

CS 178: Machine Learning & Data Mining

Homework 3: Due Friday 12 May 2023 (11:59pm)

Version 1.0 (Last Modified: 1 May 2023)

Instructions

This homework (and many 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.

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 (20 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 in Code (65 points)
 - Problem 2.1: Setting up Data (5 points)
 - Problem 2.2: Vary Amount of Data (20 points)
 - Problem 2.3: Learning Curves (10 points)
 - Problem 2.3: Tuning your Neural Network (30 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.

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

from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score

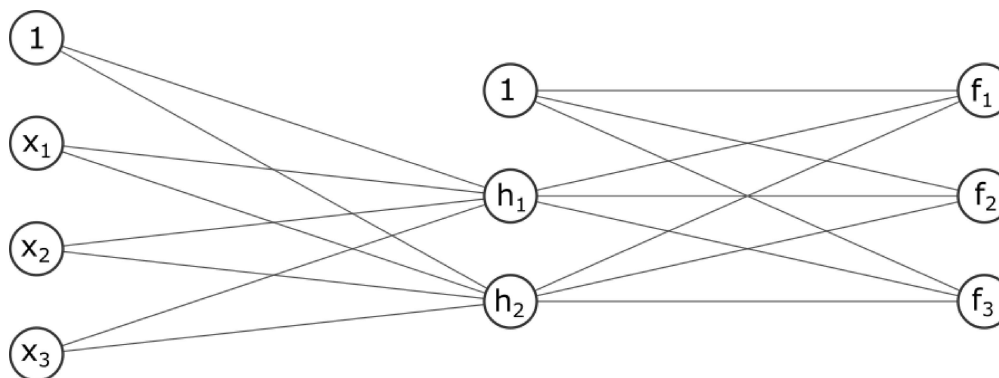
from sklearn.linear_model import LogisticRegression
from sklearn.neural_network import MLPClassifier

import warnings
warnings.filterwarnings('ignore')

# Fix the random seed for reproducibility
# !! Important !! : do not change this
seed = 1234
np.random.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} , as well as weights \mathbf{W} for the hidden layer and weights \mathbf{B} for the output layer. For example, w_{12} is the weight connecting input x_1 to hidden unit h_2 . This network uses the ReLU activation function for the hidden layer, and uses the softmax activation function for the output layer.



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

$$\mathbf{W} = \begin{bmatrix} w_{01} & w_{11} & w_{21} & w_{31} \\ w_{02} & w_{12} & w_{22} & w_{32} \end{bmatrix} = \begin{bmatrix} 1 & -1 & 0 & 6 \\ 2 & 1 & 1 & 3 \end{bmatrix}$$

$$\mathbf{B} = \begin{bmatrix} \beta_{01} & \beta_{11} & \beta_{21} \\ \beta_{02} & \beta_{12} & \beta_{22} \\ \beta_{03} & \beta_{13} & \beta_{23} \end{bmatrix} = \begin{bmatrix} 6 & -1 & 0 \\ 5 & 0 & 2 \\ 2 & 1 & 1 \end{bmatrix}$$

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 *L^AT_EX* 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} ?

$$h_1 = g(1 + (-1)x_1 + 0x_2 + 6x_3), \quad g = \text{ReLU} \\ = \max(0, 1 - x_1 + 6x_3)$$

$$h_2 = g(2 + 1x_1 + 1x_2 + 3x_3) \\ = \max(0, 2 + x_1 + x_2 + 3x_3)$$

$$f_1 = \text{softmax}(6 + (-1)h_1 + 0h_2) \\ = \text{softmax}(6 - h_1)$$

$$f_2 = \text{softmax}(5 + 0h_1 + 2h_2) \\ = \text{softmax}(5 + 2h_2)$$

$$f_3 = \text{softmax}(2 + 1h_1 + 1h_2) \\ = \text{softmax}(2 + h_1 + h_2)$$

$$\star \text{ note: } \text{softmax}(\mathbf{z}_k) = \frac{e^{z_k}}{\sum_{r=1}^K e^{z_r}}$$

