

Assignment 3 Solutions**1 Neural Network**

a) Layer 1:

$$\vec{z}_1 = W^{(0)} \vec{x} = \begin{pmatrix} 0.3 \\ 0.3 \\ 0.3 \end{pmatrix} \Rightarrow \vec{h}_1 = g(\vec{z}_1) = \begin{pmatrix} 0.5744 \\ 0.5744 \\ 0.5744 \end{pmatrix}$$

Layer 2:

$$\vec{z}_2 = W^{(1)} \vec{h}_1 = \begin{pmatrix} 0.5117 \\ 0.5117 \\ 0.5117 \end{pmatrix} \Rightarrow \vec{h}_2 = g(\vec{z}_2) = \begin{pmatrix} 0.6264 \\ 0.6264 \\ 0.6264 \end{pmatrix}$$

Layer 3 (output layer):

$$\hat{y} = W^{(2)} \vec{h}_2 = 0.9397$$

b) Step 1: Gradients calculation

$$\frac{\partial L}{\partial \hat{y}} = \hat{y} - y = -1.0603 \quad (1)$$

$$\frac{\partial L}{\partial \vec{z}_2} = \frac{\partial \hat{y}}{\partial \vec{z}_2} \frac{\partial L}{\partial \hat{y}} \quad (2a)$$

$$= \frac{\partial \vec{h}_2}{\partial \vec{z}_2} \frac{\partial \hat{y}}{\partial \vec{h}_2} \frac{\partial L}{\partial \hat{y}} \quad (2b)$$

$$= \begin{pmatrix} g'(z_{2,1}) & 0 & 0 \\ 0 & g'(z_{2,2}) & 0 \\ 0 & 0 & g'(z_{2,3}) \end{pmatrix} (W^{(2)})^\top \frac{\partial L}{\partial \hat{y}} \quad (2c)$$

$$= \begin{pmatrix} g'(z_{2,1}) w_{1,1}^{(2)} \frac{\partial L}{\partial \hat{y}} \\ g'(z_{2,2}) w_{1,2}^{(2)} \frac{\partial L}{\partial \hat{y}} \\ g'(z_{2,3}) w_{1,3}^{(2)} \frac{\partial L}{\partial \hat{y}} \end{pmatrix} \quad (2d)$$

$$= \begin{pmatrix} -0.1241 \\ -0.1241 \\ -0.1241 \end{pmatrix} \quad (2e)$$

$$\frac{\partial L}{\partial \vec{z}_1} = \frac{\partial \vec{z}_2}{\partial \vec{z}_1} \frac{\partial L}{\partial \vec{z}_2} \quad (3a)$$

$$= \frac{\partial \vec{h}_1}{\partial \vec{z}_1} \frac{\partial \vec{z}_2}{\partial \vec{h}_1} \frac{\partial L}{\partial \vec{z}_2} \quad (3b)$$

$$= \begin{pmatrix} g'(z_{1,1}) & 0 & 0 \\ 0 & g'(z_{1,2}) & 0 \\ 0 & 0 & g'(z_{1,3}) \end{pmatrix} (W^{(1)})^\top \frac{\partial L}{\partial \vec{z}_2} \quad (3c)$$

$$= \begin{pmatrix} g'(z_{1,1}) \left(\frac{\partial L}{\partial z_{2,1}} w_{1,1}^{(1)} + \frac{\partial L}{\partial z_{2,2}} w_{2,1}^{(1)} + \frac{\partial L}{\partial z_{2,3}} w_{3,1}^{(1)} \right) \\ g'(z_{1,2}) \left(\frac{\partial L}{\partial z_{2,1}} w_{1,2}^{(1)} + \frac{\partial L}{\partial z_{2,2}} w_{2,2}^{(1)} + \frac{\partial L}{\partial z_{2,3}} w_{3,2}^{(1)} \right) \\ g'(z_{1,3}) \left(\frac{\partial L}{\partial z_{2,1}} w_{1,3}^{(1)} + \frac{\partial L}{\partial z_{2,2}} w_{2,3}^{(1)} + \frac{\partial L}{\partial z_{2,3}} w_{3,3}^{(1)} \right) \end{pmatrix} \quad (3d)$$

