

Deep Learning Systems (ENGR-E 533) Homework 4

Instructions

Due date: Dec. 13, 2020, 23:59 PM (Eastern)

THERE IS NO GRACE PERIOD FOR THIS LAST HOMEWORK.

THE SUBMISSION SYSTEM WILL BE CLOSED AT THE DEADLINE.

- Start early if you're not familiar with the subject, TF or PT programming, and L^AT_EX.
- 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.
- Submit a `.ipynb` as a consolidated version of your report and code snippets. But the math should be clear with L^AT_EX symbols and the explanations should be full by using text cells. In addition, submit an `.html` version of your notebook where you embed your sound clips and images. For example, if you have a graph as a result of your code cell, it should be visible in this `.html` version before we run your code. Ditto for the sound examples.

Problem 1: Speaker Verification [6 points]

1. In this problem, we are going to build a speaker verification system. It takes two utterances as input, and predicts whether they were spoken by the same speaker (positive class) or not (negative class).
2. `trs.pkl` contains an $500 \times 16,180$ matrix, whose row is a speech signal with 16,180 samples. They are the returned vectors from the `librosa.load` function. Similarly, `tes.pkl` holds a $200 \times 22,631$ matrix.
3. The training matrix is ordered by speakers. Each speaker has 10 utterances, and there are 50 such speakers (that's why there are 500 rows). Similarly, the test set has 20 speakers, each of which is with 10 utterances.
4. Randomly sample L pairs of utterances from the ten utterance of the first speaker. In theory, there are $\binom{10}{2} = 45$ pairs you can sample from (the order of the two utterances within a pair doesn't matter). You can use all 45 of them if you want. These are the *positive* examples in your first minibatch.
5. Let's construct L *negative* pairs as well. First, randomly sample L utterances from the 49 training speakers. Second, randomly sample another L utterances from the first speaker (the speaker you sampled the positive pairs from). Using these two sets, each has L examples, form another set of L pairs. If $L > 10$, you'll need to repeatedly use the first speaker's utterance (i.e. sampling with replacement). This set is your *negative* examples, each of whose pair contains an utterance from the first speaker and a random utterance spoken by a different speaker.
6. The L positive pairs and L negative pairs form your first minibatch. You have $2L$ pairs of utterances in total.

7. Repeat this process for the other training speakers, so that each speaker is represented by L positive pairs and L negative pairs. By doing so, you can form 50 minibatches with a balanced number of positive and negative pairs.
8. Train a Siamese network that tries to predict 1 for the positive pairs and 0 for the negative ones. In a minibatch, since you have L positive and L negative pairs, respectively, your net must predict L ones and L zeros, respectively.
9. I found that STFT on the signals serves the initial feature extraction process. Therefore, your Siamese network will take as input TWO spectrograms, each of which is of size $513 \times T$. I wouldn't care too much about your choice of the network architecture this time (if it works anyway), but it has to somehow predict a fixed-length feature vector for the given sequence of spectra (consequently, TWO fixed-length vectors for the pair of input spectrograms). Using the inner product of the two latent embedding vectors as the input to the sigmoid function, you'll do a logistic regression. Use your imagination and employ whatever techniques you learned in class to design/train this network.
10. Construct similar batches from the test set, and test the verification accuracy of your network. Report your test-time speaker verification performance. I was able to get a decent result ($\sim 70\%$) with a reasonable network architecture, which converged in a reasonable amount of time (i.e. in an hour).
11. Submit your code and accuracy on the test examples.

Problem 2: Variational Autoencoders on Poor Sevens [5 points]

1. `tr7.pkl` contains 6,265 MNIST digits from its training set, but not all ten digits. I only selected 7's. Therefore, it's with a rank-3 tensor of size $6,265 \times 28 \times 28$. Similarly, `te7.pkl` contains 1,028 7's.
2. The digit images in this problem are special, because I added a special effect to them. So, they are different from the original 7's in the MNIST dataset in a way. I want you to find out what I did to the poor 7's.
3. Instead of eyeballing all those images, you need to implement a VAE that finds out a few latent dimensions, one of which should show you the effect I added.
4. Once again, I wouldn't care too much about your network architecture. This could be a good chance for you to check out the performance of the CNN encoder, followed by a decoder with deconvolution layers (or transposed convolution layers), but do something else if you feel like. I found that fully-connected networks work just fine.
5. What's important here in the VAE is, as a VAE, it needs a hidden layer that is dedicated to learn the latent embedding. In this layer, each hidden unit is governed by a standard normal distribution as its *a priori* information. Also, be careful about the re-parameterization technique and the loss function.
6. You'll need to limit the number of hidden units K in your code layer (the embedding vector) with a small number (e.g. smaller than 5) to reduce your search space. Out of K , there must be a dimension that explains the effect that I added.
7. One way to prove that you found the latent dimension of interest is to show me the digits generated by the decoder. More specifically, you may want to "generate" new 7's by feeding a few randomly generated code vectors, that are the random samples from the K normal distributions that your VAE learned. But, they won't be enough to show which dimension takes care of my added effect. Therefore, your random code vectors should be formed specially.
8. What I'd do is to generate code vectors by fixing the $K - 1$ dimensions with the same value over the codes, while varying only one of them.

