

# Programming assignment 2: shallow neural networks

```
for x,y in zip(X,Y):
    a2, (z1,a1,z2, _) = model.predict(x)
    if y == 1:
        cost -= np.log(a2)
    else:
        cost -= np.log(1-a2)

    diff = a2-y
    # layer 2
    db2 += diff
    dW2 += diff * a1

    #layer 1
    db1_tmp = diff * model.W2 * (1-a1**2)
    db1 += db1_tmp
    dW1 += db1_tmp.reshape(model.num_hidden,1) * x.reshape(1,model.num_input_features)
```

$$\frac{\partial L(a^{[2]}, y)}{\partial b^{[2]}} = a^{[2]} - y$$
$$\frac{\partial L(a^{[2]}, y)}{\partial w_i^{[2]}} = (a^{[2]} - y) a_i^{[1]}$$
$$\frac{\partial L(a^{[2]}, y)}{\partial b_i^{[1]}} = (a^{[2]} - y) w_i^{[2]} (1 - a_i^{[1]^2})$$
$$\frac{\partial L(a^{[2]}, y)}{\partial W_{ij}^{[1]}} = (a^{[2]} - y) w_i^{[2]} (1 - a_i^{[1]^2}) x_j$$

⊙: element-wise multiplication

$$\frac{\partial L(a^{[2]}, y)}{\partial b_i^{[1]}} = (a^{[2]} - y)w_i^{[2]}(1 - a_i^{[1]^2})$$

$$\frac{\partial L(a^{[2]}, y)}{\partial W_{ij}^{[1]}} = (a^{[2]} - y)w_i^{[2]}(1 - a_i^{[1]^2})x_j$$

$$\nabla_{\mathbf{b}^{[1]}} = \begin{bmatrix} \frac{\partial L(a^{[2]}, y)}{\partial b_1^{[1]}} \\ \vdots \\ \frac{\partial L(a^{[2]}, y)}{\partial b_h^{[1]}} \end{bmatrix} = \begin{bmatrix} (a^{[2]} - y)w_1^{[2]}(1 - a_1^{[1]^2}) \\ \vdots \\ (a^{[2]} - y)w_h^{[2]}(1 - a_1^{[1]^2}) \end{bmatrix} = (a^{[2]} - y)(\mathbf{w}^{[2]} \odot (1 - \mathbf{a}^{[1]^2}))$$

$$\nabla_{\mathbf{W}^{[1]}} = \begin{bmatrix} \frac{\partial L(a^{[2]}, y)}{\partial W_{11}^{[1]}} & \dots & \frac{\partial L(a^{[2]}, y)}{\partial W_{1n}^{[1]}} \\ \vdots & \ddots & \vdots \\ \frac{\partial L(a^{[2]}, y)}{\partial W_{h1}^{[1]}} & \dots & \frac{\partial L(a^{[2]}, y)}{\partial W_{hn}^{[1]}} \end{bmatrix} = \begin{bmatrix} \frac{\partial L(a^{[2]}, y)}{\partial b_1^{[1]}}x_1 & \dots & \frac{\partial L(a^{[2]}, y)}{\partial b_1^{[1]}}x_n \\ \vdots & \ddots & \vdots \\ \frac{\partial L(a^{[2]}, y)}{\partial b_h^{[1]}}x_1 & \dots & \frac{\partial L(a^{[2]}, y)}{\partial b_h^{[1]}}x_n \end{bmatrix} = \nabla_{\mathbf{b}^{[1]}} \mathbf{x}^T$$

```
#layer 1
db1_tmp = diff * model.W2 * (1-a1**2)
db1 += db1_tmp
dW1 += db1_tmp.reshape(model.num_hidden,1) * x.reshape(1,model.num_input_features)
```

# A simpler solution

```
#layer 1
db1_tmp = diff * model.W2 * (1-a1**2)
db1 += db1_tmp
dW1 += np.outer(db1_tmp, x)
```

`numpy.outer(a, b, out=None)`

[\[source\]](#)

Compute the outer product of two vectors.

