

## Hena Ghonia (20213256)

**Question 1** (5-5-6). Consider a latent variable model  $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ , where  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$  and  $\mathbf{z} \in \mathbb{R}^K$ . The encoder network (aka “recognition model”) of variational autoencoder,  $q_\phi(\mathbf{z}|\mathbf{x})$ , is used to produce an approximate (variational) posterior distribution over latent variables  $\mathbf{z}$  for any input datapoint  $\mathbf{x}$ .<sup>1</sup> This distribution is trained to match the true posterior by maximizing the evidence lower bound (ELBO):

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi}[\log p_\theta(\mathbf{x} | \mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) || p(\mathbf{z}))$$

Let  $\mathcal{Q}$  be the family of variational distributions with a feasible set of parameters  $\mathcal{P}$ ; i.e.  $\mathcal{Q} = \{q(\mathbf{z}; \pi) : \pi \in \mathcal{P}\}$ ; for example  $\pi$  can be mean and standard deviation of a normal distribution. We assume  $q_\phi$  is parameterized by a neural network (with parameters  $\phi$ ) that outputs the parameters,  $\pi_\phi(\mathbf{x})$ , of the distribution  $q \in \mathcal{Q}$ , i.e.  $q_\phi(\mathbf{z}|\mathbf{x}) := q(\mathbf{z}; \pi_\phi(\mathbf{x}))$ .

1.1 Show that maximizing the expected complete data log likelihood (ECLL)

$$\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})]$$

for a fixed  $q(\mathbf{z}|\mathbf{x})$ , wrt the model parameter  $\theta$ , is equivalent to maximizing

$$\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$$

This means the maximizer of the ECLL coincides with that of the marginal likelihood only if  $q(\mathbf{z}|\mathbf{x})$  perfectly matches  $p(\mathbf{z}|\mathbf{x})$ .

**Solution:** As we know that Variational method are based on following relationship from slide no.17 from VAE.pdf(Lecture slides):

$$\log p_\theta(\mathbf{x}) = \mathcal{L}(q, p_\theta) + D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x})) \quad (1)$$

Where  $\mathcal{L}(q, p_\theta)$  is Evidence lower bound.

Defining variational lower bound on data likelihood as:

$$\mathcal{L}(q, p_\theta) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})]$$

where  $\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z})]$  is given as expected complete data log likelihood.

Substituting  $\mathcal{L}(q, p_\theta)$  in equation (1):

$$\log p_\theta(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] + D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$$

$$\log p_\theta(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z}|\mathbf{x})] + D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$$

$$\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x})) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z}|\mathbf{x})]$$

From above equation maximizing  $\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})]$  when  $q(\mathbf{z}|\mathbf{x})$  perfectly matches  $p(\mathbf{z}|\mathbf{x})$  is equivalent to maximizing  $\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$ .

1. Using a recognition model in this way is known as “amortized inference”; this can be contrasted with traditional variational inference approaches (see, e.g., Chapter 10 of Bishop’s *Pattern Recognition and Machine Learning*), which fit a variational posterior independently for each new datapoint.

It can also be proved in following way. As we know that  $q(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{z})$  are independent of parameter  $\theta$ , we can write in following form using the condition that  $q(\mathbf{z}|\mathbf{x})$  perfectly matches  $p(\mathbf{z}|x)$ :

$$\begin{aligned}
 & \arg \max_{\theta} \{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})] \} \\
 &= \arg \max_{\theta} \left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_{\theta}(\mathbf{x}|\mathbf{z})q(\mathbf{z}|\mathbf{x})} \right] \right\} \\
 &= \arg \max_{\theta} \left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z}|\mathbf{x})} \right] \right\} \\
 &= \arg \max_{\theta} \{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \} \quad [\because p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z}) = p_{\theta}(\mathbf{x}, \mathbf{z})] \\
 &= \arg \max_{\theta} \{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] \} \\
 &= \arg \max_{\theta} \{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}|\mathbf{x})p_{\theta}(\mathbf{x}) - \log q(\mathbf{z}|\mathbf{x})] \} \quad [\because p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{z}|\mathbf{x})p_{\theta}(\mathbf{x})] \quad (2) \\
 &= \arg \max_{\theta} \{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z}|\mathbf{x})p_{\theta}(\mathbf{x}) - \log q(\mathbf{z}|\mathbf{x})] \} \quad [\because q(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}, \mathbf{z})] \\
 &= \arg \max_{\theta} \{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x})] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z}|\mathbf{x})] \} \\
 &= \arg \max_{\theta} \{ \log p_{\theta}(\mathbf{x}) - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{q(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \} \\
 &= \arg \max_{\theta} \{ \log p_{\theta}(\mathbf{x}) - D_{KL}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \}
 \end{aligned}$$

