

Intro to Image Understanding (CSC420)

Assignment 3

Due Date: November 8th, 2024, 10:59:00 pm
Total: 160 marks

General Instructions:

- You are allowed to work directly with one other person to discuss the questions. However, you are still expected to write the solutions/code/report in your own words; i.e. no copying. If you choose to work with someone else, you must indicate this in your assignment submission. For example, on the first line of your report file (after your own name and information, and before starting your answer to Q1), you should have a sentence that says *worked together with my classmate [name]* and *written the solutions/code/report in [your name]*.
- Your submission should be a single file (PDF), with the answers to the specific questions, the code, and the explanation and discussion of your results. Submit this file to MarkUs directly.
- Submit documents and code separately. Please store all your results separately. Submit the code to the folder and then submit the results to the folder. Include a **README.txt** file (inside the folder) that describes the code and the results.
- Do not worry if you submit multiple files; you can submit multiple files; you can submit multiple files.



Part I: Theoretical

[Question 1] reparameterization trick (5 marks)

Briefly (2-3 sentences) explain the purpose of the reparameterization trick in a variational autoencoder.

[Question 2] GAN (5 marks)

In a GAN we have a generator and discriminator. Calculating the loss function for which of them requires a detach()? Your answer could be either of the two, neither, or both. Briefly (1-2 lines) justify your answer.

[Question 3] VQ-VAE (5 marks)

Briefly explain what the following line of code does in a vector-quantised variational autoencoder (VQ-VAE) implementation.

```
quantized = inputs + (quantized - inputs).detach()
```

[Question 4] FID (5 marks)

The Frechet Inception Distance (FID) score is a metric used to evaluate the quality of images generated by GANs. Find and read a (short) tutorial about FID and briefly (in 3-5 sentences) explain what it measures and how it is computed.

[Question 5] Corner detection

For corner detection,

Let's denote the 2×2

1. (1 marks) Corner

2. (4 marks) Pro



[Question 6] Optical flow (5 marks)

Optical flow is problematic in which of the following conditions? Provide a Yes/No answer and a brief explanation for each case.

1. (1 mark) In homogeneous image areas.
2. (1 mark) In textured image areas.
3. (1 mark) At image edges.
4. (1 mark) At the boundaries of a moving object.
5. (1 mark) Corner of a non-moving object.

[Question 7] LSTM (10 marks)

We want to build an LSTM cell that sums its inputs over time. What should the value of the input gate and the forget gate be?

[Question 8] GAN training¹ (non-saturating generator cost) (15 marks)

Consider a GAN with generator $G(z)$ and discriminator $D(G(z))$. The figure below shows the training losses for two different generator loss functions: $J_1(G)$ and $J_2(G)$. The blue curve plots the value of $J_1(G)$ as a function of $D(G(z))$. Likewise, the red curve plots the value of $J_2(G)$ as a function of $D(G(z))$. For m generated samples, $J_1(G)$ and $J_2(G)$ are defined as follows:

$$J_1(G) = \frac{1}{m} \sum_{i=1}^m \log D(G(z_i))$$

$$J_2(G) = \frac{1}{m} \sum_{i=1}^m \|D(G(z_i)) - 1\|$$



1. **(5 mark)** Early in the training, is the value of $D(G(z))$ closer to 0 or closer to 1? Briefly explain why.
2. **(5 mark)** Which of the two cost functions would you choose to train your GAN? Briefly justify your answer.
3. **(5 mark)** "A GAN is successfully trained when $D(G(z))$ is close to 1". Is this statement TRUE or FALSE? Briefly explain your answer.

(You can use insight learned from this question in your implementation tasks.)

¹source for this question: <https://coursys.sfu.ca/2020sp-cmpt-980-g2/pages/final-questions/view>

Part II: Implementation Tasks (105 marks)

TASK I { GAN (40 marks)

In Tutorial 1, we saw a simple GAN implementation where both the generator and the discriminator used fully connected layers. Implement a GAN where the generator uses transposed convolutional layers and the discriminator uses convolutional layers. Make both of them have 5 layers. In the generator, start the first layer as follows:

```
nn.ConvTranspose2d(in_channels=64, out_channels=512, kernel_size=4, stride=1, padding=0)
```

and adjust the parameters of the following 4 transposed conv layers to map a 64-dim latent vector into a 28 × 28 grayscale image. The rest of the generator you can keep similar to that of Tutorial 1, i.e. batchnorm and ReLU after the first 4 transposed convolutional layers and sigmoid at the end. For the discriminator

```
nn.Conv2d(in_channels=1, out_channels=64, kernel_size=3, stride=1, padding=0)
```

and adjust the parameters of the following 4 conv layers to map a 28 × 28 grayscale image into a 64-dim latent vector.

Train this GAN with the same hyperparameters as the simple GAN in Tutorial 1.

Task II { For this task, you must choose between Task II.a or Task II.b. Do not do both; write only one.

Task II.a { WGAN

For this question { and only this question { you are allowed to use AI code generation as much as you want. You can also ask your favourite LLM (e.g. chatGPT) what steps you need to take to modify a GAN implementation into a WGAN. If you do use any AI tools, mention them.

The Wasserstein GAN (or WGAN) is a GAN variant from 2017 that aims to get rid of problems like mode collapse and improve the training stability of GANs.

1. Modify the code in Tutorial 1 (or write your own code) to implement a WGAN. Train this WGAN on MNIST and compare your results with that of the simple GAN in Tutorial 1.
2. Briefly explain if/how training this WGAN was different from training your conv GAN in the previous task.

Task II.b - Simple Text-guided Image Generation (40 marks)

In this task, we implement a simple text-guided image generator. To this end,

1. load a GAN model pre-trained on Imagenet. For example, you can load the `vqgan_imagenet_f16_16384` model from <https://github.com/CompVis/taming-transformers>.
2. freeze the GAN, but allow the random seed vector weights to be trained.
3. use CLIP to encode a text prompt and also the generated image from the GAN
4. Choose a loss function that tries to match the CLIP encoding of the prompt with that of the GAN-generated image. Backpropagate the loss through the GAN to update the random seed.

5. use the updated random seed to generate a new image, and backpropagate again and so on.

Using this method, we can generate images that match the text prompts:

- \a dog playing
- a prompt that y



[Task III] Corner Detection

Download two images of the city of Toronto under two different viewing directions:

- <https://commons.wikimedia.org/wiki/File:Toronto.jpg>
- https://commons.wikimedia.org/wiki/File:City_of_Toronto,_Canada.jpg

1. Calculate the eigenvalues λ_1 and λ_2 for each pixel of I_1 and I_2 .
2. Show the scatter plot of λ_1 and λ_2 (where $\lambda_1 > \lambda_2$) for all the pixels in I_1 and the same scatter plot for I_2 (**5 marks**). Each point shown at location $(x; y)$ in the scatter plot, corresponds to a pixel with eigenvalues: $\lambda_1 = x$ and $\lambda_2 = y$.
3. Based on the scatter plots, pick a threshold for $\min(\lambda_1, \lambda_2)$ to detect corners. Illustrate detected corners on each image using the chosen threshold (**2 marks**).
4. Constructing matrix M involves the choice of a window function $w(x; y)$. Often a Gaussian kernel is used. Repeat steps 1, 2, and 3 above, using a significantly different Gaussian kernel (i.e. a different σ) than the one used before. For example, choose a σ that is significantly (e.g. 5 times, or 10 times) larger than the previous one (**3 marks**). Explain how this choice influenced the corner detection in each of the images (**5 marks**).