Given two vectors,  $a = [a_0, a_1, \dots, a_M]$  and  $b = [b_0, b_1, \dots, b_N]$ , the outer product [1] is:

```
[[a0*b0  a0*b1  ...  a0*bN ]
 [a1*b0      .
 [ ...      .
 [aM*b0      aM*bN ]]
```

```
In [9]: np.outer(np.array([1,2,3]),np.array([1,2]))
```

```
Out[9]: array([[1, 2],
               [2, 4],
               [3, 6]])
```

# Deep Neural Networks



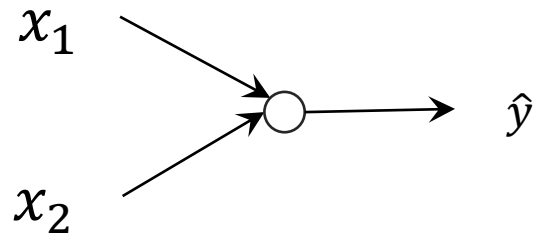
한양대학교 ERICA  
소프트웨어융합대학  
COLLEGE OF COMPUTING

인공지능학과  
Department of  
Artificial Intelligence

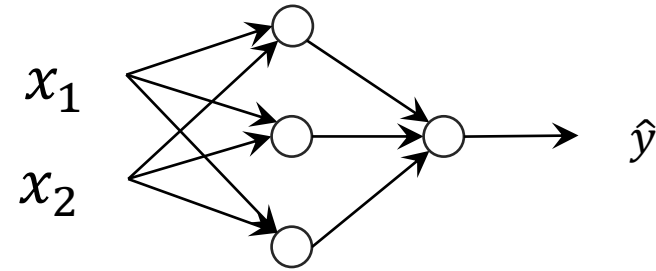
정 우 환 (whjung@hanyang.ac.kr)

Fall 2021

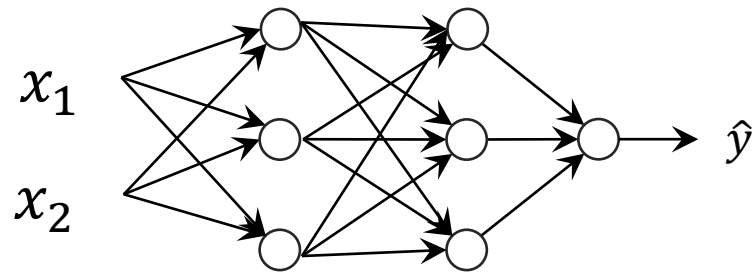
# What is a deep neural network?