$$= \begin{pmatrix} -0.0273 \\ -0.0273 \\ -0.0273 \end{pmatrix} \quad (3e)$$

Step 2: Update weights

$$W^{(2)} \leftarrow W^{(2)} - \eta (\vec{h}^{(2)})^\top \frac{\partial L}{\partial \hat{y}} = \begin{bmatrix} 1.1642 & 1.1642 & 1.1642 \end{bmatrix} \quad (4a)$$

$$W^{(1)} \leftarrow W^{(1)} - \eta \vec{h}_1 \left(\frac{\partial L}{\partial \vec{z}_2} \right)^\top = \begin{bmatrix} 0.3713 & 0.3713 & 0.3713 \\ 0.3713 & 0.3713 & 0.3713 \\ 0.3713 & 0.3713 & 0.3713 \end{bmatrix} \quad (4b)$$

$$W^{(0)} \leftarrow W^{(0)} - \eta \vec{x} \left(\frac{\partial L}{\partial \vec{z}_1} \right)^\top = \begin{bmatrix} 0.1273 & 0.1273 & 0.1273 \\ 0.1273 & 0.1273 & 0.1273 \\ 0.1273 & 0.1273 & 0.1273 \end{bmatrix} \quad (4c)$$

Note that we could have also computed the partial derivatives with respect to each component of the weight matrices, or with respect to each row of the weight matrices, to get the same result.

2 Autoencoders

2.1 part a

The attached code block shows how we implemented the Vanilla autoencoder including Encoder, Decoder, and forward passes. We also modified *Train* and *Validate* function in the .py file so as to return loss values for plotting loss curves. Although it is not required for students to implement/report, it helps to interpret results.

```

1 class Autoencoder(nn.Module):
2
3     def __init__(self, dim_latent_representation=2):
4
5         super(Autoencoder, self).__init__()
6
7         class Encoder(nn.Module):
8             def __init__(self, output_size=2):
9                 super(Encoder, self).__init__()
10
11                 self.nn = nn.Sequential(
12                     nn.Linear(28 * 28, output_size),
13                 )
14
15             def forward(self, x):
16                 return self.nn(x)
17
18         class Decoder(nn.Module):
19             def __init__(self, input_size=2):
20                 super(Decoder, self).__init__()
21
22                 self.nn = nn.Sequential(
23                     nn.Linear(input_size, 28 * 28),
24                     nn.Sigmoid(),
25                 )
26
27             def forward(self, z):
28                 return self.nn(z)
29
30         self.encoder = Encoder(output_size=dim_latent_representation)
31         self.decoder = Decoder(input_size=dim_latent_representation)
32
33     def forward(self, x):
34         x = self.encoder(x)
35         x = self.decoder(x)
36         return x

```

Listing 1: Part (a) architecture

Results from this autoencoder are as follows:

- Epoch: 20
- Training loss: 0.7271
- Validation loss (reconstruction error) : 0.7330

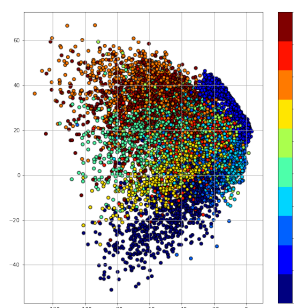


Figure 1: Scatter plot of the latent representation (Part a)

