

# SE\_ResNet的pytorch实现

残差块:

```
class Residual_block(nn.Module):
    def __init__(self, in, middle_out, out, kernel_size=3, padding=1):
        self.out_channel = middle_out
        super(Residual_block, self).__init__()

        self.shortcut = nn.Sequential(
            nn.Conv2d(nin, nout, kernel_size=1),
            nn.BatchNorm2d(nout)
        )
        self.block = nn.Sequential(
            nn.Conv2d(in, middle_out,
kernel_size=kernel_size, padding=padding),
            nn.BatchNorm2d(middle_out),
            nn.ReLU(inplace=True),
            nn.Conv2d(middle_out, out,
kernel_size=kernel_size, padding=padding),
            nn.BatchNorm2d(nout))

    def forward(self, input):
        x = self.block(input)
        return nn.ReLU(inplace=True)(x + self.shortcut(input))
```

**注意** 在input与输出相加前, 这里需要一个shortcut层来调整input的通道大小

SE-Net:

```
class SE(nn.Module):
    def __init__(self, in, middle_out, out, reduce=16):
        super(SE, self).__init__()
        self.block = Residual_block(in, middle_out, out)

        self.shortcut = nn.Sequential(
            nn.Conv2d(nin, nout, kernel_size=1),
            nn.BatchNorm2d(nout)
        )
        self.se = nn.Sequential(
            nn.Linear(out, out // reduce),
            nn.ReLU(inplace=True),
            nn.Linear(out // reduce, out),
            nn.Sigmoid())

    def forward(self, input):
        x = self.block(input)
        batch_size, channel, _, _ = x.size()
        y = nn.AvgPool2d(x.size()[2:])(x)
        y = y.view(y.shape[0], -1)
        y = self.se(y).view(batch_size, channel, 1, 1)
        y = x * y.expand_as(x)
```

```
out = y + self.shortcut(input)
return out
```

