

Part I Q1

Since each layer is linear, if the activation function is linear, then the result is just putting several linear together

Need to show that the result of linear functions is still a linear function.

Let n be the # of linear functions
prove by induction

Base case: $n=1$

$$f = wx + b$$

so f is linear

Inductive step: Assume $n-1$ linear functions can be reduced to 1 single linear function (I.H)

$$f = f_n(f_{n-1}(f_{n-2} \dots (f(x)) \dots))$$

where f_i s are linear functions

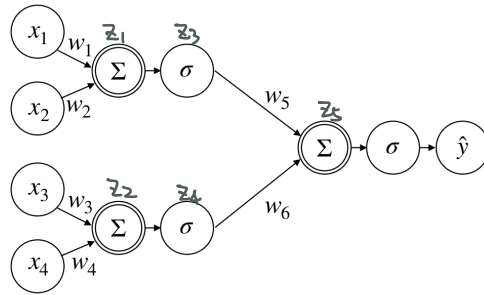
By I.H

$$\begin{aligned} f &= f_n(wx + b) \text{ for some } w \text{ and } b \\ &= w_n(wx + b) + b_n \\ &= w_nwx + (w_nb + b_n) \\ &= mx + a \text{ where } m = w_nw, a = w_nb + b_n \end{aligned}$$

so f is linear.

Thus, no matter how many layers we have, we can always reduce it to 1 single layer. So the number of layers has effectively no impact on the network.

Part 1 Q2:



$$z_1 = 1.368, \quad z_2 = -1.432, \quad z_3 = 0.7971, \quad z_4 = 0.1928, \quad z_5 = 0.5991, \quad \hat{y} = 0.6455$$

$$\frac{\partial L}{\partial \hat{y}} = 2 \| y - \hat{y} \| = 2 (0.6455 - 0.5) = 0.291$$

$$\hat{y} = \sigma(z_5)$$

$$\frac{\partial \hat{y}}{\partial z_5} = \hat{y} (1 - \hat{y}) = 0.6455 (1 - 0.6455) = 0.2288$$

$$\frac{\partial L}{\partial z_5} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_5} = 2 \| y - \hat{y} \| \hat{y} (1 - \hat{y}) = 0.2288 \cdot 0.291 = 0.0666$$

$$z_5 = w_5 z_3 + w_6 z_4$$

$$\frac{\partial z_5}{\partial z_4} = w_6$$

$$\frac{\partial L}{\partial z_4} = \frac{\partial L}{\partial z_5} \cdot \frac{\partial z_5}{\partial z_4} = 2 w_6 \| y - \hat{y} \| \hat{y} (1 - \hat{y}) = 0.0666 \cdot (-0.2) = -0.0133$$

$$z_4 = \sigma(z_2)$$

$$\frac{\partial z_4}{\partial z_2} = z_4 (1 - z_4) = 0.1556$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial z_4} \cdot \frac{\partial z_4}{\partial z_2} = z_4 (1 - z_4) 2 w_6 \| y - \hat{y} \| \hat{y} (1 - \hat{y}) = 0.1556 \cdot (-0.0133) = -0.0021$$

$$z_2 = w_3 x_3 + w_4 x_4$$

$$\frac{\partial z_2}{\partial w_3} = x_3$$

$$\begin{aligned} \frac{\partial L}{\partial w_3} &= \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_3} = x_3 z_4 (1 - z_4) 2 w_6 \| y - \hat{y} \| \hat{y} (1 - \hat{y}) \\ &= -0.3 \cdot (-0.0021) \\ &= 0.0006 \end{aligned}$$

Part 1 B3

Each filter size is $4 \times 4 \times 50$

So apply filter once required $4 \times 4 \times 50 = 800$ multiplications

and $800 - 1 = 799$ addition. In total 1599 FLOP

The output size is $(12 + 2 \cdot 1 - 4) / 2 + 1 = 6$

That means we need to apply each filter $6 \times 6 = 36$ times.

