

# CS 236 Homework 3

Instructors: Stefano Ermon and Aditya Grover

{ermon,adityag}@cs.stanford.edu

Available: 11/11/2019; Due: 23:59 PST, 12/2/2019

---

## Problem 1: Flow models (20 points)

In this problem, we will implement a Masked Autoregressive Flow (MAF) model on the Moons dataset, where we define  $p_{\text{data}}(\mathbf{x})$  over a 2-dimensional space ( $\mathbf{x} \in \mathbb{R}^n$  where  $n = 2$ ). Recall that MAF is comprised of Masked Autoregressive Distribution Estimator (MADE) blocks, which has a special masking scheme at each layer such that the autoregressive property is preserved. In particular, we consider a Gaussian autoregressive model:

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | \mathbf{x}_{<i})$$

such that the conditional Gaussians  $p(x_i | \mathbf{x}_{<i}) = \mathcal{N}(x_i | \mu_i, (\exp(\alpha_i)^2))$  are parameterized by neural networks  $\mu_i = f_{\mu_i}(\mathbf{x}_{<i})$  and  $\alpha_i = f_{\alpha_i}(\mathbf{x}_{<i})$ . Note that  $\alpha_i$  denotes the log standard deviation of the Gaussian  $p(x_i | \mathbf{x}_{<i})$ .

As seen in lecture, a normalizing flow uses a series of deterministic and invertible mappings  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  such that  $x = f(z)$  and  $z = f^{-1}(x)$  to transform a simple prior distribution  $p_z$  (e.g. isotropic Gaussian) into a more expressive one. In particular, a normalizing flow which composes  $k$  invertible transformations  $\{f_j\}_{j=1}^k$  such that  $\mathbf{x} = f^k \circ f^{k-1} \circ \dots \circ f^1(\mathbf{z}_0)$  takes advantage of the change-of-variables property:

$$\log p(\mathbf{x}) = \log p_z(f^{-1}(\mathbf{x})) + \sum_{j=1}^k \log \left| \det \left( \frac{\partial f_j^{-1}(\mathbf{z}_j)}{\partial \mathbf{z}_j} \right) \right|$$

In MAF, the forward mapping is:  $x_i = \mu_i + z_i \cdot \exp(\alpha_i)$ , and the inverse mapping is:  $z_i = (x_i - \mu_i) / \exp(\alpha_i)$ . The log of the absolute value of the Jacobian is:

$$\log \left| \det \left( \frac{\partial f^{-1}}{\partial \mathbf{x}} \right) \right| = - \sum_{i=1}^n \alpha_i$$

where  $\mu_i$  and  $\alpha_i$  are as defined above.

Your job is to implement and train a 5-layer MAF model on the Moons dataset for 100 epochs by modifying the MADE and MAF classes in the `flow_network.py` file. Note that we have provided an implementation of the sequential-ordering masking scheme for MADE.

1. [5 points] Implement the `forward` function in the MADE class in `codebase/flow_network.py`. The forward pass describes the mapping  $x_i = \mu_i + z_i \cdot \exp(\alpha_i)$ .
2. [5 points] Implement the `inverse` function in the MADE class in `codebase/flow_network.py`. The inverse pass describes the mapping  $z_i = (x_i - \mu_i) / \exp(\alpha_i)$ .
3. [7 points] Implement the `log_probs` function in the MAF class in `codebase/flow_network.py`. Then train the MAF model for 50 epochs by executing `python run_flow.py`. [Hint: you should be getting a validation/test loss of around 1.2 nats]

4. [3 points] Visualize 1000 samples drawn the model after is has been trained, which you can do by finding the figure in `maf/samples_epoch50.png`. Attach this figure in your writeup – your samples will not be perfect, but should for the most part resemble the shape of the original dataset.

Please package all of the programming parts (1.3, 2.2, 4.2, and 5.5) together using `make_submission.sh` and submit the resulting ZIP file on GradeScope.

