

[OTP8] Report on lab session #3 - Part II

BADILLO, Jose - MARTINEZ, Gabriela

jose.badillo@student-cs.fr - gabriela.martinez@student-cs.fr

Université Paris-Saclay, CentraleSupélec — December 5, 2019

1 Objective

The objective of this practical session is to build and train a Variational Autoencoder (VAE) with Pytorch. The main goal of such generative model is to provide new random samples close to the real content used to train the network.

For the training phase, we will present the training subset of MNIST data-set which will generate the output vectors of the mean and variance that properly represents the Gaussian distribution of the data.

During the test phase we will measure how close the test samples are with regards to the network distribution.

Finally we will generate new samples coming from such trained model.

2 Variational Autoencoder Architecture

A variational autoencoder keeps an encoding distribution regularized and improved during the training phase so the latent space is good enough to produce new samples. The VAE was named after the Variational Inference method used in statistics.

The data was previously accessed from <http://deeplearning.net/data/mnist/mnist.pkl.gz>. and divided into training, dev, and test sets.

2.1 Encoder

The implemented architecture starts with the *encoding* module implementing a network with one hidden layer and two different output layers representing the parameters of the Gaussian distribution which characterize the data distribution over the latent space.

```
1
2 class Encoder(nn.Module):
3
4     def __init__(self, input_dim, hidden_dim, z_dim):
5         super(Encoder, self).__init__()
6
7         self.linear = nn.Linear(input_dim, hidden_dim)
8         self.mu = nn.Linear(hidden_dim, z_dim)
9         self.var = nn.Linear(hidden_dim, z_dim)
10
11     def forward(self, x):
12         # x is of shape [batch_size, input_dim]
13         hidden = F.relu(self.linear(x))
14         # hidden is of shape [batch_size, hidden_dim]
15         z_mu = self.mu(hidden)
16         # z_mu is of shape [batch_size, latent_dim]
17         z_var = self.var(hidden)
18         # z_var is of shape [batch_size, latent_dim]
19         return z_mu, z_var
```

2.2 Decoder

A decoder is a multi-layer perceptron that will take as input a sample from the latent distribution and after passing through the network it will deliver a new representation of such sample in the original dimensions.

```
1
2 class Decoder(nn.Module):
3
4     def __init__(self, z_dim, hidden_dim, output_dim):
5         super(Decoder, self).__init__()
6
7         self.linear = nn.Linear(z_dim, hidden_dim)
8         self.out = nn.Linear(hidden_dim, output_dim)
9
10    def forward(self, x):
11        # x is of shape [batch_size, latent_dim]
12        hidden = F.relu(self.linear(x))
13        # hidden is of shape [batch_size, hidden_dim]
14        predicted = torch.sigmoid(self.out(hidden))
15        # predicted is of shape [batch_size, output_dim]
16        return predicted
```

2.3 VAE

The creation of the auto-encoder consists on the conjunction of the encoder and the decoder like any other model of such type. The particularity of this generative model is that in order to obtain the newly created output image, we obtain a random sample from the latent space z and during training we are supposed to back-propagate the error. Nevertheless due to the randomness, we cannot be sure on where such gradients should be back-propagated. In order to solve such limitation we used a reparameterization trick.

This is possible taking into account that the posterior is normally distributed and then we can make an approximation using another normal distribution with another parameters.

```
1 class VAE(nn.Module):
2     def __init__(self, enc, dec):
3         super(VAE, self).__init__()
4
5         self.enc = enc
6         self.dec = dec
7
8     def forward(self, x):
9         # encode
10        z_mu, z_var = self.enc(x)
11        # sample from the distribution having latent parameters z_mu, z_var
12        # reparameterize
13        std = torch.sqrt( torch.exp( z_var ) )
14        eps = torch.normal(0, 1., z_mu.shape)
15        x_sample = z_mu + eps * std
16        # decode
17        predicted = self.dec(x_sample)
18        return predicted, z_mu, z_var
```

3 Network training

3.1 The training loop

In order to avoid over-fitting and be sure that we could use the auto-encoder as a generative model it is necessary to add a regularization during the training process. The regularization performed for such

purposes is the encoding as a distribution on the latent space. The main steps to train a variational auto-encoder are depicted below [1]:

