

**Due Date: April 12, 2020**

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 **Christos Tsirigotis** and **Philippe Brouillard**.

This assignment covers mathematical and algorithmic techniques underlying regularization and popular families of deep generative models. Thus, we explore regularization (Question 1), variational autoencoders (VAEs, Questions 2), normalizing flows (Question 3), and generative adversarial networks (GANs, Question 4-5).

**Question 1** (7-5-5-3). The point of this question is to understand and compare the effects of different regularizers (specifically dropout and weight decay) on the weights of a network. Consider a linear regression problem with input data  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , weights  $\mathbf{w} \in \mathbb{R}^{d \times 1}$  and targets  $\mathbf{y} \in \mathbb{R}^{n \times 1}$ . Suppose that dropout is applied to the input (with probability  $1 - p$  of dropping the unit i.e. setting it to 0). Let  $\mathbf{R} \in \mathbb{R}^{n \times d}$  be the dropout mask such that  $\mathbf{R}_{ij} \sim \text{Bern}(p)$  is sampled i.i.d. from the Bernoulli distribution.

For a squared error loss function with dropout, we then have:

$$L(\mathbf{w}) = \|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|^2$$

- 1.1 Let  $\Gamma$  be a diagonal matrix with  $\Gamma_{ii} = (\mathbf{X}^\top \mathbf{X})_{ii}^{1/2}$ . Show that the *expectation* (over  $\mathbf{R}$ ) of the loss function can be rewritten as  $\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$ . *Hint: Note we are trying to find the expectation over a squared term and use  $\text{Var}(Z) = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2$ .*
- 1.2 Show that the solution  $\mathbf{w}^{\text{dropout}}$  that minimizes the expected loss from question 1.1 satisfies

$$p\mathbf{w}^{\text{dropout}} = (\mathbf{X}^\top \mathbf{X} + \lambda^{\text{dropout}} \Gamma^2)^{-1} \mathbf{X}^\top \mathbf{y}$$

where  $\lambda^{\text{dropout}}$  is a regularization coefficient depending on  $p$ . How does the value of  $p$  affect the regularization coefficient,  $\lambda^{\text{dropout}}$ ?

- 1.3 Express the loss function for a linear regression problem without dropout and with  $L^2$  regularization, with regularization coefficient  $\lambda^{L^2}$ . Derive its closed form solution  $\mathbf{w}^{L^2}$ .
- 1.4 Compare the results of 1.2 and 1.3: identify specific differences in the equations you arrive at, and discuss qualitatively what the equations tell you about the similarities and differences in the effects of weight decay and dropout (1-3 sentences).

**Question 2** (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.  $q_\phi(\mathbf{z} | \mathbf{x}) := q(\mathbf{z}; \pi_\phi(\mathbf{x}))$ .

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

2.2 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?

2.3 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)
- (b) from the computational point of view (efficiency)
- (c) in terms of memory (storage of parameters)

**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  $D_{\text{KL}}[P_X || P_{F(U)}] = D_{\text{KL}}[P_{F^{-1}(X)} || P_U]$ . 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.

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.

3.2 Suppose two scenarios: 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.

**Question 4 (3-7).** Let  $p_0$  and  $p_1$  be two probability distributions with densities  $f_0$  and  $f_1$  (respectively). We want to explore what we can do with a trained GAN discriminator. A trained discriminator is thought to be one which is "close" to the optimal one:

$$D^* := \arg \max_D \mathbb{E}_{\mathbf{x} \sim p_1} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_0} [\log(1 - D(\mathbf{x}))].$$

4.1 For the first part of this problem, derive an expression we can use to estimate the Jensen-Shannon divergence using a trained discriminator. We remind that the definition of JSD is  $\text{JSD}(p_0, p_1) = \frac{1}{2} (KL(p_0||\mu) + KL(p_1||\mu))$ , where  $\mu = \frac{1}{2}(p_0 + p_1)$ .

4.2 For the second part, we want to demonstrate that a optimal GAN Discriminator (i.e. one which is able to distinguish between examples from  $p_0$  and  $p_1$  with minimal NLL loss) can be used to express the probability density of a datapoint  $\mathbf{x}$  under  $f_1$ ,  $f_1(\mathbf{x})$  in terms of  $f_0(\mathbf{x})$ <sup>2</sup>. Assume  $f_0$  and  $f_1$  have the same support. Show that  $f_1(\mathbf{x})$  can be estimated by  $f_0(\mathbf{x})D(\mathbf{x})/(1 - D(\mathbf{x}))$  by establishing the identity  $f_1(\mathbf{x}) = f_0(\mathbf{x})D^*(\mathbf{x})/(1 - D^*(\mathbf{x}))$ .

