

**Problem 1: Optimal GAN discriminator**

Assume that the generator is defined as  $G_\theta : \mathbb{R}^{d_{\text{latent}}} \rightarrow \mathbb{R}^{d_{\text{data}}}$ , and the discriminator  $D : \mathbb{R}^{d_{\text{data}}} \rightarrow [0, 1]$  outputs a probability  $p(\text{real} | x)$ . Let  $p_{\text{data}}$  be the density of the real data distribution and  $p_\theta = (G_\theta)_\# p_{\text{latent}}$  be the density of the distribution of generated samples.

Show that for a fixed generator  $G_\theta$  the optimal discriminator  $D^*$  for the following optimisation problem

$$\max_D \left\{ \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p_{\text{latent}}} [\log(1 - D(G_\theta(z)))] \right\}$$

has the following form:  $D^*(\text{real} | x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_\theta(x)}$  (hint: for  $a, b > 0$  find the maximiser of  $a \log(x) + b \log(1 - x)$  with respect to  $x$ , then use this result to find the optimal  $D^*$ ).

**Problem 2: Alternative GAN losses**

1. Alternatively, losses for training discriminator and generator can be written as follows:

$$\mathcal{L}_{\text{MSE}}(\varphi) = \mathbb{E}_{x \sim p_{\text{data}}} [(1 - D_\varphi(x))^2] + \mathbb{E}_{z \sim p_{\text{latent}}} [(D_\varphi(G_\theta(z)))^2],$$

$$\mathcal{L}_{\text{MSE}}(\theta) = \mathbb{E}_{z \sim p_{\text{latent}}} [(1 - D_\varphi(G_\theta(z)))^2],$$

which are minimised with respect to  $\varphi$  and  $\theta$  respectively. This contrasts with the vanilla objective

$$\mathcal{L}_{\text{vanilla}}(\theta, \varphi) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D_\varphi(x)] + \mathbb{E}_{z \sim p_{\text{latent}}} [\log(1 - D_\varphi(G_\theta(z)))]$$

Show that the gradients  $\nabla_\varphi \mathcal{L}_{\text{MSE}}(\varphi)$  and  $\nabla_\theta \mathcal{L}_{\text{MSE}}(\theta)$  are proportional to  $\nabla_\varphi \mathcal{L}_{\text{vanilla}}(\varphi, \theta)$  and  $\nabla_\theta \mathcal{L}_{\text{vanilla}}(\varphi, \theta)$  respectively.

2. In order to improve the training stability one can use 'non-saturating' losses for training the generator:  $-\log D(G_\theta(z))$  in place of  $\log(1 - D(G_\theta(z)))$ :

$$\mathcal{L}_{\text{non-sat}}(\theta) = \mathbb{E}_{z \sim p_{\text{latent}}} [-\log D(G_\theta(z))].$$

Are the non-saturating GAN gradients proportional to the vanilla ones?

**Problem 3: Wasserstein GAN**

Read **Wasserstein GAN paper**. Let  $\mathbb{P}_r$  represent the distribution of real data and  $\mathbb{P}_g$  the distribution of generated data. Answer the following:

1. Explain (informally) how using  $\text{KL}(\mathbb{P}_r \| \mathbb{P}_g)$  as a measure of distance between two distributions differs from  $W_1(\mathbb{P}_r, \mathbb{P}_g)$  and why the latter is more sensible.
2. (**Hard**) Explain (informally) why:

$$\inf_{\pi \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \pi} [\|x - y\|] = \sup_{\|f\|_L \leq 1} \left\{ \mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{y \sim \mathbb{P}_g} [f(y)] \right\}$$

Answer the following questions:

- (a) Give the definition of a 1-Lipschitz function and explain (informally) why  $f$  should be 1-Lipschitz.
- (b) If  $f$  is parameterised by a neural network, what are the possible ways to make  $f$  1-Lipschitz?