There are 20 filters $\Rightarrow 20 \times 36 \times 1599$ FLOP

The output size is $6 \times 6 \times 20$

The output size of max pooling is $(6 - 3) / 1 + 1 = 4$

So # of time applying max pooling is $4 \times 4 \times 20 = 320$

↑
20 output map

Each max pooling takes $3 \times 3 - 1 = 8$ FLOP

Total # of FLOP = $20 \times 36 \times 1599 + 320 \times 8 = 1153840$ FLOP without bias

Since the output dimension is $6 \times 6 \times 20$ and each pixel can have

a bias value

So # of bias = $6 \times 6 \times 20 = 720$

Total # of FLOP = $1153840 + 720 = 1154560$ FLOP with bias

Part 1 Q4

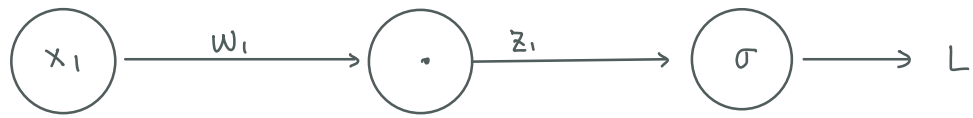
Convolution : padding=0 , stride=1 , filter size = 5×5

subsampling : padding=0 , stride=2 , filter size = 2×2

Size / operation	Filter	Depth	stride	padding	# of parameter
$32 \times 32 \times 1$					
convolution	5×5	6	1	0	$(5 \times 5 \times 1 + 1) \cdot 6 = 156$
$28 \times 28 \times 6$					
subsampling	2×2	6	2	0	
$14 \times 14 \times 6$					
convolution	5×5	16	1	0	$(5 \times 5 \times 6 + 1) \times 16 = 2416$
$10 \times 10 \times 16$					
Subsampling	2×2	16	2	0	
$5 \times 5 \times 16$					
FC					$5 \times 5 \times 16 \times 120 = 48000$
120					
FC					$120 \times 84 = 10080$
84					
GC					84×10

Total # of parameter = 61492

Part 1 Q5:



I consider only 1D case here because for each neuron
1 input variable will only associated with 1 weight parameter

Consider a input x_1 , multiplied by some weight w_1

and going through logistic function

$$\frac{\partial L}{\partial L} = 1$$

$$L = \sigma(z_1)$$

$$\frac{\partial L}{\partial z_1} = L(1-L)$$

$$z_1 = x_1 w_1$$

$$\frac{\partial z_1}{\partial x_1} = w_1$$

$$\frac{\partial L}{\partial x_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial x_1}$$

$$= w_1 L(1-L)$$

Notice that the equation does not need the value of x_1

Part 1 Q6

(a) The output range is $(-1, 1)$

where the output range for logistic function is $(0, 1)$

$$\begin{aligned} \text{(b)} \quad \frac{d}{dx} \tanh(x) &= \frac{(1+e^{-2x})(2e^{-2x}) - (1-e^{-2x})(-2e^{-2x})}{1 + 2 \cdot e^{-2x} + e^{-4x}} \\ &= \frac{2e^{-2x} + 2 \cdot e^{-4x} - 2e^{-4x} + 2e^{-2x}}{1 + 2 \cdot e^{-2x} + e^{-4x}} \\ &= \frac{4e^{-2x}}{1 + 2e^{-2x} + e^{-4x}} \\ &= 1 - \tanh^2(x) \end{aligned}$$

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

$$\sigma(2x) = \frac{1}{1 + e^{-2x}}$$

$$2\sigma(2x) = \frac{2}{1 + e^{-2x}}$$

$$2\sigma(2x) - 1 = \frac{1 - e^{-2x}}{1 + e^{-2x}} = \tanh(x)$$

$$\begin{aligned} \frac{d}{dx} \tanh(x) &= 1 - \tanh^2(x) = 1 - [2\sigma(2x) - 1]^2 \\ &= 4\sigma(2x) - 4\sigma^2(2x) \end{aligned}$$

(c) if you want to map the result to $(0, 1)$ where

each output representing the probability, then sigmoid

is a good choice

If the sign is important, we can use $\tanh(x)$

since it maps negative value to negative value and positive value to positive value.

Also if you want the training time to be shorter, $\tanh(x)$ can be used since it has higher derivatives.