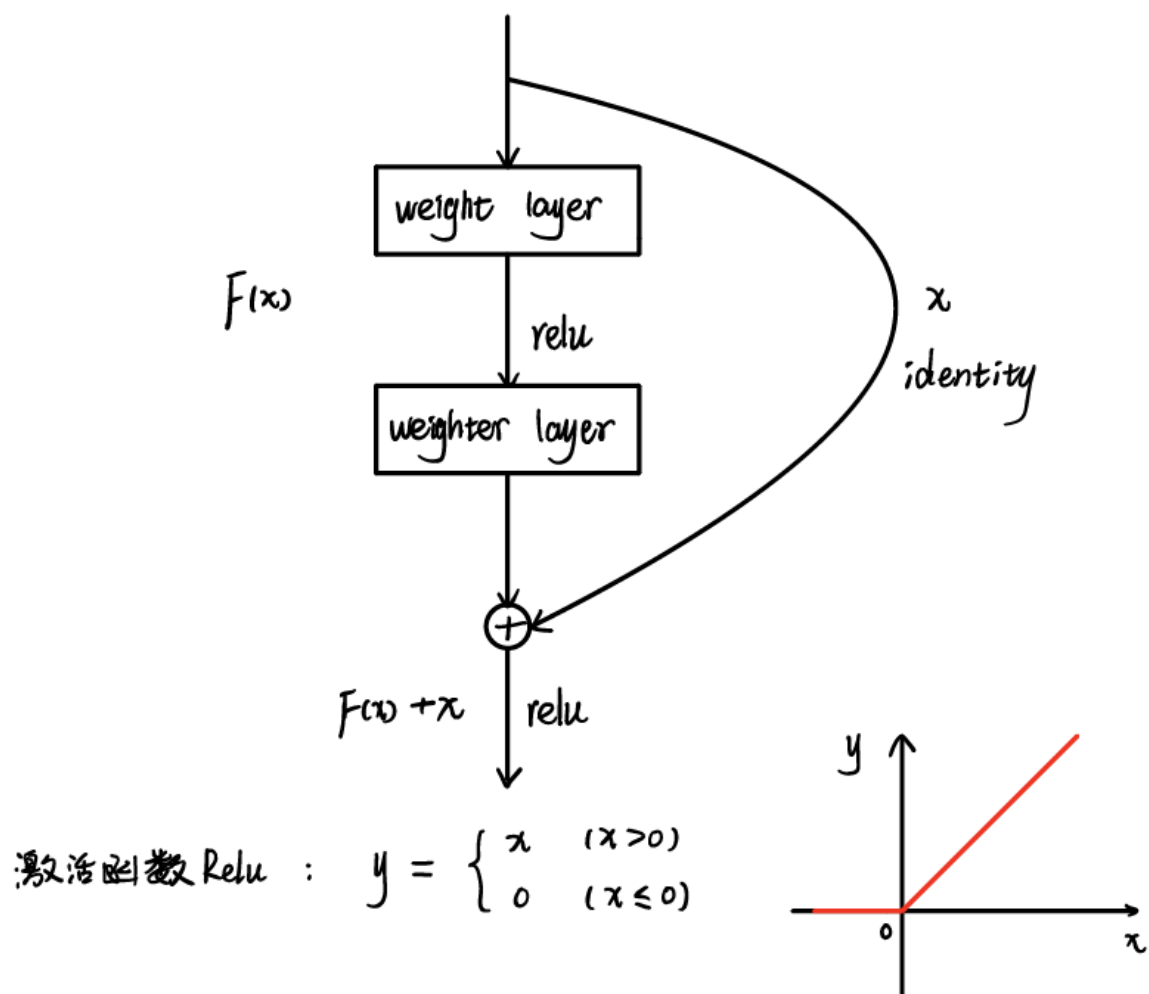
注意 我们使用平均池化层进行下采样，这里的kernel\_size是动态的

## ResNet

### 1、为什么网络层数越深，误差越大？在ResNet中作者提出了两个问题：

- 随着网络层数的加深，梯度消失和梯度爆炸以及梯度爆炸这个问题会越来越明显（假设在反向传播过程中，每一层的误差梯度是一个小于一的数，在反向传播过程中，每向前传播一次都乘以一个小于一的系数，那么当网络越来越深的时候所乘的小于一的系数就越来越多，梯度就会越来越趋近于零，造成梯度消失。相反，再反向传播过程中，如果每一层的梯度都是一个大于一的数，则在反向传播过程中梯度会越来越大，即梯度爆炸）。梯度消失和梯度爆炸问题的解决通常是通过以下三种方式解决：
  1. 对数据进行标准化处理
  2. 权重初始化
  3. Batch Normalization（放弃dropout）
- Degradation problem：即在解决了梯度消失和爆炸的问题之后，依然会存在网络层数加深，但是效果依然不好的问题。ResNet提出了残差结构，通过残差结构，解决退化问题。

### 2、ResNet网络中残差块的结构：



### 3、ResNet原论文中作者遇到的问题：

1、通过简单堆叠卷积层与池化层而形成的深层网络存在：层数越深，网络效果越差的问题。作者提出了1处的两个问题。关于梯度消失或者梯度爆炸的本质原因：网络太深，网络权值更新不稳定，本质是因为梯度反向传播过程中的连乘效应（链式法则）。关于梯度消失和梯度爆炸的直观解释：