1.2 Consider a finite training set  $\{\mathbf{x}_i : i \in \{1, \dots, n\}\}$ ,  $n$  being the size the training data. Let  $\phi^*$  be the maximizer  $\arg \max_{\phi} \sum_{i=1}^n \mathcal{L}(\theta, \phi; \mathbf{x}_i)$  with  $\theta$  fixed. In addition, for each  $\mathbf{x}_i$  let  $q_i \in \mathcal{Q}$  be an “instance-dependent” variational distribution, and denote by  $q_i^*$  the maximizer of the corresponding ELBO. Compare  $D_{KL}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_{\theta}(\mathbf{z}|\mathbf{x}_i))$  and  $D_{KL}(q_i^*(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x}_i))$ . Which one is bigger?

**Solution:** Inference gap of  $q^*$  can be written as from paper - ‘Inference Suboptimality in Variational Autoencoders’:

$$\log p_{\theta}(\mathbf{x}) - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x})] = \log p_{\theta}(\mathbf{x}_i) - \mathcal{L}[q_i^*(\mathbf{z})] + \mathcal{L}[q_i^*(\mathbf{z})] - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)]$$

where  $\log p_{\theta}(\mathbf{x}_i) - \mathcal{L}[q_i^*(\mathbf{z})]$  is Approximation gap and  $\mathcal{L}[q_i^*(\mathbf{z})] - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)]$  is Amortization gap. As we know that Variational method are based on following relationship from slide no.17 from VAE.pdf(Lecture slides):

$$\log p_{\theta}(\mathbf{x}) - D_{KL}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) = \mathcal{L}[q(\mathbf{z}|\mathbf{x})]$$

$$\log p_{\theta}(\mathbf{x}) - \mathcal{L}[q(\mathbf{z}|\mathbf{x})] = D_{KL}(q(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))$$

Replacing above form of formula in inference gap formula for  $\log p_{\theta}(\mathbf{x}) - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x})]$  and  $\log p_{\theta}(\mathbf{x}_i) - \mathcal{L}[q_i^*(\mathbf{z})]$ .

$$\log p_{\theta}(x) - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x})] = D_{KL}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_{\theta}(\mathbf{z}|\mathbf{x}_i))$$

$$\log p_{\theta}(x_i) - \mathcal{L}[q_i^*(\mathbf{z})] = D_{KL}(q_i^*(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x}_i))$$

Replacing above equations in inference gap formula:

$$D_{KL}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_{\theta}(\mathbf{z}|\mathbf{x}_i)) = D_{KL}(q^*(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x}_i)) + (\mathcal{L}[q_i^*(\mathbf{z})] - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)])$$

As mentioned in question that  $q_i^*$  is maximizer of the corresponding ELBO i.e from  $\mathcal{Q}$ ,  $(\mathcal{L}[q_i^*(\mathbf{z})] - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)])$  will be positive in above equation so following form can be written:

$$D_{KL}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_{\theta}(\mathbf{z}|\mathbf{x}_i)) \geq D_{KL}(q^*(\mathbf{z})||p_{\theta}(\mathbf{z}|\mathbf{x}_i))$$

1.3 Following the previous question, compare the two approaches in the second subquestion.

- (a) in terms of bias of estimating the marginal likelihood via the ELBO, in the best case scenario (i.e. when both approaches are optimal within the respective families)

**Solution:** There is bias since we use Variational Inference method (given in lecture slide-9 of VAE.pdf) which has KL divergence and expected data log likelihood is the term in ELBO that depends on  $\theta$ .

- (b) from the computational point of view (efficiency)

**Solution:** Approximating  $q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)$  is more computational efficient than  $q_i(\mathbf{z})$  as we don't need to update  $q_i$  for each  $x_i$  in finite training set.

- (c) in terms of memory (storage of parameters)

**Solution:** Approximating  $q_i^*(\mathbf{z})$  needs more parameter space than  $q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)$  since  $q_i^*(\mathbf{z})$  needs to store parameter at each  $i$ th iteration and hence has  $n$  times more parameter than amortized approximation.