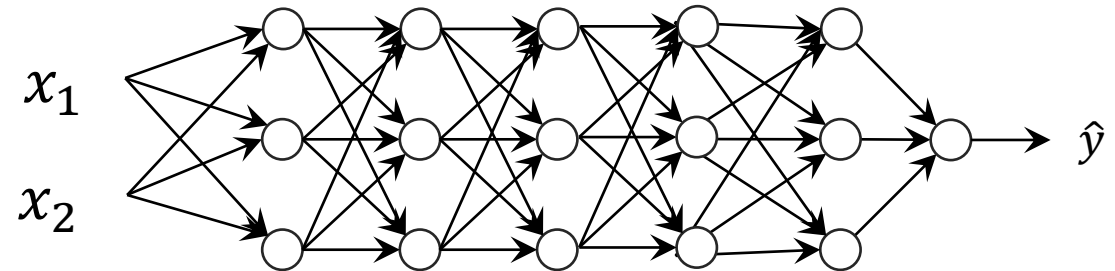
Logistic regression



1 hidden layer

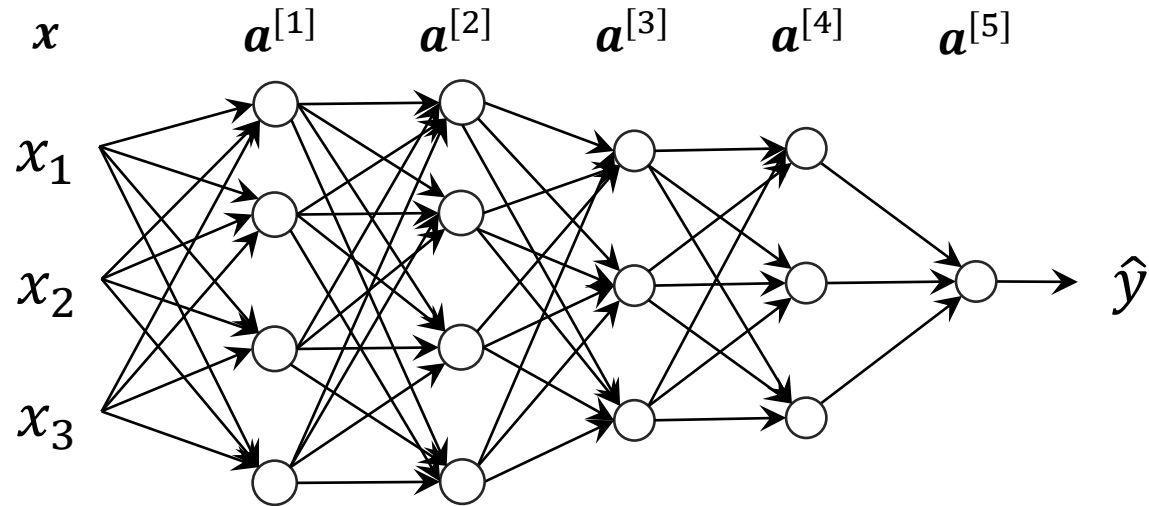


2 hidden layers



5 hidden layers

# Deep neural network



$f$ : activation function (e.g., ReLU)

$$\mathbf{a}^{[1]} = f(\mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]})$$

$$\mathbf{a}^{[i]} = f(\mathbf{W}^{[i]}\mathbf{a}^{[i-1]} + \mathbf{b}^{[i]})$$

$$\hat{y} = \mathbf{a}^{[5]} = \sigma(\mathbf{W}^{[5]}\mathbf{a}^{[4]} + \mathbf{b}^{[5]})$$

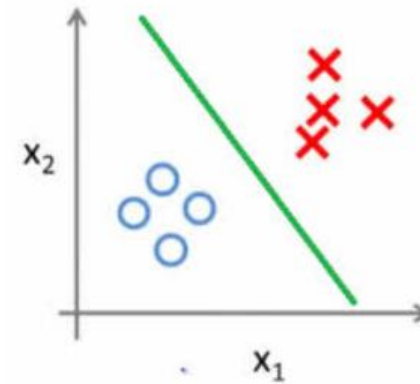
For  $i = 2, 3, 4$

# Output types

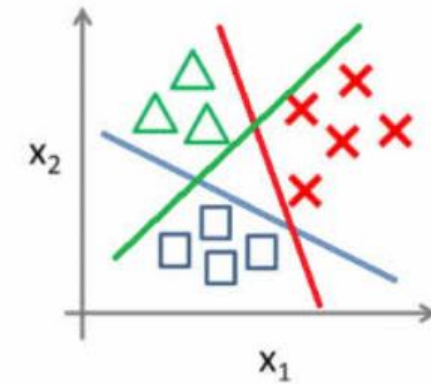
Output Type	Output Distribution	Output Layer	Cost Function
Binary	Bernoulli	Sigmoid	Binary cross-entropy
Discrete	Multinoulli	Softmax	Discrete cross-entropy
Continuous	Gaussian	Linear	Gaussian cross-entropy (MSE)
Continuous	Mixture of Gaussian	Mixture Density	Cross-entropy
Continuous	Arbitrary	See part III: GAN, VAE, FVBN	Various

# Output types

Binary classification:



Multi-class classification:



Amey band (2020)

Binary classification

(n-ary) Classification

Output classes

0/1

1/2/3/./n

Output activation  
function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{i'=1}^n e^{x_{i'}}}$$

Loss function

$$-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$-\sum_{i=1}^n \hat{y}_i \log y_i$$