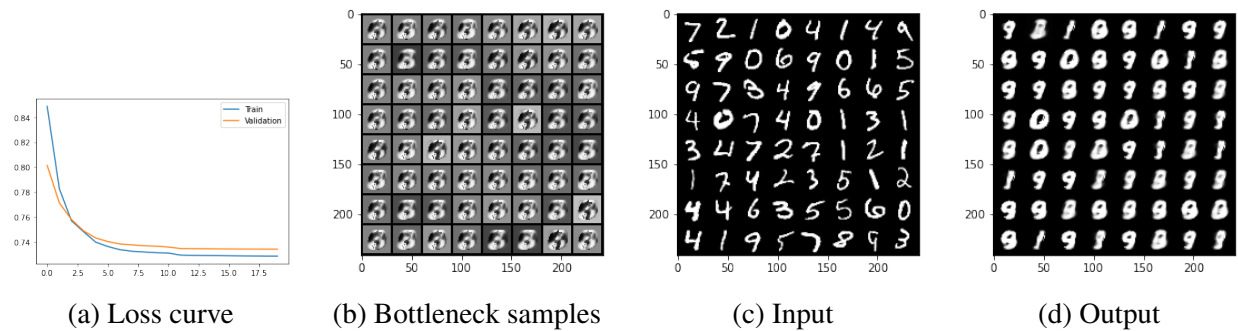


Figure 2: Results from part a

2.2 part b

```

1 class Autoencoder(nn.Module):
2
3     def __init__(self, dim_latent_representation=2):
4
5         super(Autoencoder, self).__init__()
6
7         class Encoder(nn.Module):
8             def __init__(self, output_size=2):
9                 super(Encoder, self).__init__()
10
11                 self.nn = nn.Sequential(
12                     nn.Linear(28 * 28, 1024),
13                     nn.ReLU(),
14                     nn.Linear(1024, output_size)
15                 )
16
17             def forward(self, x):
18                 return self.nn(x)
19
20         class Decoder(nn.Module):
21             def __init__(self, input_size=2):
22                 super(Decoder, self).__init__()
23
24                 self.nn = nn.Sequential(
25                     nn.Linear(input_size, 1024),
26                     nn.ReLU(),
27                     nn.Linear(1024, 28 * 28)
28                 )
29
30             def forward(self, z):
31                 return self.nn(z)
32
33         self.encoder = Encoder(output_size=dim_latent_representation)
34         self.decoder = Decoder(input_size=dim_latent_representation)
35
36     def forward(self, x):
37         x = self.encoder(x)
38         x = self.decoder(x)
39         return x

```

Listing 2: Part (b) architecture

- Epoch: 20
- Training loss: 0.6621
- Validation loss (reconstruction error) : 0.6717

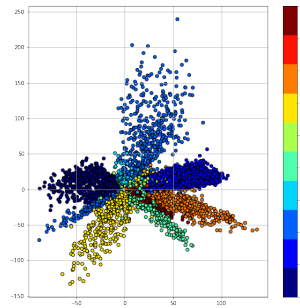


Figure 3: Scatter plot of the latent representation (Part b)

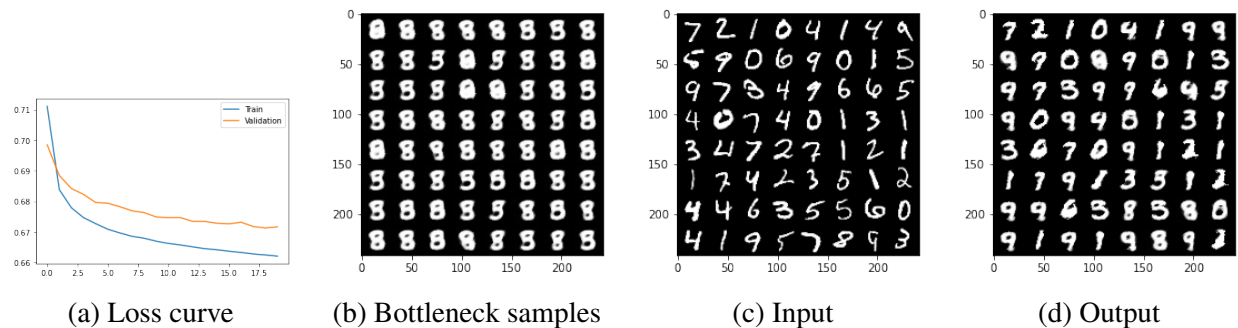


Figure 4: Results from part b

Now we compare two factors we changed here compared to the previous part. In part (a), we see undercomplete autoencoder since the 2D hidden layer is too smaller than the input (784). Using too few hidden units will result in underfitting. This constraint will also impose the model to lose part of the data information during the compression process. However, additional layers can learn more expressive model with complex representations. This can be concluded from the reconstruction loss values as well as the reconstructed sample images.

As we can see in the scatter plot, although some of the labels are encoded with similar values, this model is providing more distinctive code to the decoder than the previous model. In fact, any neural network with one fully connected layer can be used to represent linearly separable functions while adding an activation function to the encoder enables the network to capture non-linearity.