Assignment 2

Task 1:

```
15 def load_dataset():
16     # Define transform
17     transform = transforms.Compose([transforms.ToTensor(),
18                                     transforms.Normalize((0.5,), (0.5,)),
19                                     ])
20     transform_train = transforms.Compose([transforms.RandomHorizontalFlip(),
21                                           transforms.ColorJitter(),
22                                           transforms.ToTensor(),
23                                           transforms.Normalize((0.5,), (0.5,)),
24                                           ])
25     dataset = datasets.ImageFolder(root='notMNIST_small')
26     total_size = len(dataset.targets)
27     train_indices, val_indices = train_test_split(list(
28         range(total_size)), test_size=total_size - TRAINING_SIZE, stratify=dataset.targets)
29
30     # Use train_test_split twice to split dataset into train, validation and test sets.
31     train_dataset = torch.utils.data.Subset(dataset, train_indices)
32     val_dataset = torch.utils.data.Subset(dataset, val_indices)
33
34     val_indices, test_indices = train_test_split(val_indices, test_size=total_size - TRAINING_SIZE - VALIDATION_SIZE, stratify=[
35         dataset.targets[i] for i in range(len(dataset.targets)) if i in val_indices])
36
37     train_dataset = torch.utils.data.Subset(dataset, train_indices)
38     val_dataset = torch.utils.data.Subset(dataset, val_indices)
39     test_dataset = torch.utils.data.Subset(dataset, test_indices)
40
41     train_dataset.dataset = copy(train_dataset.dataset)
42     val_dataset.dataset = copy(val_dataset.dataset)
43     test_dataset.dataset = copy(test_dataset.dataset)
44
45     train_dataset.dataset.transform = transform_train
46     val_dataset.dataset.transform = transform
47     test_dataset.dataset.transform = transform
48
49     print(train_dataset.dataset.transform,
50           val_dataset.dataset.transform, test_dataset.dataset.transform)
51     # Load dataset with batch size = 32
52     trainloader = torch.utils.data.DataLoader(
53         train_dataset, batch_size=32, shuffle=True)
54     validationloader = torch.utils.data.DataLoader(
55         val_dataset, batch_size=32, shuffle=True)
56     testloader = torch.utils.data.DataLoader(
57         test_dataset, batch_size=32, shuffle=True)
58     return trainloader, validationloader, testloader
59
```

The transformation that I choose for validation set and test set are ToTensor() and Normailize(). ToTensor() is required and Normailize() generally speeds up the learning process and leads to faster convergence. For training set, in addition to ToTensor() and Normalize(), I also add RandomHorizontalFlip and ColorJitter because I want neural network to be able recognize letter even it is horizontally flipped or the brightness is different. To split data set into training, validation, test

sets, `train_test_split` is the helper function I use. After removing the corrupted images, the training set will have 15000 images, validation set has 1000 images and test set has 2724 images. A small batch size will result in a longer training time, so I choose 32 in my case which suits best on my pc.

Task 2:

```
45 def define_model(hidden_sizes: int):
46     # define the neural network model
47     # Hyperparameters for our network
48     input_size = 3 * 28 * 28
49     output_size = 10
50     # Build a feed-forward network
51     model = nn.Sequential(nn.Linear(input_size, hidden_sizes),
52                           nn.ReLU(),
53                           nn.Linear(hidden_sizes, output_size),
54                           nn.Softmax(dim=1))
55     print(model)
56     return model
57
```

The model has only one hidden layer with Relu as activation function and has a softmax function after the output layer.

```
58
59 def train_model(model, trainloader, validationloader, save_file_name, lr, n_epochs=20, max_epochs_stop=2):
60     # create model and define the loss
61
62     criterion = nn.CrossEntropyLoss()
63     # Optimizers require the parameters to optimize and a learning rate
64     optimizer = optim.SGD(model.parameters(), lr=lr)
65     overall_start = timer()
66     history = []
67     valid_loss_min = np.Inf
68
69     for epoch in range(n_epochs):
70         train_loss = 0
71         validation_loss = 0
72         train_acc = 0
```