1. The input is encoded as distribution over the latent space
2. A point from the latent space is sampled from that distribution
3. The sampled point is decoded and the reconstruction error can be computed
4. The reconstruction error is back-propagated through the network

```
1 def train():
2
3     # set the train mode
4     model.train()
5
6     # Use gradient clipping
7     torch.nn.utils.clip_grad_value_(model.parameters(), 5.)
8
9     train_loss = 0
10
11     for i, (x, _) in enumerate(train_iterator):
12
13         # reshape the data into [batch_size, 784]
14         x = x.view(-1, 28 * 28)
15         x = x.to(device)
16
17         # update the gradients to zero
18         optimizer.zero_grad()
19
20         # forward pass
21         x_sample, z_mu, z_var = model(x)
22
23         # total loss
24         loss = loss_function( x_sample, x, z_mu, z_var )
25
26         # backward pass
27         loss.backward()
28         train_loss += loss.item()
29
30         # update the weights
31         optimizer.step()
32
33     return train_loss
```

The loss function to be minimized during training comes from the Variational Inference statistics method, and it is composed of a "reconstruction" and a "regularization term". The latter is expressed as the Kullback-Leiber divergence between the returned distribution and a normal Gaussian.

```
1 # Reconstruction + KL divergence losses summed over all elements and batch
2 def loss_function(reconstruction_x, x, z_mean, z_var):
3
4     # reconstruction loss
5     BCE = F.binary_cross_entropy(reconstruction_x, x.view(-1, 784), reduction='sum')
6     # kl divergence loss
7     KLD = -0.5 * torch.sum(1 + z_var - z_mean.pow(2) - z_var.exp())
8
9     return BCE + KLD
```

3.2 The test loop

Next to the training we have developed a test loop which will assess how similar the current trained distribution of the VAE is being with regards of the images unknown for the network.

```
1 def test():
2
3     # set the evaluation mode
4     model.eval()
5
6     test_loss = 0
7
8     # we don't need to track the gradients, since we are not updating the parameters during
9     # evaluation / testing
10    with torch.no_grad():
11        for i, (x, _) in enumerate(test_iterator):
12
13            # reshape the data
14            x = x.view(-1, 28 * 28)
15            x = x.to(device)
16
17            # forward pass
18            x_sample, z_mu, z_var = model(x)
19
20            # total loss
21            loss = loss_function( x_sample, x, z_mu, z_var )
22            test_loss += loss.item()
23
24    return test_loss
```

One of the key factors in order to avoid the over-fitting is to stop the training after certain number of epochs. This is done considering as stop factor the amount of times the best_test_loss has not being improved.

```
1 best_test_loss = float('inf')
2 trainLoss = []
3 testLoss = []
4
5 for epoch in range(N_EPOCHS):
6
7     train_loss = train()
8     test_loss = test()
9     loss = []
10    train_loss /= len(train_dataset)
11    test_loss /= len(test_dataset)
12
13    if epoch % 5 == 0:
14        print('>> epoch ', epoch)
15        print('----> Train loss: >> ', train_loss)
16        print('-----> Test loss: >> ', test_loss)
17        print('')
18
19    trainLoss.append( train_loss )
20    testLoss.append( test_loss )
21
22    if best_test_loss > test_loss:
23        best_test_loss = test_loss
24        patience_counter = 1
25    else:
26        patience_counter += 1
27
28    if patience_counter > 3:
```

3.3 Gradient clipping

Also, we are implementing gradient clipping in the training process, so we can deal with numerical overflow or underflow (exploding gradients).

4 Experiments

4.1 Testing the batch-size

We trained the model several times with 4 different batch-sizes: 64, 128, 256 and 512. The hidden dimension was fixed to 784 so the encoding keeps all the pixel information. The variation of such dimension will be tested in further experiments. The latent dimension was arbitrary selected to be 2 for this experiment however in further test we are going to vary such value. The number of epochs was fixed to 50 to verify the behaviour of the `batch_size`, however as shown in the graphs of Figure 1 the test accuracy was worsen after 25-30 epochs so then the solution automatically stops the fitting phase.

