

# CS 446 / ECE 449 — Homework 5

*your NetID here*

Version 1.0

## Instructions.

- Homework is due **Wednesday, Nov. 17th, at noon CST**; you have **3** late days in total for **all Homeworks**.
- Everyone must submit individually at gradescope under **hw5** and **hw5code**.
- The “written” submission at **hw5 must be typed**, and submitted in any format gradescope accepts (to be safe, submit a PDF). You may use L<sup>A</sup>T<sub>E</sub>X, markdown, google docs, MS word, whatever you like; but it must be typed!
- When submitting at **hw5**, gradescope will ask you to **mark out boxes around each of your answers**; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full **academic integrity** information. You should cite any external reference you use.
- We reserve the right to reduce the auto-graded score for **hw5code** if we detect funny business (e.g., your solution lacks any algorithm and hard-codes answers you obtained from someone else, or simply via trial-and-error with the autograder).
- When submitting to **hw5code**, only upload **hw5.py**. Additional files will be ignored.

# 1. Expectation Maximization

In this question, we expect you to do some computation related to the EM algorithm covered in the lecture.

**Background.** On the xy-plane, we have a rigid object, and our sensor can capture  $N$  key points, whose coordinates are:  $\mathcal{P} = \{\mathbf{p}^{(1)}, \mathbf{p}^{(2)}, \dots, \mathbf{p}^{(N)}\}$ . ( $N$  is a sufficiently large integer.) An unknown force then cause a translation  $\mathbf{T}$  to this object, where  $\mathbf{T}$  is encoded  $T_x$  and  $T_y$ , meaning how long the object has moved along the x-axis and y-axis. To calculate parameter  $\mathbf{T}$ , we use our sensor to capture the key points on the rigid object one more time, acquiring a set of  $N$  key points:  $\mathcal{Q} = \{\mathbf{q}^{(1)}, \mathbf{q}^{(2)}, \dots, \mathbf{q}^{(N)}\}$ . (See Fig. 1 for an intuitive demonstration. The “L” is our rigid body,  $N = 3$ , and the blue and red dots are the key points.)

**Assumption.** During this process, we assume that an bi-jection mapping between  $\mathcal{P}$  and  $\mathcal{Q}$  exists (See Fig. 1, the key points are all the corners of “L.”) However, this bijection mapping cannot be directly got from the sensors, as the orders of the points may be shuffled during perception.

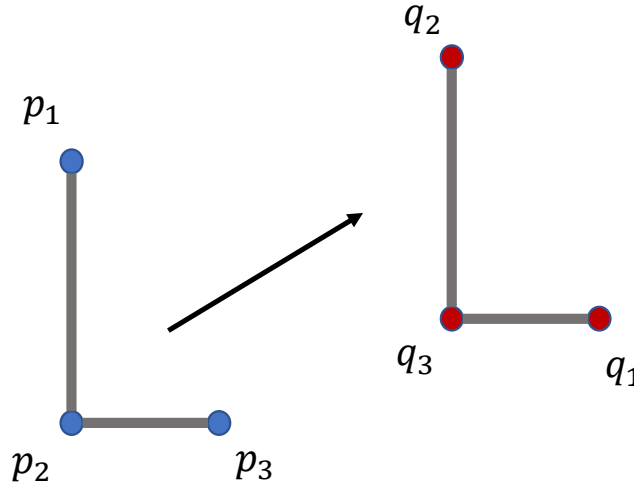


Figure 1

**Objective.** Using EM algorithm to estimate the translation parameter and the corresponding pairs of key points, by introducing binary hidden variables  $\mathbf{Z} \in \mathbb{R}^{N \times N}$ , where  $Z_{ij} = 1$  means  $\mathbf{p}^{(i)}$  and  $\mathbf{q}^{(j)}$  are a match.

**Remark.** We have the following remarks on this question.

- This questions set is for you to understand the process of EM with an intuitive example without annoying mathematical equations. However, to make it easy to understand, we use additional assumptions and simplifications. Please note the difference between this example and rigorous EM algorithm when you learn deeper machine learning courses in the future.
  - You may find EM overkill for this simple question, and you are right about it. This question originates from a classical problem in computer vision called “point cloud registration,” where this problem could be much more difficult when the sensor has noises and the bijection between  $\mathcal{P}$  and  $\mathcal{Q}$  does not exist. You may refer to “iterative closest points” (ICP) if you are interested in this.
- (a) **Joint Probability.** If we know a pair of **matched** points  $(\mathbf{p}^{(i)}, \mathbf{q}^{(j)})$ , think of it as a single data sample, with the corresponding hidden state  $Z_{ij} = 1$ . Intuitively, based on this single pair, how likely parameter  $\mathbf{T}$  is depends on  $\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(j)}\|_2$ . To make the math easier, we assume  $\|\mathbf{p} + \mathbf{T} - \mathbf{q}\|_2$  follows from the Gaussian distribution  $\mathcal{N}(0, \sigma)$  where  $\sigma$  is known, i.e.,

$$\mathbb{P}_{\mathbf{T}}((\mathbf{p}^{(i)}, \mathbf{q}^{(j)}) | Z_{ij} = 1) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(j)}\|_2^2}{2\sigma^2}\right) \quad (1)$$