For better understanding of contribution of each change (not required for students), we also trained the same model, with 1024 hidden units and no *ReLU* activation function so called ablation study. In the following, we included a side-by-side scatter plot of these three models. Training this model is beneficial since it has only one changed factor compared to part a and part b.

```

1 Autoencoder(
2   (encoder): Encoder(
3     (nn): Sequential(
4       (0): Linear(in_features=784, out_features=1024, bias=True)
5       (1): Linear(in_features=1024, out_features=2, bias=True)
6     )
7   )
8   (decoder): Decoder(
9     (nn): Sequential(
10      (0): Linear(in_features=2, out_features=1024, bias=True)
11      (1): Linear(in_features=1024, out_features=784, bias=True)
12      (2): Sigmoid()
13    )
14  )

```

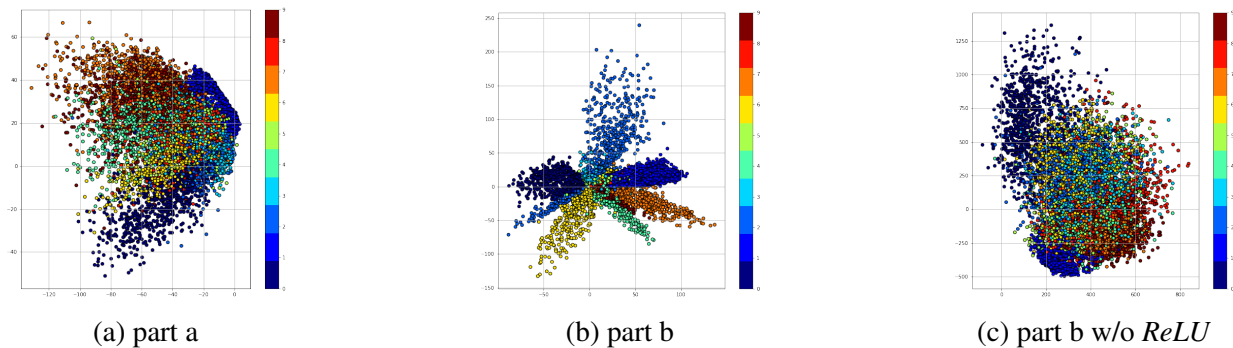


Figure 5: Latent space scatter plot for part a, part b, and part b without ReLU activation function

```
14 )
15 )
```

Listing 3: Part (b) architecture without ReLU activation

Note that we know that another advantage of ReLU could be avoiding vanishing gradients. However, since we use a shallow network, this answer would not be acceptable.

2.3 part c

In the following we provided the implementation and results for part c.

```
1 class Autoencoder(nn.Module):
2
3     def __init__(self, dim_latent_representation=2):
4
5         super(Autoencoder, self).__init__()
6
7         class Encoder(nn.Module):
8             def __init__(self, output_size=2):
9                 super(Encoder, self).__init__()
10
11                 self.nn = nn.Sequential(
12                     nn.Linear(28*28, 1024),
13                     nn.Sigmoid(),
14                     nn.Linear(1024, output_size)
15                 )
16
17             def forward(self, x):
18                 return self.nn(x)
19
20
21         class Decoder(nn.Module):
22             def __init__(self, input_size=2):
23                 super(Decoder, self).__init__()
24
25                 self.nn = nn.Sequential(
26                     nn.Linear(input_size, 1024),
27                     nn.Sigmoid(),
28                     nn.Linear(1024, 28*28),
29                     nn.Sigmoid()
30                 )
31
32             def forward(self, z):
33                 return self.nn(z)
34
35
36         self.encoder = Encoder(output_size=dim_latent_representation)
37         self.decoder = Decoder(input_size=dim_latent_representation)
38
39     def forward(self, x):
40         x = self.encoder(x)
41         x = self.decoder(x)
```

46

return x

Listing 4: Part (b) architecture without ReLU activation

- Epoch: 20
- Training loss: 0.6655
- Validation loss (reconstruction error) : 0.6770

We plotted the activation functions in fig. 6. Regarding fig. 3 and fig. 7, we can immediately observe that the benefit of ReLUs is sparsity. Sparsity arises when $input \leq 0$. The more units with negative values exist in a layer the more sparse the resulting representation. Hence, if certain neurons are less important in terms of their weights, they will be removed, and the model is sparse. On the other hand, Sigmoids shape a dense representation since they are likely to generate some non-zero value for most of the input values.

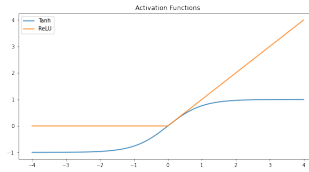


Figure 6: ReLU and Tanh activation functions.