Predict for \mathbf{x}

$$h_1(\mathbf{x}) = \max(0, 1 - 1 + 6(-2)) = \max(0, -12) = 0$$

$$h_2(\mathbf{x}) = \max(0, 2 + 1 + 3 + 3(-2)) = \max(0, 0) = 0$$

$$f_1(\mathbf{x}|\theta) = \text{softmax}(6 - 0) = \text{softmax}(6) = \frac{e^6}{e^6 + e^5 + e^2} \approx 0.72$$

$$f_2(\mathbf{x}|\theta) = \text{softmax}(5 + 2(0)) = \text{softmax}(5) = \frac{e^5}{e^6 + e^5 + e^2} \approx 0.27$$

$$f_3(\mathbf{x}|\theta) = \text{softmax}(2 + 0 + 0) = \text{softmax}(2) = \frac{e^2}{e^6 + e^5 + e^2} \approx 0.01$$

• Since f_1 is the greatest, the model would predict $y=0$ for \mathbf{x}

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 greater than zero, 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 \mathbf{x} is $y = 1$. What would be the value of our loss function based on the network's prediction for \mathbf{x} ?
- Suppose instead that the true label for the input \mathbf{x} is $y = 2$. What would be the value of our loss function based on the network's prediction for \mathbf{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 [21]: ▶ z = np.array([6, 5, 2])
exp_z = np.exp(z)
probs = exp_z / np.sum(exp_z)

from math import log

print(f"Loss for y=1: {-log(probs[1])}")
print(f"Loss for y=2: {-log(probs[2])}")

Loss for y=1: 1.3265626412674703
Loss for y=2: 4.32656264126747
```

Problem 1.3 (10 points): Network Size

- Suppose we change our network so that there are 12 hidden units instead of 2. How many total weights and biases would there be in our new network?

We would have $4 * 12 = 48$ weights for the hidden layer, then $(1 + 12) * 3 = 39$ weights for the output layer, so $39 + 48 = 87$ total weights

Problem 2: Neural Networks in Code

In the second problem of this 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 \(https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPClassifier.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPClassifier.html) for the corresponding documentation. Although there are specialized Python libraries for neural networks, like [TensorFlow \(https://www.tensorflow.org/\)](https://www.tensorflow.org/) and [PyTorch \(https://pytorch.org/\)](https://pytorch.org/), we'll stick with scikit-learn as you're already familiar with this library.

In this problem, we'll be working with the MNIST dataset, which we already saw in Homework 1. As a reminder, this is an image classification dataset, where each image is a hand-written digit. Take a look at Homework 1 to remind yourself what this dataset looks like.

Problem 2.1: Setting up the data (5 points)

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

- Use the scikit-learn class `StandardScaler` to standardize both the training and testing features. Remember that you should only fit the `StandardScaler` on the training data, and *not* the testing data.

```
In [23]: ▶ # 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 [44]: ▶ X_tr, X_te, y_tr, y_te = train_test_split(X, y, test_size=0.1, random_state=seed, shuffle=True)
```

```
In [45]: ▶ scaler = StandardScaler().fit(X_tr)
X_tr = scaler.transform(X_tr)
X_te = scaler.transform(X_te)
```

Problem 2.2: Varying the amount of training data (20 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 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 `n_tr` feature vectors in `X_tr` , where `n_tr` = `[100, 1000, 5000, 10000, 20000, 50000, 63000]` . You should use the `MLPClassifier` class from scikit-learn in your implementation.
- Train a logistic regression classifier (with the default settings in sklearn) using the first `n_tr` feature vectors in `X_tr` , where `n_tr` = `[100, 1000, 5000, 10000, 20000, 50000, 63000]` . You should use the `LogisticRegression` class from scikit-learn in your implementation. **Make sure to use the argument `random_state=seed` for reproducibility.**
- Create a plot of the training error and testing error for both the logistic regression and MLP models as a function of the number of training data points. Be sure to include an x-label, y-label, and legend in your plot. Use a log-scale on the x-axis. Give a short (one or two sentences) description of what you see in your plot.

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 `n_tr` values, and only run your code with all of the `n_tr` values given once you are certain your code is working.

