

**Due Date: April 15th, 2022, 11 p.m. EST**

Instructions

- For all questions, show your work!
- Please use a document preparation system such as LaTeX.
- Unless noted that questions are related, assume that notation and definitions for each question are self-contained and independent.
- Submit your answers electronically via Gradescope.
- TAs for this assignment are **Mohammad-Javad Bayazi** and **Naga Karthik**.

This assignment covers mathematical and algorithmic techniques in the families of deep generative models and some of the recent self-supervised learning methods. Thus, we will explore Variational autoencoders (VAEs, Question 1), Autoregressive models (Question 2), Normalizing flows (Question 3), Generative adversarial networks (GANs, Question 4), and Self-supervised models (Question 5).

**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  $p(\mathbf{z}|\mathbf{x})$ .

**Answer 1.1**

We have  $q(\mathbf{z}|\mathbf{x})$  does not depend on  $\theta$ , then :

$$\begin{aligned} \arg \max_{\theta} \left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})] \right\} &= \arg \max_{\theta} \left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \frac{p_\theta(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\mathbf{x})} \right] \right\} \\ &= \arg \max_{\theta} \left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} + \log \frac{p_\theta(\mathbf{z}|\mathbf{x})}{q(\mathbf{z}|\mathbf{x})} \right] \right\} \end{aligned}$$

---

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.

Using Bayes theorem, we have :  $\frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} = p_\theta(\mathbf{x})$ , replacing this in the equation :

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

- 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  $D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)||p_\theta(\mathbf{z}|\mathbf{x}_i))$  and  $D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x}_i))$ . Which one is bigger ?

**Answer 1.2**

We have :

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

From the lectures notes, we can write :

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

Replacing in the equality, we get :

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

Since  $q_i^*$  is the maximizer amongst  $\mathcal{Q}$ , the term  $(\mathcal{L}[q_i^*(\mathbf{z})] - \mathcal{L}[q_{\phi^*}(\mathbf{z}|\mathbf{x}_i)])$  is positive.

Thus :

$$D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x})) \leq D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$$

- 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)

**Answer (a)**

The bias due to the KL divergence is defined as :  $\text{bias}(\theta) = -D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$ , so from the question 1.2 :

$$D_{\text{KL}}(q_i^*(\mathbf{z})||p_\theta(\mathbf{z}|\mathbf{x})) \leq D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \implies \text{bias}_{x_i, \phi^*}(\theta) \leq \text{bias}_{x_i, q_i^*}(\theta) \leq 0$$

From that the bias is closer to zero when we estimate the optimal  $q_i^*$  compared to when we calculate  $\phi^*$  for recognition model. Thus the optimal  $q_i^*$  can allow us to have a more accurate estimation of the marginal likelihood via the ELBO.

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

**Answer (b)**

Calculating  $q_i^*$  for each data point  $\mathbf{x}_i$  might be computationally heavy, while the recognition model does not need to update  $q_i$  for each data point  $\mathbf{x}_i$  until getting a good approximation. Thus approximating  $q_{\phi^*}$  is more efficient.

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

**Answer (c)**

Approximating  $q_i^*$  needs more memory storage since we store the parameters for each training data point  $\mathbf{x}_i$ , which means  $n$  times more expensive than the amortization.

**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.

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

Your answer should look like Figure 1 (Right).

**Answer 2**

Since the padding size is equal to 1 and without kernel flipping, we can see from figure 3 the receptive field in feature map obtained after the first layer. We only consider pixels that actually

$k_{11}$	$k_{12}$	$k_{13}$
$k_{21}$	0	0
0	0	0

$k_{11}$	$k_{12}$	$k_{13}$
$k_{21}$	$k_{22}$	0
0	0	0

FIGURE 2 – (Left) Kernel obtained when apply the mask of type A. (Right) Kernel obtained when apply the mask of type B.

<sup>2</sup>. An example of this is the use of masking in the Transformer architecture.