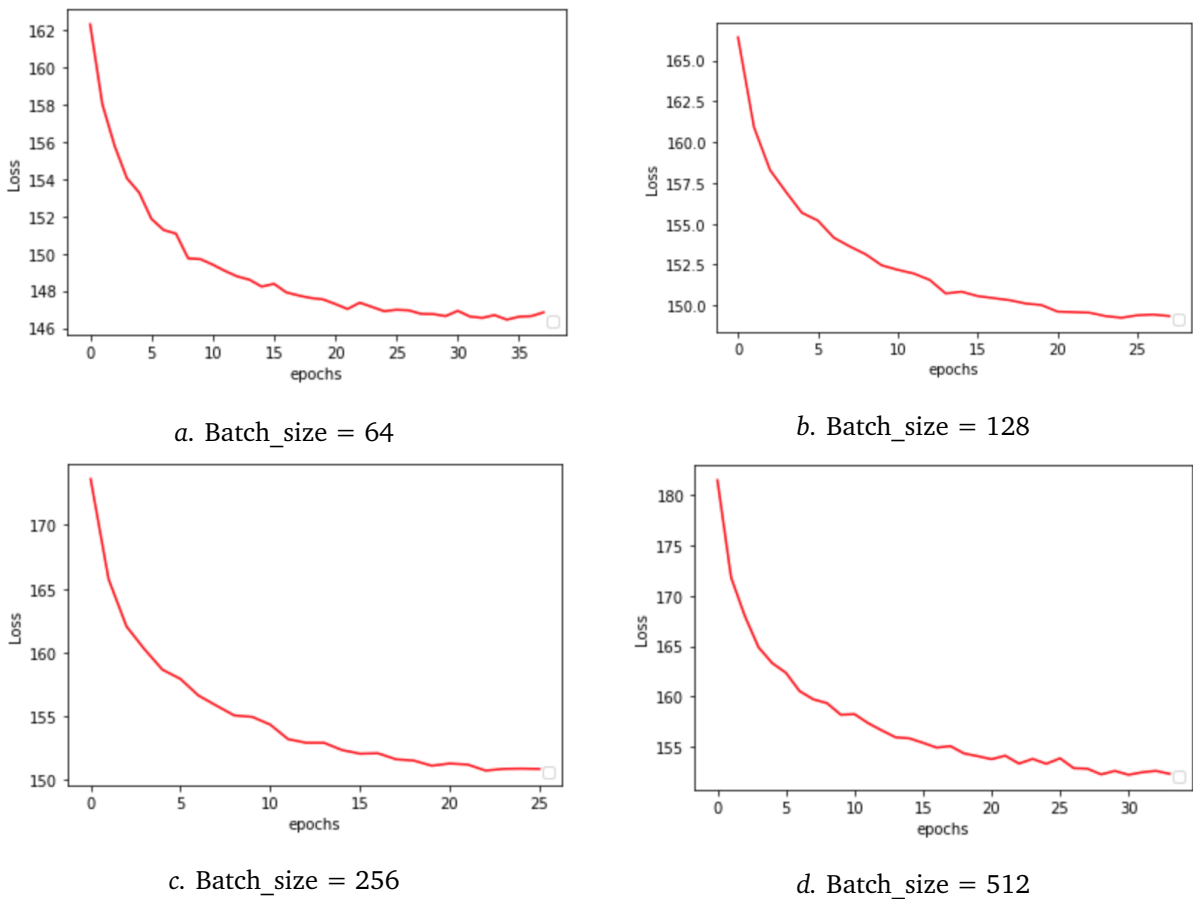


Figure 1. Testing loss of the Variational Autoencoder using different batch_sizes

From the results shown above it is clear that the lower the `batch_size` the more reduction in loss throughout the epochs. Although the variation of such parameter is not impacting in a severe way, we identified the `batch_size` with value 256 as the one presenting more stability when reducing the test loss. Moreover it is obtaining a low value of loss in lower number of epochs with a bigger amount of samples to be trained. Therefore this will help to improve the performance time during the fitting phase. The following experiments during this lab will use 256 `batch_size` value. However this is not a conclusive value and more testing would have to be done.

4.2 Testing the latent dimension

We have trained the variational auto-encoder with 3 different values for the latent dimension: 2, 10, 20. For this purpose we decided to keep a hidden dimension of 748 so we don't lose any detail of the pixels. The goal of this test is to find the most appropriate value for such latent dimension.

