

Homework 2

Start date: 29th January 2018

Due date: 11th February 2018, 11:55 pm

How to Submit

When you have completed the exercises and everything has finished running, click on 'File' in the menu-bar and then 'Download .ipynb'. This file must be submitted to Moodle named as **studentnumber_DL_hw2.ipynb** before the deadline above.

Also send a **sharable link** to the notebook at the following email: ucl.coursework.submit@gmail.com. You can also make it sharable via link to everyone, up to you.

IMPORTANT

Please make sure your submission includes **all results/plots/tables** required for grading. We will not re-run your code.

The Data

Handwritten Digit Recognition Dataset (MNIST)

In this assignment we will be using the [MNIST digit dataset \(https://yann.lecun.com/exdb/mnist/\)](https://yann.lecun.com/exdb/mnist/).

The dataset contains images of hand-written digits (0 – 9), and the corresponding labels.

The images have a resolution of 28×28 pixels.

The MNIST Dataset in TensorFlow

You can use the tensorflow build-in functionality to download and import the dataset into python (see *Setup* section below). But note that this will be the only usage of TensorFlow in this assignment.

The Assignment

Objectives

This assignment will be mirroring the first assignment (DL_hw1), but this time you are **not allowed to use any of the Tensorflow functionality for specifying nor optimizing** your neural network models. You will now use your **own implementations** of different neural network models (labelled Model 1-4, and described in the corresponding sections of the Colab).

As before, you will train these models to classify hand written digits from the Mnist dataset.

Keep in mind, the purpose of this exercise is to implement and optimize your own neural networks architectures without the toolbox/library tailored to do so. This also means, in order to train and evaluate your models, you will need to implement your own optimization procedure. You are to use the same cross-entropy loss as before and your own implementation of SGD.

Additional instruction:

Do not use any other libraries than the ones provided in the imports cell. You should be able to do everything via *numpy* (especially for the convolutional layer, rely on the in-built matrix/tensor multiplication that numpy offers).

There are a few questions at the end of the colab. **Before doing any coding, please take look at Question 1** -- this should help you with the implementations, especially the optimization part.

Variable Initialization

Initialize the variables containing the parameters using [Xavier initialization \(http://proceedings.mlr.press/v9/glorot10a.html\)](http://proceedings.mlr.press/v9/glorot10a.html).

Hyper-parameters

For each of these models you will be requested to run experiments with different hyper-parameters.

More specifically, you will be requested to try 3 sets of hyper-parameters per model, and report the resulting model accuracy.

Each combination of hyper-parameter will specify how to set each of the following:

- **num_epochs:** Number of iterations through the training section of the dataset [*a positive integer*].
- **learning_rate:** Learning rate used by the gradient descent optimizer [*a scalar between 0 and 1*]

In all experiments use a *batch_size* of 100.

Loss function

All models, should be trained as to minimize the **cross-entropy loss** function:

$$\text{loss} = - \sum_{i=1}^N \log p(y_i | x_i, \theta) = - \sum_{i=1}^N \log \left(\underbrace{\frac{\exp(z_i[y_i])}{\sum_{c=1}^{10} \exp(z_i[c])}}_{\text{softmax output}} \right) = \sum_{i=1}^N \left(-z_i[y_i] + \log \left(\sum_{c=1}^{10} \exp(z_i[c]) \right) \right)$$

where $z \in \mathbb{R}^{10}$ is the input to the softmax layer and $z[c]$ denotes the c -th entry of vector z . And i is a index for the dataset $\{(x_i, y_i)\}_{i=1}^N$.

Optimization

Use **stochastic gradient descent (SGD)** for optimizing the loss function. Sum over the batch.

Training and Evaluation

The tensorflow built-in functionality for downloading and importing the dataset into python returns a Datasets object.

This object will have three attributes:

- train
- validation
- test

Use only the **train** data in order to optimize the model.

Use `datasets.train.next_batch(100)` in order to sample mini-batches of data.

Every 20000 training samples (i.e. every 200 updates to the model), interrupt training and measure the accuracy of the model, each time evaluate the accuracy of the model both on 20% of the **train** set and on the entire **test** set.

Reporting

For each model *i*, you will collect the learning curves associated to each combination of hyper-parameters.

Use the utility function `plot_learning_curves` to plot these learning curves,

and the utility function `plot_summary_table` to generate a summary table of results.

For each run collect the train and test curves in a tuple, together with the hyper-parameters.

```
experiments_task_i = [  
    (num_epochs_1, learning_rate_1), train_accuracy_1, test_accuracy_1),  
    (num_epochs_2, learning_rate_2), train_accuracy_2, test_accuracy_2),  
    (num_epochs_3, learning_rate_3), train_accuracy_3, test_accuracy_3)]
```

Hint:

Remind yourselves of the chain rule and read through the lecture notes on back-propagation (computing the gradients by recursively applying the chain rule). This is a general procedure that applies to all model architectures you will have to code in the following steps. Thus, you are strongly encourage to implement an optimization procedure that generalizes and can be re-used to train all your models. Recall the only things that you will need for each layer are:

- (i) the gradients of layer with respect to its input
- (ii) the gradients with respect to its parameters, if any.

(See Question 1).