9. For example, if $K = 3$ and you're interested in the third dimension, your codes should look like as follows:

$$\mathbf{Z} = \begin{bmatrix} 0.23 & -0.18 & -5 \\ 0.23 & -0.18 & -4.5 \\ 0.23 & -0.18 & -4.0 \\ 0.23 & -0.18 & -3.5 \\ 0.23 & -0.18 & -3.0 \\ & \vdots & \\ 0.23 & -0.18 & 4.5 \\ 0.23 & -0.18 & 5.0 \end{bmatrix} \quad (1)$$

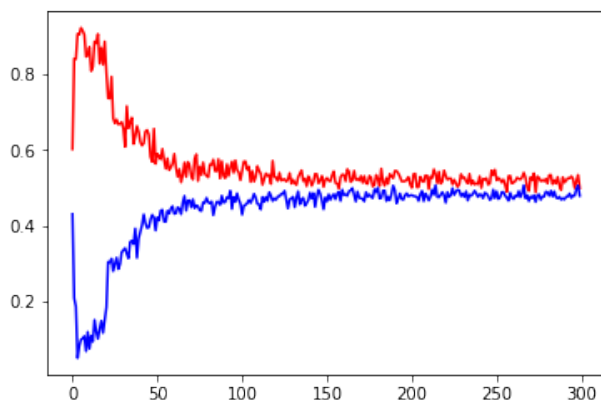
Note that the first two column vectors are once randomly sampled from the normal distributions, but then shared by all the codes so that the variation found in the decoded output relies solely on the third dimension.

10. You'll want to examine all the K dimensions by generating samples from each of them. Show me the ones you like. They should show a bunch of **similar-looking** 7's but with **gradually changing effect** on them. The generated samples that show a gradual change of the thickness of the stroke, for example, are not a good answer, because that's not the one I added, but something that was there in the dataset already.
11. Submit your notebook with figure and code.

Problem 3: Conditional GAN [5 points]

- Let's develop a GAN model that can generate MNIST digits, but based on the auxiliary input from the user indicating which digit to create.
- To this end, the generate has to be trained to receive two different kinds of input: the random vector and the class label.
- The random vector is easy to create. It should be a d -dimensional vector, sampled from a standard normal distribution $\mathcal{N}(0, 1)$. $d = 100$ worked just fine for me.
- As for the *conditioning* vector, you somehow need to inform the network of your intention. For example, if you want to generate a "0", you need to give that information to the generator. There are many different ways to condition a neural network at various stages. But, this time, let's use a simple one. We will convert the digit label into a one-hot vector. For example, if you want to generate a "7" the conditioning vector is $[0, 0, 0, 0, 0, 0, 1, 0, 0]$.
- Then, we need to combine these two different kinds of information. Again, there are many different ways, but let's just stick to a simple solution. We will concatenate the d -dimensional random vector and the 10-dimensional one-hot vector. Therefore, the input to your generator is with $d + 10$ dimensions. If your $d = 100$, the input dimension is 110.
- You are free to choose whatever network architecture you want to practice with. Here's the fully-connected one I found as a good starting point: $110 \times 200 \times 400 \times 784$. I used ReLU as the activation function, but as for the last layer, I used tanh. It means that I'll interpret -1 as the black pixel while $+1$ being the white pixel.
- The discriminator has a similar architecture: $784 \times 400 \times 200 \times 100 \times 1$. The reason why it takes a 784-dim vector is that it wants to know what the image sample is conditioned on. Also note that it does binary classification to discern whether the conditioned input image is a real or fake example, i.e., you will need to set up the last layer as a logistic regression function.
- To train this GAN model, sample a minibatch of B examples from your MNIST dataset. These are your real examples. But, instead of feeding them directly to your discriminator, you'll append their label information by turning it into the one-hot representation. Don't forget to match the scale: it has to be from -1 to $+1$ instead of $[0, 1]$ as that's how the generator defines the pixel intensity.