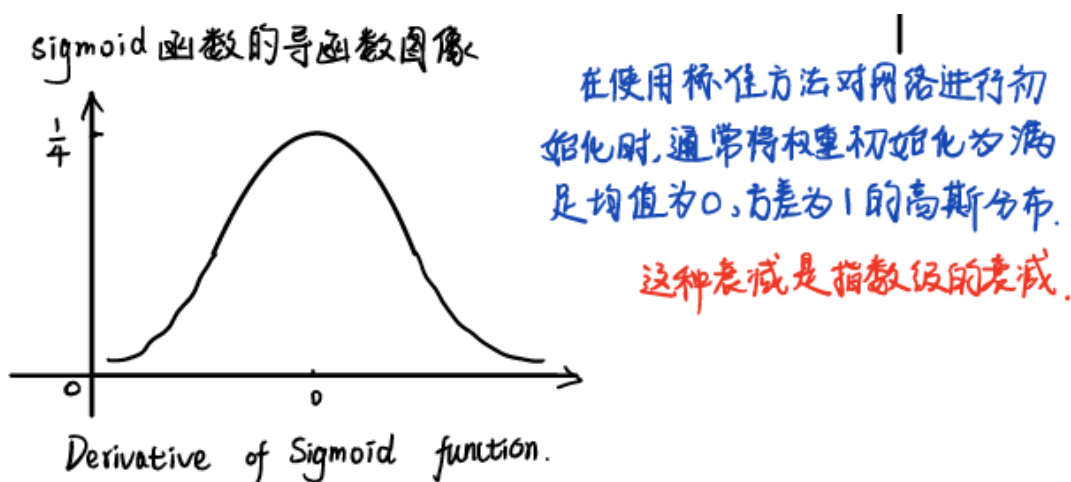
- 以三个隐藏层为例：

$y_i = \delta(z_i) = \delta(w_i x_i + b_i)$ ,  $C$  为代价函数，根据链式法则，可以得到图中表达式

$$\begin{aligned} \text{则: } \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \frac{\partial x_3}{\partial z_2} \frac{\partial z_2}{\partial x_2} \frac{\partial x_2}{\partial z_1} \frac{\partial z_1}{\partial b_1} \\ &= \frac{\partial C}{\partial y_4} \delta'(z_4) w_4 \cdot \delta'(z_3) w_3 \cdot \delta'(z_2) w_2 \cdot \delta'(z_1) \end{aligned}$$

$\delta$  为激活函数 Sigmoid:  $h(x) = \frac{1}{1 + \exp(-x)}$

- Sigmoid函数的导函数图像为：导函数最大值为1/4



可见， $\delta'(x)$  的最大值为  $\frac{1}{4}$ ，而我们初始化网络权值  $|w|$  通常都小于1。  
因此， $|\delta'(z)w| \leq \frac{1}{4}$ ，故由公式知：

$$\frac{\partial C}{\partial b_1} = \frac{\partial C}{\partial y_4} \underbrace{\delta'(z_4) w_4}_{\leq \frac{1}{4}} \cdot \underbrace{\delta'(z_3) w_3}_{\leq \frac{1}{4}} \cdot \underbrace{\delta'(z_2) w_2}_{\leq \frac{1}{4}} \cdot \underbrace{\delta'(z_1)}_{\leq \frac{1}{4}}$$

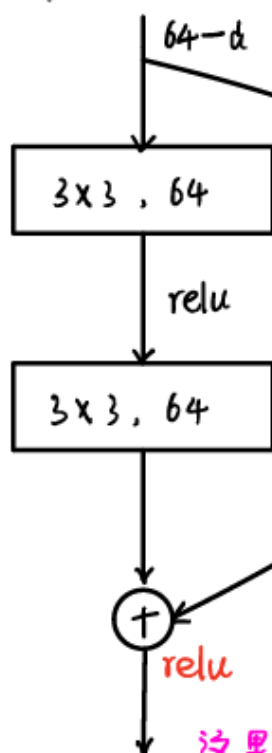
使用链式法则求导，层数越深，求导结果越小，不断趋近于零，造成梯度消失

梯度爆炸：当权重参数  $w$  被初始化为一个较大值时，当  $|\delta'(z)w| > 1$  时，由链式法则，代价函数的梯度会随网络层数的加深呈指数级增长

#### 4、ResNet残差结构:

- 主要分为两种: 针对相对浅层的网络如ResNet18, ResNet34和针对相对深层的网络ResNet50、ResNet101以及ResNet152等
- 残差结构①:

Residual 结构:



主要针对网络层数较少的网络所使用的残差结构, 如 ResNet 34

参数量: 假设输入为 256 维

$$3 \times 3 \times \underline{256} \times \underline{256} + 3 \times 3 \times 256 \times 256 = 1179648$$

一个 256 是输入的特征矩阵的深度, 另外一个 256 是卷积核的个数.

① 结构主线是两个  $3 \times 3$  的卷积层, short cut 部分将原始的输入特征矩阵相加.