**Question 2** (5-5-5-5). One way to enforce autoregressive conditioning is via masking the weight parameters.<sup>2</sup> Consider a two-hidden-layer convolutional neural network without kernel flipping, with kernel size  $3 \times 3$  and padding size 1 on each border (so that an input feature map of size  $5 \times 5$  is convolved into a  $5 \times 5$  output). Define mask of type A and mask of type B as

$$(\mathbf{M}^A)_{::ij} := \begin{cases} 1 & \text{if } i < 2 \\ 1 & \text{if } i = 2 \text{ and } j < 2 \\ 0 & \text{elsewhere} \end{cases} \quad (\mathbf{M}^B)_{::ij} := \begin{cases} 1 & \text{if } i < 2 \\ 1 & \text{if } i = 2 \text{ and } j \leq 2 \\ 0 & \text{elsewhere} \end{cases}$$

where the index starts from 1. Masking is achieved by multiplying the kernel with the binary mask (elementwise). Specify the receptive field of the output pixel that corresponds to the third row and the third column (index 33 of Figure 1 (Left)) in each of the following 4 cases:

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

FIGURE 1 – (Left)  $5 \times 5$  convolutional feature map. (Right) Template answer.

2. An example of this is the use of masking in the Transformer architecture.

1. If we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^A$  for the second layer.

**Solution:** Figure 2 (Left)) shows if we use  $\mathbf{M}^A$  for the first layer and (Right) shows result after we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^A$  for the second layer.

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

FIGURE 2 – (Left) After using  $\mathbf{M}^A$  for the first layer. (Right) After using  $\mathbf{M}^A$  for the second layer.

2. If we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^B$  for the second layer.

**Solution:** Figure 3 (Left)) shows if we use  $\mathbf{M}^A$  for the first layer and (Right) shows result after we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^B$  for the second layer.

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

FIGURE 3 – (Left) After using  $\mathbf{M}^A$  for the first layer. (Right) After using  $\mathbf{M}^B$  for the second layer.

3. If we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^A$  for the second layer.

**Solution:** Figure 4 (Left)) shows if we use  $\mathbf{M}^B$  for the first layer and (Right) shows result after we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^A$  for the second layer.

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

FIGURE 4 – (Left) After using  $\mathbf{M}^B$  for the first layer. (Right) After using  $\mathbf{M}^A$  for the second layer.

4. If we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^B$  for the second layer.

**Solution:** After using  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^B$  for the second layer, we get receptive field as Figure 3 (Right).

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

11	12	13	14	15
21	22	23	24	25
31	32	33	34	35
41	42	43	44	45
51	52	53	54	55

FIGURE 5 – (Left) After using  $\mathbf{M}^B$  for the first layer. (Right) After using  $\mathbf{M}^B$  for the second layer.

**Question 3** (6-4). In this question, we study some properties of normalizing flows. Let  $X \sim P_X$  and  $U \sim P_U$  be, respectively, the distribution of the data and a base distribution (e.g. an isotropic gaussian). We define a normalizing flow as  $F : \mathcal{U} \rightarrow \mathcal{X}$  parametrized by  $\theta$ . Starting with  $P_U$  and then applying  $F$  will induce a new distribution  $P_{F(U)}$  (used to match  $P_X$ ). Since normalizing flows are invertible, we can also consider the distribution  $P_{F^{-1}(X)}$ .

However, some flows, like planar flows, are not easily invertible in practice. If we use  $P_U$  as the base distribution, we can only sample from the flow but not evaluate the likelihood. Alternatively, if we use  $P_X$  as the base distribution, we can evaluate the likelihood, but we will not be able to sample.

3.1 Show that  $D_{KL}[P_X||P_{F(U)}] = D_{KL}[P_{F^{-1}(X)}||P_U]$ . In other words, the forward KL divergence between the data distribution and its approximation can be expressed as the reverse KL divergence between the base distribution and its approximation.

