# Report Practical(Assignment 3)

**Team members:**
Fares Ben Slimane
Parviz Haggi
Mohammed Loukili
Jorge A. Gutierrez Ortega

April 22, 2019

**Abstract**

This report includes our solutions to the problems of the 3rd practical assignment. It consists of three sections: In Section 1 we implement the original GAN and WGAN with gradient penalty along with some experimentation. Section 2 is about the implementation and experimentation of the VAEs. The last section (3) includes experimentation and evaluation of the generative models' ability to generate realistic-looking images. In each case, our code is uploaded to the following Github repository [1]

## Problem 1

In this problem we are using a GAN discriminator to estimate the Jensen Shannon divergence (JSD) and a the Wasserstein Distance (WD). We implement the discriminator/critic using a multi-layer perceptron that we can initialize to either have an output $y \in [0, 1]$ (sigmoid activation) or an output $y \in \mathbb{R}$ (no activation). The architecture consists of two hidden layers with 64 and 128 output units respectively. We use ReLU activation function in the hidden layers.

1. In this section, we implement a function to estimate the JSD. Using our MLP with sigmoid activation at the output we optimize the following objective function:

$$\arg\max_{\theta} \left\{ \log 2 + \frac{1}{2}\mathbf{E}_{x \sim p}[\log(D_{\theta}(x))] + \frac{1}{2}\mathbf{E}_{y \sim q}\log(1 - D_{\theta}(y))] \right\}$$

At its optimum, the objective function estimates the JSD between the distributions given by $p$ and $q$. An overview of the implementation can be seen below. The full code is available in the file **density_estimation.py** under our github repository [1].

```
1
2 #Implementation of the JSD
3
4 optimizer.zero_grad()
5 p = torch.cat((p1,torch.rand(batch_size,1)),1).to(device)
6 q = torch.cat((phi*torch.ones(batch_size,1),t
7                 orch.rand(batch_size,1)),1).to(device)
8 Dp = net(p)
9 Dq = net(q)
10 loss = -(math.log(2.) + (1/2.)*torch.mean(torch.log(Dp)) +
```

```
11                                                    (1/2.)*torch.mean(torch.log(1−Dq)))
12  loss.backward()
13  optimizer.step()
```

2. In this section, we implement a function to estimate the WD. In this case we use our MLP without activation function, so it has an output of a real valued scalar. Here we optimize the following objective function:

$$\arg\max_\theta \mathbf{E}_{x\sim p}[T_\theta(x)] - \mathbf{E}_{y\sim q}[T_\theta(y)] - \lambda\mathbf{E}_{z\sim r}(||\nabla_z T_\theta(z)||_2 - 1)^2.$$

$r$ is the distribution over $z = ax + (1-a)y$, where $x \sim p$, $y \sim q$ and $a \sim U[0,1]$.

At it's optimum, this objective function estimates the WD between the distributions given by $p$ and $q$. The following portion of the code shows the implementation of the optimization. The full code is also available in our repository [1]

```
1
2  #Implementation of the WD
3
4  optimizer.zero_grad()
5  p = torch.cat((p1,torch.rand(batch_size,1)),1).to(device)
6  q = torch.cat((phi*torch.ones(batch_size,1),
7                          torch.rand(batch_size,1)),1).to(device)
8  Dp = net(p)
9  Dq = net(q)
10 #  gradient penalty
11 a = torch.rand(batch_size, 1).expand(batch_size,2).to(device)
12 r = a*p + (1−a)*q
13 r.requires_grad = True
14 Dr = net(r)
15 gradients = torch.autograd.grad(outputs=Dr, inputs=r,
16                 grad_outputs=torch.ones(batch_size,1).to(device),
17                 create_graph=True, retain_graph=True, only_inputs=True)[0]
18 loss = −(torch.mean(Dp) − torch.mean(Dq) −
19         gp_coeff*torch.mean((gradients.norm(2, dim=1) − 1) ** 2))
20 loss.backward()
21 optimizer.step()
```