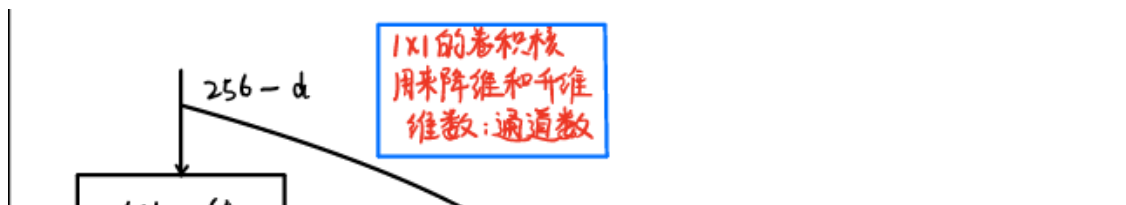
② 主线上的结果与 short cut 部分的结果相加

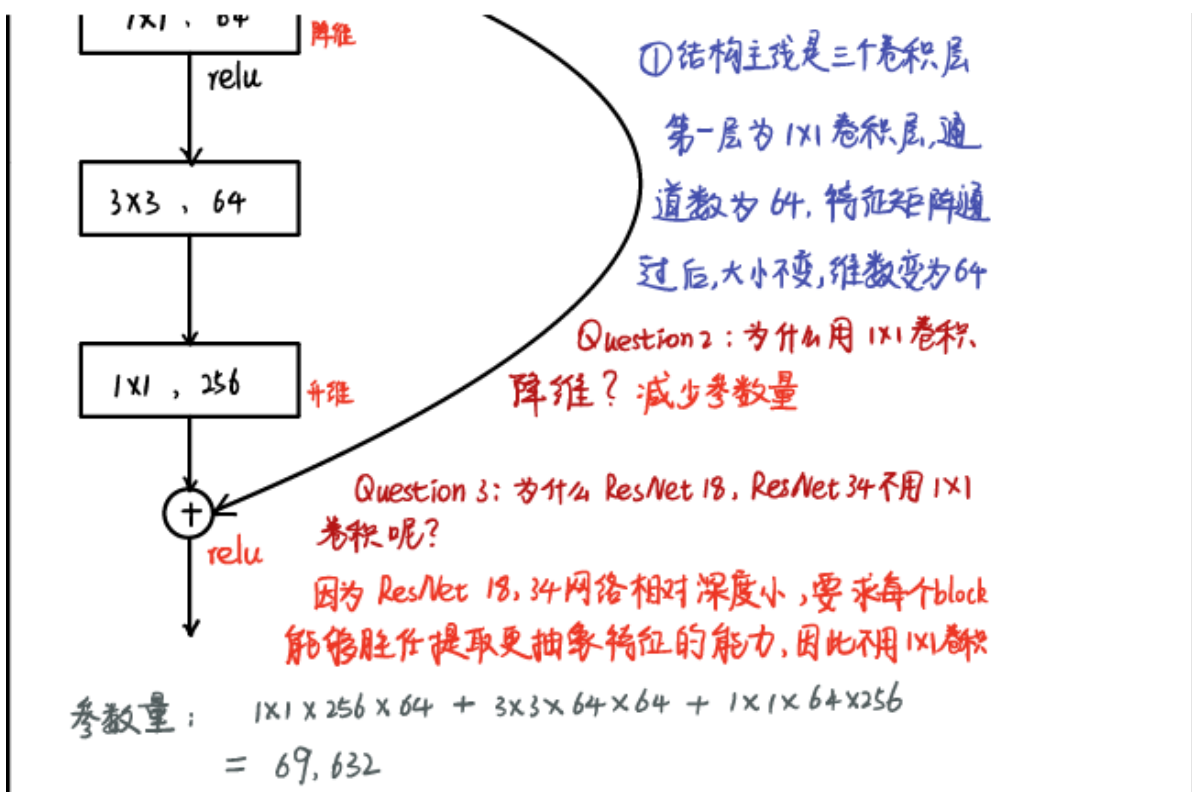
这里需注意, 主分支与 short cut 分支两个部分输出的特征矩阵的 shape 必须相同.

☆:  $\text{width} * \text{height} * \text{输入特征矩阵的维数} * \text{卷积核的个数}$

此处参数量计算的解释: 两个 256 是指, 假设输入特征矩阵的维度是 256, 主线上卷积核的个数也得对应是 256.

- 残差结构②: 此种结构主要是针对 ResNet50、Resnet101、ResNet152 等深度网络使用的。





此处参数计算量解释: 输入特征图的维度是 256, 但第一个卷积层的卷积核的数量是 64

## 5、ResNet论文中的结构图:

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

此表中的 conv2\_x、conv3\_x、conv4\_x 以及 conv5\_x 表示一系列的残差结构, 以 ResNet 34 为例.