Also from the previous assignment, you should have a good idea of what to expect, both terms of behavior and relative performance. (To keep everything comparable, we kept all the hyperparameters and reporting the same).

Imports and utility functions (do not modify!)

In [0]:

```
1 # Import useful libraries.
2 import matplotlib.pyplot as plt
3 from tensorflow.examples.tutorials.mnist import input_data
4 import numpy as np
5
6 # Global variables.
7 log_period_samples = 20000
8 batch_size = 100
9
10 # Import dataset with one-hot encoding of the class labels.
11 def get_data():
12     return input_data.read_data_sets("MNIST_data/", one_hot=True)
13
14 # Placeholders to feed train and test data into the graph.
15 # Since batch dimension is 'None', we can reuse them both for train and eval.
16 def get_placeholders():
17     x = tf.placeholder(tf.float32, [None, 784])
18     y_ = tf.placeholder(tf.float32, [None, 10])
19     return x, y_
20
21 # Plot learning curves of experiments
22 def plot_learning_curves(experiment_data):
23     # Generate figure.
24     fig, axes = plt.subplots(3, 4, figsize=(22,12))
25     st = fig.suptitle(
26         "Learning Curves for all Tasks and Hyper-parameter settings",
27         fontsize="x-large")
28     # Plot all learning curves.
29     for i, results in enumerate(experiment_data):
30         for j, (setting, train_accuracy, test_accuracy) in enumerate(results):
31             # Plot.
32             xs = [x * log_period_samples for x in range(1, len(train_accuracy)+1)]
33             axes[j, i].plot(xs, train_accuracy, label='train_accuracy')
34             axes[j, i].plot(xs, test_accuracy, label='test_accuracy')
35             # Prettyfy individual plots.
36             axes[j, i].ticklabel_format(style='sci', axis='x', scilimits=(0,0))
37             axes[j, i].set_xlabel('Number of samples processed')
38             axes[j, i].set_ylabel('Epochs: {}, Learning rate: {}'.format(*setting))
39             axes[j, i].set_title('Task {}'.format(i + 1))
40             axes[j, i].legend()
41     # Prettyfy overall figure.
42     plt.tight_layout()
43     st.set_y(0.95)
44     fig.subplots_adjust(top=0.91)
45     plt.show()
46
47 # Generate summary table of results.
48 def plot_summary_table(experiment_data):
49     # Fill Data.
50     cell_text = []
51     rows = []
52     columns = ['Setting 1', 'Setting 2', 'Setting 3']
53     for i, results in enumerate(experiment_data):
54         rows.append('Model {}'.format(i + 1))
55         cell_text.append([])
56         for j, (setting, train_accuracy, test_accuracy) in enumerate(results):
57             cell_text[i].append(test_accuracy[-1])
58     # Generate Table.
59     fig=plt.figure(frameon=False)
60     ax = plt.gca()
61     the_table = ax.table(
62         cellText=cell_text,
63         rowLabels=rows,
64         colLabels=columns,
65         loc='center')
66     the_table.scale(1, 4)
67     # Prettyfy.
68     ax.patch.set_facecolor('None')
69     ax.xaxis.set_visible(False)
70     ax.yaxis.set_visible(False)
```

Model

Train a neural network model consisting of 1 linear layer, followed by a softmax:

(input → linear layer → softmax → class probabilities)

Hyper-parameters

Train the model with three different hyper-parameter settings:

- *num_epochs*=5, *learning_rate*=0.0001
- *num_epochs*=5, *learning_rate*=0.005
- *num_epochs*=5, *learning_rate*=0.1

```
In [0]: 1 # CAREFUL: Running this CL resets the experiments_task1 dictionary where results should be stored.  
2 # Store results of runs with different configurations in a dictionary.  
3 # Use a tuple (num_epochs, learning_rate) as keys, and a tuple (training_accuracy, testing_accuracy)  
4 experiments_task1 = {}  
5 settings = [(5, 0.0001), (5, 0.005), (5, 0.1)]
```

```

In [16]: 1 print('Training Model 1')
2
3 # Train Model 1 with the different hyper-parameter settings.
4 for (num_epochs, learning_rate) in settings:
5
6     # Reset graph, recreate placeholders and dataset.
7     mnist = get_data()
8     eval_mnist = get_data()
9
10    #####
11    # Define model, loss, update and evaluation metric. #
12    #####
13
14    # define the components of the model including loss
15    graph = [linear(name='linear1',\
16                    shape=(784, 10),\
17                    initializer=xavier),\
18              cross_entropy_softmax_logits(name='loss1')]
19
20    # define the model
21    model = Model(graph)
22
23    # define the optimizer
24    optimizer = SGD(name='sgd',\
25                   learning_rate=learning_rate,\
26                   model=model)
27
28    # since evaluation metric is a function, we just need to call it directly
29
30    # Train.
31    i, train_accuracy, test_accuracy = 0, [], []
32    log_period_updates = int(log_period_samples / batch_size)
33    while mnist.train.epochs_completed < num_epochs:
34
35        # Update.
36        i += 1
37        batch_xs, batch_ys = mnist.train.next_batch(batch_size)
38
39        #####
40        # Training step #
41        #####
42        optimizer.learning(batch_xs, batch_ys)
43
44        # Periodically evaluate.
45        if i % log_period_updates == 0:
46
47            #####
48            # Compute and store train accuracy. #
49            #####
50            train_data_size = eval_mnist.train.images.shape[0]
51            eval_train_xs = eval_mnist.train.images[:int(train_data_size*0.2)]
52            eval_train_ys = eval_mnist.train.labels[:int(train_data_size*0.2)]
53            _, train_predicts = model.forward_pass(eval_train_xs)
54            train_predicts = softmax(train_predicts)
55            acc_train = accuracy_evaluation(train_predicts, eval_train_ys)
56            print ('accuracy of train data: ', acc_train)
57            train_accuracy.append(acc_train)
58            #####
59            # Compute and store test accuracy. #
60            #####
61            eval_test_xs = eval_mnist.test.images
62            eval_test_ys = eval_mnist.test.labels
63            _, test_predicts = model.forward_pass(eval_test_xs)
64            test_predicts = softmax(test_predicts)
65            acc_test = accuracy_evaluation(test_predicts, eval_test_ys)
66            print ('accuracy of test data: ', acc_test)
67            test_accuracy.append(acc_test)
68        del graph
69        del model
70        del optimizer
71        experiments_task1.append(
72            (num_epochs, learning_rate), train_accuracy, test_accuracy))

