

Homework 3

CS 5787 Deep Learning

Spring 2020

Due: See Canvas

Your homework submission must cite any references used (including articles, books, code, websites, and personal communications). All solutions must be written in your own words, and you must program the algorithms yourself. **If you do work with others, you must list the people you worked with.** Submit your solutions as a PDF to Canvas.

Your homework solution must be typed. We urge you to prepare it in \LaTeX . It must be output to PDF format. To use \LaTeX , we suggest using <http://overleaf.com>, which is free and can be accessed online.

Your programs must be written in Python. The relevant code to the problem should be in the PDF you turn in. If a problem involves programming, then the code should be shown as part of the solution to that problem. One easy way to do this in \LaTeX is to use the verbatim environment, i.e., `\begin{verbatim} YOUR CODE \end{verbatim}`. For this assignment, you may use the plotting toolbox of your choice, PyTorch, and NumPy.

If told to implement an algorithm, don't use a toolbox, or you will receive no credit.

Problem 0 - Recurrent Neural Networks (10 points)

Recurrent neural networks (RNNs) are universal Turing machines as long as they have enough hidden units. In the next homework assignment we will cover using RNNs for large-scale problems, but in this one you will find the parameters for an RNN that implements binary addition. Rather than using a toolbox, you will find them by hand.

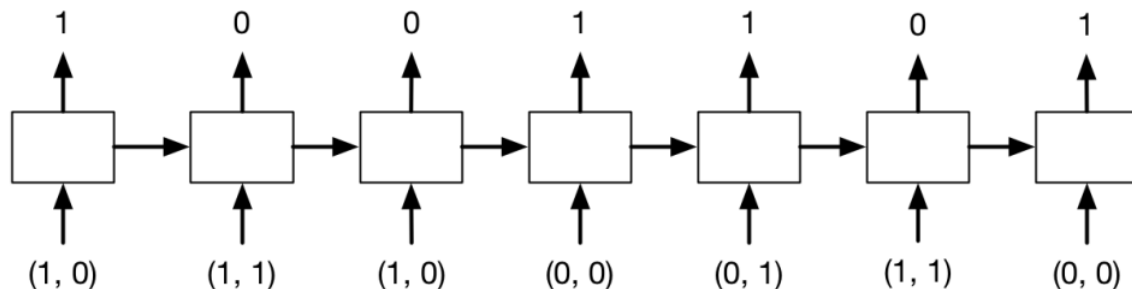
The input to your RNN will be two binary numbers, starting with the *least* significant bit. You will need to pad the largest number with an additional zero on the left side and you should make the other number the same length by padding it with zeros on the left side. For instance, the problem

$$100111 + 110010 = 1011001$$

would be input to your RNN as:

- Input 1: 1, 1, 1, 0, 0, 1, 0
- Input 2: 0, 1, 0, 0, 1, 1, 0
- Correct output: 1, 0, 0, 1, 1, 0, 1

The RNN has two input units and one output unit. In this example, the sequence of inputs and outputs would be:



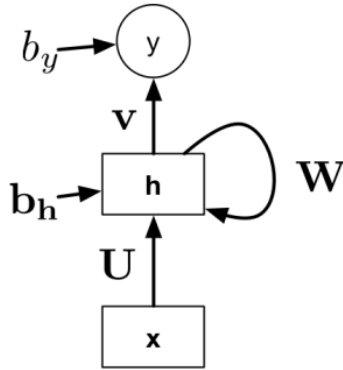
The RNN that implements binary addition has three hidden units, and all of the units use the following non-differentiable hard-threshold activation function

$$\sigma(a) = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$$

The equations for the network are given by

$$\begin{aligned} y_t &= \sigma(\mathbf{v}^T \mathbf{h}_t + b_y) \\ \mathbf{h}_t &= \sigma(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{h}_{t-1} + \mathbf{b}_h) \end{aligned}$$

where $\mathbf{x}_t \in \mathbb{R}^2$, $\mathbf{U} \in \mathbb{R}^{3 \times 2}$, $\mathbf{W} \in \mathbb{R}^{3 \times 3}$, $\mathbf{b}_h \in \mathbb{R}^3$, $\mathbf{v} \in \mathbb{R}^3$, and $b_y \in \mathbb{R}$



Part 1 - Finding Weights

Before backpropagation was invented, neural network researchers using hidden layers would set the weights by hand. Your job is to find the settings for all of the parameters by hand, including the value of \mathbf{h}_0 . Give the settings for all of the matrices, vectors, and scalars to correctly implement binary addition.

