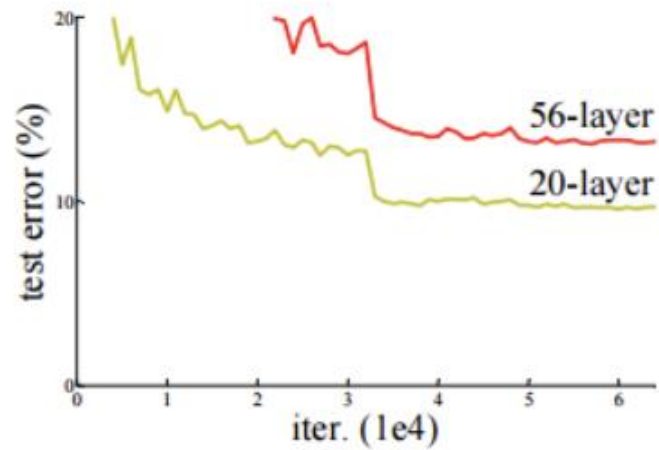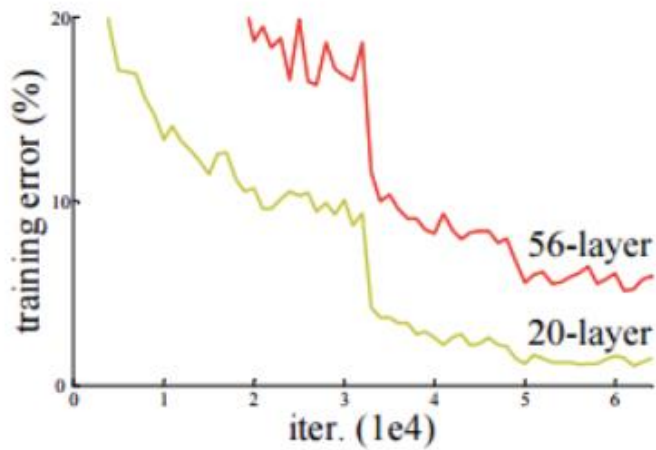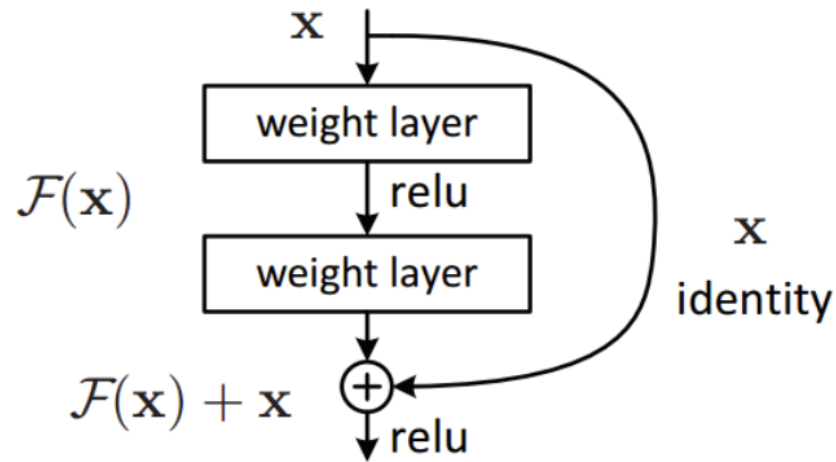# ResNet 구현