```

```

Training Model 1
Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Module : linear1 is constructed.
Module : loss1 is constructed.
accuracy of train data:  0.7436363636363637
accuracy of test data:  0.7551
accuracy of train data:  0.8075454545454546
accuracy of test data:  0.8187
accuracy of train data:  0.8293636363636364
accuracy of test data:  0.8444
accuracy of train data:  0.845
accuracy of test data:  0.8599
accuracy of train data:  0.8636363636363637

```

Model 2 (15 pts)

1 hidden layer (32 units) with a ReLU non-linearity, followed by a softmax.

(input → non-linear layer → linear layer → softmax → class probabilities)

Hyper-parameters

Train the model with three different hyper-parameter settings:

- *num_epochs*=15, *learning_rate*=0.0001
- *num_epochs*=15, *learning_rate*=0.005
- *num_epochs*=15, *learning_rate*=0.1

```
In [0]: 1 # CAREFUL: Running this CL resets the experiments_task1 dictionary where results should be stored.  
2 # Store results of runs with different configurations in a dictionary.  
3 # Use a tuple (num_epochs, learning_rate) as keys, and a tuple (training_accuracy, testing_accuracy)  
4 experiments_task2 = []  
5 settings = [(15, 0.0001), (15, 0.005), (15, 0.1)]
```

In [18]:

```
1 print('Training Model 2')
2
3 # Train Model 1 with the different hyper-parameter settings.
4 for (num_epochs, learning_rate) in settings:
5
6     mnist = get_data() # use for training.
7     eval_mnist = get_data() # use for evaluation.
8
9     #####
10    # Define model, loss, update and evaluation metric. #
11    #####
12
13    # define the components of the model including loss
14    graph = [linear(name='linear1',\
15                    shape=(784, 32),\
16                    initializer=xavier),\
17              relu(name='relu1'),\
18              linear(name='linear2',\
19                    shape=(32, 10),\
20                    initializer=xavier),\
21              cross_entropy_softmax_logits(name='loss1')]
22
23    # define the model
24    model = Model(graph)
25
26    # define the optimizer
27    optimizer = SGD(name='sgd',\
28                   learning_rate=learning_rate,\
29                   model=model)
30
31    # Train.
32    i, train_accuracy, test_accuracy = 0, [], []
33    log_period_updates = int(log_period_samples / batch_size)
34    while mnist.train.epochs_completed < num_epochs:
35
36        # Update.
37        i += 1
38        batch_xs, batch_ys = mnist.train.next_batch(batch_size)
39
40        #####
41        # Training step #
42        #####
43        optimizer.learning(batch_xs, batch_ys)
44
45        # Periodically evaluate.
46        if i % log_period_updates == 0:
47
48            #####
49            # Compute and store train accuracy. #
50            #####
51            train_data_size = eval_mnist.train.images.shape[0]
52            eval_train_xs = eval_mnist.train.images[:int(train_data_size*0.2)]
53            eval_train_ys = eval_mnist.train.labels[:int(train_data_size*0.2)]
54            _, train_predicts = model.forward_pass(eval_train_xs)
55            train_predicts = softmax(train_predicts)
56            acc_train = accuracy_evaluation(train_predicts, eval_train_ys)
57            print ('accuracy of train data: ', acc_train)
58            train_accuracy.append(acc_train)
59            #####
60            # Compute and store test accuracy. #
61            #####
62            eval_test_xs = eval_mnist.test.images
63            eval_test_ys = eval_mnist.test.labels
64            _, test_predicts = model.forward_pass(eval_test_xs)
65            test_predicts = softmax(test_predicts)
66            acc_test = accuracy_evaluation(test_predicts, eval_test_ys)
67            print ('accuracy of test data: ', acc_test)
68            test_accuracy.append(acc_test)
69        del graph
70        del model
71        del optimizer
72        experiments_task2.append(
73            ((num_epochs, learning_rate), train_accuracy, test_accuracy))
```

```
accuracy of train data: 0.9054545454545454
accuracy of test data: 0.9136
accuracy of train data: 0.9074545454545454
accuracy of test data: 0.9156
accuracy of train data: 0.9068181818181819
accuracy of test data: 0.9163
accuracy of train data: 0.9082727272727272
accuracy of test data: 0.917
accuracy of train data: 0.9094545454545454
accuracy of test data: 0.918
accuracy of train data: 0.9099090909090909
accuracy of test data: 0.9189
accuracy of train data: 0.9118181818181819
accuracy of test data: 0.9197
accuracy of train data: 0.9122727272727272
accuracy of test data: 0.9206
accuracy of train data: 0.9114545454545454
accuracy of test data: 0.9199
accuracy of train data: 0.9134545454545454
accuracy of test data: 0.9216
```

Model 3 (15 pts)

2 hidden layers (32 units) each, with ReLU non-linearity, followed by a softmax.