Hint: Have one hidden unit activate if the sum is at least 1, one hidden unit activate if the sum is at least 2, and one hidden unit if it is 3.

Solution:

Problem 1 - GRU for Sentiment Analysis

In this problem you will use a popular RNN model called the Gated Recurrent Units (GRU) to learn to predict the sentiment of a sentence. The dataset we are using is the IMDB review dataset ([link](#)). It is a binary sentiment classification (positive or negative) dataset that contains 50,000 movie reviews (50% for training and 50% for testing). We provide four text files for you to download on Canvas: `train_pos_reviews.txt`, `train_neg_reviews.txt`, `test_pos_reviews.txt`, `test_neg_reviews.txt`. Each line is an independent review for a movie.

Put your code in the appendix.

Part 1 - Preprocessing (5 points)

First you need to do proper preprocessing of the sentences so that each word is represented by a single number index in a vocabulary.

Remove all punctuation from the sentences. Build a vocabulary from the unique words

collected from text file so that each word is mapped to a number.

Now you need to convert the data to a matrix where each row is a sentence. Because sentences are of different length, you need to pad or truncate the sentences to make them same length. We are going to use 400 as the fixed length in this problem. That means any sentence that is longer than 400 words will be truncated; any sentence that is shorter than 400 words will be padded with 0s. Please note that your padded 0s should be placed *before* the sentences.

After you prepare the data, you can define a standard PyTorch dataloader directly from numpy arrays (say you have data in `train_x` and labels in `train_y`).

```
train_data = TensorDataset(torch.from_numpy(train_x), torch.from_numpy(train_y))
train_loader = DataLoader(train_data, shuffle=True, batch_size=batch_size)
```

Implement the data preprocessing procedure.

Solution:

Part 2 - Build A Binary Prediction RNN with GRU (10 points)

Your RNN module should contain the following modules: a word embedding layer, a GRU, and a prediction unit.

1. You should use `nn.Embedding` layer to convert an input word to an embedded feature vector.
2. Use `nn.GRU` module. Feel free to choose your own hidden dimension. It might be good to set the `batch_first` flag to `True` so that the GRU unit takes (batch, seq, embedding_dim) as the input shape.
3. The prediction unit should take the output from the GRU and produce a number for this binary prediction problem. Use `nn.Linear` and `nn.Sigmoid` for this unit.

At a high level, the input sequence is fed into the word embedding layer first. Then, the GRU is taking steps through each word embedding in the sequence and return output / feature at each step. The prediction unit should take the output from the final step of the GRU and make predictions.

Implement your RNN module, train the model and report accuracy on the test set.

Solution:

Part 3 - Comparison with a MLP (5 points)

Since each sentence is a fixed length input (with potentially many 0s in some samples), we can also train a standard MLP for this task.

Train a two layer MLP on the training data and report accuracy on the test set. How does it compare with the result from your RNN model?

Solution:

Problem 2 - Generative Adversarial Networks

For this problem, you will be working with Generative Adversarial Networks (GAN) on Fashion-MNIST dataset (Figure 1).

Fashion-MNIST dataset can be loaded directly in PyTorch by the following command:

```
import torchvision
fmnist = torchvision.datasets.FashionMNIST(root=".", train=True,
transform=transform, download=True)
data_loader = torch.utils.data.DataLoader(dataset=fmnist,
batch_size=batch_size, shuffle=True)
```

Similar to the well known MNIST dataset, Fashion-MNIST is designed to be a standard testbed for ML algorithms. It has the same image size and number of classes as MNIST, but is a little bit more difficult.

We are going to train GANs to generate images that looks like those in Fashion-MNIST dataset. Through the process, you will have a better understanding on GANs and their characteristics.

Training a GAN is notoriously tricky, as we shall see in this problem.

Put your code in the appendix.