9. Accordingly, generate a class-balanced set of fake examples by feeding B random vectors to your generator. Again, each of your random vectors needs to be appended by a randomly chosen one-hot vector. For example, if your $B = 100$, you may want to generate ten ones, ten twos, and so on. Although the generated images are not with any label information anymore, you know that each should belong to a particular digit class based on your conditioning vector. Therefore, when you feed these fake examples to the discriminator, you need to append the one-hot vectors once again. Of course, the one-hot vectors should match the ones you used to inform the generator as input.
10. To summarize, the input to your generator is a $d + 10$ -dim vector. The last 10 elements should be copied to augment your fake example, generated from the generator, to construct a 794-dim vector. You have B fake examples as such. The real examples are with the same size, but their first 784 elements are from the real MNIST images, accompanied by the last 10 elements representing the class, to which the image belongs.
11. Train this GAN model. I used Adam with lower-than-usual learning rates. Dropout helped the discriminator. Below is the figure that shows the change of the classification accuracy over the epochs (red for real and blue for fake examples).

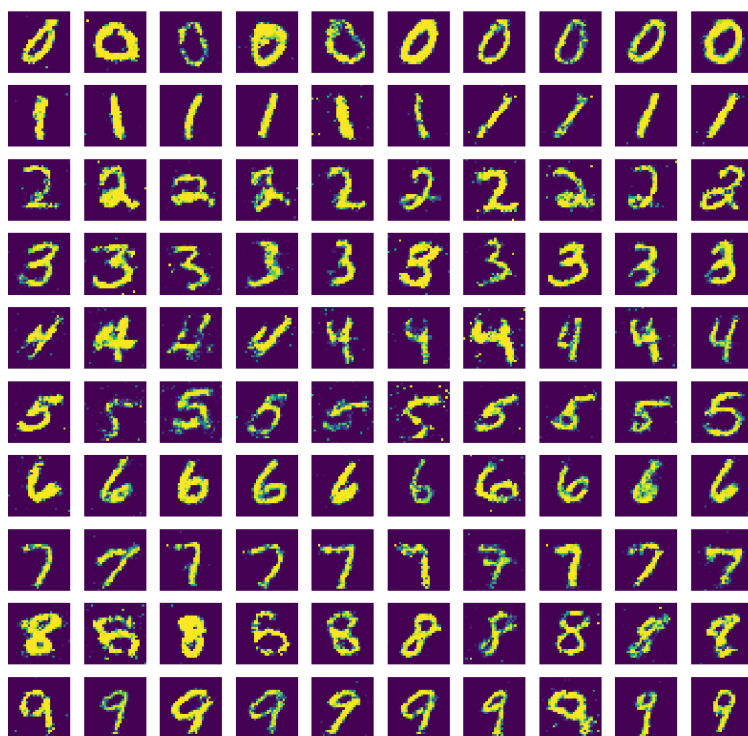


I can see that it converged to the Nash equilibrium, as the discriminator seems to be confused.

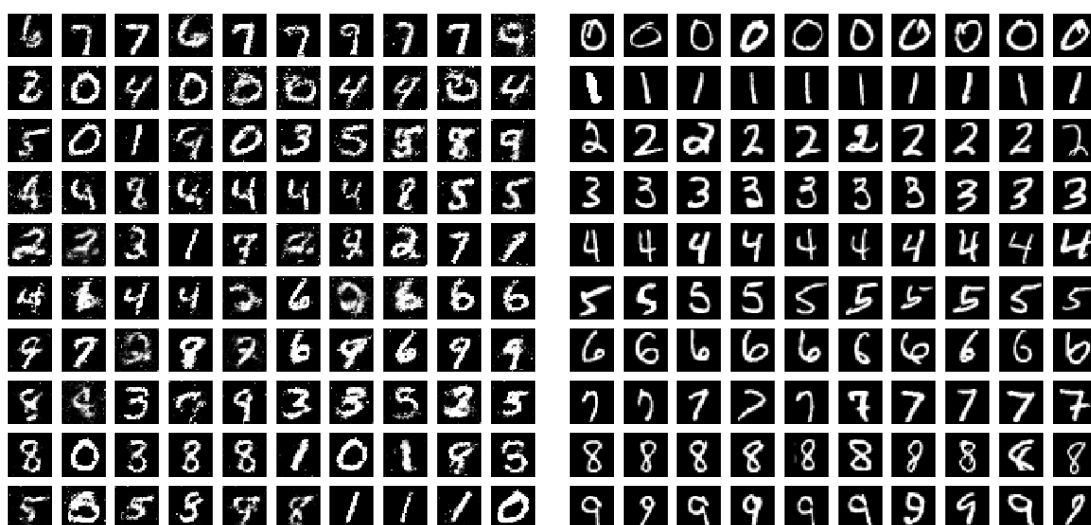
12. Below is the test examples that I generated by feeding new random vectors (plus the intended class labels). I placed ten examples per class in a row. These are of course not the best MNIST digits I can imagine, but they look fine given the simple structure and algorithm I used.
13. Please feel free to use whatever other things you want to try out, such as WGAN, but if your results are decent (like mine) we'll give away the full score.
14. Report both the convergence graph as well as the generated examples.

