

#1

$W_1$ 은  $W$ 의 first layer,  $W_2$ 은 second layer라고 할때,  $W_1$ 과  $W_2$ 은 다음과 같이 주어진다.

$$W_{1,i,j} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \in \mathbb{R}^{3 \times 3}$$

$$W_{2,i,j} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \in \mathbb{R}^{3 \times 3}$$

$$X = \begin{bmatrix} X_{1,j} & X_{1,j+1} & X_{1,j+2} \\ X_{2,j} & X_{2,j+1} & X_{2,j+2} \\ X_{3,j} & X_{3,j+1} & X_{3,j+2} \end{bmatrix}$$

$$W_1 \text{과 } X \text{의 convolution} \rightarrow -X_{i,j} + X_{i+1,j} = Y_{1,i,j}$$

$$W_2 \text{와 } X \text{의 convolution} \rightarrow -X_{i,j} + X_{i,j+1} = Y_{2,i,j}$$

가 되는 항은 0이므로 생략한다.

#2 filter  $W \in \mathbb{R}^{c \times c \times k \times k}$ . (Output의 layer가  $C$ 개일때  $c \times c \times k \times k$ 개의 filter가  $C$ 개 필요함을 알 수 있다.)

$W$ 의  $i$ 번째 layer, 즉  $i$ 번째 filter  $W_i \in \mathbb{R}^{c \times c \times k \times k}$ 를 생각하자.

$$W_{i,j} = \begin{bmatrix} 1/k^2 & 1/k^2 & \dots \\ 1/k^2 & 1/k^2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

(if  $j=i$ )

$W_i$ 의  $j$ 번째 layer를 살펴보자.  $j$ 가  $i$ 와

같은  $W_{i,j,j,j}$ 는 모든 원소가  $1/k^2$ 이여야 한다.

$j \neq i$ 인 경우,  $W_{i,j,j,j}$ 의 모든 원소는 0이다.

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \ddots \end{bmatrix}$$

(else)

이렇게  $W$ 을 구성하면  $W$ 은

Approximation:  $X \rightarrow Y$ 에 해당하는 filter이다.

따, 이 경우 convolution의 stride는  $k$ 이다.

#3  $W \in \mathbb{R}^{3 \times 1 \times 1}$ 가  $[0.299, 0.587, 0.114]$ 로 주어진다.

$X_{:,i,j}$ 과  $W$ 의 convolution은 생각해보면

$$\text{Output } Y_{i,j} = X_{1,i,j} \times W_{1,1,1} + X_{2,i,j} \times W_{2,1,1} + X_{3,i,j} \times W_{3,1,1}$$

$$= 0.299 X_{1,i,j} + 0.587 X_{2,i,j} + 0.114 X_{3,i,j} \text{ 을 얻는다.}$$

#4.  $\tau: \mathbb{R} \rightarrow \mathbb{R}$ , non decreasing activation fn.

$p: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{k \times l}$ : max pool operation.

Show for all  $X \in \mathbb{R}^{m \times n}$ ,  $\tau$  and  $p$  commute.  $X$  is  $(i,j)$  entry is  $X_{ij}$  at  $i$ th row and  $j$ th column.

Pf)  $f_1 = \frac{m}{k}$ ,  $f_2 = \frac{n}{l}$  at  $i$ th row and  $j$ th column.  $p(X) \in \mathbb{R}^{k \times l}$  is  $(i,j)$  entry is  $p_{ij}(X)$  at  $i$ th row and  $j$ th column.

$$p_{ij}(X) = \max Y_{ij}, \quad Y_{ij} := \{ X_{\alpha, \beta} \mid (i-1)f_1 < \alpha \leq if_1, (j-1)f_2 < \beta \leq jf_2 \}$$

$$\tau(p_{ij}(X)) = \tau\left(\max(Y_{ij})\right) = \tau\left(\max\left(\{X_{\alpha, \beta} \mid (i-1)f_1 < \alpha \leq if_1, (j-1)f_2 < \beta \leq jf_2\}\right)\right)$$

\* Since  $\tau$  is a non decreasing function,  
 $\max\{\tau(x_1), \tau(x_2), \dots, \tau(x_n)\} = \tau(\max\{x_1, x_2, \dots, x_n\})$  holds.

$$= \max\left(\{ \tau(X_{\alpha, \beta}) \mid (i-1)f_1 < \alpha \leq if_1, (j-1)f_2 < \beta \leq jf_2 \}\right)$$

$$= p_{ij}(\tau(X)) !$$

$$\therefore \tau(p_{ij}(X)) = p_{ij}(\tau(X)) \text{ for all } (i,j)$$

and thus,  $\tau(p(X)) = p(\tau(X))$  holds.