**Solution:** Using inverse function theorem for probability change of variable and change of variable for below integral:

$$\begin{aligned} D_{KL}[P_X(\mathbf{x})||P_{F(U)}(\mathbf{x})] &= \int_{\mathcal{X}} p_X(\mathbf{x}) \log \left( \frac{p_X(\mathbf{x})}{p_{F(U)}(\mathbf{x})} \right) d\mathbf{x} \\ &= \int_{\mathcal{X}} p_X(\mathbf{x}) \log \left( \frac{p_X(\mathbf{x})}{p_X(F^{-1}(\mathbf{x})) |\det J_{F^{-1}}(\mathbf{x})|} \right) d\mathbf{x} \\ &= \int_{F^{-1}(\mathcal{X})} p_X(F(\mathbf{u})) |\det J_F(\mathbf{u})| \log \left( \frac{p_X(F(\mathbf{u})) |\det J_F(\mathbf{u})|}{p_U(\mathbf{u})} \right) d\mathbf{u} \quad (3) \\ &= \int_{F^{-1}(\mathcal{X})} p_{F^{-1}(X)}(\mathbf{u}) \log \frac{p_{F^{-1}(X)}(\mathbf{u})}{p_U(\mathbf{u})} d\mathbf{u} \\ &= D_{KL}[P_{F^{-1}(X)}(\mathbf{u})||P_U(\mathbf{u})] \end{aligned}$$

Hence we prove  $D_{KL}[P_X||P_{F(U)}] = D_{KL}[P_{F^{-1}(X)}||P_U]$ . If  $f$  is invertible, forward divergence between the data distribution and its approximation can be expressed as the reverse KL divergence between the base distribution and its approximation i.e  $D_{KL}[P_X||P_Y] = D_{KL}[f(P_X)||f(P_Y)]$ .

3.2 Suppose two scenario: 1) you don't have samples from  $p_X(\mathbf{x})$ , but you can evaluate  $p_X(\mathbf{x})$ , 2) you have samples from  $p_X(\mathbf{x})$ , but you cannot evaluate  $p_X(\mathbf{x})$ . For each scenario, specify if you would use the forward KL divergence  $D_{KL}[P_X||P_{F(U)}]$  or the reverse KL divergence  $D_{KL}[P_{F(U)}||P_X]$  as the objective to optimize. Justify your answer.

**Solution:** For scenario 1 when we don't have samples from  $p_X(\mathbf{x})$ , but can evaluate  $p_X(\mathbf{x})$ ; reverse KL divergence is used.

Expanding reverse KL divergence:

$$\begin{aligned} \mathcal{L}(\theta) &= D_{KL}[P_{F(U)}||P_X] \\ &= \mathbb{E}_{p_{F(U)}(\mathbf{x})} [\log p_{F(U)}(\mathbf{x}; \theta) - \log p_X(\mathbf{x})] \quad (4) \\ &= \mathbb{E}_{p_U(\mathbf{u})} [\log p_U(\mathbf{u}) - \log |J_F(\mathbf{u})| - \log p_X(F(\mathbf{u}; \theta))] \end{aligned}$$

As we can see from above equation, for reverse KL divergence we only need to evaluate  $p_X(x)$  and don't need samples from  $p_X(x)$ .

For scenario 2 when we have samples from  $p_X(x)$ , but cannot evaluate  $p_X(x)$ ; forward KL divergence is used.

Expanding forward KL divergence:

$$\begin{aligned}\mathcal{L}(\theta) &= D_{KL}[P_X || P_{F(U)}] \\ &= \mathbb{E}_{p_X(\mathbf{x})} [\log p_X(\mathbf{x}) - \log p_{F(U)}(\mathbf{x}; \theta)] \\ &= \mathbb{E}_{p_X(\mathbf{x})} [-\log |J_{F^{-1}}(\mathbf{x})| - \log p_U(F^{-1}(\mathbf{x}; \theta))] + \text{const.}\end{aligned}\tag{5}$$

As we can see in above equation that expectation is w.r.t  $p_X(x)$  which means we need samples from  $p_X(x)$  but no need to evaluate  $\log p_X(\mathbf{x})$ .

**Question 4 (4-3-6).** In this question, we are concerned with analyzing the training dynamics of GANs. Consider the following value function

$$V(d, g) = dg \tag{6}$$

with  $g \in \mathbb{R}$  and  $d \in \mathbb{R}$ . We will use this simple example to study the training dynamics of GANs.

1. Consider gradient descent/ascent with learning rate  $\alpha$  as the optimization procedure to iteratively minimize  $V(d, g)$  w.r.t.  $g$  and maximize  $V(d, g)$  w.r.t.  $d$ . We will apply the gradient descent/ascent to update  $g$  and  $d$  simultaneously. What is the update rule of  $g$  and  $d$ ? Write your answer in the following form

$$[d_{k+1}, g_{k+1}]^\top = A[d_k, g_k]^\top$$

where  $A$  is a  $2 \times 2$  matrix; i.e. specify the value of  $A$ .