# 1. ResNet 설명

- **기존 Network의 문제점**
  - 층이 깊어질수록 Gradient vanishing/exploding 문제가 발생
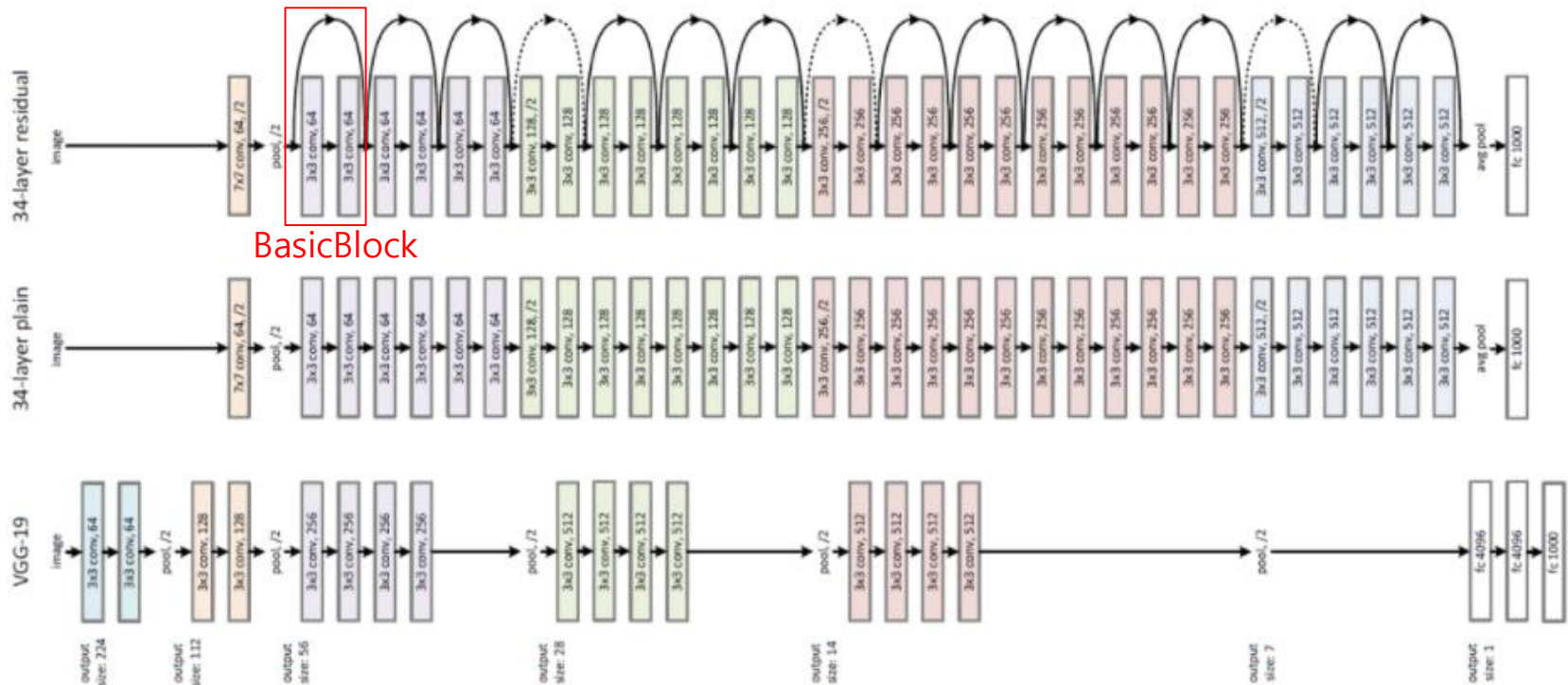
# 1. ResNet 설명

■ ResNet의 핵심 아이디어

▪ Skip connection (identity mapping)을 두어 gradient vanishing/exploding 문제를 효과적으로 해결