# Problem 5 : Non CE loss function

```
In [4]: import torch
import torch.nn as nn
from torch.optim import Optimizer
from torch.utils.data import DataLoader
from torchvision import datasets
from torchvision.transforms import transforms
import matplotlib.pyplot as plt
from random import shuffle
'''
Step 1: Data
'''
# Use data with only 4 and 9 as labels: which is hardest to classify
label_1, label_2 = 4, 9

# MNIST training data
train_set = datasets.MNIST(root='./mnist_data/', train=True, transform=transforms.ToTensor(), download=True)

# Use data with two labels
idx = (train_set.targets == label_1) + (train_set.targets == label_2)
train_set.data = train_set.data[idx]
train_set.targets = train_set.targets[idx]
train_set.targets[train_set.targets == label_1] = -1
train_set.targets[train_set.targets == label_2] = 1

# MNIST testing data
test_set = datasets.MNIST(root='./mnist_data/', train=False, transform=transforms.ToTensor())

# Use data with two labels
idx = (test_set.targets == label_1) + (test_set.targets == label_2)
test_set.data = test_set.data[idx]
test_set.targets = test_set.targets[idx]
test_set.targets[test_set.targets == label_1] = -1
test_set.targets[test_set.targets == label_2] = 1

'''
Step 2: (same step)
'''
class Linear(nn.Module) :

    def __init__(self, input_dim=28*28) :
        super().__init__()
        self.linear = nn.Linear(input_dim, 1, bias=False)

    def forward(self, x) :
        return self.linear(x.float().view(-1, 28*28))

'''
Step 3: Create the model, specify loss function and optimizer. (LOOK HERE)
'''
model_CE = Linear()
model_MSE = Linear()

def CE_loss(output, target):
    return torch.mean(-torch.nn.functional.logsigmoid(target.reshape(-1)*output.reshape(-1)))

def MSE_loss(output, target):
    output = output.reshape(-1)
    target = target.reshape(-1)
    B = len(output)
    total = 0
    for i in range(B):
        y = target[i]
        z = output[i]
        total += 0.5*(1-y)*((1-torch.sigmoid(-z))**2 + torch.sigmoid(z)**2)
        total += 0.5*(1+y)*(torch.sigmoid(-z)**2 + (1-torch.sigmoid(z))**2)
    return total/B

optimizer_CE = torch.optim.SGD(model_CE.parameters(), lr=255*1e-4)
optimizer_MSE = torch.optim.SGD(model_MSE.parameters(), lr=255*1e-4)

'''
Step 4: Train model with SGD (LOOK HERE)
'''
train_loader = DataLoader(dataset=train_set, batch_size=64, shuffle=True)

for epoch in range(3) :
    for images, labels in train_loader :
        optimizer_CE.zero_grad()
        optimizer_MSE.zero_grad()

        train_loss_CE = CE_loss(model_CE(images), labels.float())
        train_loss_CE.backward()

        train_loss_MSE = MSE_loss(model_MSE(images), labels.float())
        train_loss_MSE.backward()

        optimizer_CE.step()
        optimizer_MSE.step()

'''
Step 5: (same step)
'''
test_loss_CE, correct_CE = 0, 0
test_loss_MSE, correct_MSE = 0, 0

test_loader = DataLoader(dataset=test_set, batch_size=1, shuffle=False)

for ind, (image, label) in enumerate(test_loader) :

    output_CE = model_CE(image)
    output_MSE = model_MSE(image)
    test_loss_CE += CE_loss(output_CE, label.float()).item()
    test_loss_MSE += MSE_loss(output_MSE, label.float()).item()

    # Make a prediction
    if output_CE.item() * label.item() >= 0 :
        correct_CE += 1
    if output_MSE.item() * label.item() >= 0 :
        correct_MSE += 1

# Print out the results
print("Cross entropy loss")
print(['[Test set] Average loss: {:.4f}, Accuracy: {}/{} ({:.2f}%)\n'.format(
    test_loss_CE /len(test_loader), correct_CE, len(test_loader),
    100. * correct_CE / len(test_loader))])
print("Mean squared loss")
print(['[Test set] Average loss: {:.4f}, Accuracy: {}/{} ({:.2f}%)\n'.format(
    test_loss_MSE /len(test_loader), correct_MSE, len(test_loader),
    100. * correct_MSE / len(test_loader))])
```