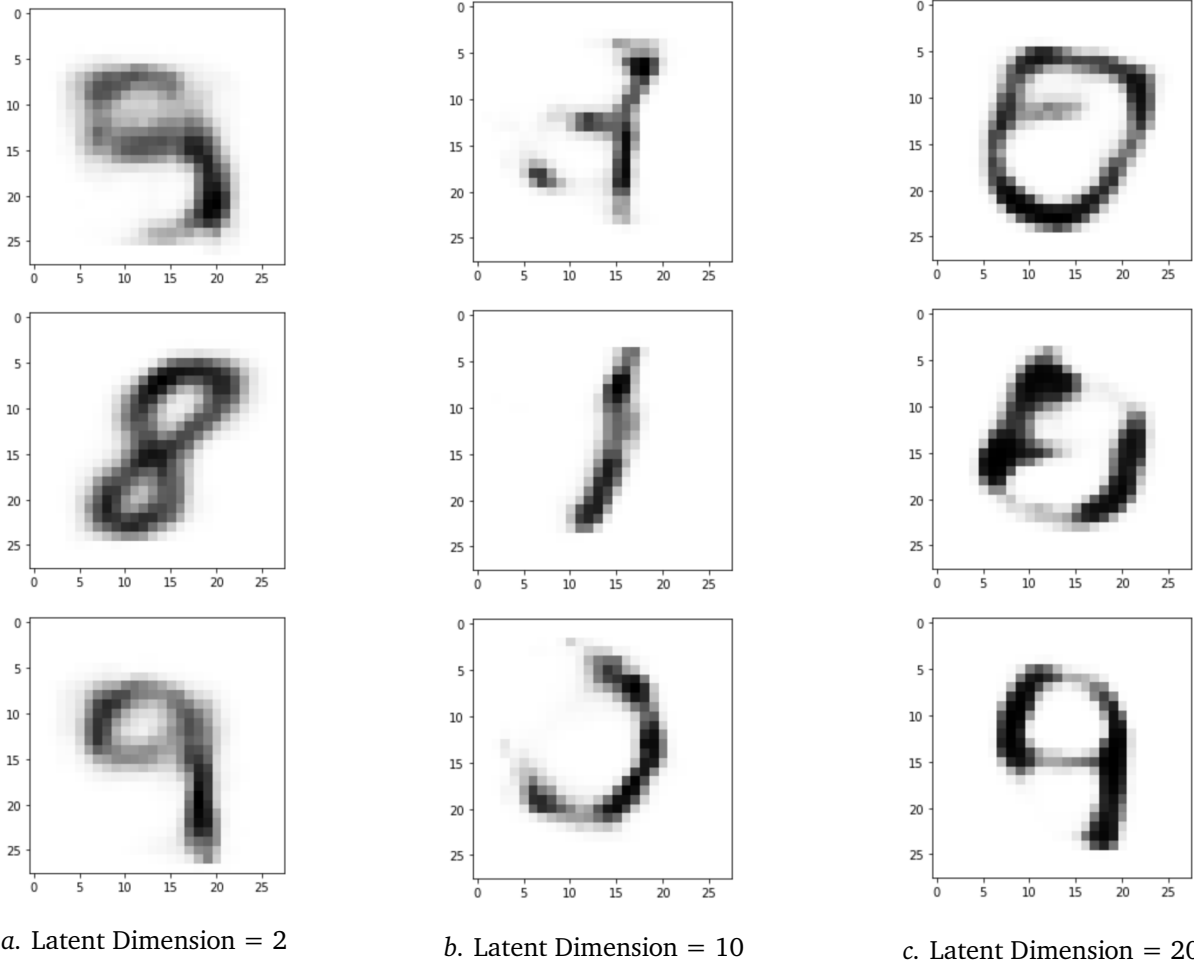


Figure 2. Generation of MNIST images from VAE with 784 nodes in hidden dimension and different values for latent dimensions

In Figure 2 we can appreciate that, when the latent dimension is bigger, the images generated from the VAE tend to produce representations that are not well defined or are far from the real representations of the original data-set. In the other hand when the latent dimension is 2 or minimal we found more meaningful images and with a better shape defined. Considering this experiment we have decided to fix the latent dimension in 2 for further experiments.

4.3 Testing the hidden dimension

We trained three different architectures of Variational Autoencoders varying the hidden layer: (1000-2, 500-2 and 100-2) For all these architectures we used a ReLU as non linear function.

Each one of these architectures was trained with a fixed batch_size of 256 samples and run with a maximum value of 50 epochs. We must recall the epochs will be different depending on the values obtained for the *best_test_loss* during the train process.