## Problem 2: Generative adversarial networks (15 points)

In this problem, we will implement a generative adversarial network (GAN) that models a high-dimensional data distribution  $p_{\text{data}}(\mathbf{x})$ , where  $\mathbf{x} \in \mathbb{R}^n$ . To do so, we will define a generator  $G_\theta : \mathbb{R}^k \rightarrow \mathbb{R}^n$ ; we obtain samples from our model by first sampling a  $k$ -dimensional random vector  $\mathbf{z} \sim \mathcal{N}(0, I)$  and then returning  $G_\theta(\mathbf{z})$ .

We will also define a discriminator  $D_\phi : \mathbb{R}^n \rightarrow (0, 1)$  that judges how realistic the generated images  $G_\theta(\mathbf{z})$  are, compared to samples from the data distribution  $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$ . Because its output is intended to be interpreted as a probability, the last layer of the discriminator is frequently the **sigmoid** function,

$$\sigma(x) = \frac{1}{1 + e^{-x}},$$

which constrains its output to fall between 0 and 1. For convenience, let  $h_\phi(\mathbf{x})$  denote the activation of the discriminator right before the sigmoid layer, i.e. let  $D_\phi(\mathbf{x}) = \sigma(h_\phi(\mathbf{x}))$ . The values  $h_\phi(\mathbf{x})$  are also called the discriminator's **logits**.

There are several common variants of the loss functions used to train GANs. They can all be described as a procedure where we alternately perform a gradient descent step on  $L_D(\phi; \theta)$  with respect to  $\phi$  to train the discriminator  $D_\phi$ , and a gradient descent step on  $L_G(\theta; \phi)$  with respect to  $\theta$  to train the generator  $G_\theta$ :

$$\min_{\phi} L_D(\phi; \theta), \quad \min_{\theta} L_G(\theta; \phi).$$

In lecture, we talked about the following losses, where the discriminator's loss is given by

$$L_D(\phi; \theta) = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}[\log D_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)}[\log(1 - D_\phi(G_\theta(\mathbf{z})))]$$

and the generator's loss is given by the **minimax loss**

$$L_G^{\text{minimax}}(\theta; \phi) = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)}[\log(1 - D_\phi(G_\theta(\mathbf{z})))]$$

1. [5 points] Unfortunately, this form of loss for  $L_G$  suffers from a *vanishing gradient* problem. In terms of the discriminator's logits, the minimax loss is

$$L_G^{\text{minimax}}(\theta; \phi) = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)}[\log(1 - \sigma(h_\phi(G_\theta(\mathbf{z})))].$$

Show that the derivative of  $L_G^{\text{minimax}}$  with respect to  $\theta$  is approximately 0 if  $D(G_\theta(\mathbf{z})) \approx 0$ , or equivalently, if  $h_\phi(G_\theta(\mathbf{z})) \ll 0$ . You may use the fact that  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$ . Why is this problematic for the training of the generator when the discriminator successfully identifies a fake sample  $G_\theta(\mathbf{z})$ ?

2. [10 points] Because of this vanishing gradient problem, in practice,  $L_G^{\text{minimax}}$  is typically replaced with the **non-saturating loss**

$$L_G^{\text{non-saturating}}(\theta; \phi) = -\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)}[\log D_\phi(G_\theta(\mathbf{z}))].$$

To turn the non-saturating loss into a concrete algorithm, we will take alternating gradient steps on Monte Carlo estimates of  $L_D$  and  $L_G^{\text{non-saturating}}$ :

$$\begin{aligned} L_D(\phi; \theta) &\approx -\frac{1}{m} \sum_{i=1}^m \log D_\phi(\mathbf{x}^{(i)}) - \frac{1}{m} \sum_{i=1}^m \log(1 - D_\phi(G_\theta(\mathbf{z}^{(i)}))), \\ L_G^{\text{non-saturating}}(\theta; \phi) &\approx -\frac{1}{m} \sum_{i=1}^m \log D_\phi(G_\theta(\mathbf{z}^{(i)})), \end{aligned}$$

where  $m$  is the batch size, and for  $i = 1, \dots, m$ , we sample  $\mathbf{x}^{(i)} \sim p_{\text{data}}(\mathbf{x})$  and  $\mathbf{z}^{(i)} \sim \mathcal{N}(0, I)$ .