3. In this example we compare the properties of JSD and WD. Given a random variable $Z \sim U[0,1]$, we compute the approximation of both metrics between the 2-dimensional distributions, $p$ given by $(0, Z)$, and $q_\phi$ given by $(\phi, Z)$, where $\phi$ is a parameter. In Figures 1 and 2 we plot the estimated JSD and WD respectively for $\phi \in [-1, 1]$ with interval of 0.1. We perform the estimation by optimizing the MLP as in points 1 and 2. We use batches of samples from the distributions of 512, and the models were trained for 5000 iterations using an SGD optimizer. For every value of $\phi$ we generate the distribution $q_\phi$ and measure its distance to $p$.
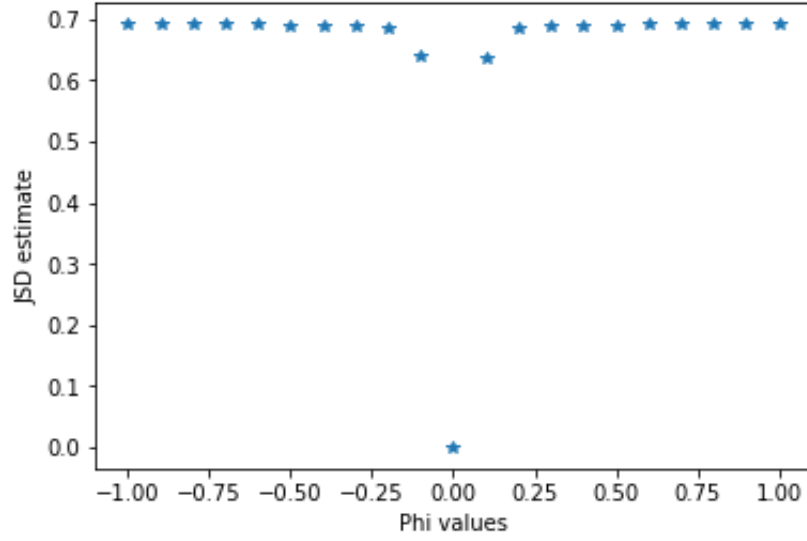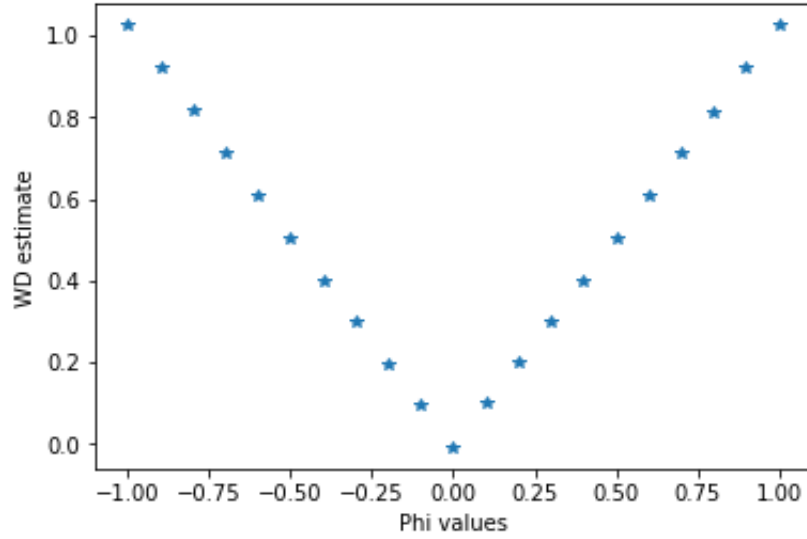
Figure 1: JSD estimation



Figure 2: WD estimation

In this experiment the distributions $p$ and $q_\phi$ have disjoint supports for all values of $\phi$ but $\phi = 0$. As expected, the JSD when the distributions are disjoint is the same (around $\log(2)$) no matter how close the distributions may be i.e. for a small value of $\phi$. This results in a difficulty to learn the distribution $p$ using JSD. In contrast, WD is continuous over the values of $\phi$ giving a better information about the *closeness* of the distributions.

The full code is given in the file **density_estimation.py** in our github repository [1].

3

4. In this section we estimate the unknown density $f_1$ using the approximation $f_0(x)D(x)/(1 - D(x))$ (proven in Question 5 from the theoretical part), where $f_0$ is a known distribution (assumed 1-dimensional standard Gaussian in this question). The full code is provided in the file **density_estimation.py** under our github repository [1].

Using the above neural network (discriminator), we minimize the following function:

$$loss = -(torch.mean(torch.log(Dx)) + torch.mean(torch.log(1 - Dy))),$$

where Dx is the feedforward of $f_1$ and Dy is the feedforward of $f_0$.