# 1. ResNet 설명

- ## ResNet의 핵심 아이디어
  - 모든 Convolution layer는 3x3 크기를 가짐
  - Residual connection (Convolution layer 2개) + Skip connection (identity mapping)으로 이루어진 BasicBlock을 만든다음, BasicBlcok을 깊게 쌓아서 높은 성능을 달성
  - Pooling을 제거하고, Stage 맨 앞 convolution layer의 stride=2로 다운샘플링
    - Stage를 거침에 따라 channel 길이는 2배, spatial size는 1/2x1/2이 됨.
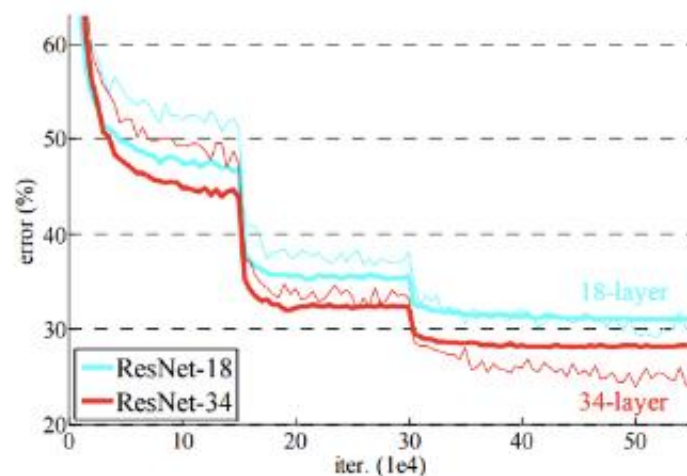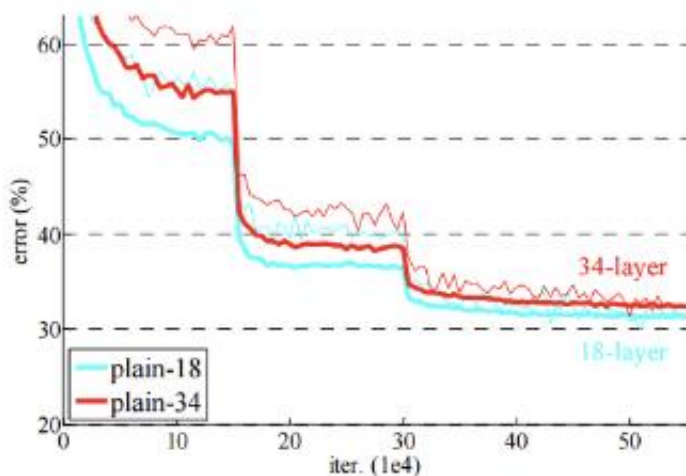  - Batch Normalization Layer(BN) 사용 **(Conv-BN-ReLU)**



BasicBlock

# 1. ResNet 설명

■ ResNet의 핵심 아이디어
  ▪ 결과적으로 깊은 층을 쌓아 성능 향상



Plain Nerwork VS ResNet

# 1. ResNet 설명

■ Batch Normalization Layer(BN)

- 공변량 시프트(Covariate shift)문제를 해결
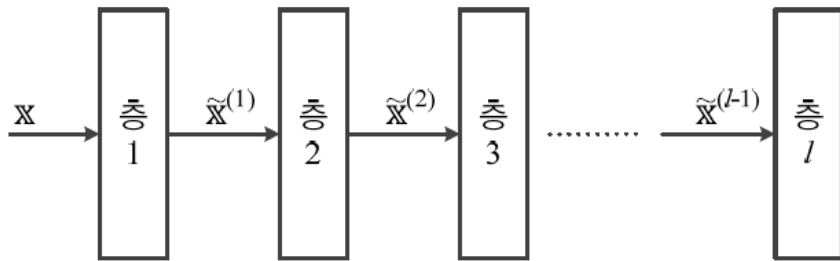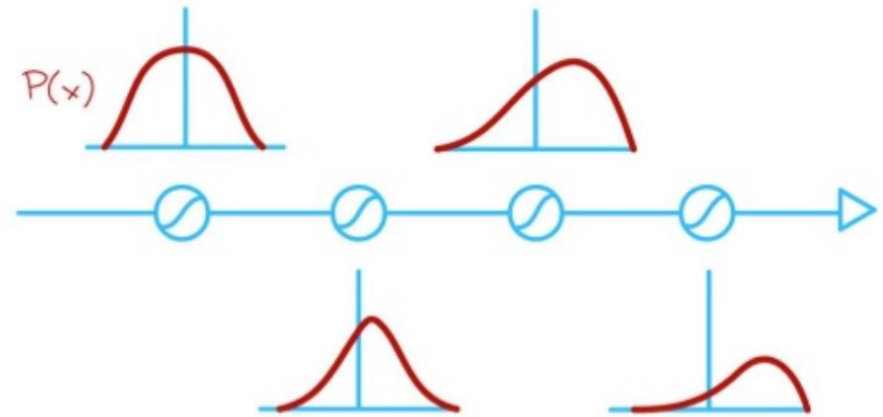- 일반적으로 Convolution Layer와 Activation Layer(e.g., ReLU) 사이에 사용