The loss function is defined as `CrossEntropyLoss`. The optimizer is `SGD`.


```

78
79     for ii, (images, labels) in enumerate(trainloader):
80         # Flatten MNIST images into a 2352 long vector
81         images = images.view(images.shape[0], -1)
82
83         # Training pass
84         optimizer.zero_grad()
85
86         output = model(images)
87         loss = criterion(output, labels)
88         loss.backward()
89         optimizer.step()
90
91         train_loss += loss.item() * images.size(0)
92
93         # Calculate accuracy by finding max log probability
94         _, pred = torch.max(output, dim=1)
95         correct_tensor = pred.eq(labels.data.view_as(pred))
96         # Need to convert correct tensor from int to float to average
97         accuracy = torch.mean(correct_tensor.type(torch.FloatTensor))
98         # Multiply average accuracy times the number of examples in batch
99         train_acc += accuracy.item() * images.size(0)
100        print(
101            f'Epoch: {epoch}\t{100 * (ii + 1) / len(trainloader):.2f}% complete. {timer()}
102            end='\r')

```

For each epoch, training loss and training accuracy has been calculated.

```

104        with torch.no_grad():
105            model.eval()
106
107            for images, labels in validationloader:
108                images = images.view(images.shape[0], -1)
109                # forward pass
110                output = model(images)
111
112                # validation loss
113                loss = criterion(output, labels)
114                validation_loss += loss.item() * images.size(0)
115
116                # Calculate validation accuracy
117                _, pred = torch.max(output, dim=1)
118                correct_tensor = pred.eq(labels.data.view_as(pred))
119                accuracy = torch.mean(
120                    correct_tensor.type(torch.FloatTensor))
121                # Multiply average accuracy times the number of examples
122                validation_acc += accuracy.item() * images.size(0)
123
124            # calculate average loss and average accuracy for training and validation
125            train_loss = train_loss / len(trainloader.dataset)
126            validation_loss = validation_loss / \
127                len(validationloader.dataset)
128
129            train_acc = train_acc / len(trainloader.dataset)
130            validation_acc = validation_acc / len(validationloader.dataset)

```

The validation accuracy and validation loss are calculated as well.

```

140         )
141         if validation_loss < valid_loss_min:
142             # Save model
143             torch.save(model.state_dict(), save_file_name)
144             # Track improvement
145             epochs_no_improve = 0
146             valid_loss_min = validation_loss
147             best_epoch = epoch
148
149         else:
150             epochs_no_improve += 1
151             # Trigger early stopping
152             if epochs_no_improve >= max_epochs_stop:
153                 print(
154                     f'\nEarly Stopping! Total epochs: {epoch}. Best epoch: {best_epoch} with loss: {valid_loss_min}'
155                 )
156                 total_time = timer() - overall_start
157                 print(
158                     f'{total_time:.2f} total seconds elapsed. {total_time / (epoch+1):.2f} seconds per epoch.'
159                 )
160
161                 # Load the best state dict
162                 model.load_state_dict(torch.load(save_file_name))
163                 # Attach the optimizer
164                 model.optimizer = optimizer
165
166                 # Format history

```

The early stop is implemented as well. By default, 2 consecutive increasing in validation loss will stop the training process.

