

CS444 Assignment 2

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

from kaggle_submission import output_submission_csv
from models.neural_net import NeuralNetwork
from utils.data_process import get_FASHION_data

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots

# For auto-reloading external modules
# See http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```

The autoreload extension is already loaded. To reload it, use:
%reload_ext autoreload

Loading Fashion-MNIST

Now that you have implemented a neural network that passes gradient checks and works on toy data, you will test your network on the Fashion-MNIST dataset.

```
In [ ]: # You can change these numbers for experimentation
# For submission be sure they are set to the default values
TRAIN_IMAGES = 50000
VAL_IMAGES = 10000
TEST_IMAGES = 10000

data = get_FASHION_data(TRAIN_IMAGES, VAL_IMAGES, TEST_IMAGES)
X_train, y_train = data['X_train'], data['y_train']
X_val, y_val = data['X_val'], data['y_val']
X_test, y_test = data['X_test'], data['y_test']

train_loss = None
train_accuracy = None
val_accuracy = None
net_2sgd, net_3sgd, net_2adam, net_3adam = None, None, None, None

sgd2_loss, sgd2_val, adam2_loss, adam2_val = None, None, None, None
```

Train using SGD

To train our network we will use SGD. In addition, we will adjust the learning rate with an exponential learning rate schedule as optimization proceeds; after each epoch, we will reduce the learning rate by multiplying it by a decay rate.

You can try different numbers of layers and other hyperparameters on the Fashion-MNIST dataset below.

```

In [ ]: def sgd_2():

    # Hyperparameters
    input_size = 28 * 28
    num_layers = 2
    hidden_size = 80
    hidden_sizes = [hidden_size] * (num_layers - 1)
    num_classes = 10
    epochs = 50
    batch_size = 50
    learning_rate = 1e-3
    learning_rate_decay = 0.9
    regularization = 0.08

    global train_loss, train_accuracy, val_accuracy, net_2sgd
    # Initialize a new neural network model
    net_2sgd = NeuralNetwork(input_size, hidden_sizes, num_classes, num_layers)

    # Variables to store performance for each epoch
    train_loss = np.zeros(epochs)
    train_accuracy = np.zeros(epochs)
    val_accuracy = np.zeros(epochs)

    # For each epoch...
    t = tqdm(range(epochs))
    for epoch in t:

        # Shuffle the dataset
        perm = np.random.permutation(X_train.shape[0])
        X, y = X_train[perm], y_train[perm]

        # Training
        # For each mini-batch...
        batches = TRAIN_IMAGES // batch_size
        for batch in range(batches):
            # Create a mini-batch of training data and labels
            X_batch = X[batch*batch_size:(batch + 1)*batch_size]
            y_batch = y[batch*batch_size:(batch + 1)*batch_size]

            # Run the forward pass of the model to get a prediction and compute the accuracy
            train_accuracy[epoch] += np.sum(
                np.argmax(net_2sgd.forward(X_batch), axis=1) == y_batch)

            # Run the backward pass of the model to compute the loss, and update the weights
            train_loss[epoch] += net_2sgd.backward(y_batch, regularization)
            net_2sgd.update(learning_rate, opt='SGD')

        # Validation
        # No need to run the backward pass here, just run the forward pass to compute accuracy
        val_accuracy[epoch] += np.sum(np.argmax(net_2sgd.forward(X_val),
                                                axis=1) == y_val) / y_val.shape[0]
        t.set_postfix({"valid": "{:.4f}".format(val_accuracy[epoch]),
                      "train": "{:.4f}".format(train_accuracy[epoch] / (batches * batch_size))})

        # Implement learning rate decay
        learning_rate = learning_rate * learning_rate_decay
        train_accuracy[epoch] = train_accuracy[epoch] / (batches * batch_size)
        train_loss[epoch] = train_loss[epoch] / (batches * batch_size)
    global sgd2_loss, sgd2_val
    sgd2_loss = train_loss.copy()
    sgd2_val = val_accuracy.copy()
sgd_2()