그림 5-17 공변량 시프트 현상

# 1. ResNet 설명

- **Batch Normalization Layer(BN)**
  - 미니배치 $\mathbb{X}_B = \{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_m\}$에 식 (5.15)를 적용하여 $\tilde{\mathbb{X}}_B = \{z_1, z_2, \cdots, z_m\}$를 얻은 후, $\tilde{\mathbb{X}}_B$를 가지고 코드 1을 수행
    - 노드마다 독립적으로 코드 1을 수행(CNN에서는 노드 단위가 아닌 특징맵 단위로 수행)
  - $\gamma$와 $\beta$: 노드마다 고유한 매개변수로서 학습으로 알아냄
  - 규제효과(overfitting 방지), 수렴속도 향상, 초기 learning-rate에 둔감성 증대 효과
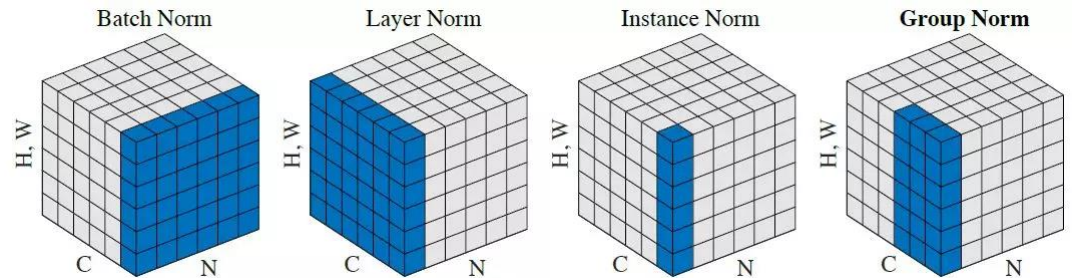
코드 1:

$$\mu_B = \frac{1}{m} \sum_{i=1}^{m} z_i$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^{m} (z_i - \mu_B)^2$$

$$\tilde{z}_i = \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}, \qquad i = 1, 2, \cdots, m$$

$$z_i' = \gamma \tilde{z}_i + \beta, \qquad i = 1, 2, \cdots, m$$

# 2. ResNet 구현

■ CIFAR-10에 동작하는 ResNet 구현
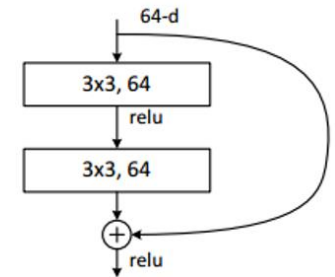
# 2. ResNet 구현