(To avoid confusion in the notations, please use  $\mathbb{P}$  to represent probability)

But we actually don't know which points are a match, or say we don't know the hidden states  $\mathbf{Z}$ , and want to use EM to help. We define that the matching is optimized from  $\mathcal{P}$  to  $\mathcal{Q}$ , and not the reverse. In this way, the matching between each point  $\mathbf{p} \in \mathcal{P}$  is independent, and two points  $\mathbf{p} \in \mathcal{P}$  can match to the same  $\mathbf{q} \in \mathcal{Q}$ , but not the reverse.

Let  $\mathbb{P}(Z_{ij} = 1) := \Pi_{ij}$  indicate the prior probability of  $\mathbf{p}^{(i)}$  and  $\mathbf{q}^{(j)}$  being a pair, where  $\mathbf{\Pi} \in \mathbb{R}^{N \times N}$  are unknown parameters. Under our assumption,  $0 \leq \Pi_{ij} \leq 1$ , and  $\sum_{k=1}^N \Pi_{ik} = 1, \forall i \in \{1, 2, \dots, N\}$ . Let us denote overall parameters  $(\mathbf{T}, \mathbf{\Pi})$  as  $\psi$  like in the class.

Given a point  $\mathbf{p}^{(i)}$ , it has  $N$  possible matches  $\mathbf{q}^{(j)}, j = 1, \dots, N$ . The probability contributed by  $\mathbf{p}^{(i)}$  in the complete data likelihood is

$$\mathbb{P}(\mathbf{p}^{(i)}, \mathcal{Q}, \mathbf{Z}_{i,:}) = \mathbb{P}(\mathbf{p}^{(i)}, \mathcal{Q} | \mathbf{Z}_{i,:}) \mathbb{P}(\mathbf{Z}_{i,:}) = \prod_{j=1}^N [\mathbb{P}_{\psi}((\mathbf{p}^{(i)}, \mathbf{q}^{(j)}) | Z_{ij} = 1) \mathbb{P}(Z_{ij} = 1)]^{Z_{ij}} \quad (2)$$

**Please write the log-likelihood  $\log \mathbb{P}_{\psi}(\mathcal{P}, \mathcal{Q}, \mathbf{Z})$  under the two following cases.**

- i. General Case: Write out the general formula for  $\log \mathbb{P}_{\psi}(\mathcal{P}, \mathcal{Q}, \mathbf{Z})$  following the “Log-likelihood of complete data” on Slides 7/27 in lecture 17.
  - ii. Special Case:  $Z_{ii} = 1, i \in \{1, 2, 3, \dots, N\}$ . To get the full credits of this question, please full expand the Gaussian distributions.
- (b) **Expectation. (1/2)** After doing the preparation work, we now start to apply EM algorithm to this question. In the E-step, our first objective is to compute  $\mathbf{R}$ , where  $\mathbf{R} = [r_{ij}] \in \mathbb{R}^{N \times N}$ , and  $r_{ij} = \mathbb{P}_{\psi^{(t)}}(Z_{ij} | \mathbf{p}^{(i)}, \mathcal{Q})$ .

Following the procedure of “the first line of the E-step on slides 8/27 in lecture 17,” answer the following questions.

