

# CS273A Homework 4

Due: Monday November 4th, 2024 (11:59pm)

---

## Instructions

This homework (and subsequent ones) will involve data analysis and reporting on methods and results using Python code. You will submit a **single PDF file** that contains everything to Gradescope. This includes any text you wish to include to describe your results, the complete code snippets of how you attempted each problem, any figures that were generated, and scans of any work on paper that you wish to include. It is important that you include enough detail that we know how you solved the problem, since otherwise we will be unable to grade it.

Your homeworks will be given to you as Jupyter notebooks containing the problem descriptions and some template code that will help you get started. You are encouraged to use these starter Jupyter notebooks to complete your assignment and to write your report. This will help you not only ensure that all of the code for the solutions is included, but also will provide an easy way to export your results to a PDF file (for example, doing *print preview* and *printing to pdf*). I recommend liberal use of Markdown cells to create headers for each problem and sub-problem, explaining your implementation/answers, and including any mathematical equations. For parts of the homework you do on paper, scan it in such that it is legible (there are a number of free Android/iOS scanning apps, if you do not have access to a scanner), and include it as an image in the Jupyter notebook.

**Double check that all of your answers are legible on Gradescope, e.g. make sure any text you have written does not get cut off.**

If you have any questions/concerns about using Jupyter notebooks, ask us on EdD. If you decide not to use Jupyter notebooks, but go with Microsoft Word or LaTeX to create your PDF file, make sure that all of the answers can be generated from the code snippets included in the document.

## Summary of Assignment: 100 total points

- Problem 1: A Small Neural Network (30 points)
  - Problem 1.1: Forward Pass (10 points)
  - Problem 1.2: Evaluate Loss (10 points)
  - Problem 1.3: Network Size (10 points)
- Problem 2: Neural Networks on MNIST (35 points)
  - Problem 2.1: Varying the Amount of Training Data (15 points)
  - Problem 2.3: Optimization Curves (10 points)
  - Problem 2.3: Tuning your Neural Network (10 points)
- Problem 3: Convolutional Networks (30 points)
  - Problem 3.1: Model structure (10 points)
  - Problem 3.2: Training (10 points)
  - Problem 3.3: Evaluation (5 points)
  - Problem 3.4: Comparing predictions (5 points)
- Statement of Collaboration (5 points)

Before we get started, let's import some libraries that you will make use of in this assignment. Make sure that you run the code cell below in order to import these libraries.

**Important: In the code block below, we set `seed=1234`. This is to ensure your code has reproducible results and is important for grading. Do not change this. If you are not using the provided Jupyter notebook, make sure to also set the random seed as below.**

**Important: Do not change any codes we give you below, except for those waiting for you to complete. This is to ensure your code has reproducible results and is important for grading.**

```
In [1]: import numpy as np
import matplotlib.pyplot as plt

import torch

from IPython import display

from sklearn.datasets import fetch_openml          # common data set access
from sklearn.preprocessing import StandardScaler   # scaling transform
from sklearn.model_selection import train_test_split # validation tools
from sklearn.metrics import accuracy_score
```

```

from sklearn.neural_network import MLPClassifier # scikit's MLP

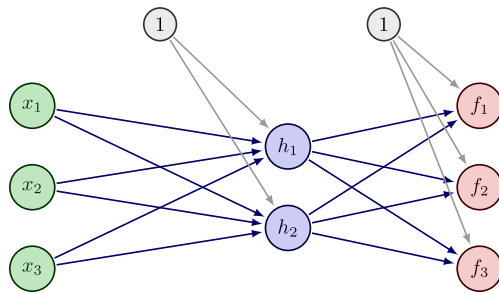
import warnings
warnings.filterwarnings('ignore')

# Fix the random seed for reproducibility
# !! Important !! : do not change this
seed = 1234
np.random.seed(seed)
torch.manual_seed(seed);

```

## Problem 1: A Small Neural Network

Consider the small neural network given in the image below, which will classify a 3-dimensional feature vector  $\mathbf{x}$  into one of three classes ( $y = 0, 1, 2$ ):



You are given an input to this network  $\mathbf{x}$ ,

$$\mathbf{x} = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} = \begin{bmatrix} 1 & 3 & -2 \end{bmatrix}$$

as well as weights  $W$  for the hidden layer and weights  $B$  for the output layer.

$$W = \begin{bmatrix} w_{01} & w_{11} & w_{21} & w_{31} \\ w_{02} & w_{12} & w_{22} & w_{32} \end{bmatrix} = \begin{bmatrix} 1 & -1 & 0 & 5 \\ 2 & 1 & 1 & 2 \end{bmatrix}$$

$$B = \begin{bmatrix} \beta_{01} & \beta_{11} & \beta_{21} \\ \beta_{02} & \beta_{12} & \beta_{22} \\ \beta_{03} & \beta_{13} & \beta_{23} \end{bmatrix} = \begin{bmatrix} 4 & -1 & 0 \\ 3 & 0 & 2 \\ 2 & 1 & 1 \end{bmatrix}$$

For example,  $w_{12}$  is the weight connecting input  $x_1$  to hidden node  $h_2$ ;  $w_{01}$  is the constant (bias) term for  $h_1$ , etc.

This network uses the ReLU activation function for the hidden layer, and uses the softmax activation function for the output layer.

Answer the following questions about this network.