Binary cross entropy

Cross entropy

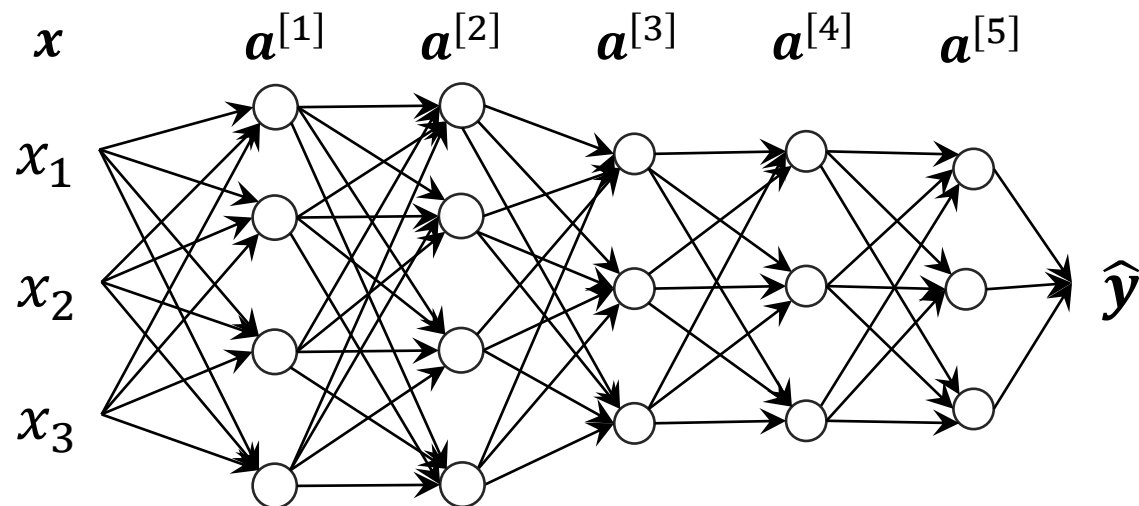


# Softmax outputs a categorical (multinoulli) distribution

$$\textit{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{i'=1}^n e^{x_{i'}}}$$

- $\textit{softmax}(\mathbf{x})_i \geq 0$
- $\sum_{i=1}^n \textit{softmax}(\mathbf{x})_i = 1$

# Deep neural network



$f$ : activation function (e.g., ReLU)

$$\mathbf{a}^{[1]} = f(\mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]})$$

$$\mathbf{a}^{[i]} = f(\mathbf{W}^{[i]}\mathbf{a}^{[i-1]} + \mathbf{b}^{[i]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[5]} = \text{softmax}(\mathbf{W}^{[5]}\mathbf{a}^{[4]} + \mathbf{b}^{[5]})$$

# Parameters vs Hyper parameters

- (Model) parameters
  - A model parameter is a configuration variable that is internal to the model and whose value can be estimated from data
  - Examples)  $\mathbf{W}^{[1]}, \mathbf{W}^{[2]}, \dots, \mathbf{b}^{[1]}, \mathbf{b}^{[2]}, \dots,$
- Hyper parameters
  - A model hyperparameter is a configuration that is external to the model and whose value cannot be estimated from data
  - Examples
    - Learning rate
    - Number of layers
    - Number of hidden units for each layer

Stochastic Gradient Descent  
Mini-Batch Gradient Descent

# Stochastic gradient descent

- A **stochastic** approximation of gradient descent optimization

Approximate  $\nabla J(\mathbf{w})$  by  $\nabla J_i(\mathbf{w}) = \nabla_{\mathbf{w}} L(y_i, \hat{y})$

## **Stochastic** gradient descent

1. Initialize parameters  $\mathbf{w}$
2. For each epoch:
3.     **Randomly shuffle** training examples
4.     For  $i = 1$  to  $n$ :
5.          $\mathbf{w} := \mathbf{w} - \eta \nabla J_i(\mathbf{w})$

Batch size: 1

Hard to parallelize

Fast update

Unstable

## Gradient descent

1. Initialize parameters  $\mathbf{w}$
2. For each epoch:
3.      $d\mathbf{w} = 0$
4.     For  $i = 1$  to  $n$ :
5.          $d\mathbf{w} := d\mathbf{w} + \nabla J_i(\mathbf{w})$
6.      $\nabla J(\mathbf{w}) = \frac{d\mathbf{w}}{n}$
7.      $\mathbf{w} := \mathbf{w} - \eta \nabla J(\mathbf{w})$