■ Dataloader 설정

```python
1  import torch
2  import torch.nn as nn
3  import torchvision
4  import torchvision.transforms as transforms
5
6
7  device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
8
9  num_epochs = 80
10 learning_rate = 0.001
11
12 transform = transforms.Compose([
13     transforms.Pad(4),
14     transforms.RandomHorizontalFlip(),
15     transforms.RandomCrop(32),
16     transforms.ToTensor()])
17
18 # CIFAR-10 dataset
19 train_dataset = torchvision.datasets.CIFAR10(root='../../data/',
20                                              train=True,
21                                              transform=transform,
22                                              download=True)
23
24 test_dataset = torchvision.datasets.CIFAR10(root='../../data/',
25                                             train=False,
26                                             transform=transforms.ToTensor())
27
28 # Data loader
29 train_loader = torch.utils.data.DataLoader(dataset=train_dataset,
30                                            batch_size=100,
31                                            shuffle=True)
32
33 test_loader = torch.utils.data.DataLoader(dataset=test_dataset,
34                                           batch_size=100,
35                                           shuffle=False)
```
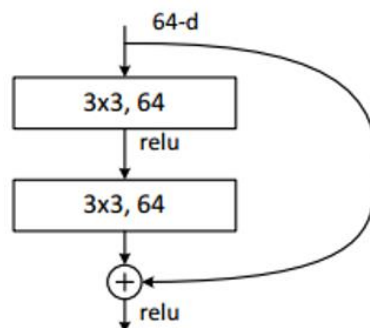
# 2. ResNet 구현

## BasicBlock (1/2)

```python
38  class BasicBlock(nn.Module):
39      def __init__(self, in_channels, out_channels, stride = 1):
40          super().__init__()
41
42          # BatchNorm에 bias가 포함되어 있으므로, conv2d는 bias=False로 설정합니다.
43          self.residual_function = nn.Sequential(
44              nn.Conv2d(in_channels, out_channels, kernel_size=3, stride=stride, padding=1, bias=False),
45              nn.BatchNorm2d(out_channels),
46              nn.ReLU(),
47              nn.Conv2d(out_channels, out_channels, kernel_size=3, stride=1, padding=1, bias=False),
48              nn.BatchNorm2d(out_channels),
49          )
50
51          # identity mapping, input과 output의 feature map size, filter 수가 동일한 경우 사용.
52          self.shortcut = nn.Sequential()
53
54          # projection mapping using 1x1conv
55          if stride != 1 or in_channels != out_channels:
56              self.shortcut = nn.Sequential(
57                  nn.Conv2d(in_channels, out_channels, kernel_size=1, stride= stride, bias=False),
58                  nn.BatchNorm2d(out_channels)
59              )
60
61          self.relu = nn.ReLU()
62          self.maxPool = nn.MaxPool2d(4)   # added for BasicBlock verificaion
63          self.fc = nn.Linear(out_channels * 8* 8, 10)  # added for BasicBlock verificaion (32/4 * 32/4)
64
```

# 2. ResNet 구현

■ BasicBlock (2/2)

```
65    def forward(self, x):
66        x = self.residual_function(x) + self.shortcut(x)
67        x = self.relu(x)
68        x = self.maxPool(x)         # added for BasicBlock verificaion
69        x = x.view(x.size(0), -1) # added for BasicBlock verificaionn
70        x = self.fc(x)              # # added for BasicBlock verificaion
71        return x
72
73 model = BasicBlock(3, 64)
74
75 criterion = nn.CrossEntropyLoss()
76 optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
```

# 2. ResNet 구현

- Training Phase

```
85  def train(epoch):
86      model.train()
87      for batch_idx, (data, target) in enumerate(train_loader):
88          data, target = data.to(device), target.to(device)
89          optimizer.zero_grad()
90          output = model(data)
91          loss = criterion(output, target)
92
93          loss.backward()
94          optimizer.step()
95
96          if (batch_idx+1) % 100 ==0:
97              print ("Epoch [{}/{}], step [{}/{}] Loss: {:4f}".format(epoch+1,
98                      num_epochs, batch_idx+1, total_step, loss.item()))
99
100     if (epoch+1) % 20 == 0:
101         curr_lr /= 3
102         update_lr(optimizer, curr_lr)
```