### Problem 1.1 (10 points): Forward Pass

- Given the inputs and weights above, compute the values of the hidden units  $h_1, h_2$  and the outputs  $f_1, f_2, f_3$ . You should do this by hand, i.e. you should not write any code to do the calculation, but feel free to use a calculator to help you do the computations.
- You can optionally use  $\texttt{! [caption] (image.png)}$  in your answer on the Jupyter notebook. Otherwise, write your answer on paper and include a picture of your answer in this notebook. In order to include an image in Jupyter notebook, save the image in the same directory as the .ipynb file and then write `! [caption] (image.png)`. Alternatively, you may go to Edit --> Insert Image at the top menu to insert an image into a Markdown cell. **Double check that your image is visible in your PDF submission.**
- What class would the network predict for the input  $\mathbf{x}$ ?

$$\mathbf{x} = \begin{bmatrix} 1 & 3 & -2 \end{bmatrix} \quad \text{classes} = 0, 1, 2$$

$$W = \begin{bmatrix} 1 & -1 & 0 & 5 \\ 2 & 1 & 1 & 2 \end{bmatrix}$$

$$B = \begin{bmatrix} 4 & -1 & 0 \\ 3 & 0 & 2 \\ 2 & 1 & 1 \end{bmatrix}$$

$$h_1 = \text{ReLU} (x_1 w_{11} + x_2 w_{21} + x_3 w_{31} + w_{01})$$

$$= \text{ReLU} (-1 + 0 - 10 + 1)$$

$$= 0$$

$$h_2 = \text{ReLU} (x_1 w_{12} + x_2 w_{22} + x_3 w_{32} + w_{02})$$

$$= \text{ReLU} (1 + 3 - 4 + 2)$$

$$= 2$$

$$f_1 = \text{softmax} (h_1 \beta_{11} + h_2 \beta_{21} + \beta_{01})$$

$$= \text{softmax} (0 + 0 + 4)$$

$$= \text{softmax} (4)$$

$$f_2 = \text{softmax} (h_1 \beta_{12} + h_2 \beta_{22} + \beta_{02})$$

$$= \text{softmax} (7)$$

$$f_3 = \text{softmax} (h_1 \beta_{13} + h_2 \beta_{23} + \beta_{03})$$

$$= \text{softmax} (4)$$

and

Output

$$\sum_k e^k = e^4 + e^7 + e^4$$

$$\therefore f_1 = \frac{e^4}{\sum_k e^k} = 0.045$$

$$f_2 = \frac{e^7}{\sum_k e^k} = 0.909$$

$$f_3 = \frac{e^4}{\sum_k e^k} = 0.045$$

So, the predicted class would be class of  $f_2$ , i.e. (1)

## Problem 1.2 (10 points): Evaluate Loss

Typically when we train neural networks for classification, we seek to minimize the log-loss function. Note that the output of the log-loss function is always nonnegative ( $\geq 0$ ), but can be arbitrarily large (you should pause for a second and make sure you understand why this is true).

- Suppose the true label for the input  $x$  is  $y = 1$ . What would be the value of our loss function based on the network's prediction for  $x$ ?
- Suppose instead that the true label for the input  $x$  is  $y = 2$ . What would be the value of our loss function based on the network's prediction for  $x$ ?

You are free to use numpy / Python to help you calculate this, but don't use any neural network libraries that will automatically calculate the loss for you.

```
In [2]: f_pred = np.array([0.045, 0.909, 0.045])

def categorical_cross_entropy_loss(f_pred, y_true):
    return -np.log(f_pred[y_true])

loss_y1 = categorical_cross_entropy_loss(f_pred, 1)
loss_y2 = categorical_cross_entropy_loss(f_pred, 2)

print(f"Loss when y=1: {loss_y1}")
print(f"Loss when y=2: {loss_y2}")
```

```
Loss when y=1: 0.09541018480465818
Loss when y=2: 3.101092789211817
```

## Problem 1.3 (10 points): Network Size

- Suppose we change our network so that there are 12 hidden nodes instead of 2. How many total parameters (weights and biases) are in our new network?

```
In [3]: def calculate_total_parameters(input_nodes, hidden_nodes, output_nodes):
    input_to_hidden_weights = input_nodes * hidden_nodes
    hidden_biases = hidden_nodes
    hidden_to_output_weights = hidden_nodes * output_nodes
    output_biases = output_nodes
    total_parameters = (input_to_hidden_weights + hidden_biases +
                       hidden_to_output_weights + output_biases)
    return total_parameters

input_nodes = 3
hidden_nodes = 12
output_nodes = 3
```

```
total_params = calculate_total_parameters(input_nodes, hidden_nodes, output_nodes)
print("Total number of parameters:", total_params)
```

Total number of parameters: 87



## Problem 2: Neural Networks on MNIST

In this part of the assignment, you will get some hands-on experience working with neural networks. We will be using the scikit-learn implementation of a multi-layer perceptron (MLP). See [here](#) for the corresponding documentation. Although there are specialized Python libraries for neural networks, like [TensorFlow](#) and [PyTorch](#), in this problem we'll just use scikit-learn since you're already familiar with it.

### Problem 2.0: Setting up the Data

First, we'll load our MNIST dataset and split it into a training set and a testing set. Here you are given code that does this for you, and you only need to run it.