Batch size:  $n$

Easy to parallelize

Inefficient

Stable

Stochastic gradient descent

(Batchsize) = 1

Mini-batch gradient descent

$1 < (\text{Batchsize}) < n$

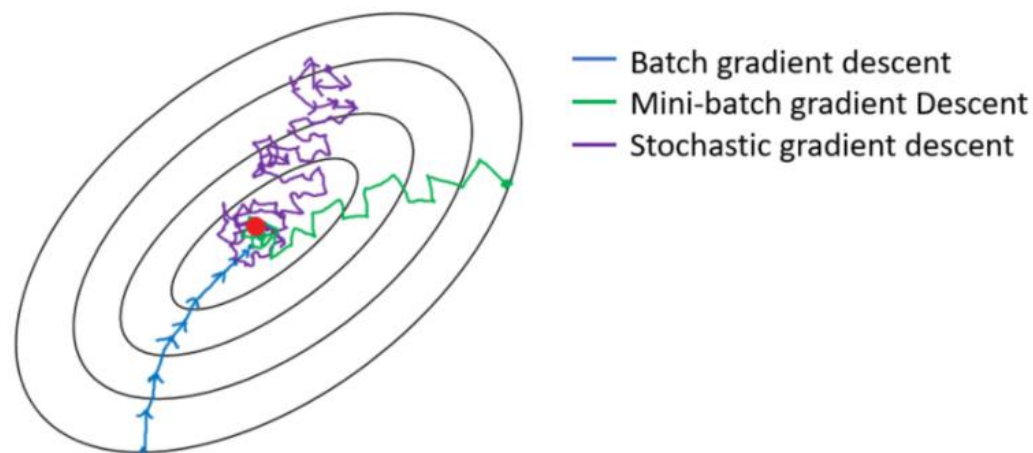
Gradient descent

(Batchsize) =  $n$

- Mini-batch gradient

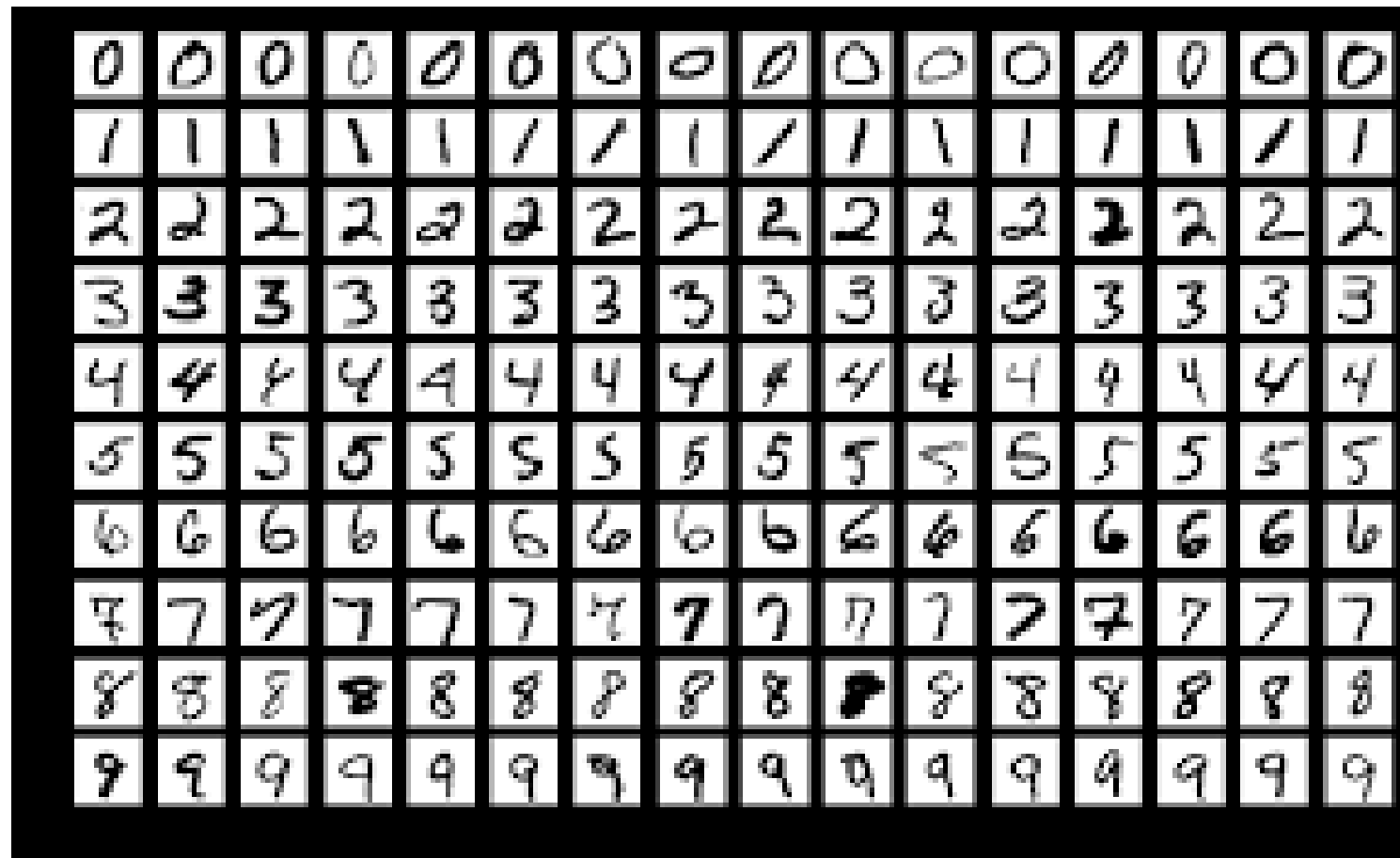
1. Initialize parameters  $\mathbf{w}$
2. For each epoch:
3.     Randomly shuffle training examples
4.     For  $b = 1$  to  $\frac{n}{\text{batchsize}}$ :
5.          $\mathbf{w} := \mathbf{w} - \eta \nabla J^b(\mathbf{w})$