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 – Receptive field in feature map 1 (Left) if we use  $\mathbf{M}^A$ . (Right) if we use  $\mathbf{M}^B$ .

participate in the calculation of pixel of index 33 in the output. The kernel masked by type A has 5 inactive pixels and 4 active pixels, while the kernel masked by type B has 4 inactive pixels and 5 active pixels. If we use mask of type A in the second layer, the value of the pixel of index 33 from the output would be calculated as :  $k_{11} * fmap_{22} + k_{12} * fmap_{23} + k_{13} * fmap_{24} + k_{21} * fmap_{32}$ .

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

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 – Receptive field if we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^A$  for the second layer.

2. If we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^B$  for the second layer 5 :

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 – Receptive field if we use  $\mathbf{M}^A$  for the first layer and  $\mathbf{M}^B$  for the second layer.

3. If we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^A$  for the second layer 6 :

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 6 – Receptive field if we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^A$  for the second layer.

4. If we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^B$  for the second layer 7 :

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 7 – Receptive field if we use  $\mathbf{M}^B$  for the first layer and  $\mathbf{M}^B$  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.

**Answer 3.1**

We have :

$$D_{KL}[p_X(\mathbf{x})||p_{F(U)}(\mathbf{x})] = \int_{\mathcal{X}} p_X(\mathbf{x}) \log \frac{p_X(\mathbf{x})}{p_{F(U)}(\mathbf{x})} d\mathbf{x}$$

Using the change of variable formula, we have :  $p_{F(U)}(\mathbf{x}) = p_U(F^{-1}(\mathbf{x}))|\det J_{F^{-1}}(\mathbf{x})|$  Replacing this in the equation :

$$\begin{aligned} D_{KL}[p_X(\mathbf{x})||p_{F(U)}(\mathbf{x})] &= \int_{\mathcal{X}} p_X(\mathbf{x}) \log \frac{p_X(\mathbf{x})}{p_{F(U)}(\mathbf{x})} d\mathbf{x} \\ &= \int_{\mathcal{X}} p_X(\mathbf{x}) \log \frac{p_X(\mathbf{x})}{p_U(F^{-1}(\mathbf{x}))|\det J_{F^{-1}}(\mathbf{x})|} d\mathbf{x} \end{aligned}$$

Let's pose  $\mathbf{u} = F^{-1}(\mathbf{x})$  and using the change of variable formula, we have :

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

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.

**Answer 3.2**

For the forward KL divergence :

$$\begin{aligned}\mathcal{L}(\theta) &= D_{KL}[p_X(\mathbf{x})||p_{F(U)}(\mathbf{x})] \\ &= \mathbb{E}_{q_X(\mathbf{x})} \left[ \log \frac{p_X(\mathbf{x})}{p_{F(U)}(\mathbf{x}; \theta)} \right] \\ &= \mathbb{E}_{q_X(\mathbf{x})} [\log p_X(\mathbf{x}) - \log p_{F(U)}(\mathbf{x}; \theta)] \\ &= \mathbb{E}_{q_X(\mathbf{x})} [-\log p_{F(U)}(\mathbf{x}; \theta)] + \text{constant} \\ &= -\mathbb{E}_{q_X(\mathbf{x})} [\log p_U(F^{-1}(\mathbf{x}; \theta)) + \log |J_{F^{-1}}(\mathbf{x})|] + \text{constant}\end{aligned}$$

For the reverse KL divergence :

$$\begin{aligned}\mathcal{L}(\theta) &= D_{KL}[p_{F(U)}(\mathbf{x})||p_X(\mathbf{x})] \\ &= \mathbb{E}_{q_{F(U)}(\mathbf{x})} \left[ \log \frac{p_{F(U)}(\mathbf{x}; \theta)}{p_X(\mathbf{x})} \right] \\ &= \mathbb{E}_{q_{F(U)}(\mathbf{x})} [\log p_{F(U)}(\mathbf{x}; \theta) - \log p_X(\mathbf{x})] \\ &= \mathbb{E}_{q_U(\mathbf{u})} [\log p_U(\mathbf{u}) - \log |J_F(\mathbf{u})| - \log p_X(F(\mathbf{u}; \theta))]\end{aligned}$$

As we can see : the forward KL divergence only require  $\mathbf{x}$  thus we can use it for the second scenario. The reverse KL divergence only require to evaluate  $p_X$  thus we can use it for the first scenario.