Cross entropy loss  
[Test set] Average loss: 0.1749, Accuracy: 1894/1991 (95.13%)

Mean squared loss  
[Test set] Average loss: 0.1007, Accuracy: 1896/1991 (95.23%)

큰 차이가 없는 것으로 보여진다.

$$\#6 (a) \quad y_L = A_L y_{L-1} + b_L, \quad A_L \in \mathbb{R}^{n_L \times n_{L-1}}, \quad b_L \in \mathbb{R}^{n_L}, \quad y_L \in \mathbb{R}^{n_L}$$

$$\Rightarrow \frac{\partial y_L}{\partial b_L} = I, \quad \frac{\partial y_L}{\partial y_{L-1}} = \left( \frac{\partial y_L}{\partial (y_{L-1})_1}, \dots, \frac{\partial y_L}{\partial (y_{L-1})_{n_{L-1}}} \right)$$

( $y_L, b_L$  is scalar)

$$\left( \frac{\partial y_L}{\partial y_{L-1}} \right)_i = \frac{\partial y_L}{\partial y_i} = \frac{\sum_{j=1}^{n_{L-1}} (A_L)_{i,j} \times (y_{L-1})_j}{\partial y_i} = (A_L)_{i,i}$$

$$\Rightarrow \frac{\partial y_L}{\partial y_{L-1}} = A_L$$

$$y_\ell = \nabla(A_\ell y_{\ell-1} + b_\ell) \quad (\ell=1, \dots, L-1) \quad A_\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}, \quad b_\ell \in \mathbb{R}^{n_\ell}$$

$$\frac{\partial y_\ell}{\partial b_\ell} = \begin{pmatrix} \frac{\partial (y_\ell)_1}{\partial (b_\ell)_1} & \dots & \frac{\partial (y_\ell)_1}{\partial (b_\ell)_{n_\ell}} \\ \vdots & & \vdots \\ \frac{\partial (y_\ell)_{n_\ell}}{\partial (b_\ell)_1} & \dots & \frac{\partial (y_\ell)_{n_\ell}}{\partial (b_\ell)_{n_\ell}} \end{pmatrix} \quad \left( \frac{\partial y_\ell}{\partial b_\ell} \right)_{i,j} = \frac{\partial (y_\ell)_i}{\partial (b_\ell)_j} = \frac{\partial \left( \nabla \left( (b_\ell)_j + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{i,k} (y_{\ell-1})_k \right) \right)}{\partial (b_\ell)_j}$$

$$= \begin{cases} \nabla'((b_\ell)_j + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{i,k} (y_{\ell-1})_k) & (i=j) \\ 0 & (\text{else}) \end{cases}$$

$$\Rightarrow \frac{\partial y_\ell}{\partial b_\ell} = \begin{pmatrix} \nabla'((b_\ell)_1 + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{1,k} (y_{\ell-1})_k) & 0 & \dots & 0 \\ 0 & \ddots & & \\ 0 & & \nabla'((b_\ell)_{n_\ell} + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{n_\ell,k} (y_{\ell-1})_k) & \\ & & & \ddots \end{pmatrix} = \text{diag}(\nabla'(b_\ell + A_\ell y_{\ell-1}))$$