$$J^b(\mathbf{w}) = \sum_{(x_i, y_i) \in \text{batch}} L(y_i, \hat{y})$$



# Implementing a DNN with PyTorch

# MNIST





# Data loading

- Import torchvision
- From torchvision import datasets

```
batch_size = 12

train_data = datasets.MNIST('D:\#datasets', train=True, download=True, transform=transforms.ToTensor())
test_data = datasets.MNIST('D:\#datasets', train=False, download=True, transform=transforms.ToTensor())

train_loader = torch.utils.data.DataLoader(train_data, batch_size = batch_size, shuffle=True)
test_loader = torch.utils.data.DataLoader(test_data, batch_size = batch_size)
```

# Model

```
class MLP(nn.Module):
    def __init__(self):
        super().__init__()

        self.in_dim = 28*28 # MNIST
        self.out_dim = 10

        self.fc1 = nn.Linear(self.in_dim, 512)
        self.fc2 = nn.Linear(512, 256)
        self.fc3 = nn.Linear(256, 128)
        self.fc4 = nn.Linear(128, 64)
        self.fc5 = nn.Linear(64, self.out_dim)

        self.relu = nn.ReLU()
        self.log_softmax = nn.LogSoftmax()

    def forward(self, x):
        a1 = self.relu(self.fc1(x.view(-1, self.in_dim)))
        a2 = self.relu(self.fc2(a1))
        a3 = self.relu(self.fc3(a2))
        a4 = self.relu(self.fc4(a3))
        logit = self.fc5(a4)
        return logit
```

# Train

```
model = MLP()
criterion = nn.CrossEntropyLoss()
optimizer = optim.SGD(model.parameters(), lr = 0.01)

for epoch in range(10): # loop over the dataset multiple times
    running_loss = 0.0
    for i, data in enumerate(train_loader, 0):
        # get the inputs; data is a list of [inputs, labels]
        inputs, labels = data

        # zero the parameter gradients
        optimizer.zero_grad()

        # forward + backward + optimize
        outputs = model(inputs)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()

        # print statistics
        running_loss += loss.item()
        if (i+1) % 2000 == 0: # print every 2000 mini-batches
            print('[%d, %5d] loss: %.3f' %
                  (epoch + 1, i + 1, running_loss / 2000))
            running_loss = 0.0

print('Finished Training')
```

---

[1,	2000]	loss: 2.209
[1,	4000]	loss: 0.739
[2,	2000]	loss: 0.316
[2,	4000]	loss: 0.230
[3,	2000]	loss: 0.154
[3,	4000]	loss: 0.144
[4,	2000]	loss: 0.112
[4,	4000]	loss: 0.101
[5,	2000]	loss: 0.075
[5,	4000]	loss: 0.082
[6,	2000]	loss: 0.061
[6,	4000]	loss: 0.063
[7,	2000]	loss: 0.046
[7,	4000]	loss: 0.054
[8,	2000]	loss: 0.033
[8,	4000]	loss: 0.041
[9,	2000]	loss: 0.029
[9,	4000]	loss: 0.033
[10,	2000]	loss: 0.021
[10,	4000]	loss: 0.025

Finished Training

# Test

```
import matplotlib.pyplot as plt
import numpy as np
```

```
def imshow(img):
    npimg = img.numpy()
    plt.imshow(np.transpose(npimg, (1, 2, 0)))
    plt.show()
```

```
dataiter = iter(test_loader)
images, labels = dataiter.next()

imshow(torchvision.utils.make_grid(images, nrow = batch_size))
print('GroundTruth')
print("    '+' '.join('%3s' % label.item() for label in labels))

outputs = model(images)
_, predicted = torch.max(outputs, 1)
print("Prediction")
print("    '+' '.join('%3s' % label.item() for label in predicted))
```