Conv2\_x: 包含 3 个残差层

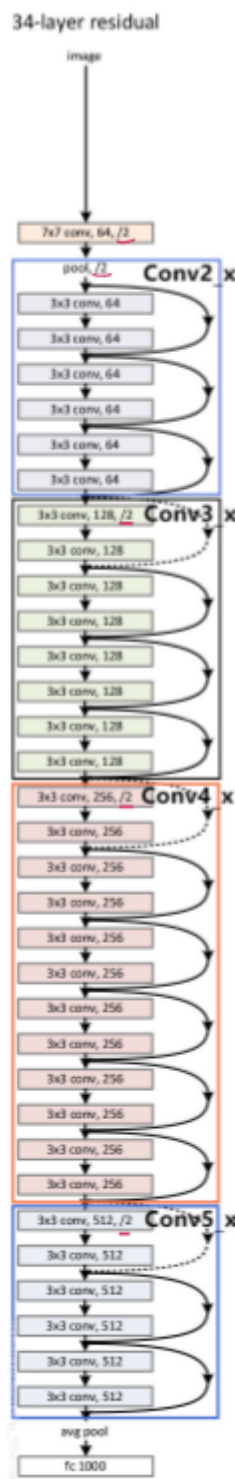
Conv3\_x: 包含 4 个残差层

Conv4\_x: 包含 6 个残差层

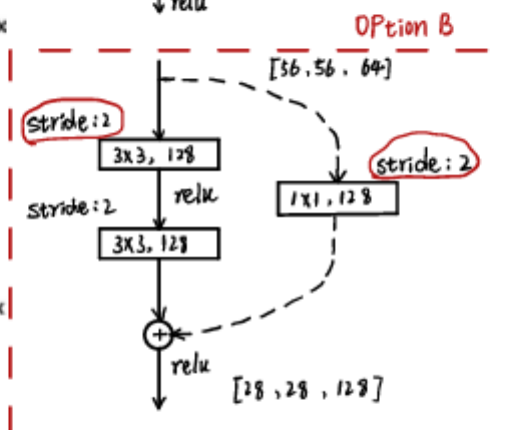
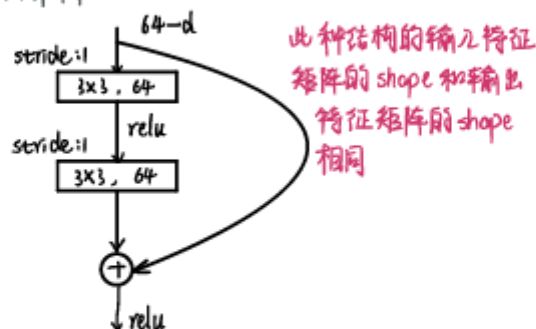
Conv5\_x: 包含 3 个残差层

layer name		output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1		112×112	7×7, 64, stride 2				
			3×3 max pool, stride 2				
layer1	conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
layer2	conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
layer3	conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
layer4	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			$1.8 \times 10^9$	$3.6 \times 10^9$	$3.8 \times 10^9$	$7.6 \times 10^9$	$11.3 \times 10^9$

## 6、以ResNet34为例：



以 ResNet 34 为例, 网络结构中即包含实线  
残差结构和虚线残差结构, 两种结构的具体  
解释: residual 结构



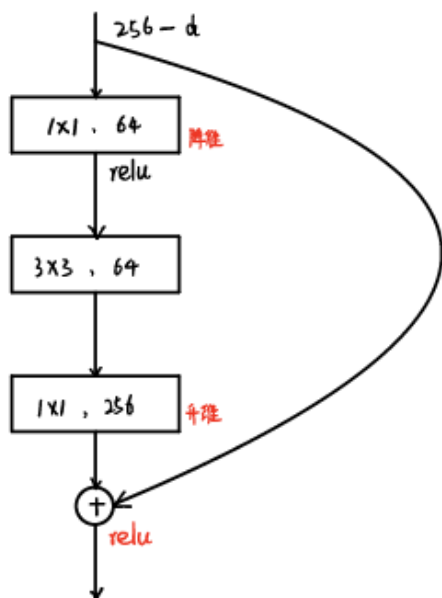
Option B 的输入特征矩阵与输出特征矩阵  
的 shape 不同. Option B 出现的位置均为两  
个连续的 Conv 连接处.

option B 中: 第一个 3x3 卷积 stride 为 2, 把  
原输入的 H 和 W 均减半. shortcut 中 stride  
也为 2 是同样的作用. 128 卷积核改变原输入  
通道数, 使最终主残部分与 shortcut 部分  
shape 相同.