We will use the scikit-learn class `StandardScaler` to standardize both the training and testing features. Notice that we **only** fit the `StandardScaler` on the training data, and *not* the testing data.

```
In [4]: # Load the features and labels for the MNIST dataset
# This might take a minute to download the images.
X, y = fetch_openml('mnist_784', as_frame=False, return_X_y=True)

# Convert labels to integer data type
y = y.astype(int)
```

```
In [5]: X_tr, X_te, y_tr, y_te = train_test_split(X, y, test_size=0.1, random_state=seed, shuffle=True)
```

```
In [6]: scaler = StandardScaler()
scaler.fit(X_tr)
X_tr = scaler.transform(X_tr)      # We can forget about the original values & work
X_te = scaler.transform(X_te)      # just with the transformed values from here
```

### Problem 2.1: Varying the amount of training data (15 points)

One reason that neural networks have become popular in recent years is that, for many problems, we now have access to very large datasets. Since neural networks are very flexible models, they are often able to take advantage of these large datasets in order to achieve high levels of accuracy. In this problem, you will vary the amount of training data available to a neural network and see what effect this has on the model's performance.

In this problem, you should use the following settings for your network:

- A single hidden layer with 64 hidden nodes
- Use the ReLU activation function
- Train the network using stochastic gradient descent (SGD) and a constant learning rate of 0.001
- Use a batch size of 256
- **Make sure to set `random_state=seed`.**

Your task is to implement the following:

- Train an MLP model (with the above hyperparameter settings) using the first `m_tr` feature vectors in `X_tr`, where `m_tr = [100, 1000, 5000, 10000, 20000, 50000, 63000]`. You should use the `MLPClassifier` class from scikit-learn in your implementation.
- Create a plot of the training error and testing error for your MLP model as a function of the number of training data points. For comparison, also plot the training and test error rates we found when we trained a logistic regression model on MNIST (these values are provided below). Again, be sure to include an x-label, y-label, and legend in your plot and use a log-scale on the x-axis.
- Give a short (one or two sentences) description of what you see in your plot. Do you think that more data (beyond these 63000 examples) would continue to improve the model's performance?

**Note** that training a neural network with a lot of data can be a **slow process**. Hence, you should be careful to implement your code such that it runs in a reasonable amount of time. One recommendation is to test your code using only a small subset of the given `m_tr` values, and only run your code with the larger values of `m_tr` once you are certain your code is working. (For reference, it took about 20 minutes to train all models on a quad-core desktop with no GPU.)

```
In [7]: import time          # helpful if you want to track execution time
tic = time.time()
```

```

train_sizes = [100, 1000, 5000, 10000, 20000, 50000, 63000]

tr_err_mlp = []
te_err_mlp = []
for m_tr in train_sizes:
    ### YOUR CODE STARTS HERE
    X_tr_subset = X_tr[:m_tr]
    y_tr_subset = y_tr[:m_tr]

    model = MLPClassifier(hidden_layer_sizes=(64,), activation='relu', solver='sgd',
                          learning_rate_init=0.001, batch_size=256, random_state=seed, max_iter=100)

    model.fit(X_tr_subset, y_tr_subset)

    tr_err_mlp.append(1 - model.score(X_tr_subset, y_tr_subset))
    te_err_mlp.append(1 - model.score(X_te, y_te))

    print(f'Total elapsed time: {time.time()-tic}')

    ### YOUR CODE ENDS HERE

```

Total elapsed time: 0.37336111068725586  
 Total elapsed time: 4.782792329788208  
 Total elapsed time: 24.73273754119873  
 Total elapsed time: 49.42767357826233  
 Total elapsed time: 92.56259107589722  
 Total elapsed time: 201.64625787734985  
 Total elapsed time: 336.4739294052124

In [8]: # When plotting, use these (rounded) values from the similar logistic regression problem solution:

```

tr_err_lr = np.array([0. , 0. , 0. , 0. , 0.024, 0.053, 0.057])
te_err_lr = np.array([0.318, 0.149, 0.142, 0.137, 0.119, 0.087, 0.083])

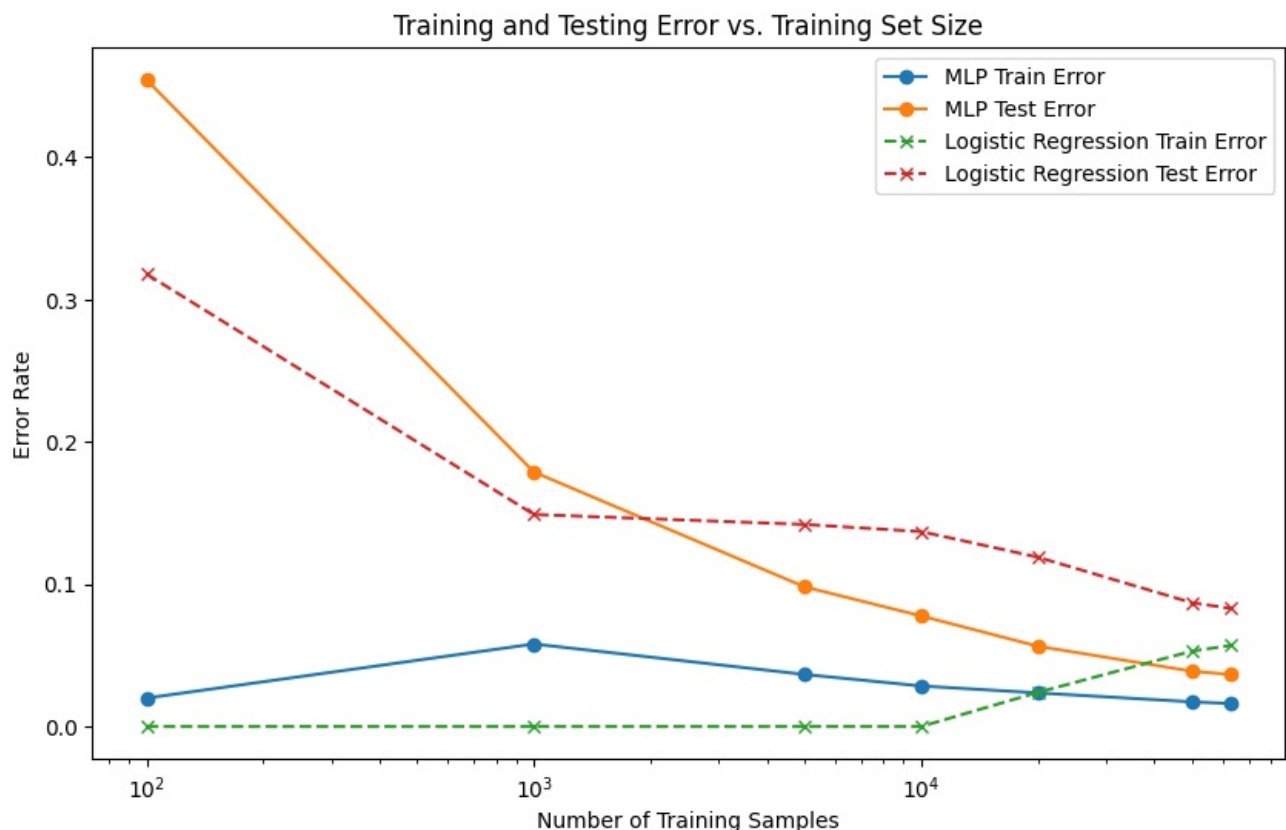
```