- i. General Case: Derive the formula for  $r_{ij}$ . For the simplicity of notations, you can also use  $\mathcal{N}(\cdot | \cdot, \cdot)$  like on the slides to denote Gaussian distribution.
  - ii. Special Case: For the  $N$  points in  $\mathcal{P}$ , the first  $\lfloor \frac{N}{2} \rfloor$  points are matched with  $\mathbf{q}^{(1)}$ , and the rest are matched to  $\mathbf{q}^{(2)}$ . Write the values of  $\mathbf{R}$ .
- (c) **Expectation. (2/2)** This question computes the values of  $Q(\psi | \psi^{(t)})$ . Please answer the following questions.
- i. General Case: Derive the formula of  $Q(\psi | \psi^{(t)})$  following the procedure on slides 8/27 in lecture 17. For the simplicity of notations, you can also use  $\mathcal{N}(\cdot | \cdot, \cdot)$  to denote Gaussian distribution.
  - ii. Special Case: Same as the special case mentioned in problem “Expectation (1/2),” fully expand the formula of  $Q(\psi | \psi^{(t)})$ . To get full credits, your answer cannot have the variables  $r_{ij}$  and the notation for Gaussian distribution  $\mathcal{N}(\cdot | \cdot, \cdot)$ . (use  $0 \times \log 0 = 0$  in your calculation)
- (d) **Maximization.** On the basis of previous derivation, complete the maximization step and the update rule. Similar to the slides 9/27 in lecture 17, write out the formulas for  $\mathbf{\Pi}^{(t+1)}$  and  $\mathbf{T}^{(t+1)}$ .

**Hint.**  $\sigma$  is fixed and you do not need to solve it.

- i. General Case: Write the formulas for  $\mathbf{\Pi}^{(t+1)}$  and  $\mathbf{T}^{(t+1)}$ . You may use  $N$ ,  $\mathbf{R}$ , and the points in  $\mathcal{P}$  and  $\mathcal{Q}$ .
- ii. Special Case: Write the formulas for  $\mathbf{T}^{(t+1)}$  for the special case in question “Expectation (1/2).” You may use  $N$ , and the points in  $\mathcal{P}$  and  $\mathcal{Q}$ .

**Solution.**

$$\begin{aligned} \text{(a)} \quad i \quad \log \mathbb{P}_{\psi}(\mathcal{P}, \mathcal{Q}, \mathbf{Z}) &= \sum_{i=1}^N \sum_{j=1}^N \log [\mathbb{P}_{\psi}((\mathbf{p}^{(i)}, \mathbf{q}^{(j)}) | Z_{ij} = 1) \mathbb{P}(Z_{ij} = 1)]^{Z_{ij}} \\ &= \sum_{i=1}^N \sum_{j=1}^N Z_{ij} \log [\mathbb{P}_{\psi}((\mathbf{p}^{(i)}, \mathbf{q}^{(j)}) | Z_{ij} = 1) \mathbb{P}(Z_{ij} = 1)] \\ &= \sum_{i=1}^N \sum_{j=1}^N Z_{ij} \log \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(j)}\|_2^2}{2\sigma^2}\right) \Pi_{ij} \right] \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^N \sum_{j=1}^N Z_{ij} (\log(\frac{1}{\sqrt{2\pi}\sigma}) + \log \Pi_{ij} - \frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(j)}\|_2^2}{2\sigma^2}) \\
\text{ii } \log \mathbb{P}_\psi(\mathcal{P}, \mathcal{Q}, \mathbf{Z}) &= \sum_{i=1}^N Z_{ii} (\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(i)}\|_2^2}{2\sigma^2}) \\
&= \sum_{i=1}^N Z_{ii} (\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(i)}\|_2^2}{2\sigma^2})
\end{aligned}$$

(b) i  $r_{ij} = \frac{\mathbb{P}(Z_{ij}=1)\mathbb{P}(\mathbf{p}^{(i)}, \mathbf{q}^j | Z_{ij}=1)}{\sum_{k=1}^N \mathbb{P}(Z_{ik}=1)\mathbb{P}(\mathbf{p}^{(i)}, \mathbf{q}^{(k)} | Z_{ik}=1)} = \frac{\Pi_{ij} \mathcal{N}(\mathbf{p}^{(i)}, \mathbf{q}^{(j)} | Z_{ij}=1)}{\sum_{k=1}^N \Pi_{ik} \mathcal{N}(\mathbf{p}^{(i)}, \mathbf{q}^{(k)} | Z_{ik}=1)}$

ii For  $i \in [1, \lfloor \frac{N}{2} \rfloor]$  and  $j = 1$ ,  $r_{ij} = 1$ . For  $i \in [\lfloor \frac{N}{2} \rfloor + 1, N]$  and  $j = 2$ ,  $r_{ij} = 1$ . For other points,  $r_{ij} = 0$ .