(input → non-linear layer → non-linear layer → linear layer → softmax → class probabilities)

Hyper-parameters

Train the model with three different hyper-parameter settings:

- *num_epochs*=5, *learning_rate*=0.003
- *num_epochs*=40, *learning_rate*=0.003
- *num_epochs*=40, *learning_rate*=0.05

```
In [0]: 1 # CAREFUL: Running this CL resets the experiments_task1 dictionary where results should be stored.
        2 # Store results of runs with different configurations in a dictionary.
        3 # Use a tuple (num_epochs, learning_rate) as keys, and a tuple (training_accuracy, testing_accuracy)
        4 experiments_task3 = []
        5 settings = [(5, 0.003), (40, 0.003), (40, 0.05)]
```


In [20]:

```
1 print('Training Model 3')
2
3 # Train Model 1 with the different hyper-parameter settings.
4 for (num_epochs, learning_rate) in settings:
5
6     mnist = get_data() # use for training.
7     eval_mnist = get_data() # use for evaluation.
8
9     #####
10    # Define model, loss, update and evaluation metric. #
11    #####
12
13    # define the components of the model including loss
14    graph = [linear(name='linear1',\
15                    shape=(784, 32),\
16                    initializer=xavier),\
17              relu(name='relu1'),\
18              linear(name='linear2',\
19                    shape=(32, 32),\
20                    initializer=xavier),\
21              relu(name='relu2'),\
22              linear(name='linear3',\
23                    shape=(32, 10),\
24                    initializer=xavier),\
25              cross_entropy_softmax_logits(name='loss1')]
26
27    # define the model
28    model = Model(graph)
29
30    # define the optimizer
31    optimizer = SGD(name='sgd',\
32                   learning_rate=learning_rate,\
33                   model=model)
34
35    # Train.
36    i, train_accuracy, test_accuracy = 0, [], []
37    log_period_updates = int(log_period_samples / batch_size)
38    while mnist.train.epochs_completed < num_epochs:
39
40        # Update.
41        i += 1
42        batch_xs, batch_ys = mnist.train.next_batch(batch_size)
43
44        #####
45        # Training step #
46        #####
47        optimizer.learning(batch_xs, batch_ys)
48
49        # Periodically evaluate.
50        if i % log_period_updates == 0:
51
52            #####
53            # Compute and store train accuracy. #
54            #####
55            train_data_size = eval_mnist.train.images.shape[0]
56            eval_train_xs = eval_mnist.train.images[:int(train_data_size*0.2)]
57            eval_train_ys = eval_mnist.train.labels[:int(train_data_size*0.2)]
58            _, train_predicts = model.forward_pass(eval_train_xs)
59            train_predicts = softmax(train_predicts)
60            acc_train = accuracy_evaluation(train_predicts, eval_train_ys)
61            print ('accuracy of train data: ', acc_train)
62            train_accuracy.append(acc_train)
63            #####
64            # Compute and store test accuracy. #
65            #####
66            eval_test_xs = eval_mnist.test.images
67            eval_test_ys = eval_mnist.test.labels
68            _, test_predicts = model.forward_pass(eval_test_xs)
69            test_predicts = softmax(test_predicts)
70            acc_test = accuracy_evaluation(test_predicts, eval_test_ys)
71            print ('accuracy of test data: ', acc_test)
72            test_accuracy.append(acc_test)
73        del graph
74        del model
75        del optimizer
76        experiments_task3.append(
77            (num_epochs, learning_rate), train_accuracy, test_accuracy))
```

Training Model 3
Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Module : linear1 is constructed.
Module : relu1 is constructed.
Module : linear2 is constructed.
Module : relu2 is constructed.
Module : linear3 is constructed.
Module : loss1 is constructed.
accuracy of train data: 0.8958181818181818
accuracy of test data: 0.9022
accuracy of train data: 0.931
accuracy of test data: 0.932
accuracy of train data: 0.9322222222222222

Model 4 (20 pts)

Model

3 layer convolutional model (2 convolutional layers followed by max pooling) + 1 non-linear layer (32 units), followed by softmax.

(input(28x28) → conv(3x3x8) + maxpool(2x2) → conv(3x3x8) + maxpool(2x2) → flatten → non-linear → linear layer → softmax → class probabilities)

- Use *padding* = 'SAME' for both the convolution and the max pooling layers.
- Employ plain convolution (no stride) and for max pooling operations use 2x2 sliding windows, with no overlapping pixels (note: this operation will down-sample the input image by 2x2).

