

# Deep Learning Systems (ENGR-E 533) Homework 1

## Instructions

**Due date: Feb. 10, 2019, 23:59 PM (Eastern)**

- Submit your pdf report and a zip file of source codes to Canvas
- No handwritten solutions
  - Go to <http://sharelatex.com> ASAP and learn how to use L<sup>A</sup>T<sub>E</sub>X
- Start early if you're not familiar with the subject, programming, and L<sup>A</sup>T<sub>E</sub>X.
- Do it yourself. Discussion is fine, but code up on your own
- Late policy
  - If the sum of the late hours (throughout the semester) < seven days (168 hours): no penalty
  - If your total late hours is larger than 168 hours, you'll get only 80% of all the late-submitted homework.
- I ask you to use either PyTorch or Tensorflow running on Python 3.
- You can submit a `.ipynb` as a consolidated version of report and code if you want. But the math should be clear with L<sup>A</sup>T<sub>E</sub>X symbols and the explanations should be full by using text cells. Your sound clips embedded in there should be playable even before running all the codes (otherwise, submit the wavefile separately). Your images in there should be visible without running the code. Don't convert your `.ipynb` into PDF or HTML. We know how to read `.ipynb`, so PDF or HTML export will slow down the grading process and make the AIs grumpy.

## Problem 1: A Detailed View to MNIST Classification [5 points]

1. Train a fully-connected net for MNIST classification (sorry, no CNN please, yet). It should be with 5 hidden layers each of which is with 1024 hidden units. Feel free to use whatever techniques you learned in class. You should be able to get the test accuracy above 98%.
2. MNIST dataset can be loaded by using an API in TF:

```
from tensorflow.examples.tutorials.mnist import input_data
mnist = input_data.read_data_sets('MNIST_data', one_hot=True)
```

In PT, you can use these lines of commands (don't worry about the batch size and normalization—you can go for your own option for them):

```
train_loader = torch.utils.data.DataLoader(
    datasets.MNIST('/opt/e533/MNIST',
        train=True,
        download=False,
        transform=transforms.Compose([
            transforms.ToTensor(),
            transforms.Normalize((0.1307,), (0.3081,))
        ])),
    batch_size=128, shuffle=True)
test_loader = torch.utils.data.DataLoader(
    datasets.MNIST('/opt/e533/MNIST',
        train=False,
        transform=transforms.Compose([
            transforms.ToTensor(),
            transforms.Normalize((0.1307,), (0.3081,))
        ])),
    batch_size=128, shuffle=True)
```

3. Once you're done with training, as a starter, do a feedforward step on your test samples, a thousand of them. Capture the output of the softmax layer, which will be a 10-dim probability vector per sample. In other words, each output dimension has 1,000 predictions corresponding to the 1,000 examples. For each 10-d output vector, find the dim with the maximum probability (which will eventually decide the class label). Plot the input image associated with that in a grid of subplots. For example, you can create a  $10 \times 10$  grid of subplots, whose first row plots first ten input images that produced the highest probabilities for the first dim (which corresponds to "0"). Eventually, if your classification was near perfect, you'll see ten 0's in the first row, ten 1's in the second, and so on.
4. Repeat the procedure in Problem 1.3 for your second to the last layer output. This time, you should have 1024-dim vector per sample. Choose 10 random dimensions of interest and repeat the procedure in 1.3 as if the 10 out of 1024 dimensions are your output vectors. Note that there can be some dimensions that are with less than 10 images associated, because they are not popular. In your  $10 \times 10$  grid, now there must be some rows that are not with enough number of images or even an empty rows. Explain your observation compared with the results from 1.3. What can you see? What would have been the ideal situation for this second-to-the-last layer? Feel free to investigate the other layers if you want, but I wouldn't care because we have a better way.
5. t-Stochastic Neighbor Embedding (tSNE) or Principal Component Analysis (PCA) are useful tools to reduce the dimension of your data and

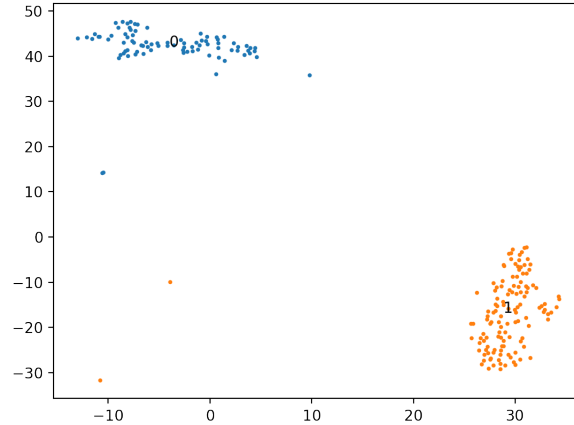


Figure 1: tSNE on two classes

visualize. By using them, you can reduce the dimension of your data, for example, down to 2D space, so that you can scatter plot your data samples. Feel free to use whatever implementation you can find for tSNE and PCA. I'd use the one in scikit-learn.