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

$$V(d, g) = dg \tag{1}$$

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  $A$ .

**Answer 4.1**

Gradient ascent on  $V(d, g)$  w.r.t  $d$  :

$$d_{k+1} = d_k + \alpha \frac{\partial}{\partial d} V(d_k, g_k) = d_k + \alpha g_k$$

Gradient descent on  $V(d, g)$  w.r.t  $g$  :

$$g_{k+1} = g_k - \alpha \frac{\partial}{\partial g} V(d_k, g_k) = g_k - \alpha d_k$$

- Do not distribute -

Rewriting in matrix :

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

With :

$$\mathbf{A} = \begin{pmatrix} 1 & \alpha \\ -\alpha & 1 \end{pmatrix}$$

2. The optimization procedure you found in 6.1 characterizes a map which has a stationary point<sup>3</sup>, what are the coordinates of the stationary points?

**Answer 4.2**

A stationary point is a point on the surface of the graph where all its partial derivatives are zero (equivalently, the gradient is zero) :

$$\begin{aligned} \nabla V(d, g) &= \begin{pmatrix} \frac{\partial}{\partial d} V(d, g) \\ \frac{\partial}{\partial g} V(d, g) \end{pmatrix} \\ &= \begin{pmatrix} g \\ d \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{cases} g = 0 \\ d = 0 \end{cases} \end{aligned}$$

Thus the stationary point is at  $(g = 0, d = 0)$ .

3. Analyze the eigenvalues of  $\mathbf{A}$  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  $k \rightarrow \infty$ .

**Answer 4.3**

Let's find the eigenvalues of  $\mathbf{A}$  :

$$\begin{aligned} \det[\mathbf{A} - \lambda \mathbf{I}] = 0 &\Rightarrow \det \begin{vmatrix} 1 - \lambda & \alpha \\ -\alpha & 1 - \lambda \end{vmatrix} = 0 \\ &\Rightarrow (1 - \lambda)^2 + \alpha^2 = 0 \\ &\Rightarrow \lambda_1 = 1 + i\alpha \quad \text{and} \quad \lambda_2 = 1 - i\alpha \end{aligned}$$

As we can see  $\mathbf{A}$  has complex eigenvalues, and we know that complex eigenvalues imply rotation. So the updates will be rotating in the parameter space and will never converge to a balance 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:

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

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

$$J(W, W_p) = \frac{1}{2} [\text{tr}(W_p^T W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(F_{12} W_p) + \text{tr}(F_2)], \quad (3)$$

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  $F_{12} = F_{21} = \mathbb{E}[\mathbf{f}_1 \mathbf{f}_2^T] = WXW^T$ . Here,  $X$  is the average augmented view of a data point  $\mathbf{x}$  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>.

### Answer 5.1

We have :

$$\begin{aligned} \|W_p \mathbf{f}_1 - \text{Stop-Grad}(\mathbf{f}_2)\|_2^2 &= (W_p \mathbf{f}_1 - \mathbf{f}_2)^T (W_p \mathbf{f}_1 - \mathbf{f}_2) \\ &= W_p^T \mathbf{f}_1^T W_p \mathbf{f}_1 - W_p^T \mathbf{f}_1^T \mathbf{f}_2 - \mathbf{f}_2^T W_p \mathbf{f}_1 + \mathbf{f}_2^T \mathbf{f}_2 \\ &= \text{tr}(W_p^T W_p \mathbf{f}_1 \mathbf{f}_1^T) - \text{tr}(W_p^T \mathbf{f}_2 \mathbf{f}_1^T) - \text{tr}(W_p \mathbf{f}_1 \mathbf{f}_2^T) + \text{tr}(\mathbf{f}_2 \mathbf{f}_2^T) \end{aligned}$$