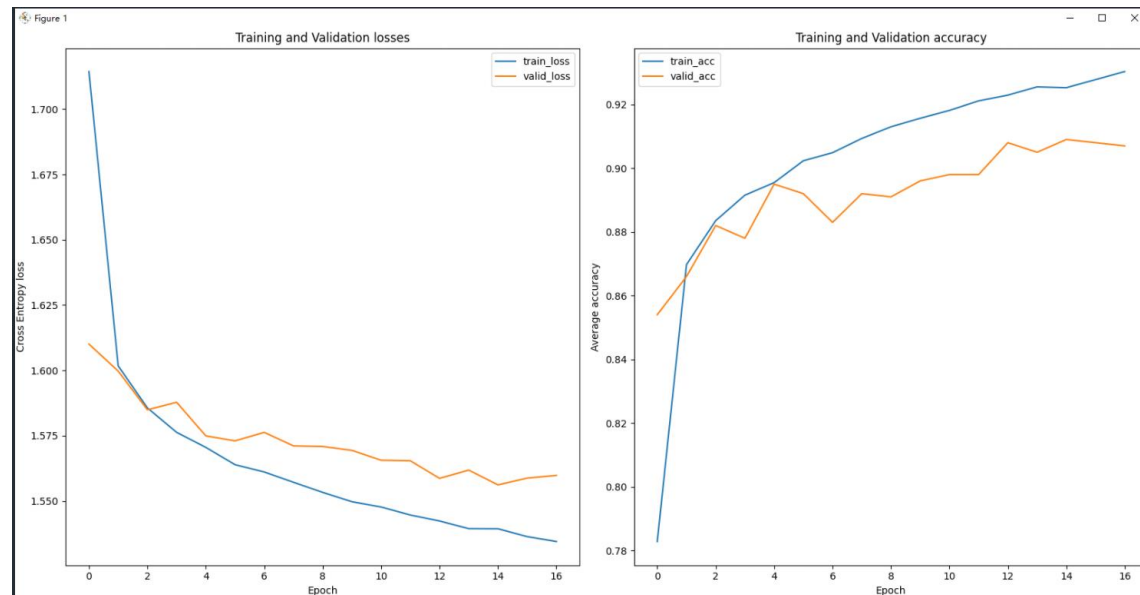
```

260     # ===== Task 2 =====
261     LR_RATE = [0.1, 0.05, 0.03, 0.01, 0.005]
262     valid_loss_min = np.Inf
263     for lr in LR_RATE:
264         model = define_model(HIDDEN_SIZE)
265         print(f'Training on learning rate: {lr}')
266         model, history, valid_loss = train_model(
267             model, trainloader, validationloader, FILE_NAME, lr, n_epochs=20)
268         if valid_loss < valid_loss_min:
269             valid_loss_min = valid_loss
270             best_model = model
271             best_history = history
272             best_lr = lr
273     print('The best learning rate is ', best_lr)
274     fig = plt.figure(figsize=(16, 8), constrained_layout=True)
275     ax1 = fig.add_subplot(121)
276     ax2 = fig.add_subplot(122)
277     for c in ['train_loss', 'valid_loss']:
278         ax1.plot(best_history[c], label=c)
279     ax1.legend()
280     ax1.set_xlabel('Epoch')
281     ax1.set_ylabel('Cross Entropy loss')
282     ax1.set_title('Training and Validation losses')
283
284     for c in ['train_acc', 'valid_acc']:
285         ax2.plot(best_history[c], label=c)
286     ax2.legend()
287     ax2.set_xlabel('Epoch')
288     ax2.set_ylabel('Average accuracy')
289     ax2.set_title('Training and Validation accuracy')
290
291     test_loss, test_acc = test_error_acc(best_model, testloader)

```

The five learning rates I choose are 0.1, 0.05, 0.03, 0.01, 0.005. The training process takes a bit longer time, so I've saved the print messages in Task 2 print result.txt file. It turns out that 0.1 works the best. The plot is shown below.

The plot has been saved as plot_result_task2.png. And for this model, the test loss is 1.5706202945345125 and the test accuracy is 89.46%. (Result can be different if training multiple times)



The best training loss and validation loss are 1.5394 and 1.5562

Task 3:

```

296 # ===== task 3 =====
297 # Based on task 2, the best learning rate is 0.1
298 try:
299     lr = best_lr
300 except:
301     lr = 0.1
302 valid_loss_min = np.Inf
303 layer_size = [100, 500, 1000]
304 valid_loss_list = []
305 best_size = 0
306 for size in layer_size:
307     model = define_model(size)
308     print(f"Training for size: {size}")
309     model, history, valid_loss = train_model(
310         model, trainloader, validationloader, FILE_NAME, lr, n_epochs=20)
311     valid_loss_list.append(valid_loss)
312     if valid_loss < valid_loss_min:
313         best_size = size
314         valid_loss_min = valid_loss
315         best_model = model
316 print(
317     f'The validation loss are {valid_loss_list[0]:.4f}, {valid_loss_list[1]:.4f}, {valid_loss_list[2]:.4f}')
318 test_loss, test_acc = test_error_acc(best_model, testloader)
319 print(
320     f'The best model has {best_size} units in hidden layer')
321 print(
322     f'The test loss is {test_loss}. The test accuracy is {100 * test_acc:.2f}%')
323

```

Training model with layer size = 100, 500, and 1000 separately. I use 0.1 as learning rate. The validation losses are 1.5982, 1.5863, 1.5828

The best model has 1000 units in hidden layer

The test loss is 1.5748151756068158. The test accuracy is 89.24%

The validation losses are extremely close to each other with hidden unit = 1000 works slightly better than 500 which works slightly better than 100. Thus, the more the number of units are, the lower the loss is. The printed messages are stored in Task 3 print result.txt file.