6. First, take a thousand test samples from your MNIST dataset. Apply tSNE and PCA on the flattened 784-dim pixels. Now you have  $2 \times 1000$  (or  $1000 \times 2$  if you transposed the data) matrix from each of the dim reduction algorithms. Scatter plot the data samples. **USE THE LABELS OF THE DATA SAMPLES SO THAT EACH SET OF SAMPLES FROM THE SAME CLASS ARE REPRESENTED WITH THE SAME COLOR. OVERLAY THE CLASS LABEL ON TOP OF THE MEAN OF THE CLASS.** By doing so, you can examine if your data is easy to classify or not. Do you think this raw image samples are easy to classify? For your information I share my scatter plot of the first two classes in Figure 1. It looks easy because there are only two classes, but with all 10 classes the situation will be different. Your plot should be similar to this but with all 10 classes.
7. Do a feedforward using your classifier. Capture the output of your first hidden layer, which will give you  $1024 \times 1000$  matrix. What that means is that you transformed your input data into a 1024-dim space. You may hope that this makes your classification easier. Check it out by doing tSNE and PCA on this matrix, which will once again give you  $2 \times 1000$  matrix. Scatter plot and check out if this layer gives you a better representation in terms of classification.

8. Repeat this procedure for all your layers including the last one. Explain your observation.

## Problem 2: Speech Denoising Using Deep Learning [5 points]

1. *If you took my MLSP class, you may think that you've seen this problem. But, it's actually somewhat different from what you did before, so read carefully. And, this time you SHOULD implement a DNN with at least two hidden layers. So, don't reuse your legacy MATLAB code for this problem.*
2. When you attended IUB, you took a course taught by Prof. K. Since you really liked his lectures, you decided to record them without the professor's permission. You felt awkward, but you did it anyway because you really wanted to review his lectures later.
3. Although you meant to review the lecture every time, it turned out that you never listened to it. After graduation, you realized that a lot of concepts you face at work were actually covered by Prof. K's class. So, you decided to revisit the lectures and study the materials once again using the recordings.
4. You should have reviewed your recordings earlier. It turned out that a fellow student who used to sit next to you always ate chips in the middle of the class right beside your microphone. So, Prof. K's beautiful deep voice was contaminated by the annoying chip eating noise.
5. But, you vaguely recall that you learned some things about speech denoising and source separation from Prof. K's class. So, you decided to build a simple deep learning-based speech denoiser that takes a noisy speech spectrum (speech plus chip eating noise) and then produces a cleaned-up speech spectrum.
6. Since you don't have Prof. K's clean speech signal, I prepared this male speech data recorded by other people. `train_dirty_male.wav` and `train_clean_male.wav` are the noisy speech and its corresponding clean speech you are going to use for training the network. Take a listen to them. Load them and covert them into spectrograms, which are the matrix representation of signals. To do so, you'll need to install `librosa` and use it by using the following codes:

```
!pip install librosa # in colab, you'll need to install this
import librosa
s, sr=librosa.load('train_clean_male.wav', sr=None)
S=librosa.stft(s, n_fft=1024, hop_length=512)
sn, sr=librosa.load('train_dirty_male.wav', sr=None)
X=librosa.stft(sn, n_fft=1024, hop_length=512)
```

which is going to give you two matrices  $\mathbf{S}$  and  $\mathbf{X}$  of size  $513 \times 2459$ . This procedure is something called Short-Time Fourier Transform.

7. Take their magnitudes by using `np.abs()` or whatever suitable method, because  $\mathbf{S}$  and  $\mathbf{X}$  are complex valued. Let's call them  $|\mathbf{S}|$  and  $|\mathbf{X}|$ .
8. Train a fully-connected deep neural network. A couple of hidden layers would work, but feel free to try out whatever structure, activation function, initialization scheme you'd like. The input to the network is a column vector of  $|\mathbf{X}|$  (a 513-dim vector) and the target is its corresponding one in  $|\mathbf{S}|$ . You may want to do some mini-batching for this. Make use of whatever functions in Tensorflow or Pytorch.
9. But, remember that your network should predict nonnegative magnitudes as output. Try to use a proper activation function in the last layer to make sure of that. I don't care which activation function you use in the middle layers.
10. `test_01_x.wav` is the test noisy signal. Load them and apply STFT as before. Feed the magnitude spectra of this test mixture  $|\mathbf{X}_{test}|$  to your network and predict their clean magnitude spectra  $|\hat{\mathbf{S}}_{test}|$ . Then, you can recover the (complex-valued) speech spectrogram of the test signal in this way:

$$\hat{\mathbf{S}} = \frac{\mathbf{X}_{test}}{|\mathbf{X}_{test}|} \odot \hat{\mathbf{S}}_{test}, \quad (1)$$

which means you take the phase information of the input noisy signal  $\frac{\mathbf{X}_{test}}{|\mathbf{X}_{test}|}$  and use that to recover the clean speech.  $\odot$  stands for the Hadamard product.

11. Recover the time domain speech signal by applying an inverse-STFT on  $\hat{\mathbf{S}}_{test}$ , which will give you a vector. Let's call this cleaned-up test speech signal  $\hat{\mathbf{s}}_{test}$ . I'll calculate something called Signal-to-Noise Ratio (SNR) by comparing it with the ground truth speech I didn't share with you. It should be reasonably good. You can actually write it out by using the following code:

```
librosa.output.write_wav('test_s_01_recons.wav', sh_test, sr)
```

12. Do the same testing procedure for `test_02_x.wav`, which actually contains Prof. K's voice along with the chip eating noise. Enjoy his enhanced voice using your DNN.