In [9]: ## YOUR CODE HERE (PLOTTING)

```

plt.figure(figsize=(10, 6))
plt.plot(train_sizes, tr_err_mlp, label="MLP Train Error", marker="o")
plt.plot(train_sizes, te_err_mlp, label="MLP Test Error", marker="o")
plt.plot(train_sizes, tr_err_lr, label="Logistic Regression Train Error", linestyle="--", marker="x")
plt.plot(train_sizes, te_err_lr, label="Logistic Regression Test Error", linestyle="--", marker="x")
plt.xscale("log")
plt.xlabel("Number of Training Samples")
plt.ylabel("Error Rate")
plt.legend()
plt.title("Training and Testing Error vs. Training Set Size")
plt.show()

```



## Discussion

The plot shows that as the number of training samples increases, both the training and test error of MLP decrease and converge. As for logistic regression, the both errors converge as well. Also, I don't think there will be much significant improvement if we increase the training samples above 63k.

## Problem 2.2: Optimization Curves (10 points)

One hyperparameter that can have a significant effect on the optimization of your model, and thus its performance, is the learning rate, which controls the step size in (stochastic) gradient descent. In this problem you will vary the learning rate to see what effect this has on how quickly training converges as well as the effect on the performance of your model.

In this problem, you should use the following settings for your network:

- A single hidden layer with 64 hidden nodes
- Use the ReLU activation function
- Train the network using stochastic gradient descent (SGD)
- Use a batch size of 256
- Set `n_iter_no_change=100` and `max_iter=100`. This ensures that all of your networks in this problem will train for 100 epochs (an *epoch* is one full pass over the training data).
- Make sure to set `random_state=seed`.

Your task is to:

- Train a neural network with the above settings, but vary the learning rate in `lr = [0.0005, 0.001, 0.005, 0.01]`.
- Create a plot showing the training loss as a function of the training epoch (i.e. the x-axis corresponds to training iterations) for each learning rate above. You should have a single plot with four curves. Make sure to include an x-label, a y-label, and a legend in your plot. (Hint: `MLPClassifier` has an attribute `loss_curve_` that you likely find useful.)
- Include a short description of what you see in your plot.

**Important: To make your code run faster, you should train all of your networks in this problem on only the first 10,000 images of `X_tr`.** In the following cell, you are provided a few lines of code that will create a small training set (with the first 10,000 images in `X_tr`) and a validation set (with the second 10,000 images in `X_tr`). You will use the validation later in Problem 3.3.

```
In [10]: # Create a smaller training set with the first 10,000 images in X_tr
#        along with a validation set from images 10,000 - 20,000 in X_tr

X_val = X_tr[10000:20000] # Validation set
y_val = y_tr[10000:20000]

X_tr = X_tr[:10000]      # From here on, we will only use these smaller sets,
y_tr = y_tr[:10000]      # so it's OK to discard the rest of the data
```