# 2. ResNet 구현

## Test Phase

```
104 def test():
105     model.eval()
106     test_loss = 0
107     correct = 0
108     for data, target in test_loader:
109         data, target = data.to(device), target.to(device)
110         output = model(data)
111         # sum up batch loss
112         test_loss +=  criterion(output, target)
113         # get the index of the max log-probability
114         pred = output.data.max(1, keepdim=True)[1]
115         correct += pred.eq(target.data.view_as(pred)).cpu().sum()
116
117     test_loss /= len(test_loader.dataset)
118     print('\nTest set: Average loss: {:.4f}, Accuracy: {}/{} ({:.0f}%)\n'.format(
119         test_loss, correct, len(test_loader.dataset),
120         100. * correct / len(test_loader.dataset)))
121
122
123 for epoch in range(0, num_epochs):
124     train(epoch)
125     test()
```

```
Files already downloaded and verified
Epoch [1/80], step [100/500] Loss: 1.752413
Epoch [1/80], step [200/500] Loss: 1.534544
Epoch [1/80], step [300/500] Loss: 1.463914
Epoch [1/80], step [400/500] Loss: 1.280996
Epoch [1/80], step [500/500] Loss: 1.088101
```

# 2. ResNet 구현

■ ResNet18 구현(1/2)

```python
73 class ResNet(nn.Module):
74     def __init__(self, block, num_block, num_classes=10, init_weights=True):
75         super().__init__()
76
77         self.in_channels=64
78
79         self.conv1 = nn.Sequential(
80             nn.Conv2d(3, 64, kernel_size=7, stride=2, padding=3, bias=False),
81             nn.BatchNorm2d(64),
82             nn.ReLU(),
83             nn.MaxPool2d(kernel_size=3, stride=2, padding=1)
84         )
85
86         self.conv2_x = self._make_layer(block, 64, num_block[0], 1)
87         self.conv3_x = self._make_layer(block, 128, num_block[1], 2)
88         self.conv4_x = self._make_layer(block, 256, num_block[2], 2)
89         self.conv5_x = self._make_layer(block, 512, num_block[3], 2)
90
91         self.avg_pool = nn.AdaptiveAvgPool2d((1,1))
92         self.fc = nn.Linear(512, num_classes)
93
```

```python
114 model = ResNet(BasicBlock, [2,2,2,2])
```

# 2. ResNet 구현

■ ResNet18 구현 (2/2)

```
94      def _make_layer(self, block, out_channels, num_blocks, stride):
95          strides = [stride] + [1] * (num_blocks - 1)
96          layers = []
97          for stride in strides:
98              layers.append(block(self.in_channels, out_channels, stride))
99              self.in_channels = out_channels
100
101         return nn.Sequential(*layers)
102
103     def forward(self,x):
104         output = self.conv1(x)
105         output = self.conv2_x(output)
106         x = self.conv3_x(output)
107         x = self.conv4_x(x)
108         x = self.conv5_x(x)
109         x = self.avg_pool(x)
110         x = x.view(x.size(0), -1)
111         x = self.fc(x)
112         return x
113
114 model = ResNet(BasicBlock, [2,2,2,2])
```

# 실습 1

■ Batch Normalization 외에 아래의 4개 이상의 Normalization 방법에 대
해 찾아보고, 구현한 결과를 보이시오

- (예시) LayerNorm: https://pytorch.org/docs/stable/generated/torch.nn.LayerNorm.html

# 실습 2

- ResNet34를 구현하시오.

| 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 | 3×3 max pool, stride 2 | | | | |
| | | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 64 \\ 3\times3, 64 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 64 \\ 3\times3, 64 \\ 1\times1, 256 \end{bmatrix}\times3$ |
| conv3_x | 28×28 | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 128 \\ 3\times3, 128 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1, 128 \\ 3\times3, 128 \\ 1\times1, 512 \end{bmatrix}\times8$ |
| conv4_x | 14×14 | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 256 \\ 3\times3, 256 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix}\times23$ | $\begin{bmatrix} 1\times1, 256 \\ 3\times3, 256 \\ 1\times1, 1024 \end{bmatrix}\times36$ |
| conv5_x | 7×7 | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3, 512 \\ 3\times3, 512 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1, 512 \\ 3\times3, 512 \\ 1\times1, 2048 \end{bmatrix}\times3$ |
| | 1×1 | average pool, 1000-d fc, softmax | | | | |
| FLOPs | | $1.8\times10^9$ | $3.6\times10^9$ | $3.8\times10^9$ | $7.6\times10^9$ | $11.3\times10^9$ |