Implement and train a non-saturating GAN on Fashion MNIST for one epoch. Read through `run_gan.py`, and in `codebase/gan.py`, implement the `loss_nonsaturating` function. To train the model, execute `python run_gan.py`. You may monitor the GAN's output in the `out_nonsaturating` directory. Note that because the GAN is only trained for one epoch, we cannot expect the model's output to produce very realistic samples, but they should be roughly recognizable as clothing items.

### Problem 3: Divergence minimization (25 points)

Now, let us analyze some theoretical properties of GANs. For convenience, we will denote  $p_{\theta}(\mathbf{x})$  to be the distribution whose samples are generated by first sampling  $\mathbf{z} \sim \mathcal{N}(0, I)$  and then returning the sample  $G_{\theta}(\mathbf{z})$ . With this notation, we may compactly express the discriminator's loss as

$$L_D(\phi; \theta) = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}[\log D_{\phi}(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x})}[\log(1 - D_{\phi}(\mathbf{x}))].$$

1. [10 points] Show that  $L_D$  is minimized when  $D_{\phi} = D^*$ , where

$$D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x}) + p_{\text{data}}(\mathbf{x})}.$$

(Hint: for a fixed  $\mathbf{x}$ , what  $t$  minimizes  $f(t) = -p_{\text{data}}(\mathbf{x}) \log t - p_{\theta}(\mathbf{x}) \log(1 - t)$ ?)

2. [5 points] Recall that  $D_{\phi}(\mathbf{x}) = \sigma(h_{\phi}(\mathbf{x}))$ . Show that the logits  $h_{\phi}(\mathbf{x})$  of the discriminator estimate the log of the likelihood ratio of  $\mathbf{x}$  under the true distribution compared to the model's distribution; that is, show that if  $D_{\phi} = D^*$ , then

$$h_{\phi}(\mathbf{x}) = \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})}.$$

3. [5 points] Consider a generator loss defined by the sum of the minimax loss and the non-saturating loss,

$$L_G(\theta; \phi) = \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x})}[\log(1 - D_{\phi}(\mathbf{x}))] - \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x})}[\log D_{\phi}(\mathbf{x})].$$

Show that if  $D_{\phi} = D^*$ , then

$$L_G(\theta; \phi) = \text{KL}(p_{\theta}(\mathbf{x}) || p_{\text{data}}(\mathbf{x})).$$

4. [5 points] Recall that when training VAEs, we minimize the negative ELBO, an upper bound to the negative log likelihood. Show that the negative log likelihood,  $-\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}[\log p_{\theta}(\mathbf{x})]$ , can be written as a KL divergence plus an additional term that is constant with respect to  $\theta$ . Does this mean that a VAE decoder trained with ELBO and a GAN generator trained with the  $L_G$  defined in the previous part are implicitly learning the same objective? Explain.

### Problem 4: Conditional GAN with projection discriminator (20 points)

So far, we have trained GANs that sample from a given dataset of images. However, many datasets come with not only images, but also labels that specify the class of that particular image. In the MNIST dataset, we have both the digit's image as well as its numerical identity. It is natural to want to generate images that correspond to a particular class.

Formally, an *unconditional* GAN is trained to produce samples  $\mathbf{x} \sim p_{\theta}(\mathbf{x})$  that mimic samples  $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$  from a data distribution. In the class-conditional setting, we instead have have labeled data  $(\mathbf{x}, y) \sim p_{\text{data}}(\mathbf{x}, y)$  and seek to train a model  $p_{\theta}(\mathbf{x}, y)$ . Since it is the class conditional generator  $p_{\theta}(\mathbf{x}|y)$  that we are interested in, we will express  $p_{\theta}(\mathbf{x}, y) = p_{\theta}(\mathbf{x}|y)p_{\theta}(y)$ . We will set  $p_{\theta}(\mathbf{x}|y)$  to be the distribution given by  $G_{\theta}(\mathbf{z}, y)$ , where  $\mathbf{z} \sim \mathcal{N}(0, I)$  as usual. For simplicity, we will assume  $p_{\text{data}}(y) = \frac{1}{m}$  and set  $p_{\theta}(y) = \frac{1}{m}$ , where  $m$  is the number of classes. In this case, the discriminator's loss becomes

$$\begin{aligned} L_D(\phi; \theta) &= -\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbf{x}, y)}[\log D_{\phi}(\mathbf{x}, y)] - \mathbb{E}_{(\mathbf{x}, y) \sim p_{\theta}(\mathbf{x}, y)}[\log(1 - D_{\phi}(\mathbf{x}, y))] \\ &= -\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbf{x}, y)}[\log D_{\phi}(\mathbf{x}, y)] - \mathbb{E}_{y \sim p_{\theta}(y)}[\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0, I)}[\log(1 - D_{\phi}(G_{\theta}(\mathbf{z}, y), y))]]. \end{aligned}$$

