

ResNet

feedforward

Goal: train DNNs (convolutional neural networks) w/ many layers.

trained on ImageNet { GoogleNet architecture → 20 layers (excluding pooling)
ResNet architecture → 150 layers

ResNet: deep residual network.

Intuition: Jacobian of a RNN $\frac{\partial \vec{h}^{(T)}}{\partial \vec{h}^{(1)}} = \frac{\partial \vec{h}^{(T)}}{\partial \vec{h}^{(T-1)}} \dots \frac{\partial \vec{h}^{(2)}}{\partial \vec{h}^{(1)}}$
product of the Jacobian of individual layers.

Multiplying all these individual matrices is why gradients explode / vanish unless we keep each matrix close to the identity.

⇒ This formula applies to feed-forward NNs (conv neural nets)

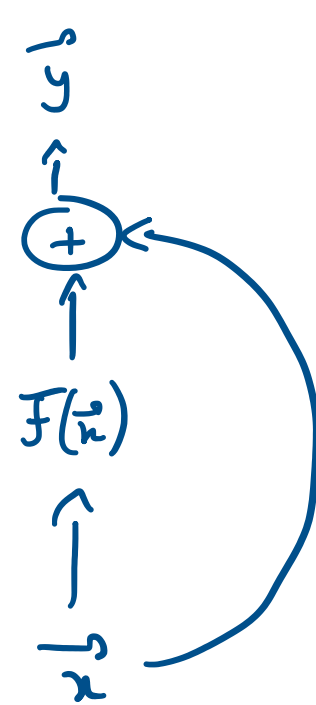
Feed forward neural network:

$\vec{h} = \phi(\vec{z})$
↑
activation function

$\vec{z} = W\vec{x} + \vec{b}$

Use this as an architectural primitive called the residual block.

Residual block $\left\{ \begin{array}{l} \vec{y} = \vec{x} + F(\vec{x}) \\ \sim \\ \text{residual function} \end{array} \right.$



Special case:

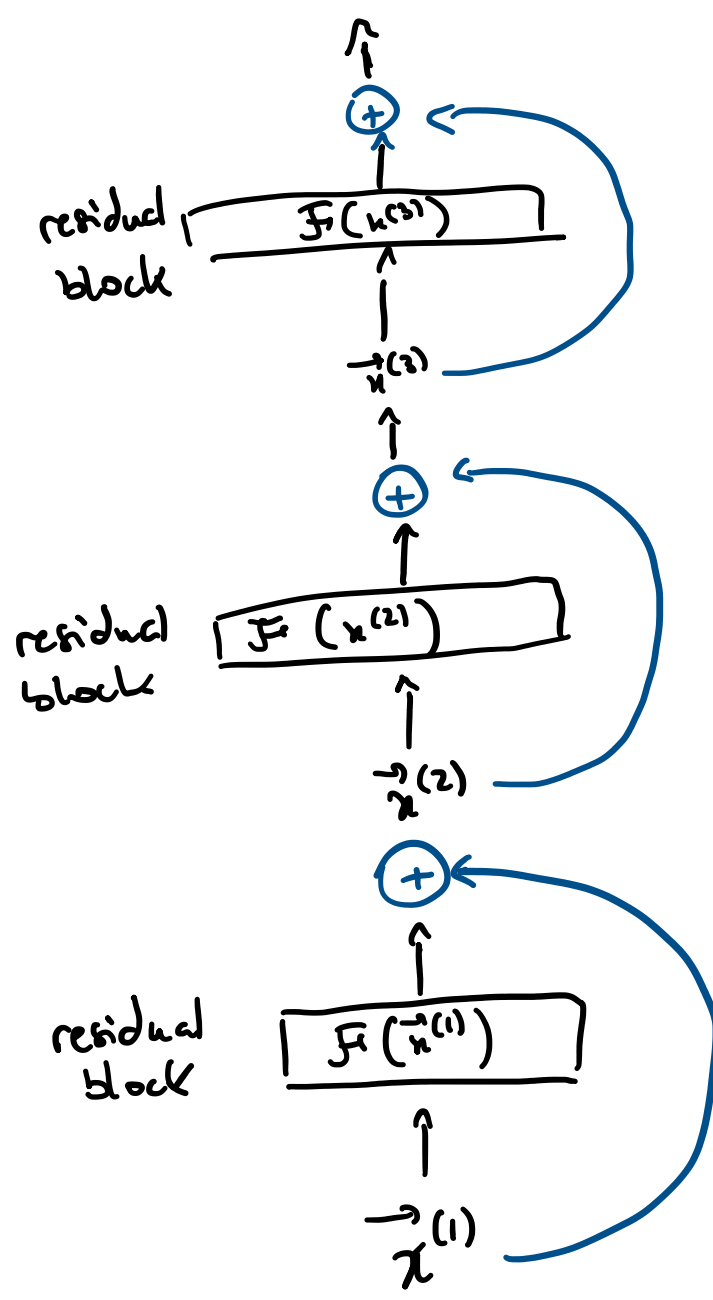
$F(\vec{x}) = \phi(\vec{z}) = W\vec{x} + b$

Typically: residual function is a small/shallow neural network w/ 1-3 hidden layers.

ResNet:

Deep residual network

Each residual block has separate trainable parameters



Why does this architecture help prevent gradients exploding / vanishing?

↳ Jacobian of a residual block should be close to the identity.

Backpropagation rules:

$\vec{\delta}^{(l)} = \vec{\delta}^{(l+1)} \frac{\partial F}{\partial \vec{x}} + \vec{\delta}^{(l+1)}$

$= \vec{\delta}^{(l+1)} \left(\frac{\partial F}{\partial \vec{x}} + I \right)$

identity matrix

if $\frac{\partial F}{\partial \vec{x}} \rightarrow 0$, the sum is close to the identity

⇒ error signal / gradients will not explode or vanish.