

Lab 8 - Neural Networks

In this lab you'll use neural networks to classify images using both [scikit-learn](#) and [PyTorch](#). PyTorch 1.9 or later is assumed to be installed. The goal is for you to see:

- that logistic regression is a special case of neural networks; and
- how to express the same type of network in both scikit-learn and in PyTorch, both shallow (logistic regression) and deep (several layers).

Run the code cell below to import the required packages.

```
In [6]: import numpy as np
import matplotlib.pyplot as plt
import sklearn
import sklearn.preprocessing # For StandardScaler
import sklearn.linear_model # For LogisticRegression
import sklearn.neural_network # For MLPClassifier
import torch
warnings.filterwarnings('ignore', category=sklearn.exceptions.ConvergenceWarning) # Annoying
np.set_printoptions(precision=3, suppress=True) # Prints as 0.001 instead of 9.976e-4
```

1. Digit classification with neural networks in scikit-learn

Exercises 1.1-1.3 ask you to load and train a model on the classic MNIST data set. It's so classic it has its own [Wikipedia page](#)! The MNIST data set contains about 60,000 training examples and 10,000 test examples. Each example consists of a 784-dimensional feature vector x_i , representing a 28x28 grayscale image of a hand-written digit ($784 = 28 \times 28$) with a label $y_i \in \{0, \dots, 9\}$.

Since there are 60,000 training cases, the matrix of training features X is provided in a 60000x784 matrix of pixel intensities. Value $X_{ij} \in \{0, \dots, 255\}$ represents the intensity (0=black, 255=white) of pixel number j in training image i . Each 784-dimensional feature vector x_i can be reshaped into a 28x28 image as depicted below.

Run the code cell below to define a function that will be useful for plotting matrices.

```
In [7]: def plot_matrix_grid(V):
    """
    Given an array V containing stacked matrices, plots them in a grid layout.
    V should have shape (K,M,N) where V[k] is a matrix of shape (M,N).
    """
    K, M, N = V.ndim == 3, "Expected V to have 3 dimensions, not %d" % V.ndim
    ncol = 8
    nrow = min(4, (K // ncol - 1) // ncol) # At most 4 rows
    V = V.reshape((-1, M*N)) # Focus on just the matrices we'll actually plot
    figsize = (2*ncol, M*N) # Guess a good figure shape based on ncol, nrow
    fig, axes = plt.subplots(nrow, ncol, sharex=True, sharey=True, figsize=figsize)
    for v, ax in zip(V, axes.flat): # Guess a good figure shape based on ncol, nrow
        img = ax.matshow(v, vmin=vmin, vmax=vmax, cmap=plt.get_cmap('gray'))
        ax.set_ticks([])
        ax.set_yticks([])
    fig.colorbar(img, cax=fig.add_axes((10.92, 0.25, 0.01, .5))) # Add a colorbar on the right
```

Exercise 1.1 – Load MNIST and plot some digits

The MNIST training data has been provided to you in a file called `mnist_train.npz`. The file is located in the same directory as this Jupyter Notebook. A `npz` file is an efficient way to store data in a file. Use Numpy's `load` function to open an `npz` file. When the file is opened, you can think of the file as being a Python dictionary where you can ask for an array by its name (its 'key'). The example below shows how to open the file and list the keys:

```
>>> with np.load("mnist_train.npz") as data:
...     print(list(data.keys()))

['x', 'y']
```

(The reason we open the file using a with-statement is because once the with-statement is complete the file ("file descriptor") is automatically closed, rather than Python trying to keep the file open. This isn't important for the lab per se, closing files when you're done with them is just good programming practice!)

Write a few lines of code to load the training data from `mnist_train.npz` and create two global variables X_{train} and y_{train} to refer to the data you loaded.

```
In [8]: # Your code here, aim for 3-4 lines.
with np.load('mnist_train.npz') as data:
    X_train = data['x']
    y_train = data['y']
```

Inspect the data by printing information about the arrays.