(c) i  $Q(\psi | \psi^{(t)}) = \sum_{i=1}^N \sum_{j=1}^N r_{ij} \log \mathbb{P}(\mathbf{p}^{(i)}, \mathbf{q}^j | Z_{ij} = 1) \mathbb{P}(Z_{ij} = 1)$   
 $= \sum_{i=1}^N \sum_{j=1}^N r_{ij} \log \mathcal{N}(\mathbf{p}^{(i)}, \mathbf{q}^j | Z_{ij} = 1) \Pi_{ij}$

ii  $Q(\psi | \psi^{(t)}) = \sum_{i=1}^{\lfloor \frac{N}{2} \rfloor} \log \mathcal{N}(\mathbf{p}^{(i)}, \mathbf{q}^1 | Z_{ij} = 1) + \sum_{i=\lfloor \frac{N}{2} \rfloor + 1}^N \log \mathcal{N}(\mathbf{p}^{(i)}, \mathbf{q}^2 | Z_{ij} = 1)$   
 $= -N \log \sqrt{2\pi}\sigma - \sum_{i=1}^{\lfloor \frac{N}{2} \rfloor} \frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(1)}\|_2^2}{2\sigma^2} - \sum_{i=\lfloor \frac{N}{2} \rfloor + 1}^N \frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(2)}\|_2^2}{2\sigma^2}$

(d) i  $\mathcal{L}(\mathbf{\Pi}, \mathbf{T}, \boldsymbol{\lambda}) = -\sum_{i=1}^N \sum_{j=1}^N r_{ij} \log \mathcal{N}(\mathbf{p}^{(i)}, \mathbf{q}^j | Z_{ij} = 1) \Pi_{ij} + \sum_{i=1}^N \lambda_i (\sum_{j=1}^N \Pi_{ij} - 1)$   
 $= \sum_{i=1}^N \sum_{j=1}^N r_{ij} (\log(\sqrt{2\pi}\sigma) - \log \Pi_{ij} + \frac{\|\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(j)}\|_2^2}{2\sigma^2}) + \sum_{i=1}^N \lambda_i (\sum_{j=1}^N \Pi_{ij} - 1)$   
Setting derivatives to zero:  
 $\frac{\partial \mathcal{L}}{\partial \Pi_{ij}} = -\frac{r_{ij}}{\Pi_{ij}} + \lambda_i = 0$   
 $\frac{\partial \mathcal{L}}{\partial \mathbf{T}} = \sum_{i=1}^N \sum_{j=1}^N r_{ij} (\mathbf{p}^{(i)} + \mathbf{T} - \mathbf{q}^{(j)}) = 0$   
 $\frac{\partial \mathcal{L}}{\partial \lambda_i} = (\sum_{j=1}^N \Pi_{ij}) - 1 = 0$   
For  $\mathbf{\Pi}^{(t+1)}$ ,  $\Pi_{ij} = \frac{r_{ij}}{\lambda_i}$ . Taking sums over rows, we have  $\sum_{j=1}^N \Pi_{ij} = \sum_{j=1}^N \frac{r_{ij}}{\lambda_i} = \frac{1}{\lambda_i} \sum_{j=1}^N r_{ij} = 1$ . Hence  $\lambda_i = \sum_{j=1}^N r_{ij} = 1$ . So  $\mathbf{\Pi}^{(t+1)} = \mathbf{R}$ .  
For  $\mathbf{T}^{(t+1)}$ ,  $\sum_{i=1}^N \sum_{j=1}^N r_{ij} (\mathbf{p}^{(i)} - \mathbf{q}^{(j)}) + \mathbf{T} \sum_{i=1}^N \sum_{j=1}^N r_{ij} = \sum_{i=1}^N \sum_{j=1}^N r_{ij} (\mathbf{p}^{(i)} - \mathbf{q}^{(j)}) + N\mathbf{T} = 0$ . Hence  $\mathbf{T}^{(t+1)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N r_{ij} (\mathbf{p}^{(i)} - \mathbf{q}^{(j)})$ .

ii  $\mathbf{T}^{(t+1)} = -\frac{1}{N} (\sum_{i=1}^{\lfloor \frac{N}{2} \rfloor} (\mathbf{p}^{(i)} - \mathbf{q}^{(1)}) + \sum_{i=\lfloor \frac{N}{2} \rfloor + 1}^N (\mathbf{p}^{(i)} - \mathbf{q}^{(2)}))$   
 $= -\frac{1}{N} (\sum_{i=1}^N \mathbf{p}^{(i)} - \sum_{i=1}^{\lfloor \frac{N}{2} \rfloor} \mathbf{q}^{(1)} - \sum_{i=\lfloor \frac{N}{2} \rfloor + 1}^N \mathbf{q}^{(2)})$

## 2. Variational Auto-Encoders

We are training a variational auto-encoder (VAE). It contains the following parts: the input are vectors  $\mathbf{x}$ , the latent vector is  $\mathbf{z}$ , the encoder models the probability of  $q_\phi(\mathbf{z}|\mathbf{x})$ , and the decoder is  $p_\theta(\mathbf{x}|\mathbf{z})$ . Based on this notation, we will first look at several problems related to the structure of variational auto-encoder.

