# Learning on shifting input distribution

$x_1$
$x_2$
$x_3$

$\hat{y}$

A second reason why batch norm works, is it makes weights, later or deeper than your network.

Cat     Non-Cat
$y = 1$     $y = 0$

$y = 1$     $y = 0$

"Covariate shift"

$x \longrightarrow y$

# Why this is a problem with neural networks?



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

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

$a_1^{[2]}$

$a_2^{[2]}$

$a_3^{[2]}$

$a_4^{[2]}$

$z_2^{[2]}$

$z_1^{[2]}$

$z_1^{[2]}$

$z_1^{[2]}$

$\hat{y}$

Mean 0
Variance 1

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

So from the perspective of the third hidden layer, these hidden unit values are changing all the time, and so it's suffering from the problem of covariate shift.
So the Batch Norm does, is it reduces the amount that the distribution of these hidden unit values shifts around.

But the take away is that batch norm means that, especially from the perspective of one of the later layers of the neural network, the earlier layers don't get to shift around as much, because they're constrained to have the same mean and variance.
And so this makes the jobs of learning on the later layers easier.

Andrew Ng

# Batch Norm as regularization

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

- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.

- This has a slight regularization effect.

$\rightarrow \tilde{z}^{[l]}$

$64, 128$   $z^{[l]}$

$X$

$X^{\{t\}}$

$\mu, \sigma^2$ are noisy

mini-batch : $64 \longrightarrow 512$   one strange property of dropout, which is that by using a bigger mini-batch size, you reduce the regularization effect.

Because by adding noise to the hidden units, it's forcing the downstream hidden units not to rely on too much on any one hidden unit.
And so similar to dropout, it adds noise to the hidden layers and therefore has a very slight regularization effect.
Because the noise added is quite small, this is not a huge regularization effect and you might choose to use batch norm together with dropout if you want the more powerful regularization effect of dropout.

Andrew Ng