# 실습 3

- 50층 이상의 ResNet은 매개변수 사용을 최소화 하기 위해 BottleNeck 을 이용한다. BottleNeck을 구현하고, 이를 이용하여 ResNet50도 구현하시오(구현 결과를 CIFAR10에 학습/검증하시오)
  - BottleNeck: Conv1x1을 이용하여 채널 차원 축소 후 Conv3x3을 사용 – 다시 Conv1x1로 채널 차원을 복원하는 방식으로 동작하는 Block
  - (참고) https://deep-learning-study.tistory.com/534

**BasicBlock**

```
64-d
  │
3x3, 64
  │ relu
3x3, 64
  │
 (+)←
  │ relu
```

**BottleNeck**

```
256-d
  │
1x1, 64
  │ relu
3x3, 64
  │ relu
1x1, 256
  │
 (+)←
  │ relu
```

| 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 | | | | |
| conv2_x | 56×56 | [3×3, 64; 3×3, 64] ×2 | [3×3, 64; 3×3, 64] ×3 | [1×1, 64; 3×3, 64; 1×1, 256] ×3 | [1×1, 64; 3×3, 64; 1×1, 256] ×3 | [1×1, 64; 3×3, 64; 1×1, 256] ×3 |
| conv3_x | 28×28 | [3×3, 128; 3×3, 128] ×2 | [3×3, 128; 3×3, 128] ×4 | [1×1, 128; 3×3, 128; 1×1, 512] ×4 | [1×1, 128; 3×3, 128; 1×1, 512] ×4 | [1×1, 128; 3×3, 128; 1×1, 512] ×8 |
| conv4_x | 14×14 | [3×3, 256; 3×3, 256] ×2 | [3×3, 256; 3×3, 256] ×6 | [1×1, 256; 3×3, 256; 1×1, 1024] ×6 | [1×1, 256; 3×3, 256; 1×1, 1024] ×23 | [1×1, 256; 3×3, 256; 1×1, 1024] ×36 |
| conv5_x | 7×7 | [3×3, 512; 3×3, 512] ×2 | [3×3, 512; 3×3, 512] ×3 | [1×1, 512; 3×3, 512; 1×1, 2048] ×3 | [1×1, 512; 3×3, 512; 1×1, 2048] ×3 | [1×1, 512; 3×3, 512; 1×1, 2048] ×3 |
| | 1×1 | average pool, 1000-d fc, softmax | | | | |
| FLOPs | | $1.8×10^9$ | $3.6×10^9$ | $3.8×10^9$ | $7.6×10^9$ | $11.3×10^9$ |

# 실습 4

- 모델의 학습 시 가중치 초기화는 성능에 큰 영향을 미친다.
  - Xavier_normal방법, Kaiming_ normal 방법이 무엇인지 인터넷을 통해 알아보고, 이 방법들로 가중치를 초기화 하여 학습 후 결과를 비교하시오.
  - (참고1) https://pytorch.org/docs/stable/nn.init.html
  - (참고2) https://deep-learning-study.tistory.com/534

# 실습 5

- 아래 ResNet 논문을 읽고, 논문과 동일한 실험 조건에서 ResNet50을 ImageNet 데이터에 대해 학습하고 실험 결과를 보이시오.
  - (논문) https://www.cv-foundation.org/openaccess/content_cvpr_2016/papers/He_Deep_Residual_Learning_CVPR_2016_paper.pdf