4.3.1 Experiment with 1000 nodes in hidden dimension

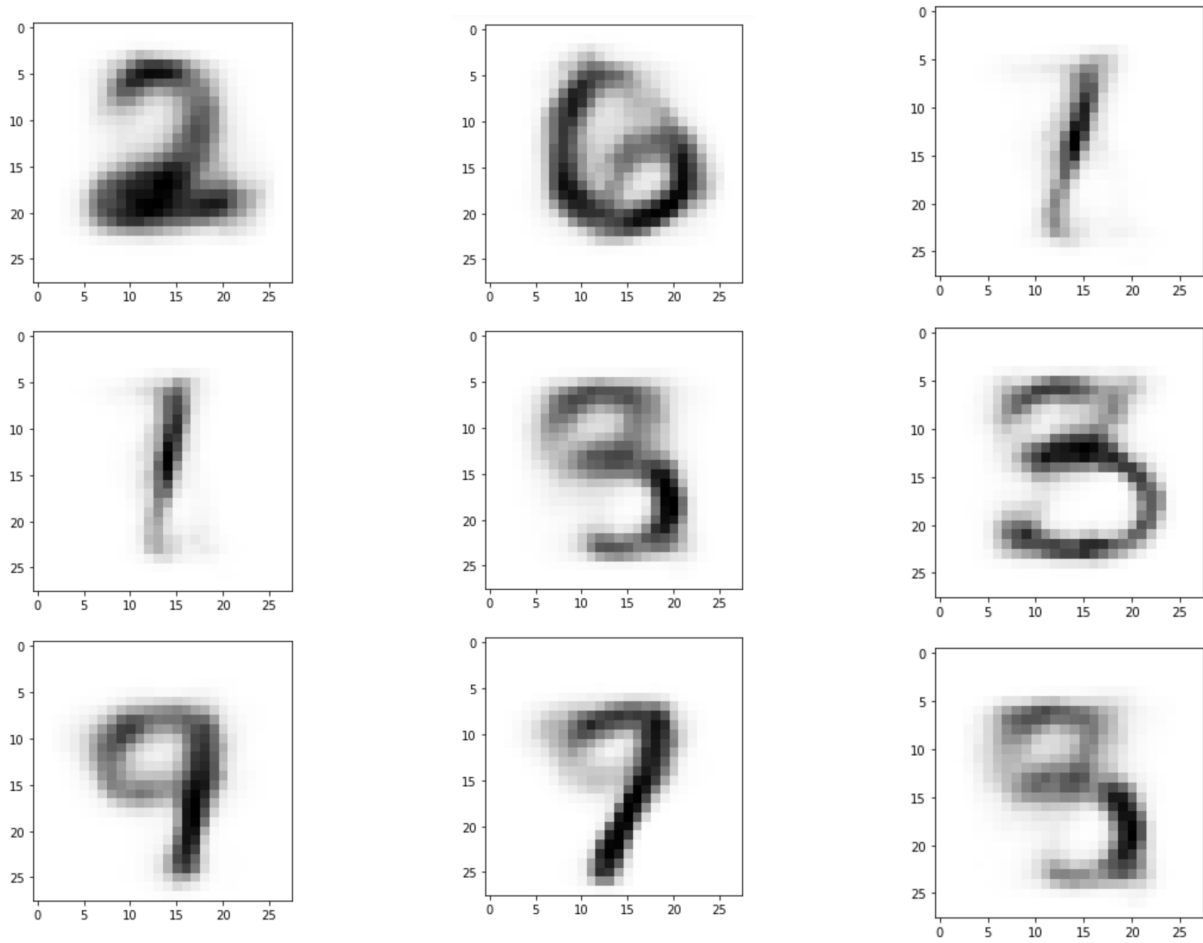


Figure 3. Generation of MNIST images from VAE with 1000 nodes in hidden dimension and 2 nodes in latent dimension

The nine images generated above from the Gaussian distribution trained within the variational auto-encoder are well defined. This satisfy the continuity (two close points in the latent space should not give two completely different contents once decoded) and completeness (for a chosen distribution, a point sampled from the latent space should give “meaningful” content once decoded) characteristics that a generative model must hold.

It is clear that our VAE is generating real values of digits. In a certain level it can be appreciated certain mixes among the digits distributions like the image with the number six where it shows a gray, blurry completion of the figure that could be trying to generate a zero, but such noise could be removed with help of other kind of models like the de-noising auto-encoder.