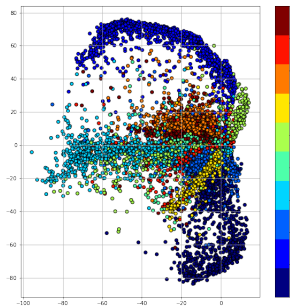


Figure 7: Scatter plot of the latent representation (Part c)

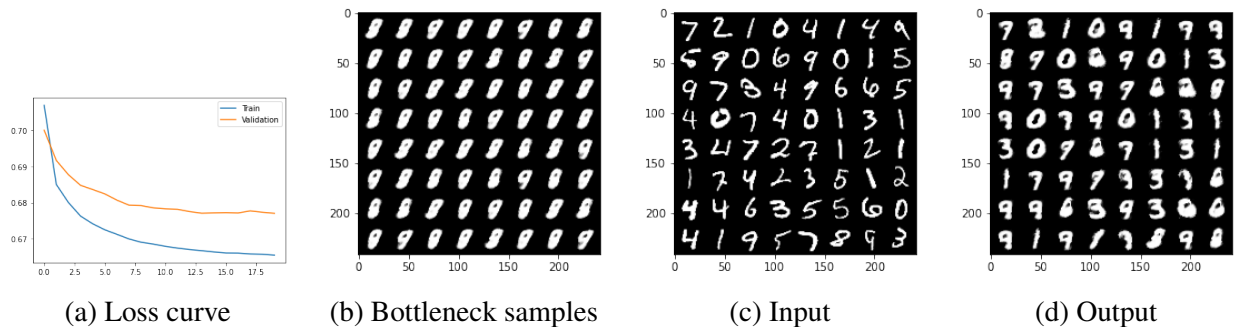


Figure 8: Results from part c

2.4 Part d

First, we implement a vanilla autoencoder similar to the part (a) with 30 dimensions as the bottle-neck size and replace *Sigmoid* to *Tanh*. The following code block shows the architecture of this step (not the implementation).

```

1 Autoencoder(
2   (encoder): Encoder(
3     (nn): Sequential(
4       (0): Linear(in_features=784, out_features=30, bias=True)
5     )
6   )
7   (decoder): Decoder(
8     (nn): Sequential(
9       (0): Linear(in_features=30, out_features=784, bias=True)
10      (1): Tanh()
11    )
12  )
13 )

```

Listing 5: Part (d) vanilla autoencoder architecture.

- Epoch: 20
- Training loss: 0.4089
- Validation loss (reconstruction error) : 0.4171

Afterward, we implemented the VAE model as shown in the code block 6. We trained this model with different KLD weights from the set 0.5, 1, 5 and loss values were as follows. Also, Fig. 9, 10, 11, and 12 show the latent space in 2D space, unseen generated samples, and reconstructed inputs for all trained models in part d.

```

1 class VAE(nn.Module):
2
3
4   def __init__(self, dim_latent_representation=2):
5
6     super(VAE, self).__init__()
7
8     class Encoder(nn.Module):
9       def __init__(self, output_size=2):
10         super(Encoder, self).__init__()
11         # needs your implementation
12         self.nn = nn.Sequential(
13           nn.Linear(28 * 28, output_size),
14         )
15
16       def forward(self, x):
17         # needs your implementation
18         return self.nn(x)
19
20     class Decoder(nn.Module):
21       def __init__(self, input_size=2):
22         super(Decoder, self).__init__()
23         # needs your implementation
24         self.nn = nn.Sequential(
25           nn.Linear(input_size, 28 * 28),
26           nn.Tanh(),
27         )
28
29       def forward(self, z):
30         # needs your implementation
31         return self.nn(z)
32
33     self.dim_latent_representation = dim_latent_representation
34     self.encoder = Encoder(output_size=dim_latent_representation)
35     self.mu_layer = nn.Linear(self.dim_latent_representation, self.dim_latent_representation)
36     self.logvar_layer = nn.Linear(self.dim_latent_representation, self.dim_latent_representation)
37     self.decoder = Decoder(input_size=dim_latent_representation)
38     # Implement this function for the VAE model
39     def reparameterise(self, mu, logvar):
40
41       if self.training:
42         std = logvar.mul(0.5).exp_()

```



```
43     eps = std.data.new(std.size()).normal_()
44     return eps.mul(std).add_(mu)
45 else:
46     return mu
47
48 def forward(self, x):
49
50     # This function should be modified for the DAE and VAE
51     x = self.encoder(x)
52     mu, logvar = self.mu_layer(x), self.logvar_layer(x)
53     z = self.reparameterise(mu, logvar)
54     return self.decoder(z), mu, logvar
```

Listing 6: Part (d) vanilla autoencoder architecture.