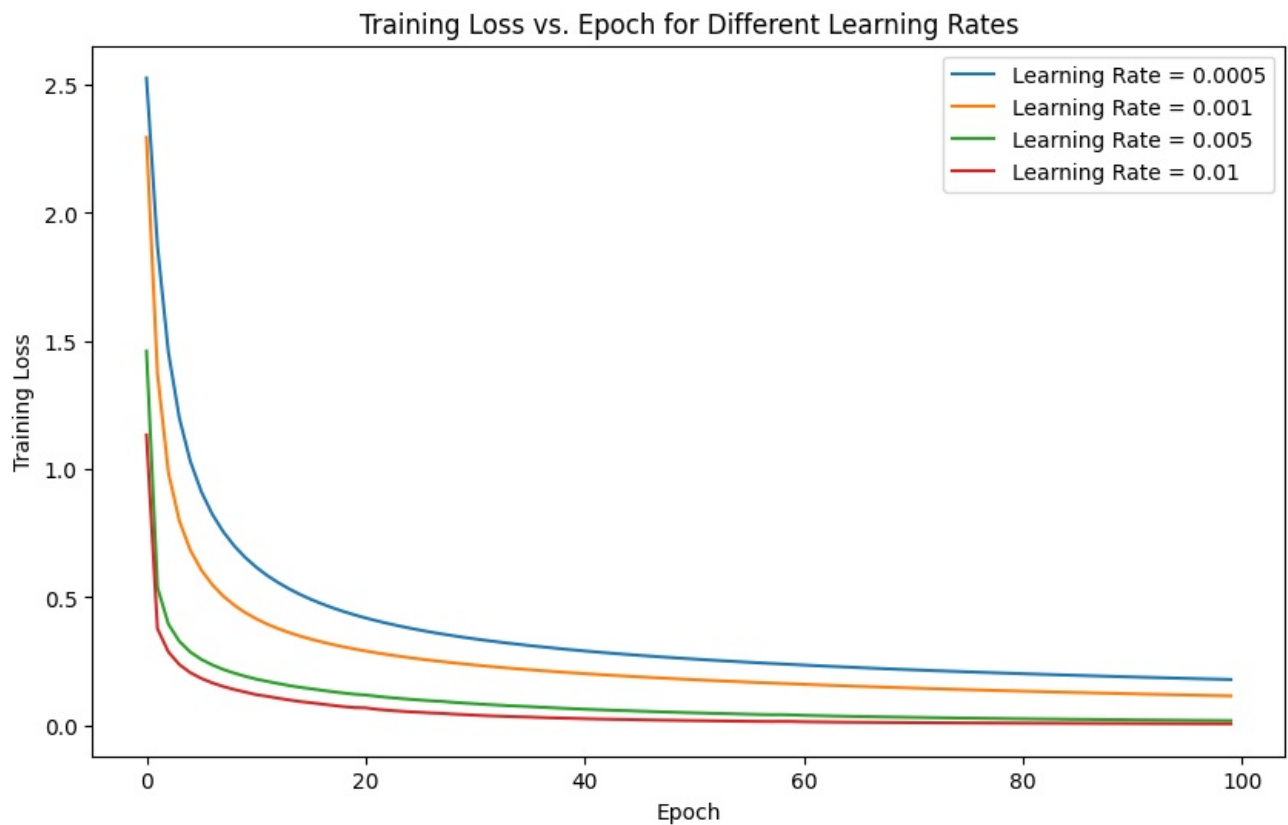
```
In [11]: learning_rates = [0.0005, 0.001, 0.005, 0.01]

err_curves = []

for lr in learning_rates:
    ### YOUR CODE STARTS HERE
    model = MLPClassifier(hidden_layer_sizes=(64,), activation='relu', solver='sgd',
                          learning_rate_init=lr, batch_size=256, max_iter=100,
                          n_iter_no_change=100, random_state=seed)
    model.fit(X_tr, y_tr)
    err_curves.append(model.loss_curve_)

plt.figure(figsize=(10, 6))
for i, lr in enumerate(learning_rates):
    plt.plot(err_curves[i], label=f"Learning Rate = {lr}")
plt.xlabel("Epoch")
plt.ylabel("Training Loss")
plt.title("Training Loss vs. Epoch for Different Learning Rates")
plt.legend()
plt.show()

### YOUR CODE ENDS HERE
```



## Discussion

Higher learning rates (like 0.01 and 0.005) lead to faster convergence, as the loss decreases more rapidly during the early epochs. However, the lowest learning rate (0.0005) converges more slowly, indicating that smaller learning rates may require more epochs to reach similar performance.

### Problem 2.3: Tuning a Neural Network (10 points)

As you saw in Problem 3.2, there are many hyperparameters of a neural network that can possibly be tuned in order to try to maximize the accuracy of your model. For the final problem of this assignment, it is your job to tune these hyperparameters.

For example, some hyperparameters you might choose to tune are:

- Learning rate
- Depth/width of the hidden layers
- Regularization strength
- Activation functions
- Batch size in stochastic optimization
- etc.

To do this, you should train a network on the training data `X_tr` and evaluate its performance on the validation set `X_val` -- your goal is to achieve the highest possible accuracy on `X_val` by changing the network hyperparameters. **Important: To make your code run faster, you should train all of your networks in this problem on only the first 10,000 images of `X_tr`.** This was already set up for you in Problem 3.2.

Try to find settings that enable you to achieve an error rate smaller than 5% on the validation data. However, tuning neural networks can be difficult; if you cannot achieve this target error rate, be sure to try at least five different neural networks (corresponding to five different settings of the hyperparameters).