```

100%|██████████| 50/50 [44:29<00:00, 53.38s/it, valid=0.8865, train=0.9296]

```

In [ ]: def sgd_3():
    # Hyperparameters
    input_size = 28 * 28
    num_layers = 3
    hidden_size = 80
    hidden_sizes = [hidden_size] * (num_layers - 1)
    num_classes = 10
    epochs = 40
    batch_size = 50
    learning_rate = 0.0005
    learning_rate_decay = 0.95
    regularization = 0.08

    global train_loss, train_accuracy, val_accuracy, net_3sgd
    # Initialize a new neural network model
    net_3sgd = NeuralNetwork(input_size, hidden_sizes, num_classes, num_layers)

    # Variables to store performance for each epoch

    train_loss = np.zeros(epochs)
    train_accuracy = np.zeros(epochs)
    val_accuracy = np.zeros(epochs)

    # For each epoch...
    t = tqdm(range(epochs))
    for epoch in t:

        # Shuffle the dataset
        perm = np.random.permutation(X_train.shape[0])
        X, y = X_train[perm], y_train[perm]

        # Training
        # For each mini-batch...
        batches = TRAIN_IMAGES // batch_size
        for batch in range(batches):
            # Create a mini-batch of training data and labels
            X_batch = X[batch*batch_size:(batch + 1)*batch_size]
            y_batch = y[batch*batch_size:(batch + 1)*batch_size]

            # Run the forward pass of the model to get a prediction and compute the accuracy
            train_accuracy[epoch] += np.sum(
                np.argmax(net_3sgd.forward(X_batch), axis=1) == y_batch)

            # Run the backward pass of the model to compute the loss, and update the weights
            train_loss[epoch] += net_3sgd.backward(y_batch, regularization)
            net_3sgd.update(learning_rate, opt='SGD')

        # Validation
        # No need to run the backward pass here, just run the forward pass to compute accuracy
        val_accuracy[epoch] += np.sum(np.argmax(net_3sgd.forward(X_val),
            axis=1) == y_val) / y_val.shape[0]
        t.set_postfix({"valid": "{:.4f}".format(val_accuracy[epoch]),
            "train": "{:.4f}".format(train_accuracy[epoch] / (batches * batch_size))})

        # Implement learning rate decay
        learning_rate = learning_rate * learning_rate_decay
        train_accuracy[epoch] = train_accuracy[epoch] / (batches * batch_size)
        train_loss[epoch] = train_loss[epoch] / (batches * batch_size)

sgd_3()

```

100%|██████████| 40/40 [41:57<00:00, 62.94s/it, valid=0.8767, train=0.8961]

Train using Adam

Next we will train the same model using the Adam optimizer. You should take the above code for SGD and modify it to use Adam instead. For implementation details, see the lecture slides. The original paper that introduced Adam is also a good reference, and contains suggestions for default values:

<https://arxiv.org/pdf/1412.6980.pdf>

```

In [ ]: def adam_2():
    # Hyperparameters
    input_size = 28 * 28
    num_layers = 2
    hidden_size = 80
    hidden_sizes = [hidden_size] * (num_layers - 1)
    num_classes = 10
    epochs = 50
    batch_size = 50
    learning_rate = 1e-4
    regularization = 0.1
    learning_rate_decay = 0.95

    global train_loss, train_accuracy, val_accuracy, net_2adam
    # Initialize a new neural network model
    net_2adam = NeuralNetwork(input_size, hidden_sizes, num_classes, num_layers)

    # Variables to store performance for each epoch
    train_loss = np.zeros(epochs)
    train_accuracy = np.zeros(epochs)
    val_accuracy = np.zeros(epochs)

    # For each epoch...
    t = tqdm(range(epochs))
    for epoch in t:

        # Shuffle the dataset
        perm = np.random.permutation(X_train.shape[0])
        X, y = X_train[perm], y_train[perm]

        # Training
        # For each mini-batch...
        batches = TRAIN_IMAGES // batch_size
        for batch in range(batches):
            # Create a mini-batch of training data and labels
            X_batch = X[batch*batch_size:(batch + 1)*batch_size]
            y_batch = y[batch*batch_size:(batch + 1)*batch_size]

            # Run the forward pass of the model to get a prediction and compute the accuracy
            train_accuracy[epoch] += np.sum(
                np.argmax(net_2adam.forward(X_batch), axis=1) == y_batch)

            # Run the backward pass of the model to compute the loss, and update the weights
            train_loss[epoch] += net_2adam.backward(y_batch, regularization)
            net_2adam.update(learning_rate, opt='Adam')

        # Validation
        # No need to run the backward pass here, just run the forward pass to compute accuracy
        val_accuracy[epoch] += np.sum(np.argmax(net_2adam.forward(X_val),
            axis=1) == y_val) / y_val.shape[0]
        t.set_postfix({"valid": "{:.4f}".format(val_accuracy[epoch]),
            "train": "{:.4f}".format(train_accuracy[epoch] / (batches * batch_size))})

        # Implement learning rate decay
        train_accuracy[epoch] = train_accuracy[epoch] / (batches * batch_size)
        train_loss[epoch] = train_loss[epoch] / (batches * batch_size)
        learning_rate = learning_rate * learning_rate_decay
    global adam2_loss, adam2_val
    adam2_loss = train_loss.copy()
    adam2_val = val_accuracy.copy()
adam_2()

