} + \boxed{\beta^{[l]}}$$

# Implementing gradient descent

for $t=1 \ldots$ num MiniBatches

    compute forward prop on $X^{\{t\}}$

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

    Use backprop to compute $dw^{[l]}$, $d\cancel{[b]}^{[l]}$, $d\beta^{[l]}$, $d\gamma^{[l]}$

    Update prometers
$$w^{[l]} := w^{[l]} - \alpha\, dw^{[l]}$$
$$\beta^{[l]} := \beta^{[l]} - \alpha\, d\beta^{[l]} \left.\phantom{\begin{matrix}1\\1\\1\\1\\1\end{matrix}}\right\}$$
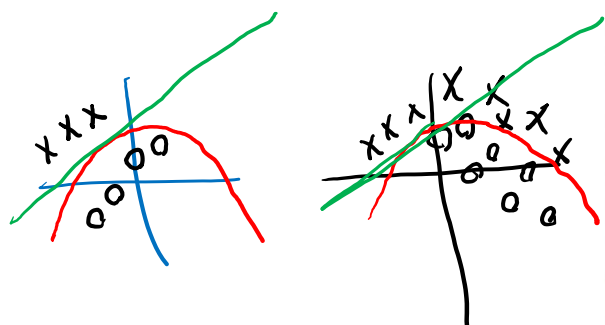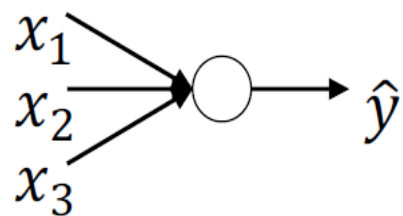$$\gamma^{[l]} := \ldots$$

RMSprob , momentum , Adam

**tf . Keras . Layers . BatchNormalization (...)**

# Batch Normalization
# Why does Batch Norm work?

# Learning on shifting input distribution

# Why this is a problem with neural networks?

$$w^{[1]}, b^{[1]}$$

$$w^{[2]}, b^{[2]}$$

$$w^{[3]}, b^{[3]}$$

$$w^{[4]}, b^{[4]}$$

$a_1^{[2]}$

$a_2^{[2]}$

$a_3^{[2]}$

$a_4^{[2]}$

$\hat{y}$

$z_2^{[2]}$

$z_2^{[2]}$

$z_1^{[2]}$

$z_1^{[2]}$

mean 0

variance 1

$$\beta^{[2]}, \gamma^{[2]}$$

# Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.

$\mu^{\{1\}}, \sigma^{\{1\}}$

$\mu^{\{2\}}, \sigma^{\{2\}}$

- This adds some noise to the values +[-] within that minibatch.

- So similar to dropout, it adds some noise to eachhidden layer's activations.

- This has a slight regularization effect.

Min-batch : 64 ⟶ 512

noise ↑        noise ↓

# Batch Normalization
# Batch Norm at Test Time

# Batch Norm at Test Time

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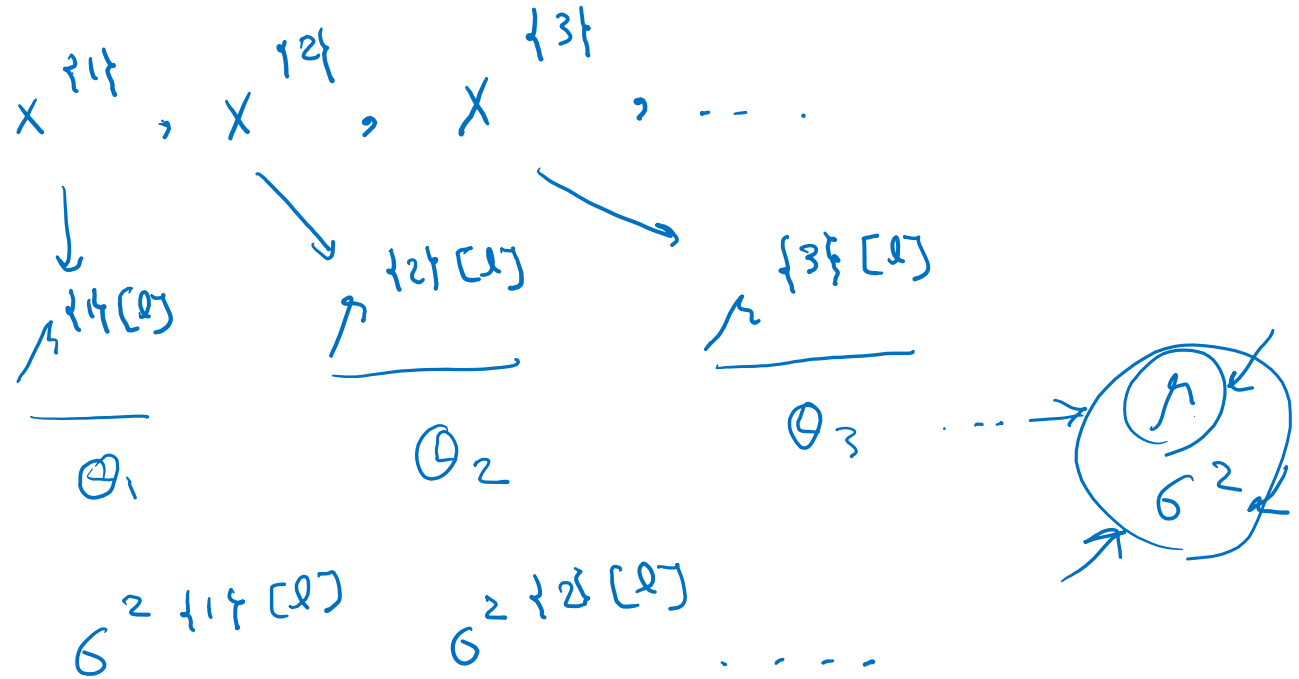
$$\mu = \frac{1}{m}\sum_i z^{(i)}$$

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

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

$$\tilde{z}^{(i)} = \gamma z_{norm}^{(i)} + \beta$$

$X^{\{1\}}$ , $X^{\{2\}}$ , $X^{\{3\}}$ , ...

$\mu^{\{1\}[\ell]}$  $\mu^{\{2\}[\ell]}$  $\mu^{\{3\}[\ell]}$

$\theta_1$   $\theta_2$   $\theta_3$ ....

$\sigma^{2\{1\}[\ell]}$   $\sigma^{2\{2\}[\ell]}$  ....

$\mu$, $\sigma^2$

$$z_{norm} = \frac{z - \mu}{\sqrt{\sigma^2 + \varepsilon}} \qquad \tilde{z} = \gamma z_{norm} + \beta$$

# Core Foundation Review

- Batch Normalization

  - Normalizing Activations in a Network

  - Fitting Batch Norm into a Neural Network

  - Why does Batch Norm work?

  - Batch Norm at Test Time