In your answer, include a table listing the different hyperparameters that you tried, along with the resulting accuracy on the training and validation sets `X_tr` and `X_val`. Indicate which of these hyperparameter settings you would choose for your final model, and report the accuracy of this final model on the testing set `X_te`.

```
In [12]: hyperparameter_settings = [
    { 'hidden_layer_sizes': (64,), 'learning_rate_init': 0.01, 'alpha': 0.001, 'batch_size': 256, 'activation': 'tanh' },
    { 'hidden_layer_sizes': (128,), 'learning_rate_init': 0.005, 'alpha': 0.001, 'batch_size': 128, 'activation': 'tanh' },
    { 'hidden_layer_sizes': (64, 64), 'learning_rate_init': 0.001, 'alpha': 0.0005, 'batch_size': 256, 'activation': 'tanh' },
    { 'hidden_layer_sizes': (128, 64), 'learning_rate_init': 0.001, 'alpha': 0.0001, 'batch_size': 512, 'activation': 'tanh' },
    { 'hidden_layer_sizes': (256,), 'learning_rate_init': 0.0005, 'alpha': 0.0001, 'batch_size': 128, 'activation': 'tanh' }
]

# Track results
```



```

results = []

for params in hyperparameter_settings:
    model = MLPClassifier(hidden_layer_sizes=params['hidden_layer_sizes'],
                          activation=params['activation'], solver='sgd',
                          learning_rate_init=params['learning_rate_init'],
                          alpha=params['alpha'],
                          batch_size=params['batch_size'],
                          max_iter=50, random_state=seed)

    model.fit(X_tr, y_tr)

    train_accuracy = accuracy_score(y_tr, model.predict(X_tr))
    val_accuracy = accuracy_score(y_val, model.predict(X_val))

    results.append({
        'Hidden Layer Sizes': params['hidden_layer_sizes'],
        'Learning Rate': params['learning_rate_init'],
        'Alpha': params['alpha'],
        'Batch Size': params['batch_size'],
        'Activation': params['activation'],
        'Training Accuracy': train_accuracy,
        'Validation Accuracy': val_accuracy
    })

import pandas as pd
results_df = pd.DataFrame(results)

results_df

```

Out[12]:

	Hidden Layer Sizes	Learning Rate	Alpha	Batch Size	Activation	Training Accuracy	Validation Accuracy
0	(64,)	0.0100	0.0010	256	logistic	0.9583	0.9200
1	(128,)	0.0050	0.0010	128	relu	0.9996	0.9423
2	(64, 64)	0.0010	0.0005	256	relu	0.9524	0.9135
3	(128, 64)	0.0010	0.0001	512	relu	0.9365	0.9030
4	(256,)	0.0005	0.0001	128	relu	0.9596	0.9214

In [13]:

```

best_params = results_df.loc[results_df['Validation Accuracy'].idxmax()]

# Train the best model on training set and evaluate on test set
best_model = MLPClassifier(hidden_layer_sizes=best_params['Hidden Layer Sizes'],
                          activation=best_params['Activation'], solver='sgd',
                          learning_rate_init=best_params['Learning Rate'],
                          alpha=best_params['Alpha'],
                          batch_size=best_params['Batch Size'],
                          max_iter=100, random_state=seed)

best_model.fit(X_tr, y_tr)

test_accuracy = accuracy_score(y_te, best_model.predict(X_te))

print("The final model's parameters are \n")

for param, value in best_params.items():
    print(f"{param}: {value}")

print(f"\nThe final model's accuracy on the test set is {test_accuracy:.4f}")

```

The final model's parameters are

```

Hidden Layer Sizes: (128,)
Learning Rate: 0.005
Alpha: 0.001
Batch Size: 128
Activation: relu
Training Accuracy: 0.9996
Validation Accuracy: 0.9423

```

The final model's accuracy on the test set is 0.9409

## Problem 3: Torch and Convolutional Networks

In this problem, we will train a small convolutional neural network and compare it to the "standard" MLP model you built in Problem 2. Since `scikit` does not support CNNs, we will implement a simple CNN model using `torch`.

The `torch` library may take a while to install if it is not yet on your system. It should be pre-installed on ICS Jupyter Hub ( `hub.ics.uci.edu` ) and Google CoLab, if you prefer to use those.

## Problem 3.0: Defining the CNN

First, we need to define a CNN model. This is done for you; it consists of one convolutional layer, a pooling layer to down-sample the hidden nodes, and a standard fully-connected or linear layer. It contains methods to calculate the 0/1 loss as well as the negative log-likelihood, and trains using the Adam variant of SGD.