```

100%|██████████| 50/50 [40:58<00:00, 49.17s/it, valid=0.8874, train=0.9412]

```

In [ ]: def adam_3():
    # Hyperparameters
    input_size = 28 * 28
    num_layers = 3
    hidden_size = 60
    hidden_sizes = [hidden_size] * (num_layers - 1)
    num_classes = 10
    epochs = 50
    batch_size = 50
    learning_rate = 0.005
    regularization = 0.03
    learning_rate_decay = 0.95

    global train_loss, train_accuracy, val_accuracy, net_3adam
    # Initialize a new neural network model
    net_3adam = NeuralNetwork(input_size, hidden_sizes, num_classes, num_layers)

    # Variables to store performance for each epoch
    train_loss = np.zeros(epochs)
    train_accuracy = np.zeros(epochs)
    val_accuracy = np.zeros(epochs)

    # For each epoch...
    t = tqdm(range(epochs))
    for epoch in t:

        # Shuffle the dataset
        perm = np.random.permutation(X_train.shape[0])
        X, y = X_train[perm], y_train[perm]

        # Training
        # For each mini-batch...
        batches = TRAIN_IMAGES // batch_size
        for batch in range(batches):
            # Create a mini-batch of training data and labels
            X_batch = X[batch*batch_size:(batch + 1)*batch_size]
            y_batch = y[batch*batch_size:(batch + 1)*batch_size]

            # Run the forward pass of the model to get a prediction and compute the accuracy
            train_accuracy[epoch] += np.sum(
                np.argmax(net_3adam.forward(X_batch), axis=1) == y_batch)

            # Run the backward pass of the model to compute the loss, and update the weights
            train_loss[epoch] += net_3adam.backward(y_batch, regularization)
            net_3adam.update(learning_rate, 0.9, 0.99, opt='SGD')

        # Validation
        # No need to run the backward pass here, just run the forward pass to compute accuracy
        val_accuracy[epoch] += np.sum(np.argmax(net_3adam.forward(X_val),
            axis=1) == y_val) / y_val.shape[0]
        t.set_postfix({"valid": "{:.4f}".format(val_accuracy[epoch]),
            "train": "{:.4f}".format(train_accuracy[epoch] / (batches * batch_size))})

        # Implement learning rate decay
        train_accuracy[epoch] = train_accuracy[epoch] / (batches * batch_size)
        train_loss[epoch] = train_loss[epoch] / (batches * batch_size)
        learning_rate = learning_rate * learning_rate_decay

    # adam_3()