Part 1 - Vanilla GAN (10 points)

A GAN is containing a Discriminator model (D) and a Generator model (G). Together they are optimized in a two player minimax game:

$$\begin{aligned}\min_D &= -\mathbb{E}_{x \in p_d} \log D(x) - \mathbb{E}_{z \in p_z} \log(1 - D(G(z))) \\ \min_G &= -\mathbb{E}_{z \in p_z} \log D(G(z))\end{aligned}$$

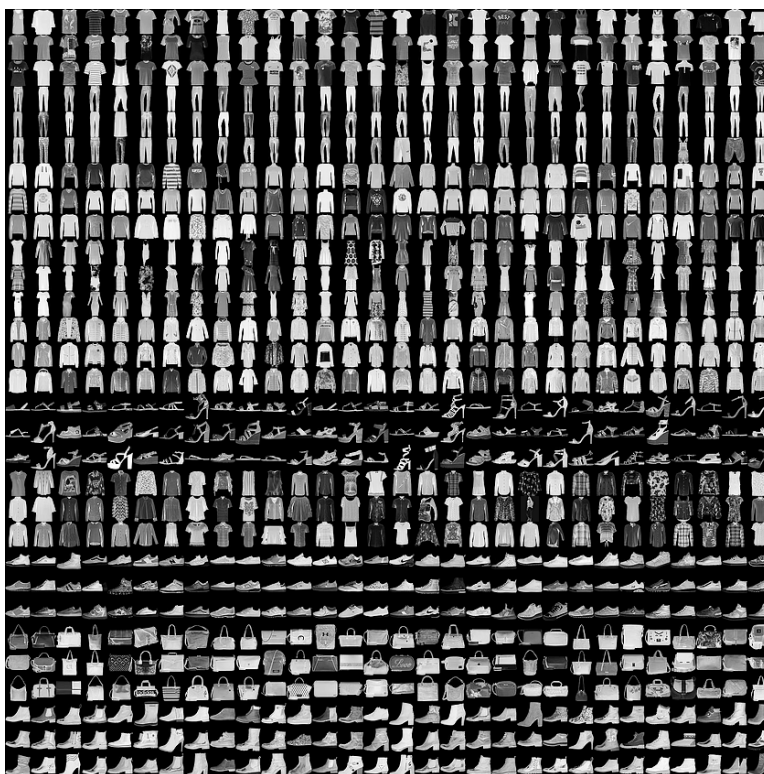


Figure 1: Fashion-Mnist dataset example images. It contains 10 classes of cloths, shoes, and bags.

In practice, a GAN is trained in an iterative fashion where we alternate between training G and training D . In pseudocode, GAN training typically looks like this:

For epoch 1:max_epochs

 Train D:

 Get a batch of real images

 Get a batch of fake samples from G

 Optimize D to correctly classify the two batches

 Train G :

 Sample a batch of random noise

 Generate fake samples using the noise

 Feed fake samples to D and get prediction scores

 Optimize G to get the scores close to 1 (means real samples)

Choice of G architecture:

Make your generator to be a simple network with three linear hidden layers with ReLU activation functions. For the output layer activation function, you should use hyperbolic tangent (tanh). This is typically used as the output for the generator because ReLU cannot output negative values.

Choice of D architecture:

Make your discriminator to be a similar network with three linear hidden layers using ReLU activation functions, but the last layer should have a logistic sigmoid as its output activation function, since it the discriminator D predicts a score between 0 and 1, where 0 means fake and 1 means real.

Train a basic GAN that can generate images from the Fashion-MNIST dataset. Plot your training loss curves for your G and D . Show the generated samples from G in 1) the beginning of the training; 2) intermediate stage of the training; and 3) after convergence.

Solution:

Part 2 - GAN Loss (10 points)

In this part, we are going to modify the model you just created in order to compare different choices of losses in GAN training.

MSE

$$\min_G \mathbb{E}_{z \in p_z, x \in p_d} (x - G(z))^2$$

You can get rid of the discriminator and directly use a MSE loss to train the generator.