- Print the shape and dtype of both your X_{train} and y_{train} arrays.
- Print the first five training samples from X_{train} and y_{train} arrays.

Since your X_{train} array is big, and because most of the first/last pixels in each image are 0 (black), to see any patterns in the features try printing a slice of values taken from the "middle" of each image. You can try, for example, pixels 400:415 are roughly from the middle row of each image (similar to blue rectangle in the diagram earlier), so try printing a slice of just those pixels. You should see `[[0 0 0 0 81 240 253 253 119 25 0 0 0 0 0]]` printed for the first row.

```
In [9]: # Your code for printing shape and dtype here. Aim for 2 lines.
print(f'Shape: {X_train.shape}, data type: {X_train.dtype}')
print(f'Shape: {y_train.shape}, data type: {y_train.dtype}')

# Your code for printing sample values. Aim for 2 lines.
print(X_train[5, 400:415])
print(y_train[5])
```

Plot a few digits to see what they look like. Use the `plot_matrix_grid` function defined earlier. To do this, you'll need to reshape the array referred to by your X_{train} variable so that the plotting code knows the images have shape 28x28 rather than being just 784-dimensional vectors.

```
In [10]: # Your code here. Aim for 1-2 lines.
x = X_train.reshape(-1, 28, 28)
plot_matrix_grid(x)
```

Look at the patterns you printed when inspecting the X_{train} variable earlier, and make sure you see where they come from in the first five images plotted above.

If you want to see more of the MNIST training digits, rather than just the first few, you can try plotting different "slices" of the X_{train} variable, such as `X_train[100:]` to start plotting at the 101st training example. (You still have to reshape the resulting array, of course.)

Finally, load the MNIST test data from the file `mnist_test.npz`, just like you did for the training data. Create global variables X_{test} and y_{test} to refer to the arrays that you loaded. These arrays will be used to evaluate test-time accuracy later on.

```
In [11]: # Your code here. Aim for 3 lines.
with np.load('mnist_test.npz') as data:
    X_test = data['x']
    y_test = data['y']
```

Exercise 1.2 – Preprocess the MNIST data

Certain models trained on MNIST work better when the features are normalized. Use `scikit-learn` to normalize the MNIST data using scaling, such as the `StandardScaler`. (You can just treat the pixels as independent features, nothing fancy.)

Write a few lines of code to normalize both your X_{train} and X_{test} variables. You can just over-write those variables with the new (normalized) feature arrays, and discard the original unnormalized data.

```
In [12]: # Your code here. Aim for 3-4 lines.
scaler = sklearn.preprocessing.StandardScaler(copy=False).fit(X_train)
X_train = scaler.transform(X_train)
X_test = scaler.transform(X_test)
```

Plot the rescaled training digits using the `plot_matrix_grid` function.

```
In [13]: # Your code here. Aim for 1-2 lines.
x = X_train.reshape(-1, 28, 28)
plot_matrix_grid(x)
```

Notice that the pixels in the center tend to be scaled down more than the pixels in the periphery. Do you understand why?

Exercise 1.3 – Train multinomial logistic regression on MNIST

Train a `LogisticRegression` object to classify MNIST digits. Use `random_state=0` and default settings otherwise.

```
In [14]: # Your code here. Aim for 2-3 lines.
lr = sklearn.linear_model.LogisticRegression(C=0.01, random_state=0)
lr.fit(X_train, y_train)
```

You can use the `score` method of the `LogisticRegression` object to compute the accuracy as a number in the range $[0.0, 1.0]$. Figure out how to convert that number (e.g., 0.934) into an error rate percentage (e.g., 6.6%).