GroundTruth

7 2 1 0 4 1 4 9 5 9 0 6

Prediction

7 2 1 0 4 1 4 9 6 9 0 6

# Test

```
n_predict = 0
n_correct = 0

for data in test_loader:
    inputs, labels = data
    outputs = model(inputs)
    _, predicted = torch.max(outputs, 1)

    n_predict += len(predicted)
    n_correct += (labels == predicted).sum()

print(f"{n_correct}/{n_predict}")
print(f"Accuracy: {n_correct/n_predict:.3f}")
```

9761/10000  
Accuracy: 0.976

# A DNN with hyper parameters

## MLP with hyper parameters

```
class MLP_h(nn.Module):
    def __init__(self, hidden_units = [512, 256, 128]):
        super().__init__()

        self.in_dim = 28*28 # MNIST
        self.out_dim = 10

        self.l_layers = []
        self.l_layers.append(nn.Linear(self.in_dim, hidden_units[0]))
        for i in range(len(hidden_units)-1):
            self.l_layers.append(nn.Linear(hidden_units[i], hidden_units[i+1]))
        self.l_layers.append(nn.Linear(hidden_units[-1], self.out_dim))

        self.relu = nn.ReLU()
        self.log_softmax = nn.LogSoftmax()

    def forward(self, x):
        a = x.view(-1, self.in_dim)
        for l in range(len(self.l_layers)):
            z = self.l_layers[l](a)
            if l == len(self.l_layers) - 1:
                logit = z
            else:
                a = self.relu(z)

        return logit
```

```
model = MLP_h([2, 3])
criterion = nn.CrossEntropyLoss()
optimizer = optim.SGD(model.parameters(), lr = 0.01)
```

```
-----
ValueError                                Traceback (most recent call last)
<ipython-input-21-e9ee64f89c5c> in <module>
      1 model = MLP_h([2, 3])
      2 criterion = nn.CrossEntropyLoss()
----> 3 optimizer = optim.SGD(model.parameters(), lr = 0.01)

~\anaconda3\lib\site-packages\torch\optim\sgd.py in __init__(self, params, defaults)
     67         if nesterov and (momentum <= 0 or dampening > 0):
     68             raise ValueError("Nesterov momentum requires momentum > 0 and dampening < 0")
----> 69         super(SGD, self).__init__(params, defaults)
     70
     71     def __setstate__(self, state):

~\anaconda3\lib\site-packages\torch\optim\optimizer.py in __init__(self, params, defaults)
     48         param_groups = list(params)
     49         if len(param_groups) == 0:
----> 50             raise ValueError("optimizer got an empty parameter list")
     51         if not isinstance(param_groups[0], dict):
     52             param_groups = [{'params': param_groups}]
```

**ValueError:** optimizer got an empty parameter list

# Sequential

- Modules will be added to it in the order they are passed in the constructor

```
# Example of using Sequential  
model = nn.Sequential(  
    nn.Conv2d(1, 20, 5),  
    nn.ReLU(),  
    nn.Conv2d(20, 64, 5),  
    nn.ReLU()  
)
```

```
# Example of using Sequential with OrderedDict  
model = nn.Sequential(OrderedDict([  
    ('conv1', nn.Conv2d(1, 20, 5)),  
    ('relu1', nn.ReLU()),  
    ('conv2', nn.Conv2d(20, 64, 5)),  
    ('relu2', nn.ReLU())  
]))
```

# ModuleList

- Holds submodules in a list
- Can be indexed like a regular Python list
- Modules it contains are properly registered, and will be visible by all Module methods

```
class MyModule(nn.Module):
    def __init__(self):
        super(MyModule, self).__init__()
        self.linears = nn.ModuleList([nn.Linear(10, 10) for i in range(10)])

    def forward(self, x):
        # ModuleList can act as an iterable, or be indexed using ints
        for i, l in enumerate(self.linears):
            x = self.linears[i // 2](x) + l(x)
        return x
```



# Programming Assignment 3: DNN

- Dataset: MNIST
- Requirement
  - Plot accuracy varying the number of layers (2,3,4,5 layers)
  - Nothing more, but, using ModuleList may save your time
- Hint:
- Due:
- Submission
  - Report (pdf or docx)
    - A figure: accuracy vs # layers
    - Source code
  - Source file