Wasserstein GAN (WGAN)

$$\begin{aligned} \min_D \quad & -\mathbb{E}_{x \in p_d} D(x) + \mathbb{E}_{z \in p_z} D(G(z)) \\ \min_G \quad & -\mathbb{E}_{z \in p_z} D(G(z)) \end{aligned}$$

WGAN is proposed to address the vanishing gradient problem in the original GAN loss when the discriminator is way ahead of the generator. One thing to change in WGAN is that the output of the discriminator should be now ‘unbounded’, namely you need to remove the sigmoid function at the output layer. And you need to clip the weights of the discriminator so that their L_1 norm is not bigger than c .

Try c from the set $\{0.1, 0.01, 0.001, 0.0001\}$ and compare their difference.

Least Square GAN

$$\begin{aligned} \min_D \quad & \mathbb{E}_{x \in p_d} (D(x) - 1)^2 + \mathbb{E}_{z \in p_z} D(G(z))^2 \\ \min_G \quad & \mathbb{E}_{z \in p_z} (D(G(z)) - 1)^2 \end{aligned}$$

The idea is to provide a smoother loss surface than the original GAN loss.

Plot training curves and show generated samples of the above mentioned losses. Discuss if you find there is any difference in training speed and generated sample's quality.

Solution:

Part 3 - Mode Collapse in GANs (10 points)

Take a copy of your vanilla GAN discriminator and change its output channel from 1 output to 10 output units. Fine-tune it as a classifier on the Fashion-MNIST training set. You should easily achieve $\sim 90\%$ accuracy on Fashion-MNIST test set.

Now generate 3000 samples using the generator you trained for Part 1. Use the classifier you just trained to predict the class labels of those samples. Plot the histogram of predicted labels.

Although the original Fashion-MNIST dataset has 10 classes equally distributed, you will find the histogram you just generated is not close to uniform (even if we consider the classifier is not perfect and 3000 samples are not too large). This is a known issue with GAN called Mode Collapse. It means the GAN is often capturing only a subset (mode) of the original data's distribution, not all of them.

Unrolled GAN is proposed to reduce the effect of mode collapse in GAN training. The intuition is that if we let G see ahead how D would change in the next k steps, G can adjust accordingly and hopefully will perform better. Its idea can be summarized in the following modified training scheme:

```
For epoch 1:max_epochs
  Train D:
    Get a batch of real images
    Get a batch of fake samples from G
```


Optimize D to correctly classify the two batches

Make a copy of D into D_{unroll}

Train D for k unrolled steps:

 Get a batch of real images

 Get a batch of fake samples from G

 Optimize D_{unroll} to correctly classify the two batches

Train G :

 Sample a batch of random noise

 Generate fake samples using the noise

 Feed fake samples to D_{unroll} and get prediction scores

 Optimize G to get the scores close to 1 (means real samples)

Note that G is trained with a copy of D at each epoch. The original D should not be updated during that part of training.

Train an unrolled GAN and re-plot the histogram from 3000 generated samples. Discuss whether unrolled GAN seems to help reduce the mode collapse problem.

WGAN is claimed to be less affected by mode collapse too. In addition to your vanilla GAN model, use the WGAN model you trained in Part 2 and plot the histogram of class distribution, compare it to unrolled GAN and vanilla GAN.

Solution:

Part 4 - Conditional GAN (10 points)

For the GANs we have been playing with, we cannot specify the class we want generated. Now, we explore adding extra information to the GAN to take more control over the generation process. Specifically, we want to generate not just *any* images from Fashion-MNIST data distribution, but images with a particular label such as shoes. This is called the Conditional GAN because now samples are drawn from a conditional distribution given a label as input.

To add the conditional input vector, we need to modify both D and G . First, we need to define the input label vector. We are going to use one-hot encoding vectors for labels: for an image sample with label k of K classes, the vector is K dimensional and has 1 at k -th element and 0 otherwise.

We then concatenate the one-hot encoding of class vector with original image pixels (flattened as a vector) and feed the augmented input to D and G . Note we need to change the number of channels in the first layer accordingly.

Train a Conditional GAN using the training script from Part 1. Plot training curves for D and G . Generate 3 samples from each of the 10 classes. Discuss differences in the generated images produced compared to the non-conditional models you built.

Solution:

Code Appendix