- Loss = MSE + 0.5 * KLD
 - Epoch: 20
 - Training loss: 0.6420
 - Validation loss (reconstruction error) : 0.6060
- Loss = MSE + 1 * KLD
 - Epoch: 20
 - Training loss: 0.6821
 - Validation loss (reconstruction error) : 0.6550
- Loss = MSE + 5 * KLD
 - Epoch: 20
 - Training loss: 0.7122
 - Validation loss (reconstruction error) : 0.7154

We observe a low loss value for the vanilla autoencoder with 30D bottleneck size which results in high quality reconstructed images. Although loss values increased in VAEs, the VAE arranged the latent space such that gaps between latent representations for different classes decreased. It means we can get better manipulated (mixed latents) output. Since VAE follows an isotropic multivariate normal distribution at the latent space, we can generate new unseen images by taking samples from the latent space with higher quality compared to the Vanilla autoencoder. However, the reconstruction quality reduced (loss values increased) since loss function is a weighted combination of MSE and KLD terms to be optimized where the KLD term forces the latent space to resembles a Gaussian distribution. As we increased the KLD weight, we achieved a more compact latent space closer to the prior distribution by sacrificing the reconstruction quality.

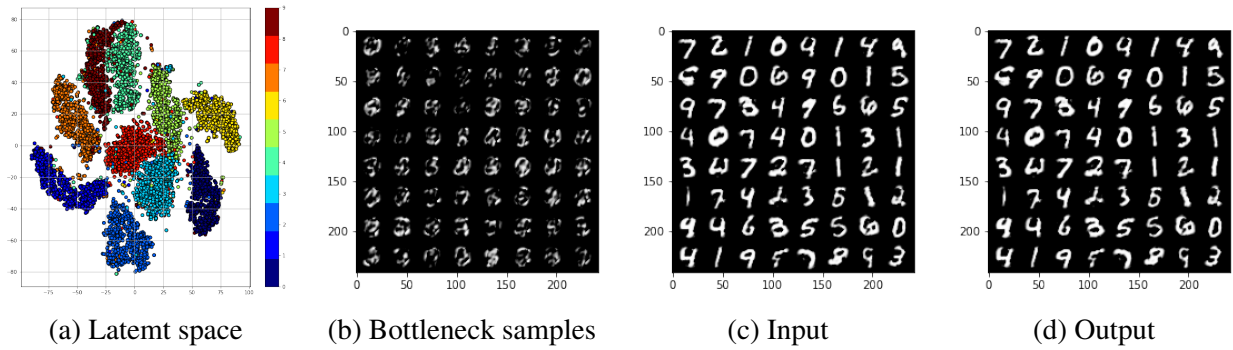


Figure 9: Results from part d Vanilla Autoencoder

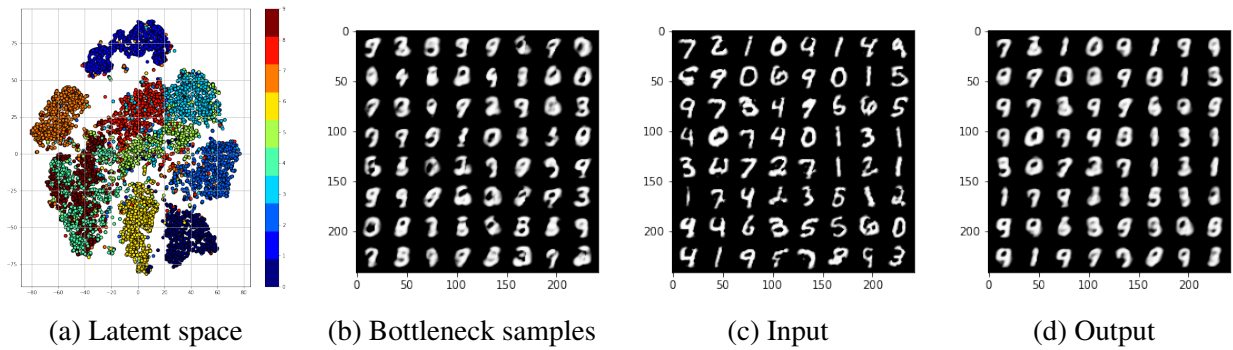


Figure 10: Results from part d Variational Autoencoder, KLD weight = 0.5

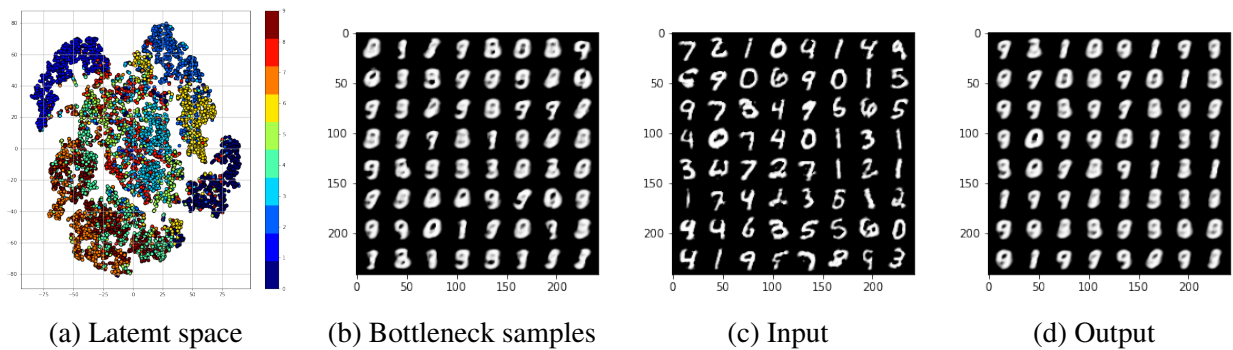


Figure 11: Results from part d Variational Autoencoder, KLD weight = 1

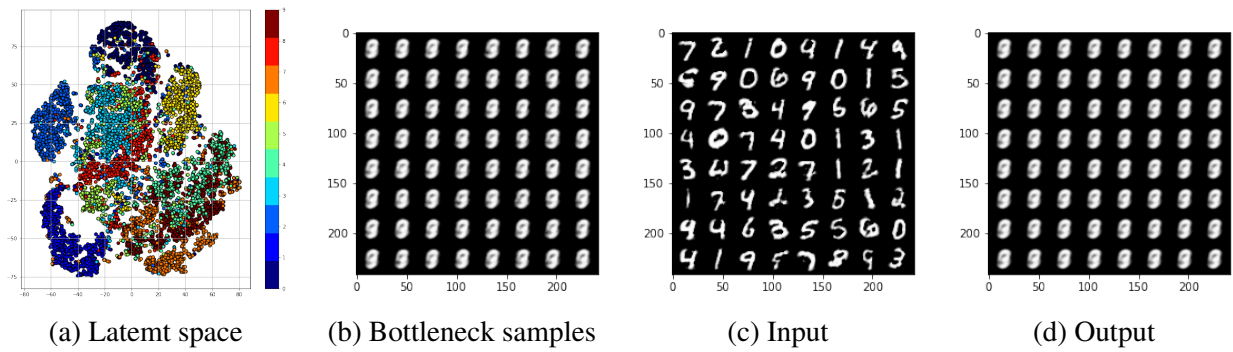


Figure 12: Results from part d Variational Autoencoder, KLD weight = 5

3 Support Vector Machines

a) We are trying to maximize $L(\vec{\lambda}) = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \phi(\vec{x}_j)^\top \phi(\vec{x}_i)$

Focusing on $\sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \phi(\vec{x}_j)^\top \phi(\vec{x}_i)$

If $i = j$, we have $9\lambda_i^2$

If $i \neq j$, we have $\lambda_i \lambda_j (y_i)(y_j)$

Therefore

$$\begin{aligned} \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \phi(\vec{x}_j)^\top \phi(\vec{x}_i) = \\ \lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 \\ - \frac{1}{2} (9\lambda_1^2 - 2\lambda_1\lambda_2 - 2\lambda_1\lambda_3 + 2\lambda_1\lambda_4 + 9\lambda_2^2 + 2\lambda_2\lambda_3 - 2\lambda_2\lambda_4 + 9\lambda_3^2 - 2\lambda_3\lambda_4 + 9\lambda_4^2) \end{aligned} \quad (5)$$