Hyper-parameters

Train the model with three different hyper-parameter settings:

- num_epochs=5, learning_rate=0.01
- num_epochs=10, learning_rate=0.001
- num_epochs=20, learning_rate=0.001

```
In [0]: 1 # CAREFUL: Running this CL resets the experiments_task1 dictionary where results should be stored.
        2 # Store results of runs with different configurations in a dictionary.
        3 # Use a tuple (num_epochs, learning_rate) as keys, and a tuple (training_accuracy, testing_accuracy)
        4 experiments_task4 = {}
        5 settings = [(5, 0.01), (10, 0.001), (20, 0.001)]
```

In [22]:

```
1 print('Training Model 4')
2
3 # Train Model 1 with the different hyper-parameter settings.
4 for (num_epochs, learning_rate) in settings:
5
6     mnist = get_data() # use for training.
7     eval_mnist = get_data() # use for evaluation.
8
9     #####
10    # Define model, loss, update and evaluation metric. #
11    #####
12
13    # define the components of the model including loss
14    graph = [convolution(name='conv1',\
15                        shape=(3,3,1,8),\
16                        initializer=xavier,\
17                        padding='SAME',\
18                        stride=(1,1)),\
19              maxpooling(name='maxpool1',\
20                        shape=(2,2),\
21                        padding='SAME',\
22                        stride=(2,2)),\
23              convolution(name='conv2',\
24                        shape=(3,3,8,8),\
25                        initializer=xavier,\
26                        padding='SAME',\
27                        stride=(1,1)),\
28              maxpooling(name='maxpool2',\
29                        shape=(2,2),\
30                        padding='SAME',\
31                        stride=(2,2)),\
32              flatten(name='flatten1',\
33                      linear(name='linear1',\
34                              shape=(7*7*8, 32),\
35                              initializer=xavier),\
36                      relu(name='relu1',\
37                          linear(name='linear2',\
38                                  shape=(32, 10),\
39                                  initializer=xavier),\
40                                  cross_entropy_softmax_logits(name='loss1'))])
41
42    # define the model
43    model = Model(graph)
44
45    # define the optimizer
46    optimizer = SGD(name='sgd',\
47                   learning_rate=learning_rate,\
48                   model=model)
49
50    # Train.
51    i, train_accuracy, test_accuracy = 0, [], []
52    log_period_updates = int(log_period_samples / batch_size)
53    while mnist.train.epochs_completed < num_epochs:
54
55        # Update.
56        i += 1
57        batch_xs, batch_ys = mnist.train.next_batch(batch_size)
58
59        #####
60        # Training step #
61        #####
62        batch_xs = batch_xs.reshape(-1, 28, 28, 1)
63        optimizer.learning(batch_xs, batch_ys)
64
65        # Periodically evaluate.
66        if i % log_period_updates == 0:
67
68            #####
69            # Compute and store train accuracy. #
70            #####
71            train_data_size = eval_mnist.train.images.shape[0]
72            eval_train_xs = eval_mnist.train.images[:int(train_data_size*0.2)]
73            eval_train_xs = eval_train_xs.reshape(-1, 28, 28, 1)
74            eval_train_ys = eval_mnist.train.labels[:int(train_data_size*0.2)]
75            _, train_predicts = model.forward_pass(eval_train_xs)
76            train_predicts = softmax(train_predicts)
77            acc_train = accuracy_evaluation(train_predicts, eval_train_ys)
78            print ('accuracy of train data: ', acc_train)
79            train_accuracy.append(acc_train)
80
81            #####
82            # Compute and store test accuracy. #
83            #####
84            eval_test_xs = eval_mnist.test.images
85            eval_test_ys = eval_mnist.test.labels
86            eval_test_xs = eval_test_xs.reshape(-1, 28, 28, 1)
87            _, test_predicts = model.forward_pass(eval_test_xs)
88            test_predicts = softmax(test_predicts)
89            acc_test = accuracy_evaluation(test_predicts, eval_test_ys)
90            print ('accuracy of test data: ', acc_test)
91            test_accuracy.append(acc_test)
92
93    del graph
94    del model
95    del optimizer
96    experiments_task4.append(
97        (num_epochs, learning_rate), train_accuracy, test_accuracy))
```