Print the training error rate and testing error rate of your logistic regression model on the MNIST data set. Your output should be in the form:

```
X.X% training error
X.X% testing error
```

How does the testing error rate you see compare to some of the error rates mentioned on the [MNIST Wikipedia page](#)?

```
In [15]: # Your code here. Aim for 2-4 lines.
test_error = 1 - lr.score(X_test, y_test)
train_error = 1 - lr.score(X_train, y_train)
print(f'train_error = {100 - 2*train_error}% training error')
print(f'test_error = {100 - 2*test_error}% testing error')
```

Print the predicted class probabilities of the first five examples in the training set. Use the `predict_proba` method of your `LogisticRegression` object. The first row of output should look something like:

```
[0.001 0. 0. 0.203 0. 0.796 0. 0. 0. ]
```

From the above probabilities we can see that the model thinks the first digit in the training set is probably digit "5" but might also be digit "3".

```
In [16]: # Your code here. Aim for 1-2 lines.
print(lr.predict_proba(X_train[5]))
```

```
[0.001 0. 0.001 0.231 0. 0.766 0. 0.001 0. 0. ]
[0. 0. 0. 0. 0. 0. 0. 0. 0. 0. ]
[0. 0.001 0.009 0.086 0.873 0. 0. 0.025 0.001 0.006]
[0. 0.965 0.019 0.002 0. 0. 0. 0. 0.013 0. ]
[0. 0. 0. 0. 0. 0.039 0. 0. 0.203 0.003 0.965]
```

Exercise 1.4 – Visualize the weights of your logistic regression model

The logistic regression model you trained in Exercise 1.3 has a `coef_` attribute. This attribute is the array of weights W we seen in Lecture 4 (e.g. slide 28). For the MNIST data, this matrix has shape $(10, 784)$, because there are 10 output classes and $784=28 \times 28$ pixels. Weight W_{ij} is the weight with which of pixel j contributes to output class i .

You are asked to visualize the weights using `plot_matrix_grid`. You may need to reshape the weight matrix to do this. The first two outputs, corresponding to predicted digit "0" and predicting digit "1" should look something like this:

Notice how the pattern for "0" is has strong negative weights in the center: that's because if there are white pixels in the center, it's unlikely that the image represents digit "0"!

If your patterns appear "noisier" than above, try repeating Exercise 1.3 but weaken "LogisticRegression"s L2 penalty by a factor of 100 from its default. Take note of any change in training/test accuracy.

Write a few lines of code to plot the weights and see what patterns they contain. You should see ten patterns. (Don't worry if the last few grid entries are just white boxes.)

```
In [17]: # Your code here. Aim for 1-2 lines.
plot_matrix_grid(lr.coef_.reshape(-1, 28, 28))
```

When an input image (of a hand-written digit) causes one of these patterns to have a large positive response (strong activation), then the corresponding class $\{0, 1, \dots, 9\}$ will be given a high probability by the final softmax operation.

Exercise 1.5 – Train a neural network on MNIST with zero hidden layers

Train a neural network on MNIST using the `sklearn.neural_network.MLPClassifier` class.

A neural network has *many* more hyperparameters to configure. Configure your neural network as follows:

- Ask for *no hidden layers*. You can do this by specifying an empty tuple `()` for the `hidden_layer_sizes` argument. This will create a neural network where the 784 input features are directly 'connected' to the 10 output predictions, which in this case corresponds to the multinomial logistic regression you did in Exercise 1.4.
- Use the `sgd` solver. This means stochastic gradient descent that we saw in Lecture 1.
- Use a batch size of 100. This means that at each step of SGD the gradient will be computed from only 100 of the 60,000 training cases. This is also called a "mini-batch". The SGD algorithm works by starting with the first 100, then the next 100, and then it continues. This is the last 100 in the training set it starts from the beginning again.
- Use `max_iter=10`. This causes the training to stop after SGD has passed over all 60,000 training cases exactly 10 times.
- Use `learning_rate_init=0.01`, which determines the step size for SGD once it has computed a gradient.
- Use `momentum=0.9`, which speeds up training.
- Use `random_state=0` for reproducibility.
- Use `verbose=True` to see progress printed out. Each time it prints "Iteration X" it means SGD has made another pass over all 60,000 training examples.

```
In [18]: # Your code here. Aim for 1-2 lines, plus whatever line wrapping you need for arguments!
mlp = sklearn.neural_network.MLPClassifier(hidden_layer_sizes=(), solver='sgd',
                                           batch_size=100, max_iter=10, learning_rate_init=0.01,
                                           momentum=0.9, random_state=0, verbose=True)
mlp.fit(X_train, y_train)
```

```
Iteration 1, loss = 0.4703096
Iteration 2, loss = 0.3895873
Iteration 3, loss = 0.3952979
Iteration 4, loss = 0.28490254
Iteration 5, loss = 0.2743585
Iteration 6, loss = 0.16956936
Iteration 7, loss = 0.1455986
Iteration 8, loss = 0.26427636
Iteration 9, loss = 0.2229657
Iteration 10, loss = 0.25639397
```