$$\left( \frac{\partial y_\ell}{\partial y_{\ell-1}} \right)_{i,j} = \frac{\partial (y_\ell)_i}{\partial (y_{\ell-1})_j} = \frac{\partial \left( \nabla \left( (b_\ell)_i + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{i,k} (y_{\ell-1})_k \right) \right)}{\partial (y_{\ell-1})_j} = \nabla'((b_\ell)_i + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{i,k} (y_{\ell-1})_k) \times (A_\ell)_{i,j}$$

$$\Rightarrow \left( \frac{\partial y_\ell}{\partial y_{\ell-1}} \right) \in \mathbb{R}^{n_\ell \times n_{\ell-1}} = \begin{pmatrix} \nabla'((b_\ell)_1 + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{1,k} (y_{\ell-1})_k) & 0 & \dots & 0 \\ 0 & \ddots & & \\ 0 & & \nabla'((b_\ell)_{n_\ell} + \sum_{k=1}^{n_{\ell-1}} (A_\ell)_{n_\ell,k} (y_{\ell-1})_k) & \\ & & & \ddots \end{pmatrix} (A_\ell)$$

$$= \text{diag}(\nabla'(b_\ell + A_\ell y_{\ell-1})) A_\ell$$

$$(b) \quad \frac{\partial y_L}{\partial A_L} \in \mathbb{R}^{n_L \times n_{L-1}} \quad \left( \frac{\partial y_L}{\partial A_L} \right)_{ij} = \frac{\partial y_L}{\partial (A_L)_{ij}}$$

$$\left( \frac{\partial y_L}{\partial A_L} \right)_{ij} = \frac{\partial y_L}{\partial (A_L)_{ij}} = \frac{\partial (b_L + \sum_{i=1}^{n_{L-1}} (A_L)_{i,j} (y_{L-1})_i)}{\partial (A_L)_{i,j}} = (y_{L-1})_i \Rightarrow \underline{\left( \frac{\partial y_L}{\partial A_L} \right) = y_{L-1}^T}$$

$$\begin{aligned} \left( \frac{\partial y_L}{\partial A_L} \right)_{i,j} &= \frac{\partial y_L}{\partial (A_L)_{i,j}} = \frac{\partial y_L}{\partial y_L} \cdot \frac{\partial y_L}{\partial (A_L)_{i,j}} = \left( \frac{\partial y_L}{\partial y_{L-1}} \cdot \frac{\partial y_{L-1}}{\partial y_{L-2}} \cdots \frac{\partial y_{L-1}}{\partial y_L} \right) \frac{\partial y_L}{\partial (A_L)_{i,j}} \\ &= \left( \frac{\partial y_L}{\partial y_L} \right) \begin{pmatrix} \frac{\partial \left( \nabla(b_L)_i + \sum_{k=1}^{n_{L-1}} (A_L)_{i,k} (y_{L-1})_k \right)}{\partial (A_L)_{i,j}} \\ \vdots \\ \frac{\partial \left( \nabla(b_L)_{n_L} + \sum_{k=1}^{n_{L-1}} (A_L)_{n_L,k} (y_{L-1})_k \right)}{\partial (A_L)_{i,j}} \end{pmatrix} \\ &= \left( \frac{\partial y_L}{\partial y_L} \right) \begin{pmatrix} \nabla'(b_L)_i + \sum_{k=1}^{n_{L-1}} (A_L)_{i,k} (y_{L-1})_k \\ \vdots \\ \nabla'(b_L)_{n_L} + \sum_{k=1}^{n_{L-1}} (A_L)_{n_L,k} (y_{L-1})_k \end{pmatrix} = \left( \frac{\partial y_L}{\partial y_L} \right)_i (y_{L-1})_j \times \nabla'(b_L)_i + \sum_{k=1}^{n_{L-1}} (A_L)_{i,k} (y_{L-1})_k \end{aligned}$$

$$\therefore \underline{\left( \frac{\partial y_L}{\partial A_L} \right)} = \text{diag} \left( \nabla'(b_L + A_L y_{L-1}) \right) \begin{pmatrix} \left( \frac{\partial y_L}{\partial y_L} \right)_1 (y_{L-1})_1, \left( \frac{\partial y_L}{\partial y_L} \right)_1 (y_{L-1})_2, \dots, \left( \frac{\partial y_L}{\partial y_L} \right)_1 (y_{L-1})_{n_{L-1}} \\ \vdots \\ \left( \frac{\partial y_L}{\partial y_L} \right)_{n_L} (y_{L-1})_1, \dots, \left( \frac{\partial y_L}{\partial y_L} \right)_{n_L} (y_{L-1})_{n_{L-1}} \end{pmatrix}$$

$$= \underline{\text{diag} \left( \nabla'(b_L + A_L y_{L-1}) \right) \left( \frac{\partial y_L}{\partial y_L} \right)^T y_{L-1}^T}$$