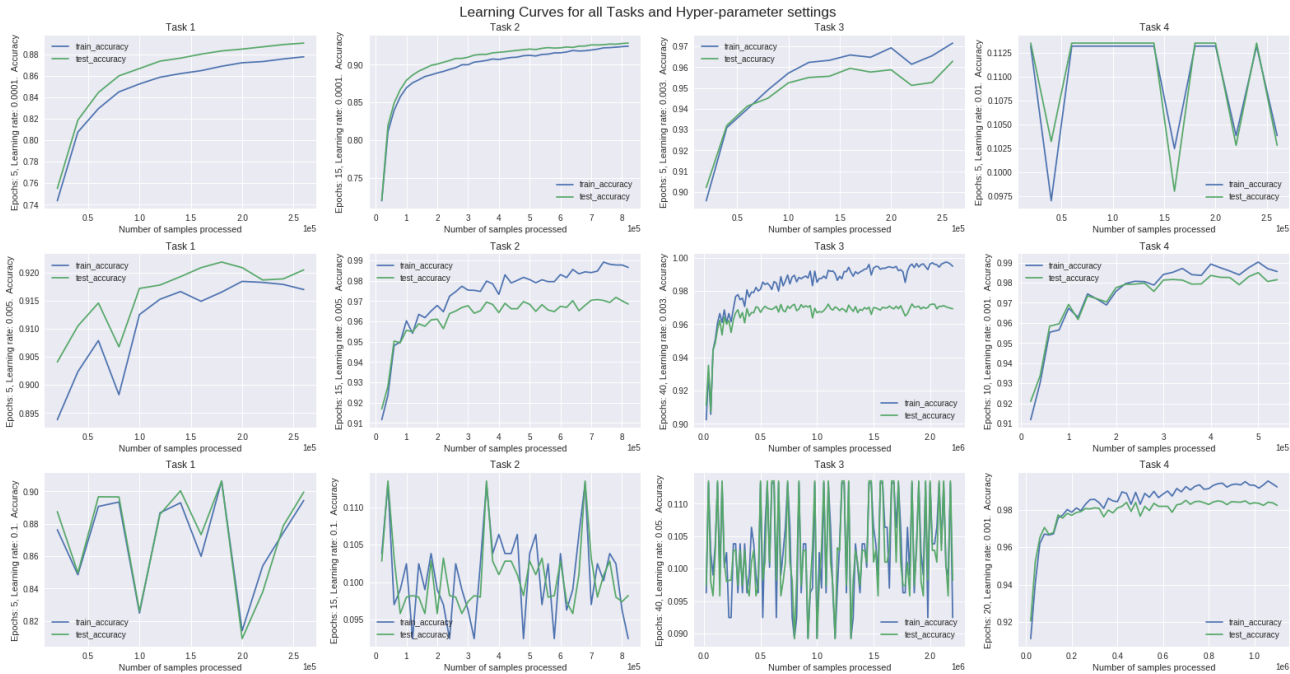
Training Model 4

Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz

```
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Module : conv1 is constructed.
Module : maxpool1 is constructed.
Module : conv2 is constructed.
Module : maxpool2 is constructed.
Module : flatten1 is constructed.
Module : linear1 is constructed.
Module : relu1 is constructed.
Module : linear2 is constructed.
Module : loss1 is constructed.
accuracy of train data: 0.11318181818181818
----- accuracy of test data: 0.1135
```

Results

```
In [23]: 1 plot_learning_curves([experiments_task1, experiments_task2, experiments_task3, experiments_task4])
```



```
In [24]: 1 plot_summary_table([experiments_task1, experiments_task2, experiments_task3, experiments_task4])
```

	Setting 1	Setting 2	Setting 3
Model 1	0.8904	0.9205	0.8996
Model 2	0.9285	0.9685	0.0982
Model 3	0.9629	0.9693	0.0982
Model 4	0.1028	0.9815	0.9824

Questions

###Q1 (32 pts): Compute the following derivatives Show all intermediate steps in the derivation (in markdown below). Provide the final results in vector/matrix/tensor form whenever appropriate.

1. [5 pts] Give the cross-entropy loss above, compute the derivative of the loss function with respect to the scores z (the input to the softmax layer).

$\frac{\partial loss}{\partial z} = ?$

2. [12 pts] Consider the first model (M1: linear + softmax). Compute the derivative of the loss with respect to

- the input x

$\frac{\partial loss}{\partial x} = ?$

- the parameters of the linear layer: weights W and bias b

$\frac{\partial loss}{\partial W} = ?$

$\frac{\partial loss}{\partial b} = ?$

3. [10 pts] Compute the derivative of a convolution layer wrt. to its parameters W and wrt. to its input (4-dim tensor). Assume a filter of size $H \times W \times D$, and stride 1.

$$\frac{\partial \text{loss}}{\partial W} = ?$$

A1: (Your answer here)

1.

$$\text{loss} = - \sum_{i=1}^N \log p(y_i | x_i, \theta) = - \sum_{i=1}^N \log \left(\underbrace{\frac{\exp(z_i[y_i])}{\sum_{c=1}^{10} \exp(z_i[c])}}_{\text{softmax output}} \right) = \sum_{i=1}^N \left(-z_i[y_i] + \log \left(\sum_{c=1}^{10} \exp(z_i[c]) \right) \right)$$

where $z \in \mathbb{R}^{10}$ is the input to the softmax layer and $z[c]$ denotes the c -th entry of vector z . And i is a index for the dataset $\{(x_i, y_i)\}_{i=1}^N$.

$$\frac{\partial \text{loss}}{\partial z_i} = \begin{cases} -1 + \frac{\exp(z_i[k])}{\sum_{c=1}^{10} \exp(z_i[c])} & \text{if } k = y \\ \frac{\exp(z_i[k])}{\sum_{c=1}^{10} \exp(z_i[c])} & \text{if } k \neq y \end{cases}$$

To reform it more compact, we can define a $\delta_{x,y}$ function such that

$$\delta_{x,y} = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$$

So, we now can write the derivative as

$$\frac{\partial \text{loss}}{\partial z_i} = \left[-\delta_{1,y} + \frac{\exp(z_i[1])}{\sum_{c=1}^{10} \exp(z_i[c])}, -\delta_{2,y} + \frac{\exp(z_i[2])}{\sum_{c=1}^{10} \exp(z_i[c])}, \dots, -\delta_{9,y} + \frac{\exp(z_i[9])}{\sum_{c=1}^{10} \exp(z_i[c])}, -\delta_{10,y} + \frac{\exp(z_i[10])}{\sum_{c=1}^{10} \exp(z_i[c])} \right] \in \mathbb{R}^{1 \times 10}$$

So, we can get

$$\frac{\partial \text{loss}}{\partial z} = \begin{bmatrix} \frac{\partial \text{loss}}{\partial z_1} \\ \frac{\partial \text{loss}}{\partial z_2} \\ \dots \\ \dots \\ \dots \\ \frac{\partial \text{loss}}{\partial z_N} \end{bmatrix} \in \mathbb{R}^{N \times 10}$$

2.