```

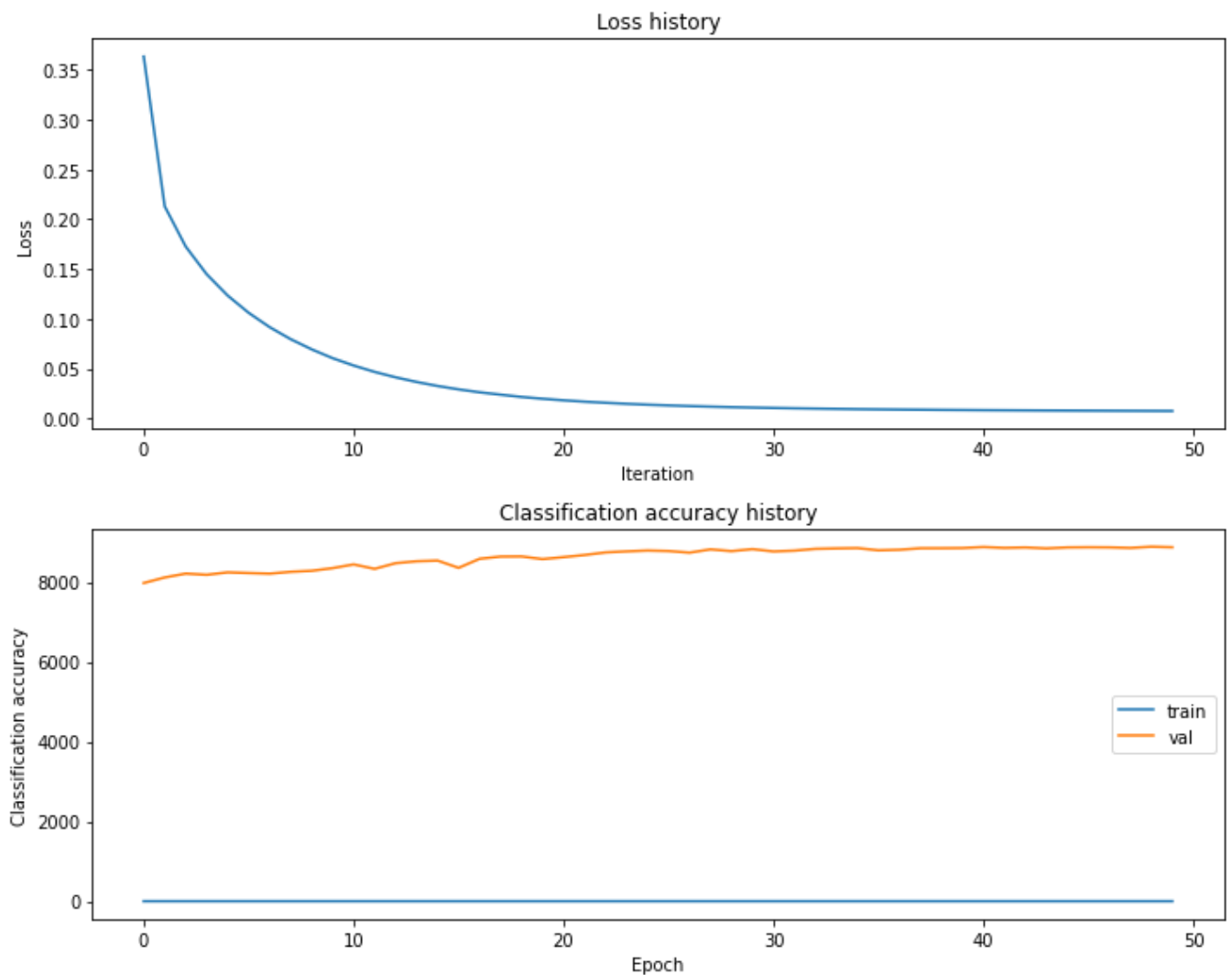
Graph loss and train/val accuracies

Examining the loss graph along with the train and val accuracy graphs should help you gain some intuition for the hyperparameters you should try in the hyperparameter tuning below. It should also help with debugging any issues you might have with your network.

```
In [ ]: # Plot the loss function and train / validation accuracies
plt.subplot(2, 1, 1)
plt.plot(train_loss)
plt.title('Loss history')
plt.xlabel('Iteration')
plt.ylabel('Loss')

plt.subplot(2, 1, 2)
plt.plot(train_accuracy, label='train')
plt.plot(val_accuracy, label='val')
plt.title('Classification accuracy history')
plt.xlabel('Epoch')
plt.ylabel('Classification accuracy')
plt.legend()

plt.tight_layout()
plt.show()
```



Hyperparameter tuning

Once you have successfully trained a network you can tune your hyperparameters to increase your accuracy.

