

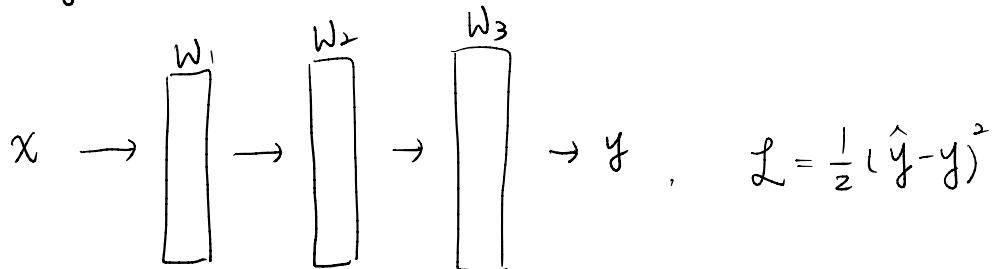


ResNet

1. What is the main problems of general network?

- Gradient Vanishing / Explosion

Assuming a single model



If we want to calculate gradient related to each param...

$$a = w_1 x$$

$$b = w_2 a \Rightarrow y = c = w_3 (w_2 (w_1 x))$$

$$c = w_3 b$$

$$L = \frac{1}{2} (\hat{y} - c)^2$$

$$\frac{dL}{dw_3} = \frac{dL}{dc} \frac{dc}{dw_3} = (\hat{y} - c) \cdot b$$

$$\frac{dL}{dw_2} = \frac{dL}{dc} \frac{dc}{db} \frac{db}{dw_2} = (\hat{y} - c) \cdot w_3 \cdot a$$

$$\frac{dL}{dw_1} = \frac{dL}{dc} \dots = (\hat{y} - c) \cdot w_3 \cdot w_2 \cdot x$$

If we have 1000 layers, the params in first

layer will be $(\hat{y} - c) \cdot \underbrace{w_{1000} \dots w_2}_{1000-1}$

a. What if $w_2 \sim w_{1000}$ are small?

$$\frac{\partial L}{\partial w_1} \sim 0, \quad w'_1 \leftarrow w_1 - \eta \frac{\partial L}{\partial w_1} \sim w_1$$

The weight will not update!

G vanishing

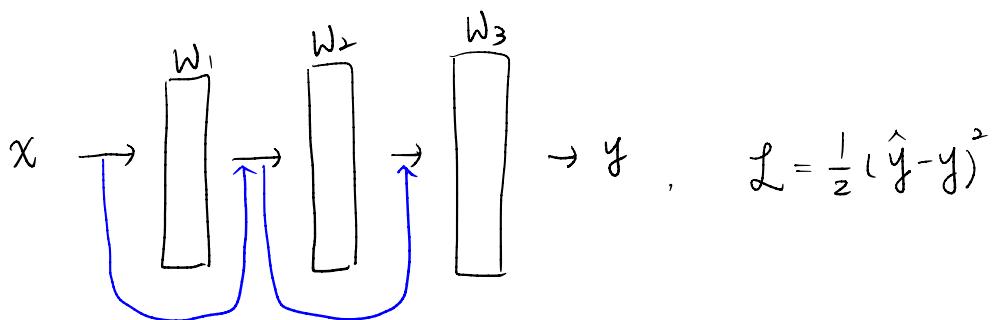
b. What if $w_2 \sim w_{1000}$ are large?

$$\frac{\partial L}{\partial w_1} \sim \infty, \quad w'_1 \rightarrow -\infty$$

G explosion

2. Skip Connection

* Dealing with gradient vanishing problem.



Now we change to ...

$$a = w_1 x + x = x(w_1 + 1)$$

$$b = w_2 a + a = a(w_2 + 1)$$

$$c = w_3 b + b = b(w_3 + 1)$$

$$\frac{\partial L}{\partial w_3} = (\hat{y} - c) \cdot b$$

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= (\hat{y} - c) \cdot (w_3 + 1) \cdot (w_2 + 1) \cdot x \\ &= (\hat{y} - c) w_3 w_2 x + (\hat{y} - c) w_2 x + (\hat{y} - c) w_3 x + (\hat{y} - c) x \end{aligned}$$

↙
 skip w_3
↓
 skip
 ↓
 skip both w_2, w_3

which can prevent gradient vanishing.

For gradient explosion, please check

- He initialisation / Xavier initialisation
- gradient clipping.

3. Architecture of ResNet

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

- For every convolution, add padding to retain the $H \times W$.
- Notice the output block of each conv-block will decrease $\frac{1}{2}$. so we need to scale the figure at the very beginning of conv-block
- * Except conv2-X, since it already exists max pool.

The reason is in traditional CNN, we like to use less resource to enhance computation efficiency. So at the very first of each layer, we will decrease the dimension while increase the channels.

4. Output dimension (Supplement)

The main param to control output dim is **stride**. Kernel size and padding are just adjusting or maintaining the output shape.

$$\text{Output size} = \text{floor}\left(\frac{(\text{Input} + 2 \times \text{padding} - \text{kernel})}{\text{stride}}\right) + 1$$

5. Deep Dive into Conv block.

$$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$$

- The first layer is used to decrease image shape, increase computational efficiency.
- 3×3 is used to extract image pattern.

- Final $|x|$ is used to retain image recognition ability, add more channels to make sure the ability.

5. Some important tips.

a. We need to make sure the ^{size} channel of x and $f(x)$ are same so that they can add together.

To downsample (modify) the size, we use $| \times |$ kernel to adjust.

b. There are two situation we need to modify the size.

① Channel

$$(B_x, C_x, H_x, W_x) \leftrightarrow (B_x, C_{fx}, H_x, W_x)$$

② Figure shape

$$(B_x, C_x, H_x, W_x) \leftrightarrow (B_x, C_x, H_{fx}, W_{fx})$$

① will happen when ResNet are 50/101/152

If 18/34, then we don't need downsampling.

② Will happen at each layer, since every new conv-block we need to adjust both channel size and image shape.

C. How to decide padding?

For conv-1, input = 224, output = 112, stride = 2

$$112 = \text{floor} \left(\frac{224 + 2p - 7}{2} \right) + 1$$

$$\Rightarrow p = 3.$$

For conv2 ~ conv5,

1×1 are used to decrease the shape $\Rightarrow p = 0$

3×3 are pattern recognition, to maintain image shape $\Rightarrow p = 1$