```

In [41]: ► n_tr = [100, 1000, 5000, 10000, 20000, 50000, 63000]

MLP_tr_errors = []
MLP_te_errors = []
LR_tr_errors = []
LR_te_errors = []

MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(64,), activation='relu',
                    solver='sgd', learning_rate_init=0.001, batch_size=256)
LR = LogisticRegression(random_state=seed)

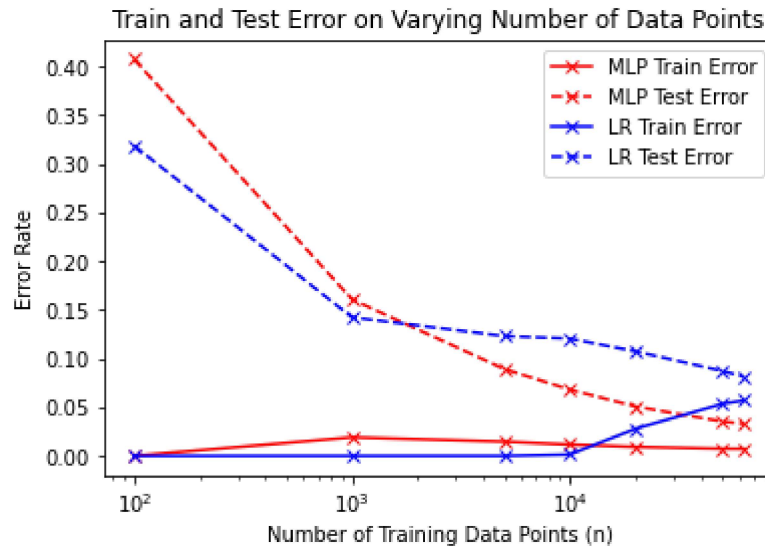
for n in n_tr:
    MLP.fit(X_tr[:n], y_tr[:n])
    MLP_tr_errors.append(1 - MLP.score(X_tr[:n], y_tr[:n]))
    MLP_te_errors.append(1 - MLP.score(X_te, y_te))

    LR.fit(X_tr[:n], y_tr[:n])
    LR_tr_errors.append(1 - LR.score(X_tr[:n], y_tr[:n]))
    LR_te_errors.append(1 - LR.score(X_te, y_te))

fig, axes = plt.subplots()
plt.xscale("log")
axes.set_title("Train and Test Error on Varying Number of Data Points")
axes.set_xlabel("Number of Training Data Points (n)")
axes.set_ylabel("Error Rate")
axes.plot(n_tr, MLP_tr_errors, 'rx-', label="MLP Train Error")
axes.plot(n_tr, MLP_te_errors, 'rx--', label="MLP Test Error")
axes.plot(n_tr, LR_tr_errors, 'bx-', label="LR Train Error")
axes.plot(n_tr, LR_te_errors, 'bx--', label="LR Test Error")
axes.legend()

```

Out[41]: <matplotlib.legend.Legend at 0x2011cc6bc10>



In the plot above, the training data of course has fairly low error rate until the logistic regression shoots up a bit with large amounts of data. For the testing data, the logistic regression has lower error rates than the neural network at small amounts of data, but is beaten by the neural network eventually with large amounts of data to accurately train it.

Problem 2.3: Learning Curves (10 points)

One hyperparameter that can have a significant effect on the performance of your model 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)
- 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).
- Use a batch size of 256
- **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 loss on the training set 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 2.4.

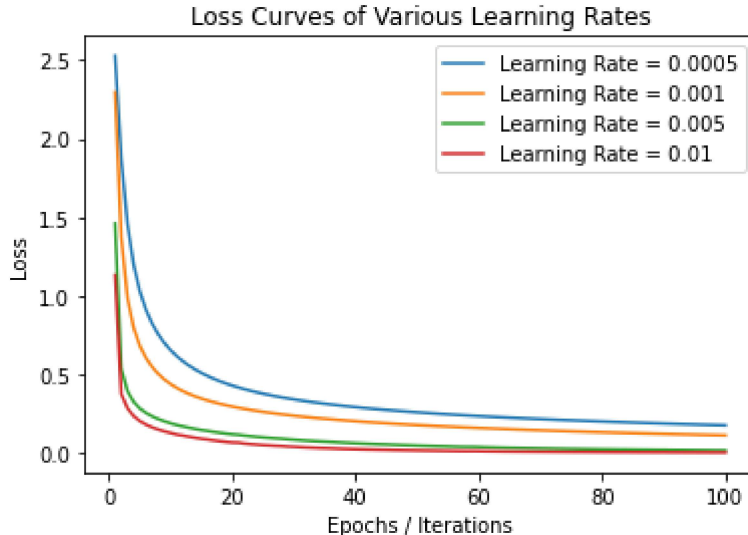
```
In [46]: ▶ # Create validation sets from the second 10k images in X_tr
X_val = X_tr[10000:20000]
y_val = y_tr[10000:20000]