- (a) We assume the latent vector  $\mathbf{z} \in \mathbb{R}^2$  follows a multi-variate Gaussian distribution  $\mathcal{N}$ . Please compute the output dimension of the encoder  $q_\phi(\cdot)$  under the following cases and briefly explain why. (If “output dimension” is not clear enough for you, think of it as “how many real numbers  $r \in \mathbb{R}$  are needed to output for the sampling of latent vectors.”)
- We assume  $\mathcal{N}$  follows a multi-variate Gaussian distribution with an **identity matrix** as the covariance matrix.
  - We assume  $\mathcal{N}$  follows a multi-variate Gaussian distribution with an **diagonal matrix** as the covariance matrix.
- (b) We then consider the problems related to the understanding of KL-Divergence.
- Using the inequality of  $\log(x) \leq x - 1$ , prove that  $D_{KL}(p(x), q(x)) \geq 0$  holds for two arbitrary distributions  $p(x)$  and  $q(x)$ .
  - Consider a binary classification problem with input vectors  $\mathbf{x}$  and labels  $y \in \{0, 1\}$ . The distribution of the ground truth label is denoted as  $P(y)$ . The expression of  $P(y)$  is as Eq 3, where  $y_{gt}$  is the ground truth label.

$$P(y = y_{gt}) = 1, P(y = 1 - y_{gt}) = 0 \quad (3)$$

Suppose we are trying to predict the label of  $\mathbf{x}$  with a linear model  $\mathbf{w}$  and sigmoid function, then the distribution of  $y$  is denoted as  $Q(y)$  and computed as Eq. 4.

$$Q(y = 0|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}, \quad Q(y = 1|\mathbf{x}) = \frac{\exp(-\mathbf{w}^\top \mathbf{x})}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \quad (4)$$

With the above information, compute the KL Divergence between the distributions of  $P(y)$  and  $Q(y|\mathbf{x})$ , specifically  $D_{KL}(P(y), Q(y|\mathbf{x})) = \mathbf{E}_{y \sim P(y)}[\log \frac{P(y)}{Q(y|\mathbf{x})}]$ .

Expand your solution to the clearest form. To get full credits, you may only use  $y_{gt}$ ,  $\mathbf{w}$ ,  $\mathbf{x}$  and related constants in your expression.

- (c) VAE is a special branch of generative method in sampling the latent vectors  $\tilde{\mathbf{z}}$  from  $q_\phi(\mathbf{z}|\mathbf{x})$  instead of directly regressing the values of  $\mathbf{z}$ . Read an example implementation of VAE at [https://github.com/AntixK/PyTorch-VAE/blob/master/models/vanilla\\_vae.py](https://github.com/AntixK/PyTorch-VAE/blob/master/models/vanilla_vae.py) and answer the following questions:
- Find the functions and lines related to the sampling of  $\tilde{\mathbf{z}}$  from  $q_\phi(\mathbf{z}|\mathbf{x})$ . Specifying the names of the functions and the related lines can lead to full credits. Please note that if your range is too broad (in the extreme case, covering every line in the file) we cannot give your full credit.
  - Suppose our latent variable is  $\mathbf{z} \in \mathbb{R}^2$  sampled from a Gaussian distribution with mean  $\boldsymbol{\mu} \in \mathbb{R}^2$  and a diagonal covariance matrix  $\boldsymbol{\Sigma} = \text{Diag}\{\sigma_1^2, \sigma_2^2\}$ . Then another random variable  $\mathbf{v} \in \mathbb{R}^2$  is sampled from a Gaussian distribution  $\mathcal{N}(0, \mathbf{I})$ . Show that  $\mathbf{V} = [\sigma_1, \sigma_2]^\top \circ \mathbf{v} + \boldsymbol{\mu}$  follows the same distribution as  $\mathbf{z}$ . ( $\circ$  denotes Hadamard product, which means element-wise product;  $\mathcal{N}(0, \mathbf{I})$  denotes the multi-variate Gaussian with zero mean and identity matrix as covariance.)
  - Under the same setting of the Question ii, we can sample the latent vector  $\tilde{\mathbf{z}}$  by the process  $\tilde{\mathbf{z}} = [\sigma_1, \sigma_2]^\top \circ \tilde{\mathbf{v}} + \boldsymbol{\mu}$ , where  $\tilde{\mathbf{v}}$  is a sampled random variable from  $\mathcal{N}(0, \mathbf{I})$ . Consider the process of training, where we apply back-propagation to train the neural networks. Given the gradient on  $\tilde{\mathbf{z}}$  as  $\tilde{\mathbf{g}} \in \mathbb{R}^2$ , which can be written as  $[\tilde{g}_1, \tilde{g}_2]$ . **What are the gradients of the output of the encoder:  $\boldsymbol{\mu}, \sigma_1, \sigma_2$ ?** (Assume the KL-Divergence loss is not considered in this part.)  
**Note:** To get full credit, you can use any constants and the variables of  $\tilde{\mathbf{v}} = [\tilde{v}_1, \tilde{v}_2]$ ,  $\tilde{\mathbf{g}} = [\tilde{g}_1, \tilde{g}_2]$ , and  $\boldsymbol{\mu}, \sigma_1, \sigma_2$ .