(1)

First, we define the definition of the linear function as

$$Z = W \cdot X + b \text{ where } W \in \mathbb{R}^{10 \times K}, b \in \mathbb{R}^{10}, X \in \mathbb{R}^{K \times N}, Z \in \mathbb{R}^{10 \times N}$$

Now, we calculate the derivative of Z w.r.t X such that

$$\begin{aligned} \frac{\partial Z[i,j]}{\partial X[p,q]} &= \frac{\partial}{\partial X[p,q]} \left(\sum_l W[i,l] \cdot X[l,j] + b[j] \right) = W[i,q] \cdot \delta_{i,p} \cdot \delta_{j,q} \\ \frac{\partial \text{loss}}{\partial X} &= \frac{\partial \text{loss}}{\partial Z} \cdot \frac{\partial Z}{\partial X} = \frac{\partial \text{loss}}{\partial Z} \cdot W \in \mathbb{R}^{N \times K} \end{aligned}$$

(2)

Then, we derive the derivative of loss w.r.t W such that

$$\begin{aligned} \frac{\partial Z[i,j]}{\partial W[p,q]} &= \frac{\partial}{\partial W[p,q]} \left(\sum_l W[i,l] \cdot X[l,j] + b[j] \right) = X[i,q] \cdot \delta_{i,p} \cdot \delta_{j,q} \\ \frac{\partial \text{loss}}{\partial W} &= \frac{\partial Z}{\partial W} \cdot \frac{\partial \text{loss}}{\partial Z} = X \cdot \frac{\partial \text{loss}}{\partial Z} \in \mathbb{R}^{K \times 10} \end{aligned}$$

Finally, we derive the derivative of loss w.r.t b such that

$$\begin{aligned} \frac{\partial Z[i,j]}{\partial b[q]} &= \frac{\partial}{\partial b[q]} \left(\sum_l W[i,l] \cdot X[l,j] + b[j] \right) = \delta_{q,j} \\ \frac{\partial \text{loss}}{\partial b} &= \sum_i \frac{\partial \text{loss}}{\partial Z}[i, :] = \mathbf{1}^T \cdot \frac{\partial \text{loss}}{\partial Z} \end{aligned}$$

where $\mathbf{1} \in \mathbb{R}^N$ is a vector of ones

3.

Now, we can define $Y = X_{pad} *_{\mathcal{S}} W + b$, where the convolution of X_{pad} and W (filter) is along the last 3 axes of X_{pad} and the first 3 axes of W .

$W \in \mathbb{R}^{W_h \times W_w \times D \times O}$, $X_{pad} \in \mathbb{R}^{N \times X_h \times X_w \times D}$, $Y \in \mathbb{R}^{N \times Y_h \times Y_w \times O}$, $b \in \mathbb{R}^O$, where W_h is the height of the kernel, W_w is the width of the kernel, D is the input channel length, O is the output channel length, N is the batch size, X_h is the height of the input, X_w is the width of the input, Y_h is the height of the output, Y_w is the width of the output, and s is stride where $s_h = 1$ is the stride length for height and $s_w = 1$ is the stride length for width.

So, we can write the derivative of $loss$ w.r.t $W[i, j, d, o]$ and $b[o]$

$$\frac{\partial Y[n, k, l, o]}{\partial W[i, j, d, o]} = X_{pad}[n, k + i, l + j, d, o]$$

$$\frac{\partial Y[n, h, w]}{\partial b} = 1$$

$$\frac{\partial loss}{\partial W[i, j, d, o]} = \sum_{n=0}^{N-1} \sum_{k=0}^{Y_h-1} \sum_{l=0}^{Y_w-1} \frac{\partial loss}{\partial Y[n, k, l, o]} \cdot \frac{\partial Y[n, k, l, o]}{\partial W[i, j, d, o]} = \sum_{n=0}^{N-1} \sum_{k=0}^{Y_h-1} \sum_{l=0}^{Y_w-1} \frac{\partial loss}{\partial Y[n, k, l, o]} \cdot X_{pad}[n, k + i, l + j, d]$$

$$\frac{\partial loss}{\partial b[o]} = \sum_{n=0}^{N-1} \sum_{k=0}^{Y_h-1} \sum_{l=0}^{Y_w-1} \frac{\partial loss}{\partial Y[n, k, l, o]}$$

We can also write $\frac{\partial loss}{\partial W}$ as

$$\frac{\partial loss}{\partial W} = X_{pad} * \frac{\partial loss}{\partial Y}$$

where the convolution between X_{pad} and $\frac{\partial loss}{\partial Y}$ is along the first 3 axes of X_{pad} and the first 3 axes of $\frac{\partial loss}{\partial Y}$.