# Create a smaller training set with the first 10k images in X_tr
X_tr = X_tr[:10000]
y_tr = y_tr[:10000]
```

```
In [52]: fig, axes = plt.subplots()
axes.set_title("Loss Curves of Various Learning Rates")
axes.set_xlabel("Epochs / Iterations")
axes.set_ylabel("Loss")

lrs = [0.0005, 0.001, 0.005, 0.01]
for lr in lrs:
    MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(64,), activation='relu',
                        solver='sgd', n_iter_no_change=100, max_iter=100, batch_size=256, learning_rate=lr)
    MLP.fit(X_tr, y_tr)
    axes.plot(np.linspace(1, 100, 100), MLP.loss_curve_, label=f"Learning Rate = {lr}")
axes.legend()
```

Out[52]: <matplotlib.legend.Legend at 0x2011f28ff40>



Problem 2.4: Tuning a Neural Network (30 points)

As you saw in Problem 2.2, 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

- Regularization strength
- Activation functions
- Batch size
- 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 2.3.

To receive full credit for this problem, you will need to tune your network hyperparameters until you achieve an error rate smaller than 5% on the validation data. However, tuning neural networks can be a difficult task, and you may not be able to achieve this target error rate. Hence, you will receive most of the credit for this problem as long as you train at least ten different neural networks with 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 [62]: ▶ MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(200,), activation='relu',
                                solver='adam', batch_size=256, learning_rate_init=0.001, max_iter=300, n_iter_no_
MLP.fit(X_tr, y_tr)
tr_error = 1 - MLP.score(X_tr, y_tr)
te_error = 1 - MLP.score(X_val, y_val)
print(f"    --> training error = {tr_error}")
print(f"    --> testing error  = {te_error}")

--> training error = 0.0
--> testing error  = 0.048200000000000002
```

```
In [64]: ▶ """
1) MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(64,), activation='relu',
                        solver='adam', batch_size=256, learning_rate_init=0.001)
   --> training error = 0.0
   --> testing error  = 0.0544

2) MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(32,), activation='relu',
                        solver='adam', batch_size=256, learning_rate_init=0.001)
   --> training error = 0.0
   --> testing error  = 0.064

3) MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(100,), activation='relu',
                        solver='adam', batch_size=256, learning_rate_init=0.001)
   --> training error = 0.0
   --> testing error  = 0.052

4) MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(128,), activation='relu',
                        solver='adam', batch_size=256, learning_rate_init=0.001, max_iter=300, n_iter_no_
   --> training error = 0.0
   --> testing error  = 0.05069999999999997

5) MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(150,), activation='relu',
                        solver='adam', batch_size=256, learning_rate_init=0.001, max_iter=300, n_iter_no_
   --> training error = 0.0
   --> testing error  = 0.048300000000000001

6) MLP = MLPClassifier(random_state=seed, hidden_layer_sizes=(200,), activation='relu',
                        solver='adam', batch_size=256, learning_rate_init=0.001, max_iter=300, n_iter_no_
   --> training error = 0.0
   --> testing error  = 0.0482
"""
```

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. If you did not collaborate with anyone, you should write something like "I completed this assignment without any collaboration."

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 completed this assignment without any collaboration.