## #Problem 7

먼저 trainable parameter의 수에 대해 계산을 하고 시작하자.

C3\_layer\_full, 즉 일반적인 convolutional net을 사용한다면,

$$6 \times (1 \times 5 \times 5 + 1) + 16 \times (6 \times 5 \times 5 + 1) + 120 \times (1 + 5 \times 5 \times 16) + 84 \times (1 + 120) + 10 \times (1 + 84) = 61706$$

(실제로 돌려보고 print로 확인해도 같은 결과가 나온다..)

Original LeNet의 C3 layer를 사용한다면,

$$\begin{aligned} &6 \times (1 \times 5 \times 5 + 1) + \\ &6 \times (3 \times 5 \times 5 + 1) + 9 \times (4 \times 5 \times 5 + 1) + (6 \times 5 \times 5 + 1) + \\ &120 \times (1 + 5 \times 5 \times 16) + 84 \times (1 + 120) + 10 \times (1 + 84) = 60806 \end{aligned}$$

개의 parameter가 있을 것이다. 총 900개의 parameter reduction이 발생한다.

코드는 다음과 같이 짜주면 된다.

```
lenet_original.py x
1 import torch
2 import torch.nn as nn
3 from torch.optim import Optimizer
4 from torch.utils.data import DataLoader
5 from torchvision import datasets
6 from torchvision.transforms import transforms
7
8 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
9
10
11 Step 1:
12
13
14 # MNIST dataset
15 train_dataset = datasets.MNIST(root='./mnist_data/',
16                               train=True,
17                               transform=transforms.ToTensor(),
18                               download=True)
19
20 test_dataset = datasets.MNIST(root='./mnist_data/',
21                              train=False,
22                              transform=transforms.ToTensor())
23
24
25
26 Step 2: LeNet5
27
28
29 # Modern LeNet uses this layer for C3
30 class C3_layer_full(nn.Module):
31     def __init__(self):
32         super(C3_layer_full, self).__init__()
33         self.conv_layer = nn.Conv2d(6, 16, kernel_size=5)
34
35     def forward(self, x):
36         return self.conv_layer(x)
```



```

38 # Original LeNet uses this layer for C3
39 class C3_layer(nn.Module):
40     def __init__(self):
41         super(C3_layer, self).__init__()
42         self.ch_in_3 = [[0, 1, 2],
43                         [1, 2, 3],
44                         [2, 3, 4],
45                         [3, 4, 5],
46                         [0, 4, 5],
47                         [0, 1, 5]] # filter with 3 subset of input channels
48         self.ch_in_4 = [[0, 1, 2, 3],
49                         [1, 2, 3, 4],
50                         [2, 3, 4, 5],
51                         [0, 3, 4, 5],
52                         [0, 1, 4, 5],
53                         [0, 1, 2, 5],
54                         [0, 1, 3, 4],
55                         [1, 2, 4, 5],
56                         [0, 2, 3, 5]] # filter with 4 subset of input channels
57
58         # put implementation here
59         self.conv3layers = nn.ModuleList([nn.Conv2d(3,1, kernel_size = 5) for i in range(6)])
60         self.conv4layers = nn.ModuleList([nn.Conv2d(4,1, kernel_size = 5) for i in range(9)])
61         self.conv6layer = nn.Conv2d(6,1, kernel_size = 5)
62
63
64     def forward(self, x):
65         L = []
66         for i, layer in enumerate(self.conv3layers):
67             L.append(layer(x[:, self.ch_in_3[i], :, :]))
68         for i, layer in enumerate(self.conv4layers):
69             L.append(layer(x[:, self.ch_in_4[i], :, :]))
70         L.append(self.conv6layer(x))
71         return torch.cat(L, dim=1) #batch
72