Therefore, the main difference for the conditional GAN is that we must structure our generator  $G_\theta(\mathbf{z}, y)$  and discriminator  $D_\phi(\mathbf{x}, y)$  to accept the class label  $y$  as well. For the generator, one simple way to do so is to encode  $y$  as a one-hot vector  $\mathbf{y}$  and concatenate it to  $\mathbf{z}$ , and then apply neural network layers normally. (A one-hot representation of a class label  $y$  is an  $m$ -dimensional vector  $\mathbf{y}$  that is 1 in the  $y$ th entry and 0 everywhere else.)

In practice, the effectiveness of the model is strongly dependent on the way the discriminator depends on  $y$ . One heuristic with which to design the discriminator is to mimic the form of the theoretically optimal discriminator. That is, we can structure the neural network used to model  $D_\phi$  based on the form of  $D^*$ , where  $D^*$  minimizes  $L_D$ . To calculate the theoretically optimal discriminator, though, it is necessary to make some assumptions.

1. **[10 points]** Suppose that when  $(\mathbf{x}, y) \sim p_{\text{data}}(\mathbf{x}, y)$ , there exists a feature mapping  $\varphi$  under which  $\varphi(\mathbf{x})$  becomes a mixture of  $m$  unit Gaussians, with one Gaussian per class label  $y$ . Assume that when  $(\mathbf{x}, y) \sim p_\theta(\mathbf{x}, y)$ ,  $\varphi(\mathbf{x})$  also becomes a mixture of  $m$  unit Gaussians, again with one Gaussian per class label  $y$ . Concretely, we assume that the ratio of the conditional probabilities can be written as

$$\frac{p_{\text{data}}(\mathbf{x}|y)}{p_\theta(\mathbf{x}|y)} = \frac{\mathcal{N}(\varphi(\mathbf{x})|\boldsymbol{\mu}_y, I)}{\mathcal{N}(\varphi(\mathbf{x})|\hat{\boldsymbol{\mu}}_y, I)},$$

where  $\boldsymbol{\mu}_y$  and  $\hat{\boldsymbol{\mu}}_y$  are the means of the Gaussians for  $p_{\text{data}}$  and  $p_\theta$  respectively.

Show that under this simplifying assumption, the optimal discriminator's logits  $h^*(\mathbf{x}, y)$  can be written in the form

$$h^*(\mathbf{x}, y) = \mathbf{y}^T (A\varphi(\mathbf{x}) + \mathbf{b})$$

for some matrix  $A$  and vector  $\mathbf{b}$ , where  $\mathbf{y}$  is a one-hot vector denoting the class  $y$ . In this problem, the discriminator's output and logits are related by  $D_\phi(\mathbf{x}, y) = \sigma(h_\phi(\mathbf{x}, y))$ . (Hint: use the result from problem 2.2.)

2. **[10 points]** Implement and train a conditional GAN on Fashion MNIST for one epoch. The discriminator has the structure described in part 1, with  $\varphi$ ,  $A$  and  $\mathbf{b}$  parameterized by a neural network with a final linear layer, and the generator accepts a one-hot encoding of the class. In `codebase/gan.py`, implement the `conditional_loss_nonsaturating` function. To train the model, execute `python run_conditional_gan.py`. You may monitor the GAN's output in the `out_nonsaturating_conditional` directory. You should be able to roughly recognize the categories that correspond to each column.

### Problem 5: Wasserstein GAN (35 points)

In many cases, the GAN algorithm can be thought of as minimizing a divergence between a data distribution  $p_{\text{data}}(\mathbf{x})$  and the model distribution  $p_\theta(\mathbf{x})$ . For example, the minimax GAN discussed in the lectures minimizes the Jensen-Shannon divergence, and the loss in problem 2.3 minimizes the KL divergence. In this problem, we will explore an issue with these divergences and one potential way to fix it.