Print the training error rate and test error rate of your neural network classifier, just like you did for logistic regression. How does your error rate compare to multinomial logistic regression? (Exercises 1.3 and 1.5)

```
In [19]: # Your code here. Aim for 3-4 lines.
test_error = 1 - mlp.score(X_test, y_test)
train_error = 1 - mlp.score(X_train, y_train)
print(f'train_error = {100 - 2*train_error}% training error')
print(f'test_error = {100 - 2*test_error}% testing error')
```

```
4.18% training error
7.68% testing error
```

Exercise 1.6 – Visualize the weights of a neural network (no hidden layers)

The `MLPClassifier` object has a `coefs_` attribute that works just like the `coef_` attribute that contained coefficient matrix W of `LogisticRegression`, except that for a neural network there are two differences:

- `coefs_` is a list of coefficient matrices, so `coefs_[0]` is $W^{(0)}$, the coefficient matrix of the first layer. Since the neural network you trained in Exercise 1.5 has no hidden layers, this $W^{(0)}$ matrix holds the same weights as the W matrix for `LogisticRegression`.
- The weight matrix for `MLPClassifier` has a different layout: it's 784×10 rather than 10×784 . Do you now how to account for this?

Write a few lines of code to repeat Exercise 1.4 but this time with the neural network weights.

```
In [20]: # Your code here. Aim for 1-2 lines.
plot_matrix_grid(np.array(mlp.coefs_[0]).T.reshape(-1, 28, 28))
```

If your patterns look sketchy then you may need to try transposing your weight matrix to account for the different layout.

Exercise 1.7 – Train and visualize the weights of a neural network with 1 hidden layer

Here you're asked to train a neural network with just one hidden layer in Exercise 1.5, but this time add a hidden layer with 16 "tanh" hidden units to your neural network. Then you'll visualize the weights of this network.

Read the documentation for `MLPClassifier` to learn how to do specify a hidden layer. (Note: in Python if you want to create a tuple object with only one item in it, you can use `(item,)` with an extra comma, rather than `item`; Python interprets to just be regular parentheses.) All the other hyperparameters can stay the same as Exercise 1.5.

Write a few lines of code to train a neural network, this time with 16 tanh hidden units. In other words, this will be a 784-16-10 neural network where the hidden layer uses tanh activations.

```
In [21]: # Your code here. Aim for 1-2 lines, plus whatever line wrapping you need for arguments!
mlp2 = sklearn.neural_network.MLPClassifier(hidden_layer_sizes=(16), activation='tanh', solver='sgd',
                                           batch_size=100, max_iter=10, learning_rate_init=0.01,
                                           momentum=0.9, random_state=0, verbose=True)
mlp2.fit(X_train, y_train)
```

```
Iteration 1, loss = 0.4703096
Iteration 2, loss = 0.3895873
Iteration 3, loss = 0.3952979
Iteration 4, loss = 0.07193538
Iteration 5, loss = 0.05528783
Iteration 6, loss = 0.04298461
Iteration 7, loss = 0.03308290
Iteration 8, loss = 0.02760081
Iteration 9, loss = 0.0229657
Iteration 10, loss = 0.0193338
```