We can also write $\frac{\partial loss}{\partial b}$ as

$$\frac{\partial loss}{\partial b} = \frac{\partial loss}{\partial Y} * \mathbf{1}$$

where the convolution between $\frac{\partial loss}{\partial Y}$ and $\mathbf{1}$ is along all of axes of $\mathbf{1}$ and the first 3 axes of $\frac{\partial loss}{\partial Y}$, and $\mathbf{1} \in \mathbb{R}^{N \times Y_h \times Y_w}$ is a 3D tensor of ones.

We can also write the derivative of $loss$ w.r.t $X_{pad}[n, i, j, d]$

$$\frac{\partial Y[n, k, l, o]}{\partial X_{pad}[n, i, j, d]} = W[i - k, j - l, d, o] \cdot (0 \leq i - k \leq W_h - 1) \cdot (0 \leq j - l \leq W_w - 1)$$

$$\begin{aligned} \frac{\partial loss}{\partial X_{pad}[n, i, j, d]} &= \sum_{k=i-W_h+1}^i \sum_{l=j-W_w+1}^j \sum_{o=0}^{O-1} \frac{\partial loss}{\partial Y[n, k, l, o]} \cdot \frac{\partial Y[n, k, l, o]}{\partial X_{pad}[n, i, j, d]} \\ &= \sum_{k=i-W_h+1}^i \sum_{l=j-W_w+1}^j \sum_{o=0}^{O-1} \frac{\partial loss}{\partial Y[n, k, l, o]} \cdot W[i - k, j - l, d, o] \end{aligned}$$

We can also write $\frac{\partial loss}{\partial X_{pad}}$ as

$$\frac{\partial loss}{\partial X_{pad}} = \frac{\partial loss}{\partial Y} * W$$

where the convolution between $\frac{\partial loss}{\partial Y}$ and W is along the last 3 axes of $\frac{\partial loss}{\partial Y}$, and the first 2 and the last axis of W .

Now, we can crop $\frac{\partial loss}{\partial X_{pad}}$ so as to get $\frac{\partial loss}{\partial X}$

$$\frac{\partial loss}{\partial X} = \frac{\partial loss}{\partial X_{pad}}[:, p_h : -p_h, p_w : -p_w, :]$$

where p_h is the padding of height of X and p_w is the padding of width of X .

Q2 (8 pts): How do the results compare to the ones you got when implementing these models in TensorFlow?

- [4 pts] For each of the models, please comment on any differences or discrepancies in results -- runtime, performance and stability in training and final performance. (This would be the place to justify design decisions in the implementation and their effects).
- [2 pts] Which of the models show under-fitting?
- [2 pts] Which of the models show over-fitting?

A2: (Your answer here)

1.

For model 1 with all of settings, the performance, stability and runtime of my model written in numpy are almost the same as the model written in tensorflow.

For model 2 with all of settings, the performance, stability and runtime of my model written in numpy are almost the same as the model written in tensorflow.

For model 3 with setting 1, the stability of my model written in numpy is a bit better than the model written in tensorflow, and the performance and runtime are almost the same as tensorflow. The reason may be caused by the different initialization values in our model compared with the ones in the tensorflow model (since the initialization is xavier, which can cause a various selection of initialization values). For model 3 with setting 2 and setting 3, the performance, runtime and stability of my model written in numpy are almost the same as the model written in tensorflow.

For model 4 with all of settings, the stability and performance of my model written in numpy is almost the same as the model written in tensorflow. However, the runtime of my model written in numpy is longer than the model written in tensorflow.

For the whole model, the basic programming logic is to build up each component of the model as a module, which includes a function **forward_pass**, a function **params_gradient** and a function **backward_pass**. Meanwhile, the updated variables and calculated gradients are stored in each module. Then, several components of modules defined before are composed together using a list called **graph**. After that, we also implement a module called **Model** including the list **graph**, which has a function called

forward_pass and a function called **backprop** so as to complete the procedures of forward passing and backpropagation. Finally, we define an optimizer module called **SGD** including the instance of **Model**, which has a function called **learning** that can operate the update of parameters. So, for training we just need to call this **learning** function, which is convenient.

About, the implementation of the loss function, we combine the softmax and cross entropy together and divide the maximal value to each output value, so that the numerical problem can be solved out. For linear module, the implementation should be simple, we just use the matrix multiplication to implement it. For convolution module and maxpooling module, we use two **for loop** to control the sliding window to finish the forward pass, backward pass and parameter gradients calculation, which is the same as what we derive in the above section. The implementation using two **for loop** may be the reason why our model is worse than tensorflow on runtime on CNN. To increase the performance, we set every variable to be **float64** so that the problem of overflow can be mitigated.

2.

Model 1 with setting 1, setting 2 and setting 3 are under-fitting.

Model 2 with setting 1 and setting 3 are under-fitting.

Model 3 with setting 3 is under-fitting.

Model 4 with setting 1 is under-fitting.

3.

Model 2 with setting 2 is over-fitting.

Model 3 with setting 2 is over-fitting.

Model 4 with setting 3 is over-fitting.