Task 4:

```
76 def define_model_two_layers(hidden_sizes: list):
77     # define the neural network model
78     # Hyperparameters for our network
79     input_size = 3 * 28 * 28
80     output_size = 10
81     # Build a feed-forward network
82     model = nn.Sequential(nn.Linear(input_size, hidden_sizes[0]),
83                           nn.ReLU(),
84                           nn.Linear(hidden_sizes[0], hidden_sizes[1]),
85                           nn.ReLU(),
86                           nn.Linear(hidden_sizes[1], output_size),
87                           nn.Softmax(dim=1))
88     print(model)
89     return model
90
```

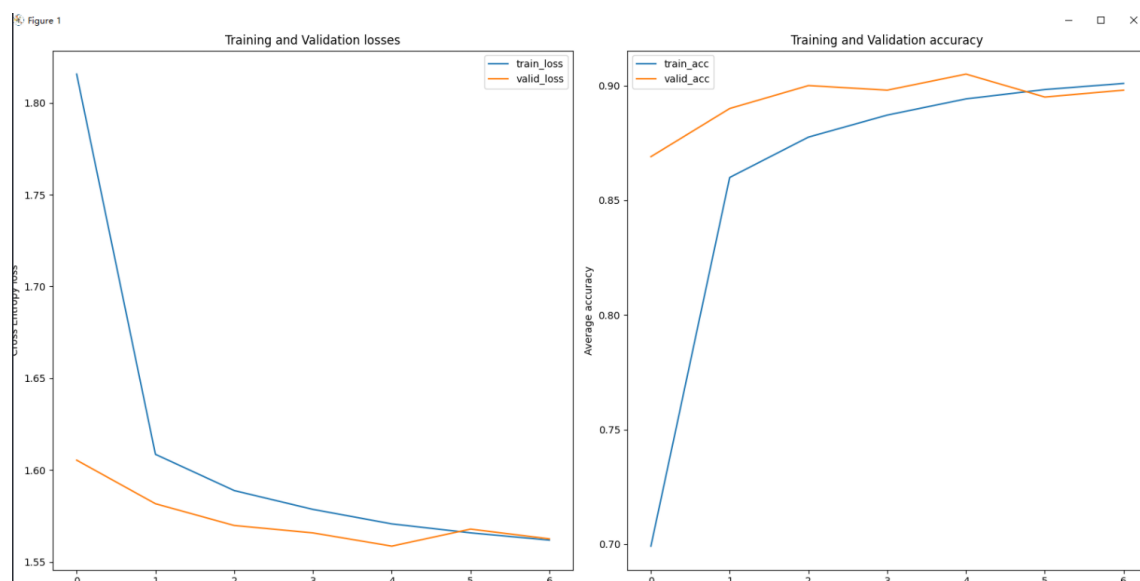
The model is defined above

```

324 # ===== task 4 =====
325 try:
326     lr = best_lr
327 except:
328     lr = 0.1
329 model = define_model_two_layers([500, 500])
330 model, history, valid_loss = train_model(
331     model, trainloader, validationloader, FILE_NAME, lr, n_epochs=20)
332 try:
333     fig.clear()
334 except:
335     fig = plt.figure(figsize=(16, 8), constrained_layout=True)
336 ax3 = fig.add_subplot(121)
337 ax4 = fig.add_subplot(122)
338 for c in ['train_loss', 'valid_loss']:
339     ax3.plot(history[c], label=c)
340 ax3.legend()
341 ax3.set_xlabel('Epoch')
342 ax3.set_ylabel('Cross Entropy loss')
343 ax3.set_title('Training and Validation losses')
344
345 for c in ['train_acc', 'valid_acc']:
346     ax4.plot(history[c], label=c)
347 ax4.legend()
348 ax4.set_xlabel('Epoch')
349 ax4.set_ylabel('Average accuracy')
350 ax4.set_title('Training and Validation accuracy')
351
352 test_loss, test_acc = test_error_acc(model, testloader)
353 print(
354     f'The test loss is {test_loss}. The test accuracy is {100 * test_acc:.2f}%')
355 plt.show()
356

```

The learning rate I use is 0.1 since it works the best in task 2. The validation loss and training loss are shown below. The plot has been saved as [plot_result_task4.png](#).