Print the training error rate and test error rate of your neural network classifier, just like you did for logistic regression. How does your error rate compare to multinomial logistic regression? (Exercises 1.3 and 1.5)

```
In [22]: # Your code here. Aim for 2-4 lines.
test_error = 1 - mlp2.score(X_test, y_test)
train_error = 1 - mlp2.score(X_train, y_train)
print(f'train_error = {100 - 2*train_error}% training error')
print(f'test_error = {100 - 2*test_error}% testing error')
```

```
4.18% training error
6.42% testing error
```

Plot the first-layer weights $W^{(0)}$ of your neural network using the `plot_matrix_grid` function, just in Exercise 1.6.

```
In [23]: # Your code here. Aim for 1-2 lines.
plot_matrix_grid(np.array(mlp2.coefs_[0]).T.reshape(-1, 28, 28))
```

Notice that there are now 16 patterns, not 10, and they no longer seem to correspond to the digits $\{0, 1, \dots, 9\}$ in any particular order. Do you understand why?

Plot the second-layer weights $W^{(2)}$ of your neural network using the `plot_matrix_grid` function.

However, this time if you inspect the shape of the second weight matrix, `coefs_[1]`, you'll see that it has shape $(16, 10)$, and so it cannot be reshaped into a 28×28 pattern. In fact the second layer has only dimension: the "hidden layer" is just a vector of 16 values (the 16 tanh-transformed activations of the first-layer patterns). Each of the 10 output units has 16 weights contributing to it, rather than 784 weights like in Exercise 1.6.

Figure out how to reshape the weight matrix so that when you call `plot_matrix_grid` you see a grid of 16×16 weight vectors, like the two examples below.

```
In [24]: # Your code here. Aim for 1-2 lines.
plot_matrix_grid(np.array(mlp2.coefs_[1]).T.reshape(-1, 4, 4))
```

Exercise 1.8 – Train a neural network with lots of hidden units

Repeat Exercise 1.7 but with two hidden layers having 100 and 50 hidden units respectively. This time use `relu` activations. All other hyperparameters can stay the same.

Write a few lines of code to train the model here.

```
In [25]: # Your code here. Aim for 1-2 lines, plus whatever line wrapping you need for arguments!
mlp3 = sklearn.neural_network.MLPClassifier(hidden_layer_sizes=(100, 50), activation='relu', solver='sgd',
                                           batch_size=100, max_iter=10, learning_rate_init=0.01,
                                           momentum=0.9, random_state=0, verbose=True)
mlp3.fit(X_train, y_train)
```

```
Iteration 1, loss = 0.33306221
Iteration 2, loss = 0.13727616
Iteration 3, loss = 0.05239719
Iteration 4, loss = 0.07193538
Iteration 5, loss = 0.05528783
Iteration 6, loss = 0.04298461
Iteration 7, loss = 0.03308290
Iteration 8, loss = 0.02760081
Iteration 9, loss = 0.0229657
Iteration 10, loss = 0.0193338
```

Print the training and testing error rates here. How do they compare to earlier models?

```
In [26]: # Your code here. Aim for 2-4 lines.
test_error = 1 - mlp3.score(X_test, y_test)
train_error = 1 - mlp3.score(X_train, y_train)
print(f'train_error = {100 - 2*train_error}% training error')
print(f'test_error = {100 - 2*test_error}% testing error')
```

```
0.28% training error
2.22% testing error
```

Plot the first-layer weights $W^{(1)}$ of your neural network here. Are the pattern detectors here qualitatively different than for earlier models? For example, do you see any digits?

```
In [27]: # Your code here. Aim for 1-2 lines.
plot_matrix_grid(np.array(mlp3.coefs_[0]).T.reshape(-1, 28, 28))
```

Don't bother plotting the 2nd and 3rd layer weights, they are high-dimensional and hard to interpret.