**Solution:** To perform gradient descent on  $V(d, g)$  w.r.t.  $g$ (generator) since we need to minimize  $V(d, g)$ . Update rule for  $g$ :

$$g_{t+1} = g_t - \alpha \frac{\partial V(d, g)}{\partial g} = g_k - \alpha \frac{\partial (dg)}{\partial g} = g_k - \alpha \cdot d_k$$

To perform gradient ascent on  $V(d, g)$  w.r.t.  $d$ (discriminator) since we need to maximize  $V(d, g)$ . Update rule for  $d$ :

$$d_{t+1} = d_t + \alpha \frac{\partial V(d, g)}{\partial d} = d_k + \alpha \frac{\partial (dg)}{\partial d} = d_k + \alpha \cdot g_k$$

. To update rule of  $d$  and  $g$  in following form:

$$\begin{aligned}[d_{k+1}, g_{k+1}]^\top &= A[d_k, g_k]^\top \\ \begin{bmatrix} d_{k+1} \\ g_{k+1} \end{bmatrix} &= \begin{bmatrix} 1 & \alpha \\ -\alpha & 1 \end{bmatrix} \begin{bmatrix} d_k \\ g_k \end{bmatrix} \\ [d_{k+1}, g_{k+1}]^\top &= \begin{bmatrix} 1 & \alpha \\ -\alpha & 1 \end{bmatrix} [d_k, g_k]^\top \\ \text{where } A &= \begin{bmatrix} 1 & \alpha \\ -\alpha & 1 \end{bmatrix}\end{aligned}$$

2. The optimization procedure you found in 4.1 characterizes a map which has a stationary point<sup>3</sup>, what are the coordinates of the stationary points?

**Solution:** Stationary points are those where gradients of  $\nabla V(d, g)$  is zero.

$$\nabla V(d, g) = \begin{bmatrix} \frac{\partial V(d, g)}{\partial d} \\ \frac{\partial V(d, g)}{\partial g} \end{bmatrix} = \begin{bmatrix} g \\ d \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

The coordinates of stationary points  $(g, d) = (0, 0)$

3. Analyze the eigenvalues of A and predict what will happen to  $d$  and  $g$  as you update them jointly. In other word, predict the behaviour of  $d_k$  and  $g_k$  as  $k \rightarrow \infty$ .

**Solution:** To analyze eigenvalues of A:

$$\det(A - \lambda I) = 0$$

$$\det \left( \begin{bmatrix} 1 & \alpha \\ -\alpha & 1 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) = 0$$

$$\det \left( \begin{bmatrix} 1 - \lambda & \alpha \\ -\alpha & 1 - \lambda \end{bmatrix} \right) = 0$$

$$(1 - \lambda)^2 + \alpha^2 = 0$$

$$(1 - \lambda)^2 = -\alpha^2$$

$$(\lambda - 1)^2 = -\alpha^2 \Rightarrow (\lambda - 1) = \pm \sqrt{-\alpha^2}$$

As  $(1 - \lambda)^2 = (\lambda - 1)^2$  and suppose  $i = \sqrt{-1}$

$$\lambda - 1 = \pm i\alpha$$

$$\lambda = 1 \pm i\alpha$$

So let  $\lambda_1 = 1 - i\alpha$  and  $\lambda_2 = 1 + i\alpha$ . Eigenvalues with complex number has rotation so updates will be performed by  $\lambda_1$  and  $\lambda_2$  alternatively in parameter space and will never converge to a point as  $k \rightarrow \infty$ .

**Question 5** (4-2-8-4-2). In this question, we will see why stop-gradient is critical for non-contrastive SSL methods like SimSiam and BYOL. We will show that removing stop-gradient results in collapsed representations, using the dynamics of SimSiam as our running example.

Consider a two-layer linear SimSiam model with the time-varying weight matrices given by  $W(t) \in \mathbb{R}^{n_2 \times n_1}$  and  $W_p(t) \in \mathbb{R}^{n_2 \times n_2}$ . Note that  $W(t)$  corresponds to the weights of the online **and** the target network, while  $W_p(t)$  denotes the weights of the predictor. Let  $\mathbf{x} \in \mathbb{R}^{n_1}$  be an input datapoint and  $\mathbf{x}_1, \mathbf{x}_2$  be the two augmented versions of the input  $\mathbf{x}$ . Also note that in some instances, the dependence on time ( $t$ ) is omitted for notational simplicity, and the weight matrices are referred to as  $W$  and  $W_p$ .