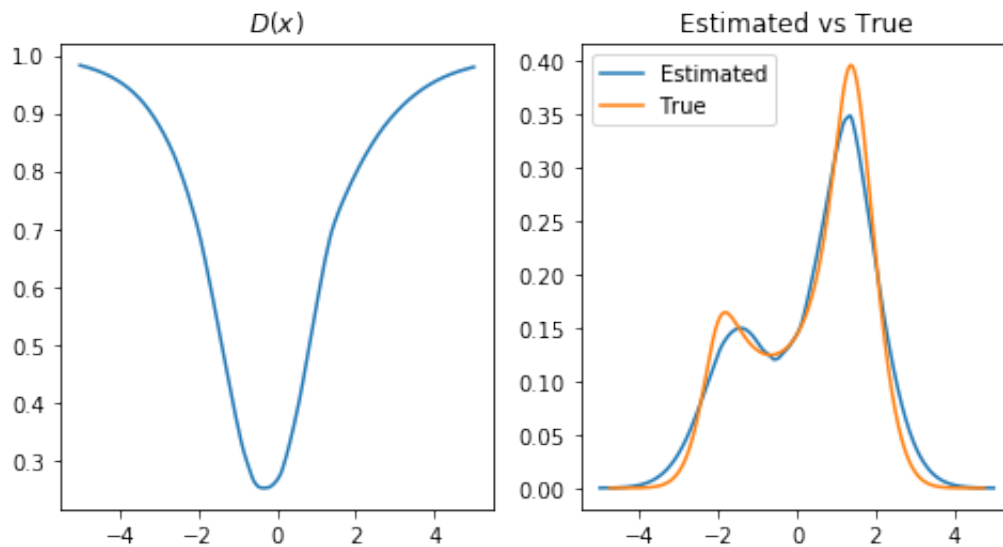The following figures show the discriminator's output and the estimated density:



Figure 3: (left) Discriminator output (Right)Estimated $f_1$

# Problem 2

A. Training VAE

We used the given architecture and ADAM with the provided learning rate. After training the model for 20 Epochs, we achieved an average value of ELBO of $-95$ on validation set. It is clearly higher than the reference value provided in question. The code for training the VAE can be found in script vae.py under the folder Problem2.

B.1 Evaluating log-likelihood with VAE

Here we implement the Importance Sampling procedure that takes as parameters the trained model, an array of $x_i$ and an array of samples $z_{ik}$ from the distribution $q(z|x_i)$. The procedure returns an array of log-likelihood $\log p(x_i)$ of the size of the mini-batches. The code snippet below demonstrates our implementation.

```
1  ##Evaluate VAE using importance sampling
2  #Q2.1
3
4  def loss_IS(model, true_x, z):
```

```
5
6    #Loop over the elements i of batch
7    M = true_x.shape[0]
8
9    #Save logp(x)
10   logp_x = np.zeros([M])
11
12   #Get mean and std from encoder
13   #2 Vectors of 100
14   mu, logvar = model.encode(true_x.to(device))
15   std = torch.exp(0.5*logvar)
16
17   K = 200
18
19   #Loop over tha batch
20    for i in range(M):
21      #z_ik
22      samples = z[i,:,:]
23
24      #Compute the reconstructed x's from sampled z's
25      x = model.decode(samples.to(device))
26
27      #Compute the p(x_i|z_ik) of x sampled from z_ik
28      #Bernoulli dist = Apply BCE
29      #Output an array of losses
30
31      true_xi = true_x[i,:,:].view(-1, 784)
32      x = x.view(-1, 784)
33
34      #P(x/z) is shape (200, 784)
35      p_xz = -(true_xi * torch.log(x) + (1.0-true_xi) * torch.log(1.0-x))
36
37      #Multiply the independent probabilities (784)
38      #Apply logsumexp to avoid small image prob
39      #P(x/z) is shape (200, 784)
40      p_xz = logsumexp(p_xz.cpu().numpy(), axis=1)
41
42
43      ##q(z_ik/x_i) follows a normal dist
44      #q_z = mgd(samples, mu, std)
45      s = std[i, :].view([std.shape[1]])
46      m = mu[i, :].view([std.shape[1]])
47
48      q_z = multivariate_normal.pdf(samples.cpu().numpy(),mean=m.cpu().numpy(),
       cov=np.diag(s.cpu().numpy()**2))
49
50
51      ##p(z_ik) follows a normal dist with mean 0/variance 1
52      #Normally distributed with loc=0 and scale=1
53      std_1 = torch.ones(samples.shape[1])
54      mu_0 = torch.zeros(samples.shape[1])
55
56      p_z = multivariate_normal.pdf(samples.cpu().numpy(),mean=mu_0.cpu().numpy
       (), cov=np.diag(std_1.cpu().numpy()**2))
57
58      #Multiply the probablities
59
60      #logp_x[i] = np.log((1.0/K) * np.sum(np.exp(np.log(p_xz) + np.log(p_z) -
       np.log(q_z))))
61      #logp_x[i] = np.log((1.0/K) * np.sum(np.exp(p_xz.cpu().numpy() + np.log(
       p_z) - np.log(q_z))))
62      logp_x[i] = np.log((1.0/K) * np.sum((p_xz * p_z)/q_z))
```