Problem 4: Missing Value Imputation Using Conditional GAN [6 points]

1. We've already seen that LSTM can act like generative model that can "predict future patches" given the "past patches" in HW3 P5.
2. This time, we'll do something similar but by using GAN. This time, it works like a missing value imputation system. We assume that only the center part of the image is known, while the generator has to predict what the other surrounding pixels are.
3. We'll formulate this as a conditional GAN. First, take a batch of MNIST images. Take their center 10×10 patches, and then flatten it. This is your 100-dimensional conditioning vector. Since there are 28×28 pixels in each image, you'll do something like this, `X[:,9:19,9:19]`, to take the center patch. This will form a $B \times 100$ matrix, for your batch of B conditioning vectors.

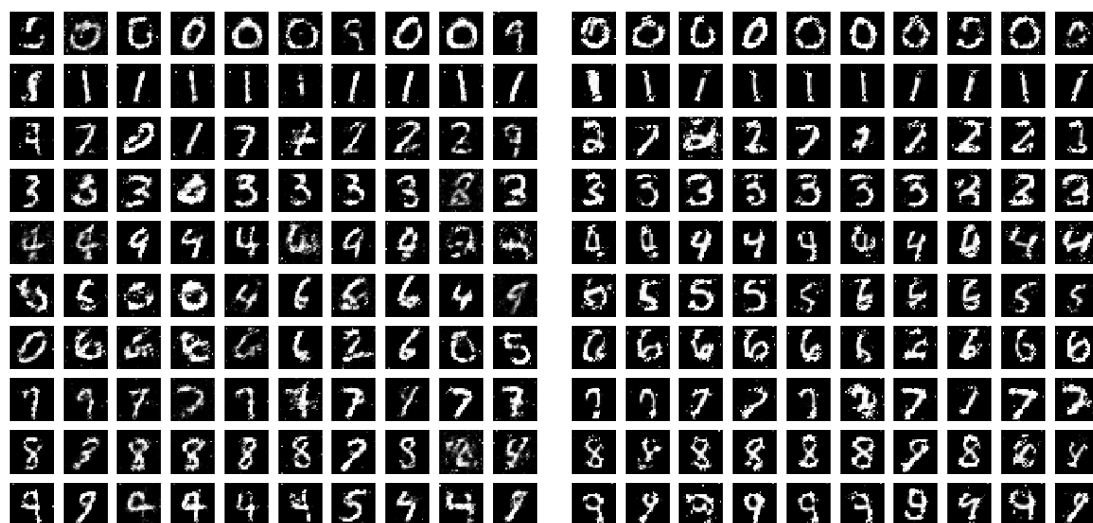


4. Append this matrix to your random vectors of 100 dimensions drawn from the standard normal distribution. This $B \times 200$ matrix is the input to your generator.
5. The generator takes this 200 dimensional vectors and synthesizes MNIST-looking digits. You will need to prepare another set of B real examples. Eventually, you feed $2B$ examples in total to your discriminator as a minibatch.
6. If both discriminator and generator are trained properly, you can see that the results are some MNIST-looking digits. But, I found that the generator simply ignores the conditioning vector and generate whatever it wants to generate. They all certainly look like MNIST digits, but the conditioning part doesn't work. Below is the generated images (left) and the ground-truth images that I extracted the center patches from (right).



They are completely different from each other.

7. So, even though I did feed the center patch as the conditioning vector to the generator, it ignores it and generate something totally different. It's because, I think, the generator has no way to know the conditioning vector is actually the center patch of the digit that it must generate. In other words, the generator is generating the whole image, although it doesn't have to generate the center patch, which is known to me. Instead, I wanted it to generate the surrounding pixels, that are the missing values.
8. As a remedy, I added another regularizer to my generator so that it functions as an autoencoder at least for the center pixels. You know, in an ordinary GAN setup, the generator loss has to penalize the discriminator's decision that classifies the fake examples into the fake class (i.e., when the generator fails to fool the discriminator). On top of this ordinary generator loss, I add a simple mean squared error term that penalizes the difference between the conditioning vector and the center patch of the generated image, as they have to be the same, essentially.
9. Since it's a regularizer, I needed to investigate different λ values to control its contribution to the total loss of the generator. It turned out that the generator is not too sensitive to this choice, although it does generate a "less conditioned" example when it comes to a too small λ . Below is the two sets of examples when I set $\lambda = 0.1$ (left) and $\lambda = 10$ (right).



10. Replicate what I did and submit your code and generated examples. Once again, you can try some other fancy models and different ways to condition the model. But we'll give you a full score if your results are as good as mine.