---

3. A stationary point is a point on the surface of the graph (of the function) where all its partial derivatives are zero (equivalently, the gradient is zero). Source: [https://en.wikipedia.org/wiki/Stationary\\_point](https://en.wikipedia.org/wiki/Stationary_point)

Let  $\mathbf{f}_1 = W\mathbf{x}_1$  be the online representation of  $\mathbf{x}_1$  and  $\mathbf{f}_2 = W\mathbf{x}_2$  be the target representation of  $\mathbf{x}_2$ . The learning dynamics of  $W$  and  $W_p$  can be obtained by minimizing SimSiam's objective function as shown below:

$$J(W, W_p) = \frac{1}{2} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\|W_p \mathbf{f}_1 - \text{Stop-Grad}(\mathbf{f}_2)\|_2^2]. \quad (7)$$

5.1 Show (with proof) that the above objective can be simplified to:

$$J(W, W_p) = \frac{1}{2} [\text{tr}(W_p^T W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(F_{12} W_p) + \text{tr}(F_2)], \quad (8)$$

where  $F_1 = \mathbb{E}[\mathbf{f}_1 \mathbf{f}_1^T] = W(X + X')W^T$ ,  $F_2 = \mathbb{E}[\mathbf{f}_2 \mathbf{f}_2^T] = W(X + X')W^T$ , and  $F_{12} = F_{21} = \mathbb{E}[\mathbf{f}_1 \mathbf{f}_2^T] = W X W^T$ . Here,  $X$  is the average augmented view of a datapoint  $\mathbf{x}$  and  $X'$  is the covariance matrix of augmented views  $\mathbf{x}'$  conditioned on  $\mathbf{x}$  and then averaged over the data  $\mathbf{x}$ , and  $\text{tr}$  is the Trace operation<sup>4</sup>.

**Solution:**

Consider that :

$$\begin{aligned} (W_p \mathbf{f}_1 - \mathbf{f}_2)^T (W_p \mathbf{f}_1 - \mathbf{f}_2) &= \mathbf{f}_1^T W_p^T W_p \mathbf{f}_1 - \mathbf{f}_1^T W_p^T \mathbf{f}_2 - \mathbf{f}_2^T W_p \mathbf{f}_1 + \mathbf{f}_2^T \mathbf{f}_2 \\ &= \text{tr}(W_p^T W_p \mathbf{f}_1 \mathbf{f}_1^T) - \text{tr}(W_p \mathbf{f}_1 \mathbf{f}_2^T) - \text{tr}(W_p^T \mathbf{f}_2 \mathbf{f}_1^T) + \text{tr}(\mathbf{f}_2 \mathbf{f}_2^T) \end{aligned} \quad (9)$$

$$J(W, W_p) = \frac{1}{2} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\|W_p \mathbf{f}_1 - \text{Stop-Grad}(\mathbf{f}_2)\|_2^2]$$

$$\begin{aligned} J(W, W_p) &= \frac{1}{2} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [(W_p \mathbf{f}_1 - \mathbf{f}_2)^T (W_p \mathbf{f}_1 - \mathbf{f}_2)] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\text{tr}(W_p^T W_p \mathbf{f}_1 \mathbf{f}_1^T) - \text{tr}(W_p \mathbf{f}_1 \mathbf{f}_2^T) - \text{tr}(W_p^T \mathbf{f}_2 \mathbf{f}_1^T) + \text{tr}(\mathbf{f}_2 \mathbf{f}_2^T)] \end{aligned} \quad (10)$$

Let  $F_1 = \mathbb{E}[\mathbf{f}_1 \mathbf{f}_1^T] = W(X + X')W^T$ ,  $F_2 = \mathbb{E}[\mathbf{f}_2 \mathbf{f}_2^T] = W(X + X')W^T$ , and  $F_{12} = F_{21} = \mathbb{E}[\mathbf{f}_1 \mathbf{f}_2^T] = W X W^T$  and  $X$  is the average augmented view of a datapoint  $\mathbf{x}$  and  $X'$  is the covariance matrix of augmented views  $\mathbf{x}'$  conditioned on  $\mathbf{x}$  and then averaged over the data  $\mathbf{x}$ , and  $\text{tr}$  is the Trace.