1. **[5 points]** Let  $p_\theta(x) = \mathcal{N}(x|\theta, \epsilon^2)$  and  $p_{\text{data}}(x) = \mathcal{N}(x|\theta_0, \epsilon^2)$  be normal distributions with standard deviation  $\epsilon$  centered at  $\theta \in \mathbb{R}$  and  $\theta_0 \in \mathbb{R}$  respectively. Show that

$$\text{KL}(p_\theta(x)||p_{\text{data}}(x)) = \frac{(\theta - \theta_0)^2}{2\epsilon^2}.$$

2. **[5 points]** Suppose  $p_\theta(x)$  and  $p_{\text{data}}(x)$  both place probability mass in only a very small part of the domain; that is, consider the limit  $\epsilon \rightarrow 0$ . What happens to  $\text{KL}(p_\theta(x)||p_{\text{data}}(x))$  and its derivative with respect to  $\theta$ , assuming that  $\theta \neq \theta_0$ ? Why is this problematic for a GAN trained with the loss function  $L_G$  defined in problem 2.3?
3. **[5 points]** To avoid this problem, we'll propose an alternative objective for the discriminator and generator. Consider the following alternative objectives:

$$\begin{aligned} L_D(\phi; \theta) &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[D_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}[D_\phi(\mathbf{x})] \\ L_G(\theta; \phi) &= -\mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[D_\phi(\mathbf{x})], \end{aligned}$$

where  $D_\phi$  is no longer constrained to functions that output a probability; instead  $D_\phi$  can be a function that outputs any real number. As defined, however, these losses are still problematic. Again consider the limit  $\epsilon \rightarrow 0$ ; that is, let  $p_\theta(x)$  be the distribution that outputs  $\theta \in \mathbb{R}$  with probability 1, and let  $p_{\text{data}}(x)$  be the distribution that outputs  $\theta_0 \in \mathbb{R}$  with probability 1. Why is there no discriminator  $D_\phi$  that minimizes this new objective  $L_D$ ?

4. **[5 points]** Let's tweak the alternate objective so that an optimal discriminator exists. Consider the same objective  $L_D$  and the same limit  $\epsilon \rightarrow 0$ . Now, suppose that  $D_\phi$  is restricted to differentiable functions whose derivative is always between  $-1$  and  $1$ . It can still output any real number. Is there now a discriminator  $D_\phi$  out of this class of functions that minimizes  $L_D$ ? Briefly describe what the optimal  $D_\phi$  looks like as a function of  $x$ .
5. **[15 points]** The Wasserstein GAN with gradient penalty (WGAN-GP) enables stable training by penalizing functions whose derivatives are too large. It achieves this by adding a penalty on the 2-norm of the gradient of the discriminator at various points in the domain. It is defined by

$$\begin{aligned} L_D(\phi; \theta) &= \mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[D_\phi(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})}[D_\phi(\mathbf{x})] + \lambda \mathbb{E}_{\mathbf{x} \sim r_\theta(\mathbf{x})}[(\|\nabla D_\phi(\mathbf{x})\|_2 - 1)^2] \\ L_G(\theta; \phi) &= -\mathbb{E}_{\mathbf{x} \sim p_\theta(\mathbf{x})}[D_\phi(\mathbf{x})], \end{aligned}$$

where  $r_\theta(\mathbf{x})$  is defined by sampling  $\alpha \sim \text{Uniform}([0, 1])$ ,  $\mathbf{x}_1 \sim p_\theta(\mathbf{x})$ , and  $\mathbf{x}_2 \sim p_{\text{data}}(\mathbf{x})$ , and returning  $\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2$ . The hyperparameter  $\lambda$  controls the strength of the penalty; a setting that usually works is  $\lambda = 10$ .

Implement and train WGAN-GP for one epoch on Fashion MNIST. In `codebase/gan.py`, implement the `loss_wasserstein_gp` function. To train the model, execute `python run_gan.py --loss.type wasserstein_gp`. You may monitor the GAN's output in the `out_wasserstein_gp` directory.