2. Neural networks in PyTorch

Exercises 2.1–2.3 ask you to train a simple neural network in `PyTorch`. Here you'll use `PyTorch` to train an MNIST classifier using the same MNIST data that you already preprocess in Part 1. The goal is just to get you familiar with PyTorch basics and how they compare to scikit-learn.

`PyTorch` is a deep learning framework like `TensorFlow`. `PyTorch` tends to be popular with deep learning researchers because it's very flexible for trying new ideas. `TensorFlow` is also flexible but is designed in such a way that it's more popular for companies trying to deploy high-performance models (in the cloud etc). Both can be used for research, of course!

Exercise 2.1 – Convert MNIST from Numpy arrays to PyTorch tensors

`PyTorch` has its own Numpy-like array class, called `Tensor`. In order to train a `PyTorch` model, you must first convert the Numpy arrays. `PyTorch` understands Numpy arrays, so this is easy. The only tricky part is that, in order to be fast and not waste memory, `PyTorch` tends to be more picky about the dtype of the arrays you give it.

Write a few lines of code to create four global variables: $X_{\text{train_torch}}$, $y_{\text{train_torch}}$, $X_{\text{test_torch}}$, $y_{\text{test_torch}}$ that are `PyTorch` versions of your preprocessed MNIST training data from Part 1. The X tensors should have dtype `float32`, and the y tensors should have dtype `int64`.

```
In [28]: # Your code here. Aim for 2-4 lines.
X_train_torch = torch.tensor(X_train, dtype=torch.float32)
X_test_torch = torch.tensor(X_test, dtype=torch.float32)
y_train_torch = torch.tensor(y_train, dtype=torch.int64)
y_test_torch = torch.tensor(y_test, dtype=torch.int64)
```

Run the code cell below to check your answer.

```
In [29]: assert 'X_train_torch' in globals(), "You didn't declare a X_train_torch variable!"
assert 'y_train_torch' in globals(), "You didn't declare a y_train_torch variable!"
assert 'X_test_torch' in globals(), "You didn't declare a X_test_torch variable!"
assert 'y_test_torch' in globals(), "You didn't declare a y_test_torch variable!"
assert isinstance(X_train_torch, torch.Tensor)
assert isinstance(y_train_torch, torch.Tensor)
assert isinstance(X_test_torch, torch.Tensor)
assert isinstance(y_test_torch, torch.Tensor)
assert X_train_torch.dtype == torch.float32
assert y_train_torch.dtype == torch.int64
assert X_test_torch.dtype == torch.float32
assert y_test_torch.dtype == torch.int64
assert X_train_torch.shape == (60000, 784)
assert y_train_torch.shape == (60000,)
assert X_test_torch.shape == (10000, 784)
assert y_test_torch.shape == (10000,)
print("Correct!")
```

Correct!

Exercise 2.2 – Train a PyTorch neural network without hidden layers

This exercise only asks you to run existing code so that you learn how `PyTorch` works. The code in this cell defines a simple logistic model, and then you're asked to modify the code to add hidden layers to the network.

Useful documentation for understanding the code that you see:

- `torch.nn` (neural network)
- `torch.optim` (optimizers such as SGD)

Here are some comments to help you understand the "starter code" below:

- A neural network is a sequence of non-linear transformations, so `PyTorch` provides a `Sequential` class that accepts a list of desired transformations.
- In a standard neural network, the transformations are just linear, i.e. $Wx + b$, and in `PyTorch` this is implemented by a `Linear` class where constructing one of these objects with `Linear(D, M)` tells the new object that it should be expecting an D -dimensional input and transform it into a M -dimensional output. To do this, the `Linear` object will create its own parameter matrix $W \in \mathbb{R}^{M \times D}$ and bias vector $b \in \mathbb{R}^M$.
- In `PyTorch`, the `CrossEntropyLoss` class conveniently combines applying a softmax and then computing the negative log-likelihood, so you don't explicitly apply softmax while training. Once you have a `CrossEntropyLoss` object, you can call it with your predictions and targets (both vectors), and it will compute the negative log likelihood, which is just one number (a scalar).