$$J(W, W_p) = \frac{1}{2} [\text{tr}(W_p^T W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(F_{12} W_p) + \text{tr}(F_2)]$$

Reference: <https://arxiv.org/pdf/2102.06810.pdf>

5.2 Based on the above expression for  $J(W, W_p)$ , find the gradient update for  $W_p$  (the predictor network), denoting it as  $\dot{W}_p$ . In other words, obtain an expression for  $\dot{W}_p = -\frac{\partial J}{\partial W_p}$  (the derivative of the objective function w.r.t the parameters  $W_p$ ).

**Solution:** Taking partial derivative of the objective function w.r.t to parameter  $W_p$ :

$$\dot{W}_p = -\frac{\partial J}{\partial W_p}$$

4. [https://en.wikipedia.org/wiki/Trace\\_\(linear\\_algebra\)](https://en.wikipedia.org/wiki/Trace_(linear_algebra)).



$$\begin{aligned}
 \dot{W}_p &= -\frac{1}{2} \frac{\partial [\text{tr}(W_p^\top W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(F_{12} W_p) + \text{tr}(F_2)]}{\partial W_p} \\
 &= -\frac{1}{2} [W_p F_1 + W_p F_1^\top - F_{12}^\top - F_{12}] \\
 &= -W_p F_1 + F_{12}^\top
 \end{aligned} \tag{11}$$

5.3 Consider the case when the Stop-Grad is removed. The gradient of the objective function  $J(W, W_p)$  w.r.t the parameters  $W$  i.e.  $\dot{W}(t) = -\frac{\partial J}{\partial W(t)}$ , is given by:

$$\dot{W}(t) = \frac{d}{dt} \text{vec}(W(t)) = -H(t) \text{vec}(W(t)),$$

where  $H(t)$  is a time-varying positive semi-definite matrix defined as

$$H(t) = X' \otimes (W_p(t)^\top W_p(t) + I_{n_2}) + X \otimes (\tilde{W}_p(t)^\top \tilde{W}_p(t)).$$

Here,  $\otimes$  is the Kronecker product<sup>5</sup>,  $\tilde{W}_p(t) = (W_p(t) - I_{n_2})$ , and "vec(W)" refers to the *vectorization* of a matrix W<sup>6</sup>. For simplicity, we are not taking weight decay into account here<sup>7</sup>.

If the minimal eigenvalue  $\lambda_{\min}(H(t))$  is bounded away from zero, i.e.  $\inf_{t \geq 0} \lambda_{\min}(H(t)) \geq \lambda_0 > 0$ , then **prove that**  $W(t) \rightarrow 0$ .

**Note:** In order to prove the above question, the following property must be used:

For a time-varying positive definite matrix  $H(t)$  whose minimal eigenvalues are bounded away from 0, the dynamics shown below:

$$\frac{d}{dt} \mathbf{w}(t) = -H(t) \mathbf{w}(t),$$

satisfies the constraint  $\|\mathbf{w}(t)\|_2 = e^{-\lambda_0 t} \|\mathbf{w}(0)\|_2$ , implying that  $\mathbf{w}(t) \rightarrow 0$ .

**Solution:** If we don't have stop gradient then we need to calculate  $\frac{\partial J}{\partial F_2}$  for obtaining gradient of the objective function.

Replacing value of  $F_1 = W(X + X')W^\top$ ,  $F_2 = W(X + X')W^\top$  and  $F_{12} = W X W^\top$  in objective function and calculating gradient:

$$\begin{aligned}
 \dot{W}(t) &= -\frac{dJ}{d(W(t))} \\
 &= -\frac{1/2 * [\text{tr}(W_p^\top W_p W(X + X')W^\top) - \text{tr}(W_p W X W^\top) - \text{tr}(W X W^\top) + \text{tr}(W(X + X')W^\top)]}{d(W(t))} \\
 &= [-W_p^\top W_p W(X + X') + (W_p^\top + W_p) W X - W(X + X') - W] \\
 &\quad - (W_p^\top W_p + I) W X' - (W_p^\top W_p - W_p^\top - W_p + I) W X - W \\
 &= -(W_p^\top W_p + I) W X' - (W_p - I)^\top (W_p - I) W X - W \\
 &= -(W_p^\top W_p + I) W X' - \tilde{W}_p^\top \tilde{W}_p W X - W \left[ \because \tilde{W}_p = W_p - I_{n_2} \right]
 \end{aligned} \tag{12}$$

5. For more information, see [https://en.wikipedia.org/wiki/Kronecker\\_product#Matrix\\_equations](https://en.wikipedia.org/wiki/Kronecker_product#Matrix_equations)

6. Also known as the "vec trick", it is obtained by stacking all the columns of a matrix A into a single vector.

7. Although omitted here, it must be noted that having weight decay is important. It has also been shown that, in practice, weight decay leads to stable learning.