It can (optionally) output a real-time plot of the training process at each epoch, if you would like to assess how it is doing.

```
In [14]: import torch
torch.set_default_dtype(torch.float64)

class myConvNet(object):
    def __init__(self):
        # Initialize parameters: assumes data size! 28x28 and 10 classes
        self.conv_ = torch.nn.Conv2d(1, 16, (5,5), stride=2) # Be careful when declaring sizes;
        self.pool_ = torch.nn.MaxPool2d(3, stride=2) # inconsistent sizes will give you
        self.lin_ = torch.nn.Linear(400,10) # hard-to-read error messages.

    def forward_(self,X):
        """Compute NN forward pass and output class probabilities (as tensor) """
        r1 = self.conv_(X) # X is (m,1,28,28); R is (m,16,24,24)/2 = (m,16,12,12)
        h1 = torch.relu(r1) #
        h1_pooled = self.pool_(h1) # H1 is (m,16,12,12), so H1p is (m,16,10,10)/2 = (m,16,5,5)
        h1_flat = torch.nn.Flatten()(h1_pooled) # and H1f is (m,400)
        r2 = self.lin_(h1_flat)
        f = torch.softmax(r2,axis=1) # Output is (m,10)
        return f

    def parameters(self):
        return list(self.conv_.parameters())+list(self.pool_.parameters())+list(self.lin_.parameters())

    def predict(self,X):
        """Compute NN class predictions (as array) """
        m,n = X.shape
        Xtorch = torch.tensor(X).reshape(m,1,int(np.sqrt(n)),int(np.sqrt(n)))
        return self.classes_[np.argmax(self.forward_(Xtorch).detach().numpy(),axis=1)] # pick the most probab

    def J01(self,X,y): return (y != self.predict(X)).mean()
    def JNLL_(self,X,y): return -torch.log(self.forward_(X)[range(len(y)),y.astype(int)]).mean()

    def fit(self, X,y, batch_size=256, max_iter=100, learning_rate_init=.005, momentum=0.9, alpha=.001, plot=False):
        self.classes_ = np.unique(y)
        m,n = X.shape
        Xtorch = torch.tensor(X).reshape(m,1,int(np.sqrt(n)),int(np.sqrt(n)))
        self.loss01, self.lossNLL = [self.J01(X,y)], [float(self.JNLL_(Xtorch,y))]

        optimizer = torch.optim.Adam(self.parameters(), lr=learning_rate_init)
        for epoch in range(max_iter):
            # 1 epoch = pass through all data
            pi = np.random.permutation(m) # per epoch: permute data order
            for ii,i in enumerate(range(0,m,batch_size)): # & split into mini-batches
                ival = pi[i:i+batch_size]
                optimizer.zero_grad() # Reset the gradient computation
                Ji = self.JNLL_(Xtorch[ival,:,:,],y[ival])
                Ji.backward()
                optimizer.step()
            self.loss01.append(self.J01(X,y)) # track 0/1 and NLL losses
            self.lossNLL.append(float(self.JNLL_(Xtorch,y)))

            if plot: # optionally visualize progress
                display.clear_output(wait=True)
                plt.plot(range(epoch+2),self.loss01,'b-',range(epoch+2),self.lossNLL,'c-')
                plt.title(f'J01: {self.loss01[-1]}, NLL: {self.lossNLL[-1]}')
                plt.draw(); plt.pause(.01);
```

## Problem 3.1: CNN model structure (10 points)

How many (trainable) parameters are specified in the convolutional network? (If you like, you can count them -- you can access them through each trainable element, e.g., `myConvNet().conv_.parameters['weights']` and `...parameters['bias']`). List how many from each layer, and the total.

```
In [15]: model = myConvNet()

conv_params = sum(p.numel() for p in model.conv_.parameters())
pool_params = sum(p.numel() for p in model.pool_.parameters())
lin_params = sum(p.numel() for p in model.lin_.parameters())

total_params = conv_params + pool_params + lin_params

print(f"Convolutional Layer Parameters: {conv_params}")
print(f"Pooling Layer Parameters: {pool_params}")
```

```
print(f"Fully Connected Layer Parameters: {lin_params}")
print(f"Total Trainable Parameters: {total_params}")
```

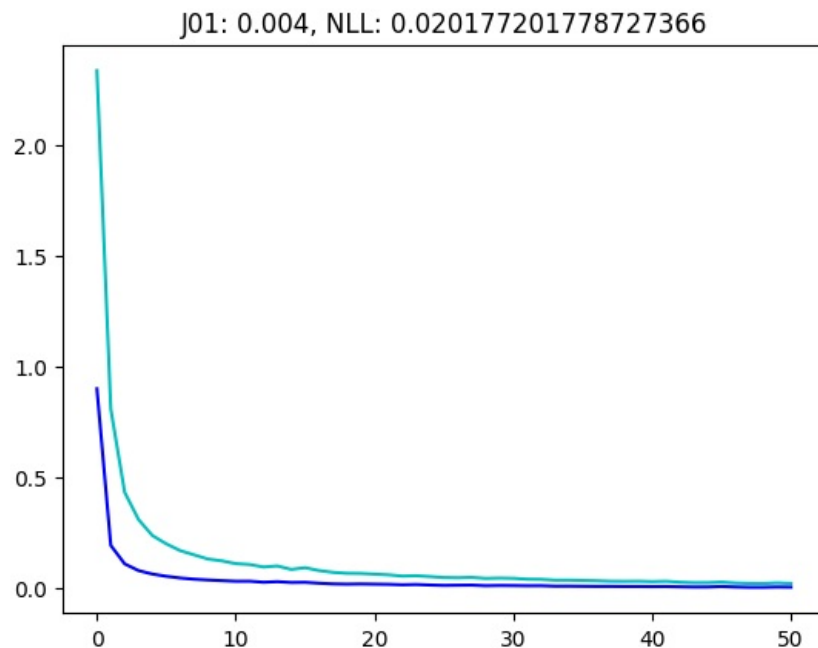
Convolutional Layer Parameters: 416  
 Pooling Layer Parameters: 0  
 Fully Connected Layer Parameters: 4010  
 Total Trainable Parameters: 4426

### Problem 3.2: Training the model (10 points)

Now train your model on `X_tr`. (Note that this should now include only 10k data points.) If you like you can plot while training to monitor its progress. Train for 50 epochs, using a learning rate of .001.