Based on the graphs of the loss function above you should be able to develop some intuition about what hyperparameter adjustments may be necessary. A very noisy loss implies that the learning rate might be too high, while a linearly decreasing loss would suggest that the learning rate may be too low. A large gap between training and validation accuracy would suggest overfitting due to a large model without much regularization. No gap between training and validation accuracy would indicate low model capacity.

You will compare networks of two and three layers using the different optimization methods you implemented.

The different hyperparameters you can experiment with are:

- **Batch size:** We recommend you leave this at 200 initially which is the batch size we used.
- **Number of iterations:** You can gain an intuition for how many iterations to run by checking when the validation accuracy plateaus in your train/val accuracy graph.
- **Initialization** Weight initialization is very important for neural networks. We used the initialization $W = \text{np.random.randn}(n) / \text{sqrt}(n)$ where n is the input dimension for layer corresponding to W . We recommend you stick with the given initializations, but you may explore modifying these. Typical initialization practices: <http://cs231n.github.io/neural-networks-2/#init>
- **Learning rate:** Generally from around $1e-4$ to $1e-1$ is a good range to explore according to our implementation.
- **Learning rate decay:** We recommend a 0.95 decay to start.
- **Hidden layer size:** You should explore up to around 120 units per layer. For three-layer network, we fixed the two hidden layers to be the same size when obtaining the target numbers. However, you may experiment with having different size hidden layers.
- **Regularization coefficient:** We recommend trying values in the range 0 to 0.1.

Hints:

- After getting a sense of the parameters by trying a few values yourself, you will likely want to write a few for-loops to traverse over a set of hyperparameters.
- If you find that your train loss is decreasing, but your train and val accuracy start to decrease rather than increase, your model likely started minimizing the regularization term. To prevent this you will need to decrease the regularization coefficient.

Run on the test set

When you are done experimenting, you should evaluate your final trained networks on the test set.

```
In [ ]: best_2layer_sgd_prediction = np.argmax(net_2sgd.forward(X_test), axis=1)
best_3layer_sgd_prediction = np.argmax(net_3sgd.forward(X_test), axis=1)
best_2layer_adam_prediction = np.argmax(net_2adam.forward(X_test), axis=1)
best_3layer_adam_prediction = np.argmax(net_3adam.forward(X_test), axis=1)
```


Kaggle output

Once you are satisfied with your solution and test accuracy, output a file to submit your test set predictions to the Kaggle for Assignment 2 Neural Network. Use the following code to do so:

```
In [ ]: output_submission_csv('./nn_2layer_sgd_submission.csv', best_2layer_sgd_prediction)
output_submission_csv('./nn_3layer_sgd_submission.csv', best_3layer_sgd_prediction)
output_submission_csv('./nn_2layer_adam_submission.csv', best_2layer_adam_prediction)
output_submission_csv('./nn_3layer_adam_submission.csv', best_3layer_adam_prediction)
```

Compare SGD and Adam

Create graphs to compare training loss and validation accuracy between SGD and Adam. The code is similar to the above code, but instead of comparing train and validation, we are comparing SGD and Adam.

```
In [ ]: # TODO: implement me

# Plot the loss function and train / validation accuracies
plt.subplot(2, 1, 1)
plt.plot(adam2_loss, label="adam")
plt.plot(sgd2_loss, label="sgd")
plt.title('Loss history')
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.legend()

plt.subplot(2, 1, 2)
plt.plot(adam2_val, label='adam')
plt.plot(sgd2_val, label='sgd')
plt.title('Classification accuracy history')
plt.xlabel('Epoch')
plt.ylabel('Classification accuracy')
plt.legend()

plt.tight_layout()
plt.show()
```