- iv. During reading the code, you might feel confused about why we are sampling  $\tilde{\mathbf{z}}$  in such a way, instead of generating a random value directly. But now, you could have some clues. Please briefly explain “Why we are sampling  $\tilde{\mathbf{z}}$  with  $\mathcal{N}(0, 1)$ , instead of directly generating the values.”

**Solution.**

- (a) i 2. Only  $\mu_1$  and  $\mu_2$  are needed because the variations of  $z$  is 1 for both entries. And the encoder output should have the same distribution as the prior  $z$  distribution.  
 ii 4.  $\mu_1, \mu_2, \sigma_1^2$  and  $\sigma_2^2$  are needed for the description of the output distribution.
- (b) i  $-D_{KL}(p||q) = -\sum_x p(x) \ln(\frac{p(x)}{q(x)}) = \sum_x p(x) \ln(\frac{q(x)}{p(x)})$   
 $\leq \sum_x p(x) (\frac{q(x)}{p(x)} - 1) = \sum_x q(x) - \sum_p(x) = 1 - 1 = 0$   
 Hence  $-D_{KL}(p||q) \geq 0$ .  
 ii  $D_{KL}(P(y), Q(y|\mathbf{x})) = \mathbf{E}_{y \sim P(y)} [\log \frac{P(y)}{Q(y|\mathbf{x})}] = P(y_{gt}) \log \frac{P(y_{gt})}{Q(y=y_{gt}|\mathbf{x})}$   
 $= -\log(Q(y=y_{gt}|\mathbf{x})) = -\log \frac{(\exp(-\mathbf{w}^T \mathbf{x}))^{y_{gt}}}{1 + \exp(-\mathbf{w}^T \mathbf{x})} = -\log \frac{(\exp(-y_{gt} \mathbf{w}^T \mathbf{x}))}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$   
 $= y_{gt} \mathbf{w}^T \mathbf{x} + \log(1 + \exp(-\mathbf{w}^T \mathbf{x}))$
- (c) i Lines 120 and 121 are related to the sampling of  $\tilde{\mathbf{z}}$ , where the function *reparameterize()* is used to move the  $\mathcal{N}(0, 1)$  distribution to  $\mathcal{N}(\mu, \sigma^2)$ .  
 ii  $\mathbf{V}$  has the same distribution as  $\mathbf{z}$  due to the following facts:
- 1 The linear combinations of Gaussian distributions are still Gaussian distributions.
  - 2  $E[\sigma_1 v_1 + \mu_1] = \sigma E[v_1] + \mu_1 = \mu_1$ .
  - 3  $Var(\sigma_1 v_1 + \mu_1) = E[(\sigma_1 v_1 + \mu_1 - \mu_1)^2] = E[\sigma^2 v_1^2] = \sigma^2 E[(v_1 - 0)^2] = \sigma^2$ .
  - 4  $Cov(\sigma_1 v_1 + \mu_1, \sigma_2 v_2 + \mu_2) = E[(\sigma_1 v_1 + \mu_1)(\sigma_2 v_2 + \mu_2)] - E[\sigma_1 v_1 + \mu_1] E[\sigma_2 v_2 + \mu_2]$   
 $= \sigma_1 \sigma_2 E[v_1 v_2] = \sigma_1 \sigma_2 E[v_1 v_2] = \sigma_1 \sigma_2 E[v_1 v_2] - \sigma_1 \sigma_2 E[v_1] E[v_2]$   
 $= \sigma_1 \sigma_2 Cov(v_1, v_2) = 0$
- Therefore,  $\mathbf{V}$  and  $\mathbf{z}$  are two multivariate Gaussian distributions with exactly the same parameters.
- iii  $\frac{\partial L}{\partial \mathbf{z}} = \tilde{\mathbf{g}}$ .  
 $\frac{\partial L}{\partial \boldsymbol{\mu}} = \frac{\partial L}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \boldsymbol{\mu}} = \tilde{\mathbf{g}}$   
 $\frac{\partial L}{\partial \sigma_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial \sigma_1} = \tilde{\mathbf{g}} v_1$   
 $\frac{\partial L}{\partial \sigma_2} = \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial \sigma_2} = \tilde{\mathbf{g}} v_2$
- iv Generation of Gaussian random variables is mostly easily done by using the function *randlike()* to generate a value that follows  $\mathcal{N}(0, 1)$ . Sampling this way enables us to avoid building complicated self-defined random variables.