5

Figure 4: A sample of generated images

```
63
64    return logp_x
```

B.2 The evaluation of the training model using the ELBO:

    a. Validation: $-95$

    b. Test: $-94.43$

The evaluation of the training model using the log-likelihood:

    a. Validation: $-18.41$

    b. Test: $-18.27$

Below (Figure 4) is a sample of the obtained images generated by the trained model:

# Problem 3

We have used in this problem a similar architecture for the VAE's decoder and the GAN's generator which is an MLP with 6 layers, as shown in this code snippet:

```
1  class Generator(nn.Module):
2      def __init__(self):
3          super(Generator, self).__init__()
4
5          self.model = nn.Sequential(
6              nn.Linear(latent_dim, 128),
7              nn.ReLU()
8              nn.Linear(128, 256),
```

```
 9              nn.ReLU()
10              nn.Linear(256, 512),
11              nn.ReLU()
12              nn.Linear(512, 1024),
13              nn.ReLU()
14              nn.Linear(1024, 2048),
15              nn.ReLU()
16              nn.Linear(2048, int(np.prod(img_shape))),
17              nn.Tanh()
18          )
19
20      def forward(self, z):
21          img = self.model(z)
22          img = img.view(img.shape[0], *img_shape)
23          return img
```

We have decided to go with this architecture, after having tested several architectures for both models, including convolutionnal neural network. We noticed that the VAE model works fine with a convolutional architecture whereas the GAN have not got good result with that kind of architecture. Both models have been trained with a latent variable dimension of 100.

A. **Qualitative Evaluations**

1. Visual samples
   The scripts used to generate the samples for the Qualitative evaluation can be found in the QualitativeVAE.py and QualitativeGAN.py under folder Problem 3 for both VAE and GAN models. We have generated different samples from both models (Figures 5 and 6). We notice that the images generated by the VAE seem realistic albeit blurry, whereas the images generated by GAN are more diversified, and seem less realistic.



Figure 5: Samples generated with VAE

Figure 6: Samples generated with GAN.

2. Learning the disentangled representation in the latent space. Figures 7 and 8 show how the models learned a disentangled representation. Starting from one point in the latent space we apply a perturbation to two arbitrary dimensions. We can observe how in both models the resulting images vary smoothly. The value of perturbation ($\epsilon$) is however different for each model, as they require different levels of perturbation to result in visible changes.
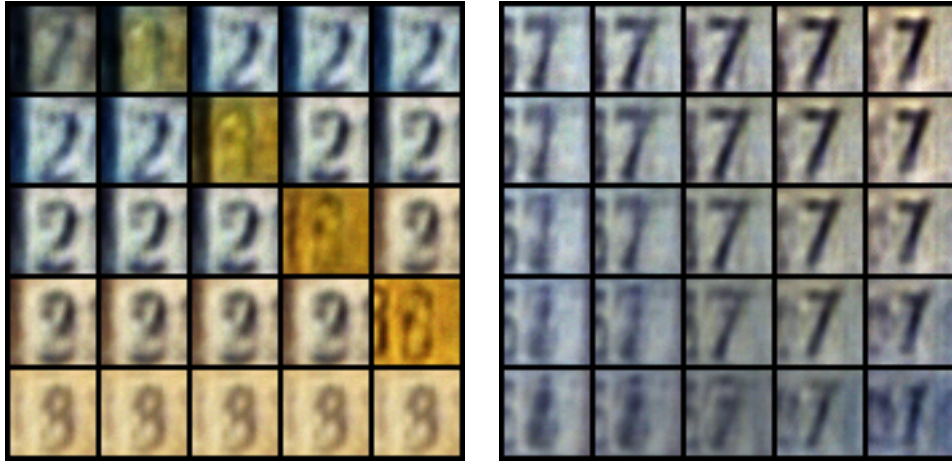


Figure 7: Samples by perturbating the dimensions, left: 6 and 96, right: 25 and 75

Figure 8: Samples by perturbating the dimensions, left: 65 and 90, right: 8 and 60

3. Interpolation in the data space and in the latent space: Figures 9, 10, 11 and12 show the interpolations in the latent and image space for VAE and GAN respectively. For both models we can observe how the interpolation in the latent space produce images that have more *meaning* compared to the interpolation in the image space, where the transition images are only superposition of the patterns (numbers) in each of the initial images. This can be explained by the latent space forming a kind of *manifold* along which the interpolation makes more sense than in the pixel space.