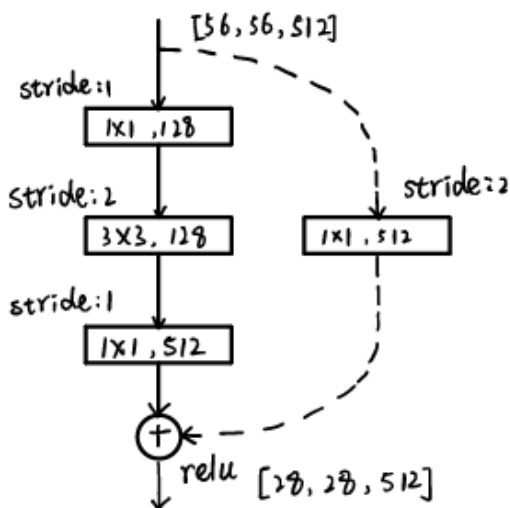
注意: Conv3\_x、Conv4\_x、Conv5\_x 各个部分之间相连接的第一个部分均为 option B 这种残差结构,  
因为通道数不匹配, 但 Conv2\_x 处直接使用第一种结构, 因为在这一层输入与输出的通道数是相同.

## 7、ResNet50、ResNet101、ResNet152 基础结构:





原论文中的, 此处可能有A, B, C  
等多种结构, 但B这种作为  
认为最好



option B

☆ 主分支与 shortcut 的输出  
特征矩阵 shape 必须相同.

注: 原论文中, 左侧虚线残差  
结构的主分支上, 第一个 1x1 卷  
积层的步距为 2, 第二个 3x3  
卷积层的步距为 1

但 pytorch 官方的实现过程中  
第一个 1x1 卷积层的步距为 1,  
第二个 3x3 卷积层步距是 2

这样在 imagenet 的 top1

上提升大概 0.5% 的准确率.

- 虚线的残差结构键输入图片的高和宽以及深度等都做变化, 而实线部分的输入特征矩阵的和输出特征矩阵是一样的。
- Conv3\_x、Conv4\_x、Conv5\_x对应的一系列残差结构的第一层均为虚线结构, 因为第一层必须将上一层的特征矩阵的高和宽以及深度调整为当前层特征矩阵所需的高和宽以及深度 (原文中: Downsampling is performed by Conv3\_1, Conv4\_1, Conv5\_1 with a stride of 2)。

## Batch Normalization

1、Batch Normalization的目的是将一批 (Batch) feature map满足均值为0, 方差为1的分布规律, 即将一批数据对应的feature map (即特征矩阵) 的每一个维度也就是每一个channel对应的维度满足均值为0, 方差为一的分布规律。

2、通常情况下, 在将图像输入到网络之前, 会对图像进行预处理。即将输入图像数据调整到满足某一分布规律, 这样可以加速网络训练。但原始图像在通过卷积层之后, 得到的feature map却不一定满足某一分布规律, 因此考虑使用Batch Normalization调整feature map, 使每一层的feature map都能满足均值为0和方差为1的分布。

3、



3. 原论文 " For a layer with  $d$ -dimensional input  $x = (x^{(1)}, \dots, x^{(d)})$ , we will normalize each dimension

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

即对于一个拥有  $d$  维的输入  $x$ , 我们将对它的每一个维度进行处理. 假设我们输入的  $x$  是 RGB 三通道的彩色图像, 那么  $d$  就是输入图像的 channel 即  $d=3$ ,  $x = (x^{(1)}, x^{(2)}, x^{(3)})$ , 其中  $x^{(1)}$  就代表 R 通道对应的特征矩阵. 依此类推, 标准化处理也就是分别对我们的 R 通道, G 通道, B 通道进行处理.

计算详解:

**Input :** Values of  $x$  over a mini-batch :  $B = \{x_1, \dots, x_m\}$ ;  
 Parameters to be learned :  $\gamma, \beta$

**Output :**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \underbrace{\gamma}_{\text{调方差}} \hat{x}_i + \underbrace{\beta}_{\text{调均值}} \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

满足 0-1 分布的效果不一定是最好的, 因此设置两个参数  $\gamma$  和  $\beta$  进行调整,  $\gamma$  和  $\beta$  是通过反向传播学习得到的. 而均值和方差是通过一批批数据计算统计得到的.

$\left\{ \begin{array}{l} \mu, \sigma^2 \text{ 在正向传播过程中统计得到} \\ \gamma, \beta \text{ 在反向传播过程中训练得到} \end{array} \right.$

注: 此处  $\beta$  和  $\gamma$  的初始值分别为 0 和 1.