### 3. Generative Adversarial Networks

Let's implement a Generative Adversarial Network(GAN) to create images of hand-written digits!

GAN consists of two parts: a generator network  $G$  and a discriminator network  $D$ .  $G$  is expected to generate a fake image from a random latent variable  $z$ , and  $D$  is expected to distinguish fake images and real images.  $G$  and  $D$  are trained jointly with a minimax objective. In this question, we will use training data from MNIST to train our GAN, and let it produce some fake images that look like hand-written digits.

- (a) First, let's implement the **Discriminator** network. It should take  $32 \times 32$  gray-scale images as input, and output the probability of each image being a real one. Its architecture is summarized in Table 1.

Table 1: **Discriminator Architecture**

Layer Index	Layer Type	Input Channels	Output Channels	Kernel Size	Stride	Padding
1	Conv2d	1	16	3	1	1
2	LeakyReLU					
3	MaxPool			2	2	0
4	Conv2d	16	32	3	1	1
5	LeakyReLU					
6	MaxPool			2	2	0
7	Conv2d	32	64	3	1	1
8	LeakyReLU					
9	MaxPool			2	2	0
10	Conv2d	64	128	3	1	1
11	LeakyReLU					
12	MaxPool			4	4	0
13	Linear	128	1			
14	Sigmoid					

A few notes:

- All Conv2d and Linear layers have bias terms. You do not have to explicitly set `Conv2d(..., bias=True)`, since it is default in PyTorch.
  - Also, you do not need to explicitly initialize the weights in Conv2d and Linear layers. The default initialization by PyTorch is good enough.
  - LeakyReLU is a variant of ReLU activation, which has a smaller gradient for negative inputs. Set `negative_slope=0.2` for all LeakyReLU layers. More info about LeakyReLU at <https://pytorch.org/docs/stable/generated/torch.nn.LeakyReLU.html>.
  - You need to reshape the tensor sometimes in the forward pass.
  - Given a batch of images with shape `(batch_size, 1, 32, 32)`, the output of this network should be a tensor with shape `(batch_size)`, and the values in it are float numbers in  $(0, 1)$ . Our autograder will only be able to check the shape and range of the output, so be careful even if you have passed the test.
- (b) Next, we can implement the **Generator** network. It should take 128-d vectors (sampled from a Gaussian distribution) as input, and output fake images. Its architecture is summarized in Table 2. We will make use of transposed convolutional layers. Given an input, a transposed convolutional layer can produce an output with a higher resolution. Thus, we can generate a  $32 \times 32$  image from a vector by stacking such layers. A visualization of how transposed convolutional layers work can be found at [https://github.com/vdumoulin/conv\\_arithmetic/blob/master/README.md](https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md).

Table 2: **Generator Architecture**

Layer Index	Layer Type	Input Channels	Output Channels	Kernel Size	Stride	Padding
1	ConvTranspose2d	128	64	4	1	0
2	LeakyReLU					
3	ConvTranspose2d	64	32	4	2	1
4	LeakyReLU					
5	ConvTranspose2d	32	16	4	2	1
6	LeakyReLU					
7	ConvTranspose2d	16	1	4	2	1
8	Tanh					

A few notes:

- Again, all Conv2d and Linear layers have bias terms and are initialized by the default setup.
  - Same LeakyReLU as above, with `negative_slope=0.2` for all LeakyReLU layers.
  - You need to reshape the tensor sometimes in the forward pass.
  - Given a batch of latent vectors with shape `(batch_size, 128)`, the output of this network should be a tensor with shape `(batch_size, 1, 32, 32)`, and the values in it are float numbers in  $(-1, 1)$ . Our autograder will only be able to check the shape and range of the output, so be careful even if you have passed the test.
- (c) In class we have learned that to jointly train the generator and discriminator, we optimize them with a minimax objective:

$$V(G, D) := \frac{1}{N} \sum_{i=1}^N \log D(\mathbf{x}_i) + \frac{1}{N} \sum_{j=1}^N \log(1 - D(G(\mathbf{z}_j)))$$

$$\min_G \max_D V(G, D)$$

Here  $N$  is the batch size (set to 64 in our implementation),  $\mathbf{x}_i$  is a real image,  $\mathbf{z}_j$  is a random latent variable sampled from a Gaussian distribution, and  $G(\mathbf{z}_j)$  is a fake image generated from it. Note that we are taking average to approximate the expectation, since we are using SGD to optimize  $G$  and  $D$ .

