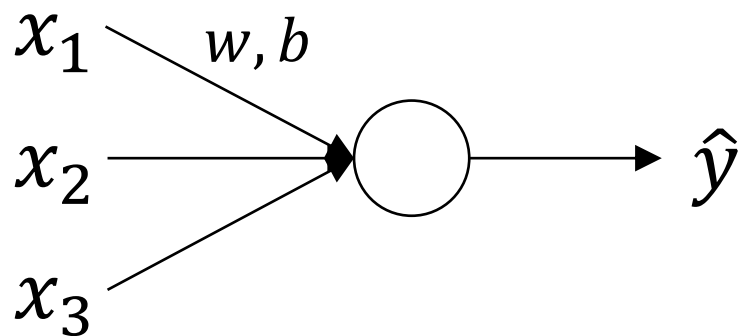# Batch Normalization

Normalizing activations in a network

(Allows you to set hyper params more easily - compared to manual selection)

deeplearning.ai

# Normalizing inputs to speed up learning

$x_1$

$w, b$

$x_2$

$x_3$

$\hat{y}$

Till now, we were only
normalizing the i/p layer,

ie, $\mu = \frac{1}{M} \sum_{i=1}^{M} x^{(i)}$

$X = X - \mu$

$\sigma^2 = \frac{1}{M} \sum_{i=1}^{M} \left( x^{(i)} \right)^2$
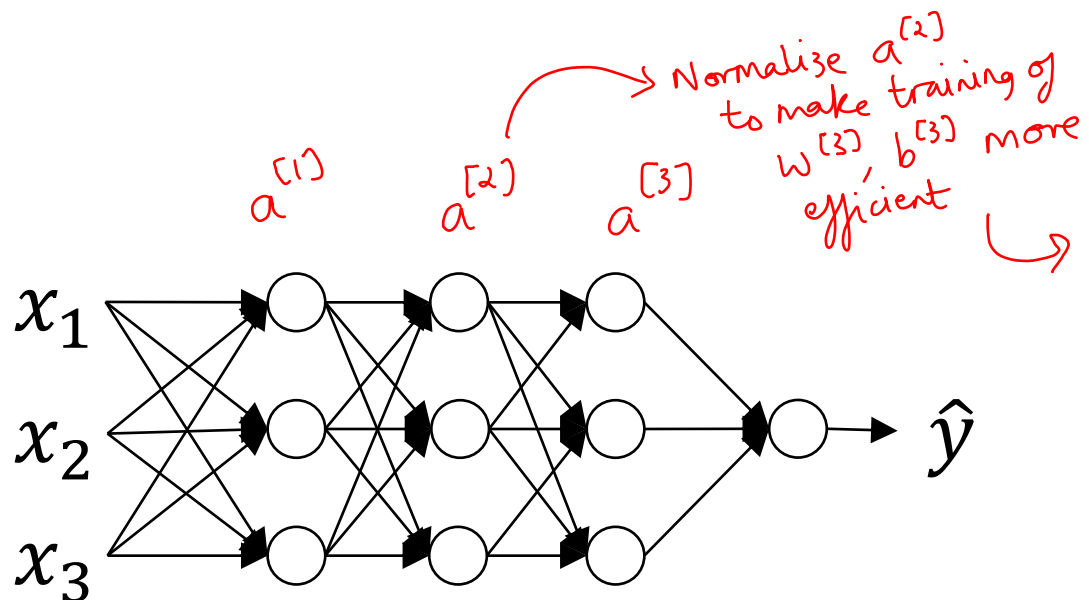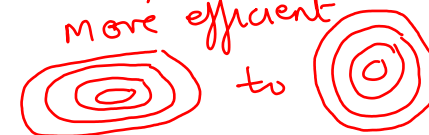
$X = X / \sigma^2$

What if we could
normalize the
Activation inputs $(a^{[1]}, a^{[2]} \ldots)$
as well?

Why?
↳ So that the training
of inner w & bs, say
$W^{[3]}, b^{[3]}$ would be
more efficient
to

$a^{[1]}$

$a^{[2]}$

$a^{[3]}$

Normalize $a^{[2]}$
to make training of
$W^{[3]}, b^{[3]}$ more
efficient

$Z^{[3]} = W^{[3]} \underbrace{a^{[2]}}_{} + b^{[3]}$

↓
Normalize this

Normalizing $a^{[2]} \approx$ Normalizing $Z^{[2]}$

$x_1$

$x_2$

$x_3$

$\hat{y}$

So lets Normalize $Z^{[2]}$

Andrew Ng

# Implementing Batch Norm

Given some intermediate hidden unit values in the NN in layer $l$
$$\hookrightarrow Z^{(1)}, Z^{(2)} \ldots Z$$

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

$$\sigma^2 = \frac{1}{M} \sum_i (Z_i - \mu)^2$$

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

$\hookrightarrow \varepsilon$ to avoid $\sqrt{0}$

Why?

Imagine you have sigmoid Activation func

for $Z^{(i)}$

If all of $Z$'s in all layers had mean=0, var=1, the model's prediction may become limited as you will only capture activations from A to B ~ why? Because all $Z$s lie b/w A & B

$\Rightarrow g(z) \in (\cdot2, \cdot8)$ for example which limits predictability, y this rule is applied to all layers

- This gives mean=0, var=1 for Z
- But maybe we don't want mean=0, var=1 for Z

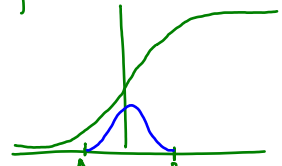$$\therefore \quad \tilde{Z}^{(i)} = \gamma \cdot Z_{norm}^{(i)} + \beta$$

Learnable params of the model (just like the weight/bias)

Why write $Z$ like this?
— It allows $Z$ to have whatever mean we want it to have

If $\gamma = \sqrt{\sigma^2 + \varepsilon}$
& $\beta = \mu$

then, we would get back
$\tilde{Z}^{(i)} = $ original $Z^{(i)}$
with mean $\mu$, var$= \sigma^2$

Now, because $\tilde{Z}^{(i)}$ is more flexible than $Z^{(i)}$ $\therefore$ we will use $\tilde{Z}^{(i)}$ in ALL FUTURE COMPUTATIONS of $Z$ in layer $L$

Andrew Ng