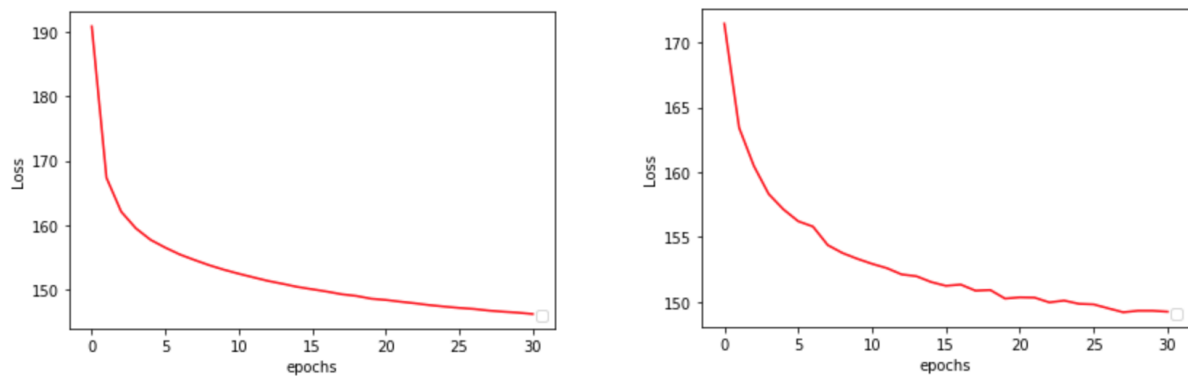


Figure 4. Training and Testing Loss functions for VAE 1000-2

4.3.2 Experiment with 500 nodes in hidden dimension

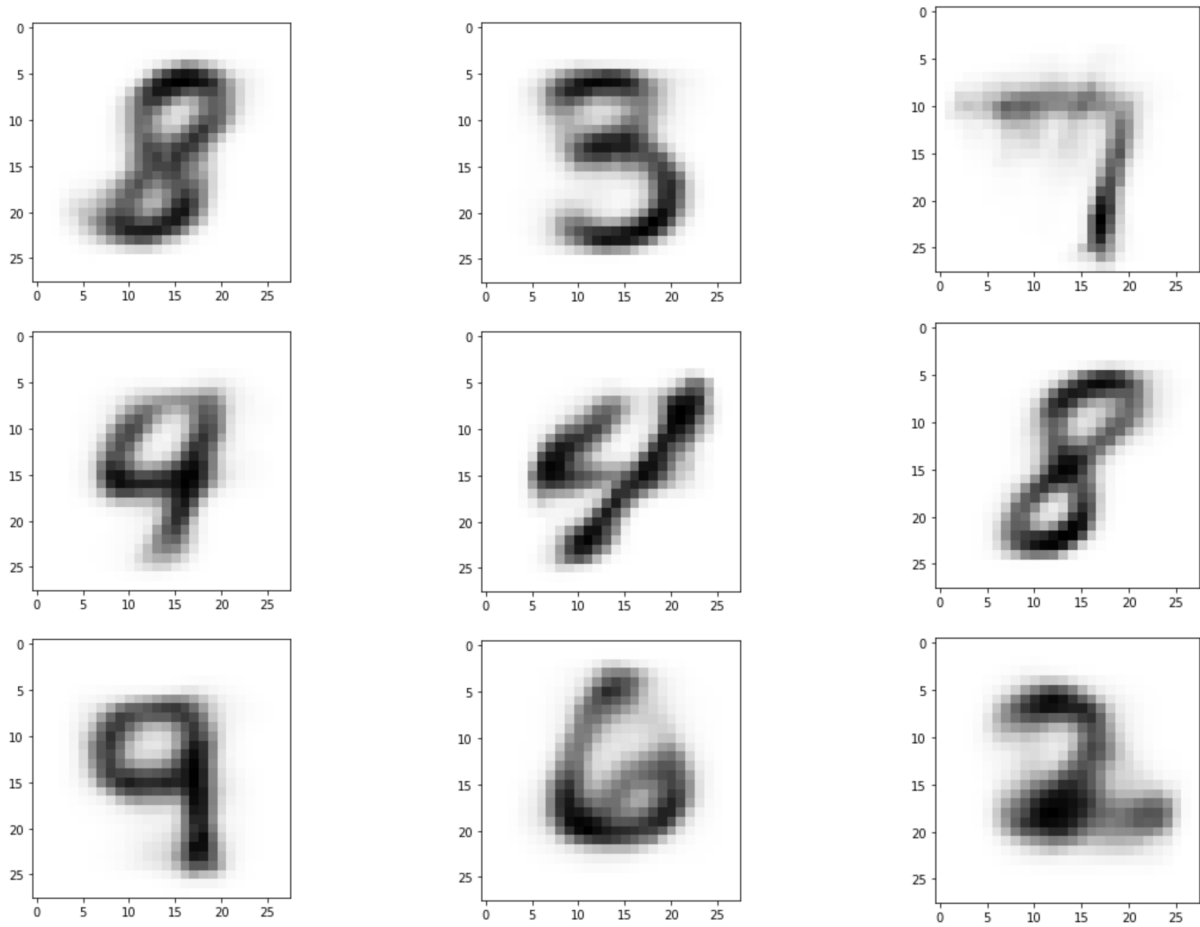


Figure 5. Generation of MNIST images from VAE with 500 nodes in hidden dimension and 2 nodes in latent dimension

The Figure 5 above show a meaningful generation of samples coming from the trained VAE. The hidden dimension of 500 was lower than the original representation of 784 but still close enough to the original dimensions in order to be able to keep the main characteristics of the data and generate proper samples than can be clearly identified.