Please complete the function `calculate_V()` in `GAN`. You may (but not required to) use the binary cross entropy loss (see <https://pytorch.org/docs/stable/generated/torch.nn.BCELoss.html>) to simplify the implementation, but be careful about the sign and reduction method of BCELoss.

- (d) We are ready to start training our GAN. The training pipeline is already provided in `train()`, and there is a `visualize()` function for your convenience. Train our GAN for 10 epochs, and **include the generated images after training in your PDF submission**.

Notes from TA:

- Training 10 epochs takes me about an hour on my laptop without GPU support. I can see interesting images after two or three epochs.
- You can make use of Google Colab(<https://colab.research.google.com/>), where you can access GPUs freely and accelerate the training. Remember to set `Runtime->Change runtime type->Hardware accelerator`.
- Some random seeds may lead to degenerated results. It's OK to try a few and manually set your random seed (`torch.manual.seed()`).

**Solution.**



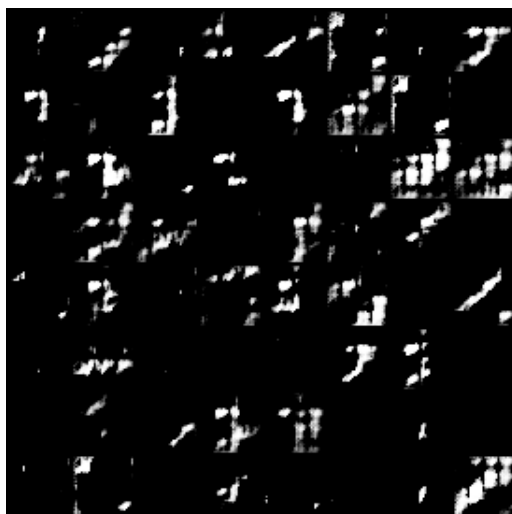


Figure 2: Q3(d):Epoch 0

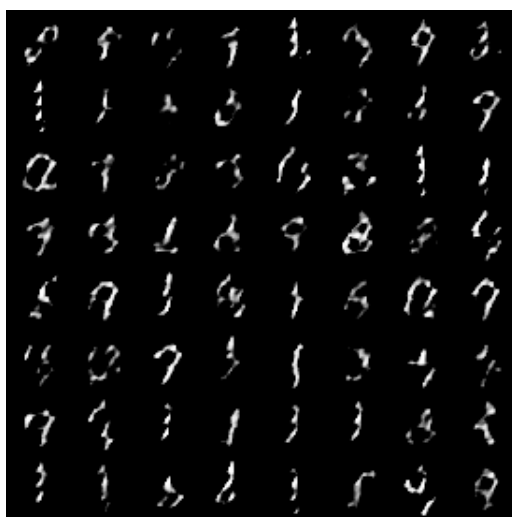


Figure 3: Q3(d):Epoch 1



Figure 4: Q3(d):Epoch 2

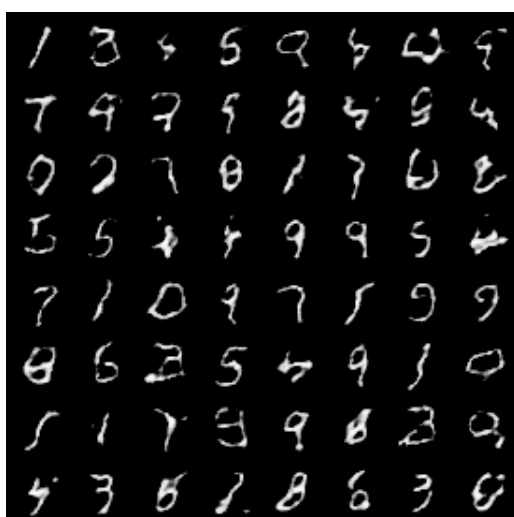


Figure 5: Q3(d):Epoch 3



Figure 6: Q3(d):Epoch 4



Figure 7: Q3(d):Epoch 5



Figure 8: Q3(d):Epoch 6



Figure 9: Q3(d):Epoch 7



Figure 10: Q3(d):Epoch 8



Figure 11: Q3(d):Epoch 9