Run the code cell below to define a simple 784-10 neural network (i.e. logistic regression).

```
In [30]: torch.manual_seed(0) # Ensure model weights initialized with same random numbers
model = torch.nn.Sequential(
    torch.nn.Linear(28*28, 100),
    torch.nn.ReLU(inplace=True),
    torch.nn.Linear(100, 50),
    torch.nn.ReLU(inplace=True),
    torch.nn.Linear(50, 10),
)
```

Run the code cell below to define some objects and variables needed for training the neural network.

```
In [31]: # Create an object that can compute "negative log likelihood of a softmax"
loss = torch.nn.CrossEntropyLoss()

# Use stochastic gradient descent to train the model
optimizer = torch.optim.SGD(model.parameters(), lr=0.01, momentum=0.9)

# Use 100 training samples at a time to compute the gradient.
batch_size = 100

# Make 10 passes over the training data, each time using batch_size samples to compute gradient
num_epoch = 10
next_epoch = 1
```

Run the code cell below to train the neural network using stochastic gradient descent (SGD). Note that if you re-run this code cell multiple times it will "continue" training from the current parameters, and if you want to "reset" the model you need to re-run the earlier code cell that defined the model!

```
In [32]: for epoch in range(next_epoch, next_epoch+num_epoch):
    # Make an entire pass (an "epoch") over the training data in batch_size chunks
    for i in range(0, len(X_train), batch_size):
        X = X_train_torch[i:i+batch_size]
        y = y_train_torch[i:i+batch_size]
        # Make predictions (final-layer activations)
        y_pred = model(X)
        l = loss(y_pred, y)
        # Compute all gradient accumulators to zero (PyTorch thing)
        model.zero_grad_()
        l.backward()
        optimizer.step()
        # Compute gradient of loss wrt all parameters (backprop!)
        # Use the gradients to take a step with SGD.
        print("Epoch %2d: loss on final training batch: %4f" % (epoch, l.item()))
    print("Epoch %2d: loss on test set: %4f" % (epoch, loss(model(X_test_torch), y_test_torch)))
    next_epoch = epoch+1
```

```
Epoch 1: loss on final training batch: 0.5178
Epoch 2: loss on final training batch: 0.2210
Epoch 3: loss on final training batch: 0.1507
Epoch 4: loss on final training batch: 0.1352
Epoch 5: loss on final training batch: 0.1076
Epoch 6: loss on final training batch: 0.3239
Epoch 7: loss on final training batch: 0.3168
Epoch 8: loss on final training batch: 0.2128
Epoch 9: loss on final training batch: 0.3074
Epoch 10: loss on final training batch: 0.3044
Epoch 11: loss on final training batch: 0.2990
Epoch 12: loss on test set: 0.3211
```

Run the code cell below to retrieve the `PyTorch` model's parameters, convert them back to Numpy, and plot them like before.

```
In [33]: W, b, *_ = model.parameters()
W = W.detach().numpy().reshape(-1, 28, 28)
plot_matrix_grid(W)
```

Exercise 2.3 – Train a PyTorch neural network with hidden layers

Using Exercise 2.2 as a starting point, write new code to implement a 784-100-50-10 neural network with `relu` activations just like you did in Exercise 1.8, but now implemented with `PyTorch`.

To do this, you will need to:

- Create a new model object that has more sequential steps to it, including the `Linear` and `ReLU` objects.
- Create a new optimizer object that knows about your new model's parameters.

If you succeed, you should be able to get


```
0.weight torch.Size([100, 784])
0.bias torch.Size([100])
2.weight torch.Size([50, 100])
2.bias torch.Size([50])
4.weight torch.Size([10, 50])
4.bias torch.Size([10])
```

In []:



```
0.weight torch.Size([100, 784])
0.bias torch.Size([100])
2.weight torch.Size([50, 100])
2.bias torch.Size([50])
4.weight torch.Size([10, 50])
4.bias torch.Size([10])
```