*Hint: Find the closed form solution for  $D^*$ .*

**Question 5 (1-2-1-1-2-3).** In this question, we are concerned with analyzing the training dynamics of GANs under gradient ascent-descent. We denote the parameters of the critic and the generator by  $\psi$  and  $\theta$  respectively. The objective function under consideration is the Jensen-Shannon (standard) GAN one:

$$\mathcal{L}(\psi, \theta) = \mathbb{E}_{p_D} \log(\sigma(C_\psi(x))) + \mathbb{E}_{p_\theta} \log(\sigma(-C_\psi(x)))$$

where  $\sigma$  is the logistic function. For ease of exposition, we will study the continuous-time system which results from the (alternating) discrete-time system when learning rate,  $\eta > 0$ , approaches zero:

$$\begin{aligned} \psi^{(k+1)} &= \psi^{(k)} + \eta v_\psi(\psi^{(k)}, \theta^{(k)}) \\ \theta^{(k+1)} &= \theta^{(k)} + \eta v_\theta(\psi^{(k+1)}, \theta^{(k)}) \end{aligned} \xrightarrow{\eta \rightarrow 0^+} \begin{aligned} \dot{\psi} &= v_\psi(\psi, \theta) \\ \dot{\theta} &= v_\theta(\psi, \theta) \end{aligned} \quad \begin{aligned} v_\psi(\psi, \theta) &:= \nabla_\psi \mathcal{L}(\psi, \theta) \\ v_\theta(\psi, \theta) &:= -\nabla_\theta \mathcal{L}(\psi, \theta) \end{aligned}$$

The purpose is to initiate a study on the stability of the training algorithm. For this reason, we will utilize the following simple setting: Both training and generated data have support on  $\mathbb{R}$ . In addition,  $p_D = \delta_0$  and  $p_\theta = \delta_\theta$ . This means that both of them are Dirac distributions<sup>3</sup> which are centered at  $x = 0$ , for the real data, and at  $x = \theta$ , for the generated. The critic,  $C_\psi : \mathbb{R} \rightarrow \mathbb{R}$ , is  $C_\psi(x) = \psi_0 x + \psi_1$ .

5.1 Derive the expressions for the "velocity" field,  $v$ , of the dynamical system in the joint parameter space  $(\psi_0, \psi_1, \theta)$ , and find its stationary points  $(\psi_0^*, \psi_1^*, \theta^*)$ .<sup>4</sup>

2. You might need to use the "functional derivative" to solve this problem. See "19.4.2 Calculus of Variations" of the Deep Learning book or "Appendix D Calculus of Variations" of Bishop's Pattern Recognition and Machine Learning for more information.

3. If  $p_X = \delta_z$ , then  $p(X = z) = 1$ .

4. To find the stationary points, set  $v = 0$  and solve for each of the parameters.

5.2 Derive  $J^*$ , the  $(3 \times 3)$  Jacobian of  $v$  at  $(\psi_0^*, \psi_1^*, \theta^*)$ .

For a continuous-time system to be locally asymptotically stable it suffices that all eigenvalues of  $J^*$  have negative real part. Otherwise, further study is needed to conclude. However, this case is not great news since the fastest achievable convergence is sublinear.

5.3 Find the eigenvalues of  $J^*$  and comment on the system's local stability around the stationary points.

Now we will include a gradient penalty,  $\mathcal{R}_1(\psi) = \mathbb{E}_{p_D} \|\nabla_x C_\psi(x)\|^2$ , to the critic's loss, so the regularized system becomes:

$$\begin{aligned}\dot{\psi} &= \bar{v}_\psi(\psi, \theta) & \bar{v}_\psi(\psi, \theta) &:= \nabla_\psi \mathcal{L}(\psi, \theta) - \frac{\gamma}{2} \nabla_\psi \mathcal{R}_1(\psi) \\ \dot{\theta} &= \bar{v}_\theta(\psi, \theta) & \bar{v}_\theta(\psi, \theta) &:= -\nabla_\theta \mathcal{L}(\psi, \theta)\end{aligned}$$

for  $\gamma > 0$ . Repeat 1-2-3 for the modified system and compare the stability of the two.

5.4 Derive the expressions for the "velocity" field,  $\bar{v}$ , of the dynamical system in the joint parameter space  $(\psi_0, \psi_1, \theta)$ , and find its stationary points  $(\psi_0^*, \psi_1^*, \theta^*)$ .<sup>5</sup>

5.5 Derive  $\bar{J}^*$ , the  $(3 \times 3)$  Jacobian of  $\bar{v}$  at  $(\psi_0^*, \psi_1^*, \theta^*)$ .

5.6 Find the eigenvalues of  $\bar{J}^*$  and comment on the system's local stability around the stationary points.

In Problem 2 of the programming assignment, you will verify empirically your claims.

---

5. To find the stationary points, set  $v = 0$  and solve for each of the parameters.