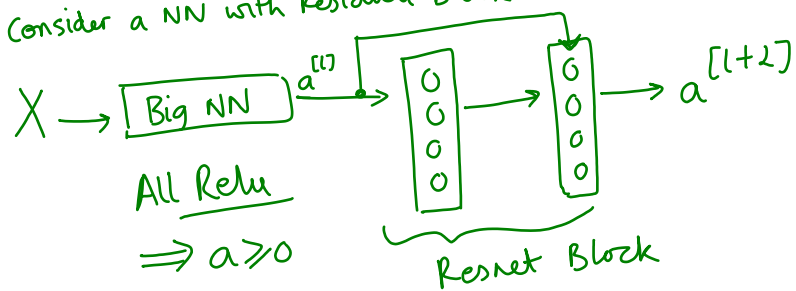


Why Resnets Work?

- consider a normal N/w
- Consider a NN with Residual Block



$$a^{[L+2]} = g(z^{[L+2]} + a^{[L]})$$

$$= g(W^{[L+2]} \cdot a^{[L+1]} + b^{[L+2]} + a^{[L]})$$

$$\text{If } W^{[L+2]} = b^{[L+2]} = 0$$

$$\text{then } a^{[L+2]} = g(a^{[L]}) = \text{Relu}(a^{[L]}) = a^{[L]} \quad (\text{since } a \geq 0 \text{ \& Relu}(a \geq 0) = a)$$

\Rightarrow Identity func. is easy for Residual block to learn
 ie, $a^{[L+2]} = a^{[L]}$ is possible
 Because of Residual block
 \Rightarrow Residual NN performs atleast as good as normal NN

\Rightarrow If we add a few Residual blocks in between/end of NN, we will achieve performance \gg performance to normal NN. The other issue with normal NN is that, if the layers are large, then it becomes harder to even learn params for the Identity func in later layers, in Resnet, you just pass the prev state of the N/w & w/o learning new params you can achieve the Identity matrix

- Because in normal NN, they can't even predict the Identity Activation, performance worsens (see graph on prev page), while in Resnet, they predict Identity + other learning using Residual Block \Rightarrow predictive performance \uparrow as layers \uparrow

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

$$\Rightarrow \dim(z^{[l+2]}) = \dim(a^{[l]})$$

How to ensure this?

- Make the convolutions the same in size
- If they were diff in dim, say $a^{[l]} = (128 \times 1)$

$$\text{and } z^{[l+2]} = (256 \times 1)$$

then Add W_5 in between

$$z^{[l+2]} + W_5 \cdot a^{[l]}$$

$$\hookrightarrow \dim = (256 \times 128)$$

$$\Rightarrow (256 \times 128) \times (128 \times 1) \\ = (256 \times 1)$$

\Rightarrow Now you can Add the 2.

Resnet example given In Video $\rightarrow W2 \mid 4$