```

```

74 class LeNet(nn.Module):
75     def __init__(self):
76         super(LeNet, self).__init__()
77         #padding=2 makes 28x28 image into 32x32
78         self.C1_layer = nn.Sequential(
79             nn.Conv2d(1, 6, kernel_size=5, padding=2),
80             nn.Tanh()
81         )
82         self.P2_layer = nn.Sequential(
83             nn.AvgPool2d(kernel_size=2, stride=2),
84             nn.Tanh()
85         )
86         self.C3_layer = nn.Sequential(
87             #C3_layer_full(),
88             C3_layer(),
89             nn.Tanh()
90         )
91         self.P4_layer = nn.Sequential(
92             nn.AvgPool2d(kernel_size=2, stride=2),
93             nn.Tanh()
94         )
95         self.C5_layer = nn.Sequential(
96             nn.Linear(5*5*16, 120),
97             nn.Tanh()
98         )
99         self.F6_layer = nn.Sequential(
100             nn.Linear(120, 84),
101             nn.Tanh()
102         )
103         self.F7_layer = nn.Linear(84, 10)
104         self.tanh = nn.Tanh()
105
106     def forward(self, x):
107         output = self.C1_layer(x)
108         output = self.P2_layer(output)
109         output = self.C3_layer(output)
110         output = self.P4_layer(output)
111         output = output.view(-1, 5*5*16)
112         output = self.C5_layer(output)
113         output = self.F6_layer(output)
114         output = self.F7_layer(output)
115         return output
116
117
118 """
119 Step 3
120 """
121 model = LeNet().to(device)
122 loss_function = torch.nn.CrossEntropyLoss()
123 optimizer = torch.optim.SGD(model.parameters(), lr=1e-1)
124
125 # print total number of trainable parameters
126 param_ct = sum([p.numel() for p in model.parameters()])
127 print(f"Total number of trainable parameters: {param_ct}")
128
129 """

```

```

129     '''
130     Step 4
131     '''
132     train_loader = torch.utils.data.DataLoader(dataset=train_dataset, batch_size=100, shuffle=True)
133
134     import time
135     start = time.time()
136     for epoch in range(10) :
137         print("{}th epoch starting.".format(epoch))
138         for images, labels in train_loader :
139             images, labels = images.to(device), labels.to(device)
140
141             optimizer.zero_grad()
142             train_loss = loss_function(model(images), labels)
143             train_loss.backward()
144
145             optimizer.step()
146     end = time.time()
147     print("Time ellapsed in training is: {}".format(end - start))
148
149     '''
150     Step 5
151     '''
152     test_loss, correct, total = 0, 0, 0
153
154     test_loader = torch.utils.data.DataLoader(dataset=test_dataset, batch_size=100, shuffle=False)
155
156     for images, labels in test_loader :
157         images, labels = images.to(device), labels.to(device)
158
159         output = model(images)
160         test_loss += loss_function(output, labels).item()
161
162         pred = output.max(1, keepdim=True)[1]
163         correct += pred.eq(labels.view_as(pred)).sum().item()
164
165         total += labels.size(0)
166
167     print('[Test set] Average loss: {:.4f}, Accuracy: {}/{} ({:.2f}%)\n'.format(
168         test_loss / total, correct, total,
169         100. * correct / total))
170
171

```

```

In [3]: runfile('C:/Users/sylee/OneDrive/바탕 화면/실신개/HW4/
lenet_original.py', wdir='C:/Users/sylee/OneDrive/바탕 화면/실신개/HW4')
Total number of trainable parameters: 60806
0th epoch starting.
1th epoch starting.
2th epoch starting.
3th epoch starting.
4th epoch starting.
5th epoch starting.
6th epoch starting.
7th epoch starting.
8th epoch starting.
9th epoch starting.
Time ellapsed in training is: 189.81498003005981
[Test set] Average loss: 0.0004, Accuracy: 9867/10000 (98.67%)

```

위와 같이 잘 train 된 것을 볼 수 있다. 계산한 파라미터 개수도 맞는다.