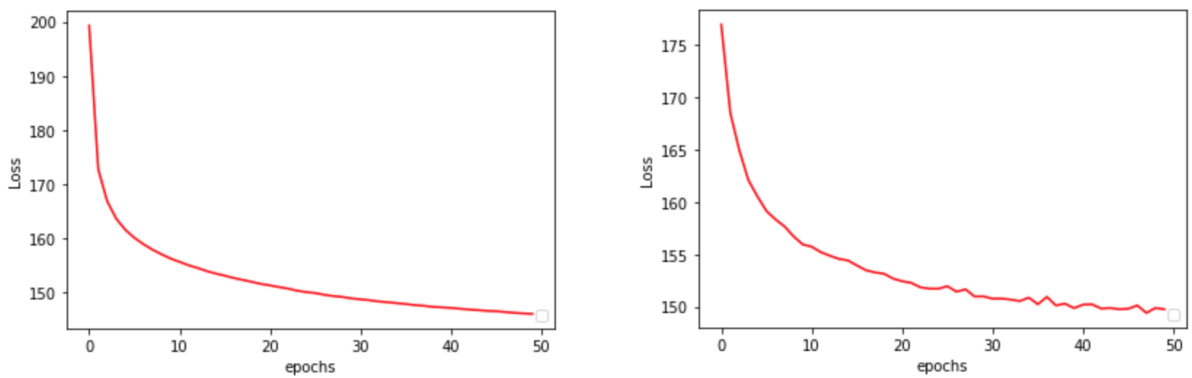


Figure 6. Training and Testing Loss functions for VAE 500-2

4.3.3 Experiment with 100 nodes in hidden dimension

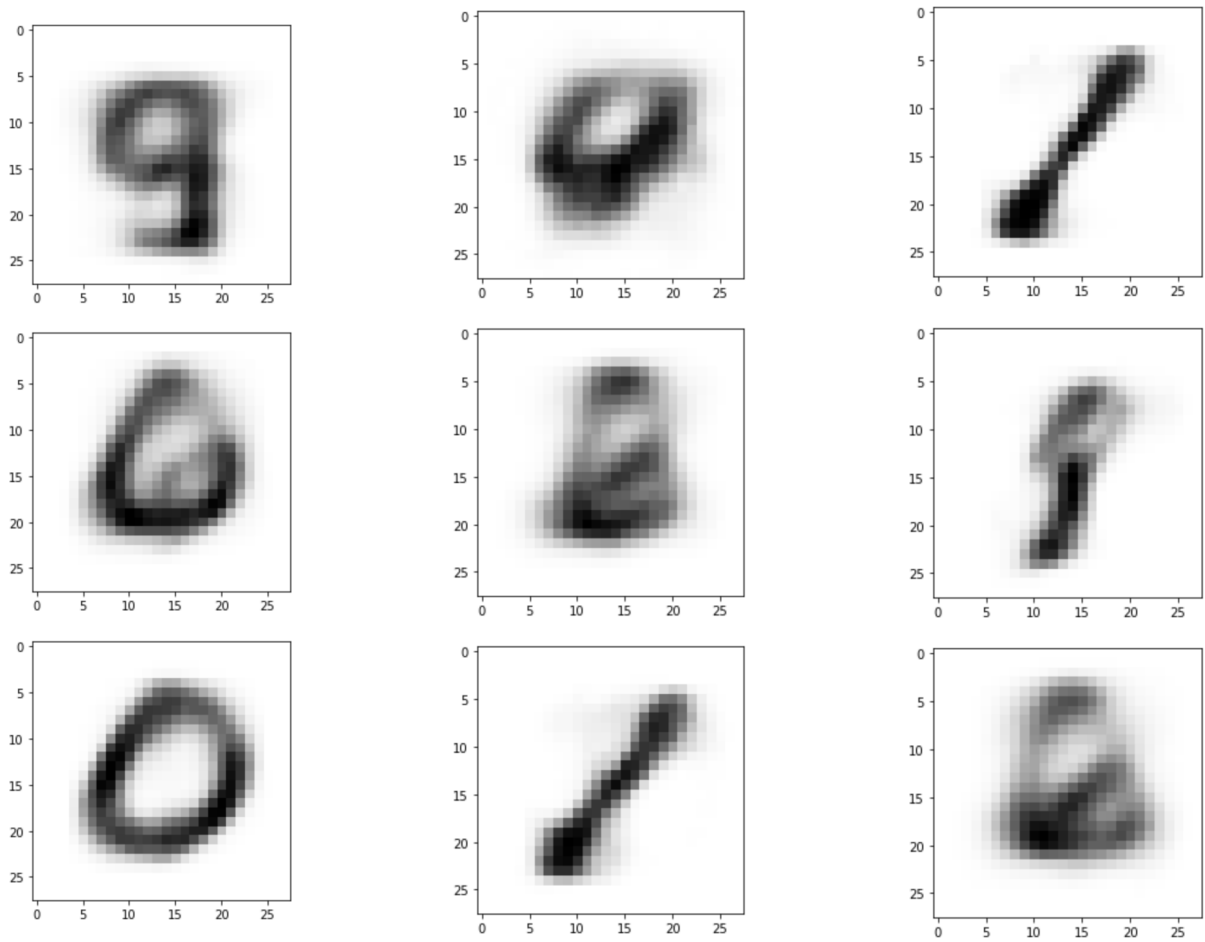


Figure 7. Generation of MNIST images from VAE with 100 nodes in hidden dimension and 2 nodes in latent dimension

The Figure 7 shows different samples generated from the VAE with 100 nodes of hidden dimension. We can clearly see that we have lost a big amount of information in the encoding phase as the decoder is not able to generate meaningful content. Certain classes like zeros or ones are still useful, however it is not easy to visualize the nines or eights as clear as we did with the previous configurations.

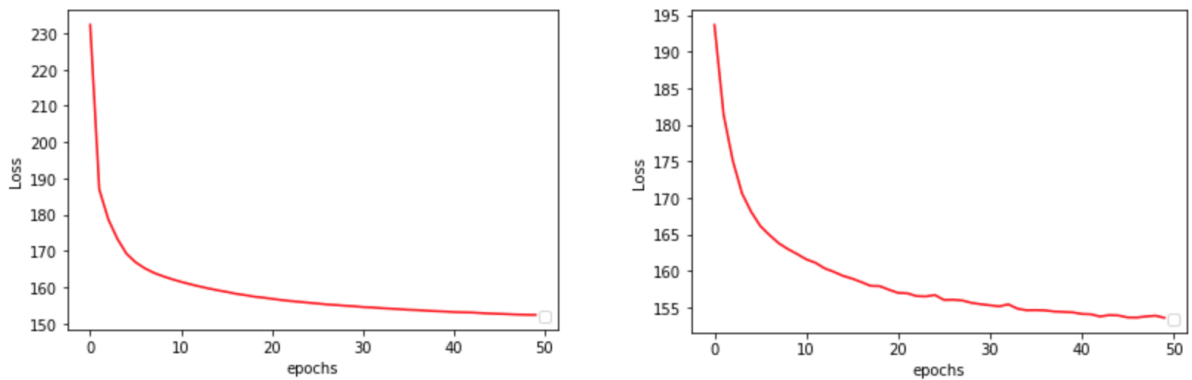


Figure 8. Training and Testing Loss functions for VAE 100-2

5 Final notes

- For the MNIST data-set when training deep learning models, a small `batch_size` is important to avoid impacts on the accuracy of the model. For variational auto-encoder the results show a similar behaviour: the lower the `batch_size` the better performance the model will have when reducing the loss function. However, it seems the impact in the performance of this generative model is not huge it is possible to use big batches like 128 or 256 without losing generative capabilities.
- For the MNIST data-set, if we consider a hidden dimension with dimensions similar than the original dimensions of an image, the fitting phase of the VAE converge after 25-30 epochs. The latter when using a very small `step_size` value of 0.001. However when we vary the number of nodes and generally we apply dimensionality reduction, then more epochs are necessary in order to converge (as expected). The number of epochs can be reduced with the variation of the `step_size` to a bigger one.
- The hidden dimension in the variational auto-encoder can be lower to the original if the purpose of the model is to reduce dimensions and then provide a distribution representation of the data. We must be aware that, the lower the dimension of the hidden layer, the less accurate are going to be the samples generated by the auto-encoder.

References

- [1] ROCCA, J. Understanding variational autoencoders (vae).