Thus we have :

$$\begin{aligned} J(W, W_p) &= \frac{1}{2} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\|W_p \mathbf{f}_1 - \text{Stop-Grad}(\mathbf{f}_2)\|_2^2] \\ &= \frac{1}{2} [\text{tr}(W_p^T W_p \mathbb{E}[\mathbf{f}_1 \mathbf{f}_1^T]) - \text{tr}(W_p \mathbb{E}[\mathbf{f}_1 \mathbf{f}_2^T]) - \text{tr}(\mathbb{E}[\mathbf{f}_2 \mathbf{f}_1^T] W_p) + \text{tr}(\mathbb{E}[\mathbf{f}_2 \mathbf{f}_2^T])] \\ &= \frac{1}{2} [\text{tr}(W_p^T W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(F_{12} W_p) + \text{tr}(F_2)] \end{aligned}$$

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  $W_p$ ).

### Answer 5.2

Taking the partial derivative w.r.t  $W_p$  and we get the gradient update rule :

$$\begin{aligned} \dot{W}_p &= -\frac{\partial J}{\partial W_p} \\ &= -W_p F_1 + F_{12} \end{aligned}$$

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:

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

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



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

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

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  $\mathbf{w}(t) \rightarrow 0$ .

### **Answer 5.3**

With  $\tilde{W}_p(t) = (W_p(t) - I_{n_2})$ , let's take the partial derivative w.r.t  $W$  :

$$\begin{aligned} \dot{W} &= -\frac{\partial J}{\partial W(t)} = -W_p^T W_p W (X + X') + (W_p^T + W_p) W X - W (X + X') \\ &= -(W_p^T W_p + I) W X' - (W_p^T W_p - W_p^T - W_p + I) W X \\ &= -(W_p^T W_p + I) W X' - (W_p - I)^T (W_p - I) W X \\ &= -(W_p^T W_p + I) W X' - \tilde{W}_p^T \tilde{W}_p W X \end{aligned}$$

Using :  $\text{vect}(AXB) = (B^T \otimes A) \text{vect}(X)$ , a property that connects vectorization and the Kronecker product, we have :

$$\begin{aligned} \dot{W} &= \frac{d}{dt} \text{vec}(W(t)) \\ &= - \left[ X' \otimes (W_p(t)^T W_p(t) + I_{n_2}) + X \otimes (\tilde{W}_p(t)^T \tilde{W}_p(t)) \right] \text{vec}(W(t)) \\ &= -H(t) \text{vec}(W(t)) \end{aligned}$$

So if the minimal eigenvalue  $\lambda_{\min}(H(t))$  is bounded away from zero :  $\inf_{t \geq 0} \lambda_{\min}(H(t)) \geq \lambda_0 > 0$ , then by using the property above, we have :  $\|W(t)\|_2 = e^{-\lambda_0 t} \|W(0)\|_2$ . Thus  $W(t) \rightarrow 0$ .

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.

### **Answer 5.4**

When both the Stop-Grad and the predictor are removed :  $W_p = I$ , replacing this in the partial derivative w.r.t  $W$ , we obtain a similar form as question 5.3 :

$$\dot{W} = -X' W$$

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.

By assuming that  $X'$  is a positive definite matrix, with a similar arguments as the question above, if  $\inf_{t \geq 0} \lambda_{\min}(X') \geq \lambda_0 > 0$ , we have :  $\|W(t)\|_2 = e^{-\lambda_0 t} \|W(0)\|_2$ . Thus  $W(t) \rightarrow 0$ .

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

**Answer 5.5**

Removing the stop-gradient and the predictor yields a gradient update for the parameters  $W(t)$ , and as we show in the last questions : if the minimal eigenvalue over time is bounded away from zero then  $W(t) \rightarrow 0$ . Thus there is no chance for  $W(t)$  to learn any meaningful features.