   (a) VAE:



Figure 9: Interpolation in the latent space using VAE



Figure 10: Interpolation in the data space using VAE

   (b) GAN:

9

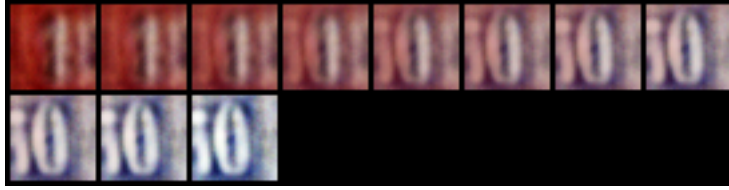Figure 11: Interpolation in the latent space using GAN



Figure 12: Interpolation in the data space using GAN

## B. Quantitative Evaluations

1. We have used the provided functions to extract the representations of the images. We compute the Frechet Inception Distance by estimating the mean and covariance of the generator's/decoder's distribution. The calculation steps are explained in the following code snippet:

```python
def calculate_fid_score(sample_feature_iterator,
                        testset_feature_iterator):
    gen_features = np.array([])
    test_features = np.array([])

    print("Extracting the features ...")
    #sample_feature_iterator is a generator of a minibatch of features images,
    #that is the last conv2d layer of the classifier of 512 features

    #For generated images
    gen_size=0
    for i in sample_feature_iterator: #iterate over minibatch images
        #Now let's get the activation of the images
        gen_features = np.vstack([gen_features,i.reshape(1,512)])  \
        if gen_features.size else i.reshape(1,512)
        gen_size+=1
        if gen_size==1000: break
    gen_features = gen_features.T

    #For test images
    test_size=0
    for i in testset_feature_iterator: #iterate over test images
        test_features= np.vstack([test_features,i.reshape(1,512)])
        test_size+=1
        if test_size==1000: break
    test_features = test_features.T

    print("Estimating the mean ...")
    #Estimating the mean of the generated images
```

10

```
30        mu_gen = np.mean(gen_features, axis=1).reshape(512,1)
31
32        #Estimating the mean of the test images
33        mu_test = np.mean(test_features, axis=1).reshape(512,1)
34
35        print("Estimating the variance ...")
36        # We use the unbiased variance estimate which
37        #is given by (X–mu)(X–mu)^T/(n–1)
38        gen_centered  = gen_features − mu_gen
39        test_centered = test_features − mu_test
40
41        sigma_gen = np.matmul(gen_centered, gen_centered.T) / (gen_size − 1)
42        sigma_test = np.matmul(test_centered, test_centered.T) / (test_size
          − 1)
43
44        print("Calculating the sqrt of cov matrices product ...")
45        # The sqrt of a matrix A needs A to be symetric, but if A, and B
46        # are sysmetric A.B is not symeyric necessarly.
47        # To solve that we use this trick:
48        # sqrt(sigma1 sigma2) = sqrt(A sigma2 A), where A = sqrt(sigma1)
49        # the covariance matrix are by definition symetric
50
51        # to prevent negative values in the cov product
52        eps = np.eye(512) * 1e−5
53
54        root_sigma_gen = linalg.sqrtm(sigma_gen + eps)
55        sigmas_prod = np.matmul(root_sigma_gen,np.matmul(sigma_test,
          root_sigma_gen))
56        # given np.matmul(root_sigma_gen,np.matmul(sigma_test,
          root_sigma_gen)) is symetric:
57        root_sigmas_prod = linalg.sqrtm(sigmas_prod + eps)
58
59        print("Calculating the FID score ...")
60        # Calculating the trace
61        trace = np.trace(sigma_test + sigma_gen − 2.0 * root_sigmas_prod)
62
63        # Calculate the squared norm between means
64        squared_norm = np.sum((mu_test − mu_gen)**2)
65
66        # Calculate the fid score
67        fid = squared_norm + trace
68
69        return fid
```

2. We sampled 1000 images from each generative models and calculate the FID-score as instructed. The results are:

   – For the GAN, the FID score is: 29526.37
   – For the VAE, the FID score is: 51355.12

   This metric confirme our ascertainment that the GAN is more realistic than the VAE, given the ground truth given by the provided classifier.

# References

[1] Github repository for assignment 3
    https://github.com/faresbs/Representation-Learning.git