(Note that my simple training process takes these arguments directly into `fit`, rather than being part of the model properties as is typical in `scikit`.)

```
In [16]: # model = myConvNet()
model.fit(X_tr, y_tr, batch_size=256, max_iter=50, learning_rate_init=0.001, plot=True)
```



### Problem 3.3: Evaluation and Discussion (5 points)

Evaluate your CNN model's training, validation, and test error. Compare these to the values you got after optimizing your model's training process in Problem 2.3 (Tuning). Why do you think these differences occur? (Note that your answer may depend on how well your model in P2.3 did, of course.)

```
In [17]: train_error = model.J01(X_tr, y_tr)

val_error = model.J01(X_val, y_val)

test_error = model.J01(X_te, y_te)

print("Training Error (CNN):", train_error)
print("Validation Error (CNN):", val_error)
print("Test Error (CNN):", test_error)

mlp_train_error = 1 - best_params['Training Accuracy']
mlp_val_error = 1 - best_params['Validation Accuracy']
mlp_test_error = 1 - test_accuracy

print("\nComparison with Optimized MLP Model (Problem 2.3)")
print("Training Error (MLP):", mlp_train_error)
print("Validation Error (MLP):", mlp_val_error)
print("Test Error (MLP):", mlp_test_error)
```

Training Error (CNN): 0.004  
 Validation Error (CNN): 0.0268  
 Test Error (CNN): 0.032285714285714286

Comparison with Optimized MLP Model (Problem 2.3)  
 Training Error (MLP): 0.000399999999999995595  
 Validation Error (MLP): 0.057699999999999974  
 Test Error (MLP): 0.059142857142857164

# Discussion

1. CNNs generally perform better on image data compared to MLPs due to their ability to capture spatial features through convolutional layers.
2. Differences in model architecture, such as convolutional filters in CNNs, make them better suited for structured data like images.

## Problem 3.4: Comparing Predictions (5 points)

Consider the "somewhat ambiguous" data point `X_val[592]`. Display the data point (it will look a bit weird since it is already normalized). Then, use your trained `MLPClassifier` model to predict the class probabilities. If there are other classes with non-negligible probability, are they plausible? Similarly, find the class probabilities predicted by your CNN model. Compare the two models' uncertainty.

```
In [18]: plt.imshow(X_val[592].reshape(28, 28), cmap='gray')
plt.title("Ambiguous Data Point (X_val[592])")
plt.axis('off')

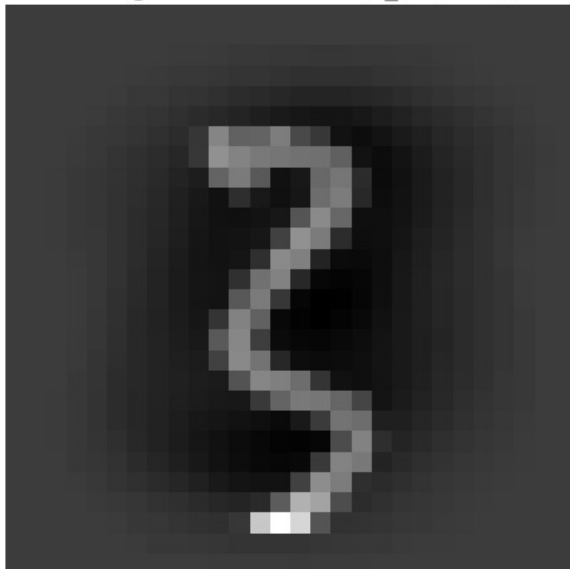
mlp_probabilities = best_model.predict_proba(X_val[592].reshape(1, -1))[0]

X_val_tensor = torch.tensor(X_val[592].reshape(1, 1, 28, 28), dtype=torch.float64)
cnn_probabilities = model.forward(X_val_tensor).detach().numpy()[0]

print("\nComparison of MLP and CNN Probabilities:")
print(f"{'Class':<10} {'MLP Probability':<15} {'CNN Probability':<15}")
for i in range(10):
    print(f"Class {i:<4} {mlp_probabilities[i]:<15.4f} {cnn_probabilities[i]:<15.4f}")
```

```
Comparison of MLP and CNN Probabilities:
Class      MLP Probability CNN Probability
Class 0    0.0042         0.0000
Class 1    0.1226         0.0000
Class 2    0.0539         0.0152
Class 3    0.0570         0.9769
Class 4    0.0000         0.0000
Class 5    0.0302         0.0000
Class 6    0.0154         0.0000
Class 7    0.2324         0.0032
Class 8    0.3082         0.0046
Class 9    0.1761         0.0000
```

Ambiguous Data Point (X\_val[592])



# Discussion

Although CNN is confident that it's 3, the MLP has confused it with digit 8 or 7 which is plausible due to the fact that in this image, the shape of 3 resembles some features of digit 8 and 7



## Statement of Collaboration (5 points)

It is **mandatory** to include a Statement of Collaboration in each submission, with respect to the guidelines below. Include the names of everyone involved in the discussions (especially in-person ones), and what was discussed.

All students are required to follow the academic honesty guidelines posted on the course website. For programming assignments, in particular, I encourage the students to organize (perhaps using EdD) to discuss the task descriptions, requirements, bugs in my code, and the relevant technical content before they start working on it. However, you should not discuss the specific solutions, and, as a guiding principle, you are not allowed to take anything written or drawn away from these discussions (i.e. no photographs of the blackboard, written notes, referring to EdD, etc.). Especially after you have started working on the assignment, try to restrict the discussion to EdD as much as possible, so that there is no doubt as to the extent of your collaboration.

I took help from scikit & torch documentation and some blogpost. Also, I discussed problem 1 with another student named Mahbub.

Processing math: 100%