For  $\text{vec}(A \times B) = (B^\top \otimes A)\text{vec}(X)$ . Using this we write for  $\frac{d}{dt}\text{vec}(W)$ :

$$\begin{aligned}\frac{d}{dt}\text{vec}(W(t)) &= -\left(X' \otimes (W_p^\top W_p + I_{n_2}) + X \otimes \tilde{W}_p^\top \tilde{W}_p + I_{n_1 n_2}\right) \text{vec}(W) \\ &= -\left(X' \otimes (W_p^\top W_p + I_{n_2}) + X \otimes \tilde{W}_p^\top \tilde{W}_p\right) \text{vec}(W)\end{aligned}\quad (13)$$

$$\frac{d}{dt}\text{vec}(W(t)) = H(t)\text{vec}(W(t)) \left[ \cdot : H(t) = \left(X' \otimes (W_p^\top W_p + I_{n_2}) + X \otimes \tilde{W}_p^\top \tilde{W}_p\right) \right]$$

Constraint given in question  $\|\mathbf{w}(t)\|_2 = e^{-\lambda_0 t} \|\mathbf{w}(0)\|_2$ , implying that  $\mathbf{w}(t) \rightarrow 0$  for dynamics  $\frac{d}{dt}\mathbf{w}(t) = -H(t)\mathbf{w}(t)$ . Now using this property for  $\frac{d}{dt}W(t) = -H(t)W(t)$  where  $\inf_{t \geq 0} \lambda_{\min}(H(t)) \geq \lambda_0 > 0$  we have:

$$\|\text{vec}(W(t))\|_2^2 \leq e^{-\lambda_0 t} \|\text{vec}(W(0))\|_2^2 \rightarrow 0$$

which proves  $W(t) \rightarrow 0$  and  $W$  has no possibility of learning any useful features.

5.4 Consider the case when both the Stop-Grad **and** the predictor are removed. Show that the representations collapse i.e.  $W(t) \rightarrow 0$ . You may assume that  $X'$  is a positive definite matrix.

**Solution:** When both the Stop - Grad and predictor are removed i.e  $W_p = I$ ,  $W(t)$  always tends to zero. Previous question already proved without Stop-gradient it the case. But if there is Stop gradient:

$$\begin{aligned}\dot{W} &= -(W_p^\top W_p + I)W X' - \tilde{W}_p^\top \tilde{W}_p W X - W \\ &= -W X' - W [\cdot : W_p = I] \\ &= -W (X' + I)\end{aligned}\quad (14)$$

$X'$  is positive definite matrix and using property from previous question:

$$\|\text{vec}(W(t))\|_2^2 \leq e^{-\lambda_0 t} \|\text{vec}(W(0))\|_2^2 \rightarrow 0$$

Hence  $W(t) \rightarrow 0$

5.5 Speculate (in 1-2 sentences) as to why the stop-gradient and the predictor are necessary for avoiding representational collapse.

**Solution:** As seen in previous two question, we have proven analytically that removing either stop-gradient and the predictor leads to representational collapse. Without stopgradient, the optimizer quickly finds a degenerated solution and reaches the minimum possible loss of -1.

Recent research says that its due to asymmetrical structure of SimSiam where the role of stop gradient is to only allow the path with predictor to be optimized with the encoder output as the target, not vice versa.

Reference-

Chen, Xinlei, and Kaiming He. "Exploring simple siamese representation learning." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.

-Zhang, Chaoning, et al. "How Does SimSiam Avoid Collapse Without Negative Samples? A Unified Understanding with Self-supervised Contrastive Learning." arXiv preprint arXiv:2203.16262 (2022).

-Tian, Yuandong, Xinlei Chen, and Surya Ganguli. "Understanding self-supervised learning dynamics without contrastive pairs." International Conference on Machine Learning. PMLR, 2021.