

Batch normalization

can we normalize $a^{[l]}$ so as to train $w^{[l]}, b^{[l]}$ faster?

we normally normalize $z^{[l]}$

Given some intermediate values in NN

$$z^{[l]}(1), \dots, z^{[l]}(m)$$

$$\mu = \frac{1}{m} \sum_i z^{[l]}(i)$$

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

$$z_{\text{norm}}^{[l]}(i) = \frac{z^{[l]}(i) - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

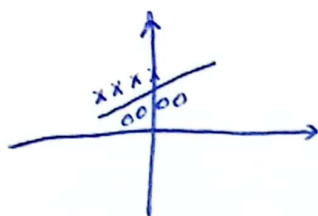
$$\tilde{z}_{\text{norm}}^{[l]}(i) = \gamma \cdot z_{\text{norm}}^{[l]}(i) + \beta$$

↑ ↑
learnable parameters of
the model

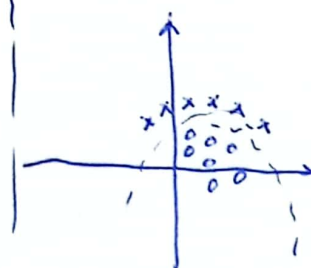
instead $\tilde{z}^{[l]}(i)$ instead of $z^{[l]}(i)$

if we have:

training set



test set



data distribution changing from train to test

$$w^{[l]}, b^{[l]} \rightarrow z^{[l]} \xrightarrow{\text{batch norm}} \tilde{z}^{[l]} \rightsquigarrow a^{[l]} = g^{(l)}(\tilde{z}^{[l]})$$

for $l = 1, \dots, \text{num Mini batches}$

- compute forward pass on $x^{(i)}$
- In each hidden layer use Batch norm to replace $z^{[l]}$ with $\tilde{z}^{[l]}$
- use backprop to compute $dw^{[l]}, db^{[l]}, d\gamma, d\beta$
- update parameters: $w^{[l]}, b^{[l]}, \gamma^{[l]}, \beta^{[l]}$

Batch norm as a regularization

- each mini batch is scaled by its mean/variance computed on just that mini-batch
- this adds some noise to the values $z^{[p]}$ within that minibatch.
 - ⇒ similar to dropout, it adds some noise to each hidden layer's activations
 - ⇒ this has a slight regularization effect.

Batch norm at test time

we don't have mini batches to add noise when calculating $z^{[i]}$

##

⇒ we use exponentially weighted average to estimate μ and σ^2