It early stops at epoch 6.

The best training Loss is 1.5708 and the best validation Loss: 1.5587 which happens at epoch4.

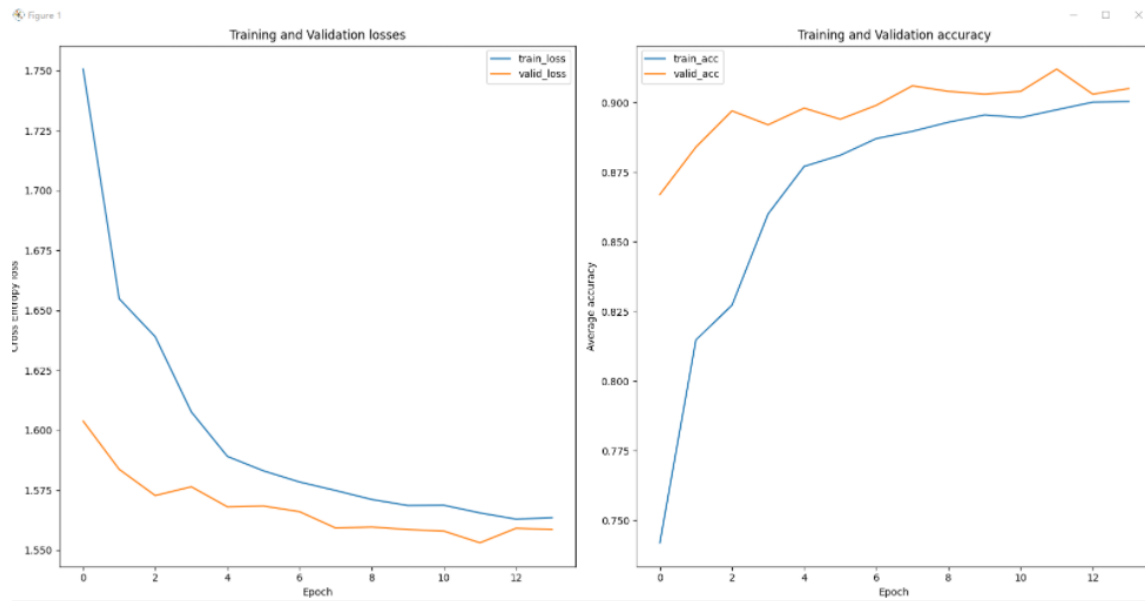
The test loss is 1.5820134526307363. The test accuracy is 88.40%

Compare with model that has only 1 layer with 1000 hidden units, the validation loss of two layers neural network is slightly lower. However, the test results are quite closed. The print messages are stored in Task 4 print result.txt file.

Task 5:

```
91
92 def define_model_with_dropout(hidden_sizes: int):
93     # define the neural network model
94     # Hyperparameters for our network
95     input_size = 3 * 28 * 28
96     output_size = 10
97     # Build a feed-forward network
98     model = nn.Sequential(nn.Linear(input_size, hidden_sizes),
99                           nn.ReLU(),
100                          nn.Dropout(0.5),
101                          nn.Linear(hidden_sizes, output_size),
102                          nn.Softmax(dim=1))
103     print(model)
104     return model
105
```

The model is defined above. The final training loss and validation loss for this model is 1.5654 and 1.5530. The training loss and validation loss in task 2 are 1.5394 and 1.5562. As we can see the training loss is higher when applying dropout. That's because the dropout prevent model from overfitting, so the training loss is expected to be higher. The print messages are stored in Task 5 print result.txt file.



The test loss is 1.581805121355995. The test accuracy is 87.89%

The plot has been saved as plot_result_task5.png.