For each λ_i , we need to maximize the following objective:

$$L_1(\lambda_1) = \lambda_1 - \frac{1}{2} (9\lambda_1^2 - 2\lambda_1\lambda_2 - 2\lambda_1\lambda_3 + 2\lambda_1\lambda_4) \quad (6)$$

$$L_2(\lambda_2) = \lambda_2 - \frac{1}{2} (9\lambda_2^2 - 2\lambda_2\lambda_1 + 2\lambda_2\lambda_3 - 2\lambda_2\lambda_4) \quad (7)$$

$$L_3(\lambda_3) = \lambda_3 - \frac{1}{2} (9\lambda_3^2 - 2\lambda_3\lambda_1 + 2\lambda_3\lambda_2 - 2\lambda_3\lambda_4) \quad (8)$$

$$L_4(\lambda_4) = \lambda_4 - \frac{1}{2} (9\lambda_4^2 + 2\lambda_4\lambda_1 - 2\lambda_4\lambda_2 - 2\lambda_4\lambda_3) \quad (9)$$

Since each $L_i(\lambda_i)$ is concave, finding λ_i such that $\frac{\partial L_i}{\partial \lambda_i} = 0$ for each i will maximize the dual of SVM. Solving the following linear equation:

$$\begin{bmatrix} 9 & -1 & -1 & 1 \\ -1 & 9 & 1 & -1 \\ -1 & 1 & 9 & -1 \\ 1 & -1 & -1 & 9 \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \lambda_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

gives: $\lambda_1 = \lambda_2 = \lambda_3 = \lambda_4 = \frac{1}{8}$

Note that $\sum_{i=1}^4 \lambda_i y_i = 0$ so $\vec{\lambda}$ satisfies the constrained dual form.

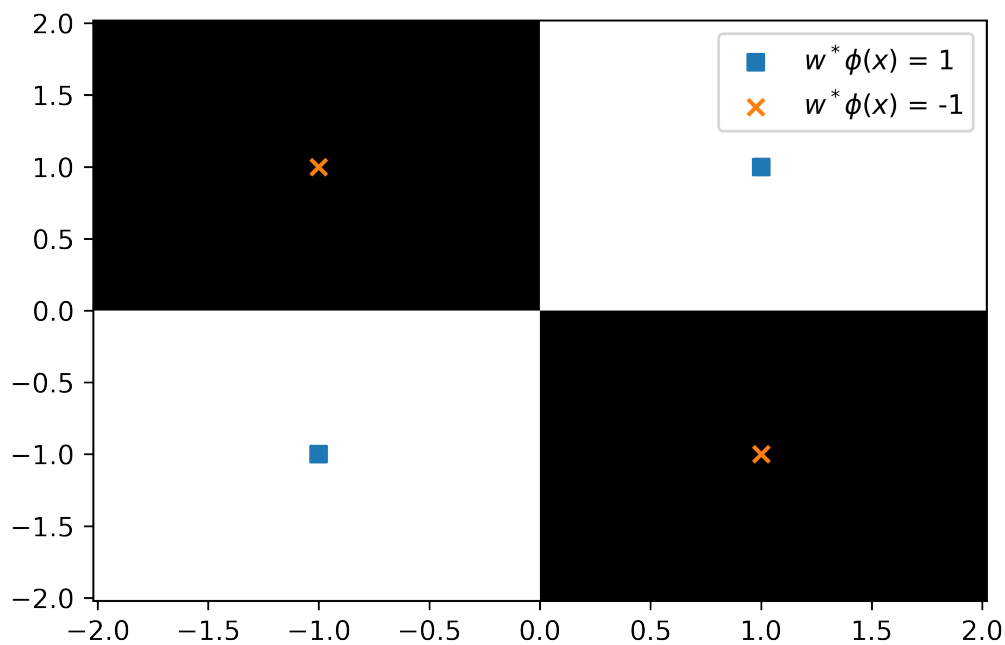
b)

$$\vec{w}^* = \sum_{i=1}^N \lambda_i^* y_i \phi(\vec{x}_i) = \frac{1}{8} \left[- \begin{bmatrix} 1 \\ -\sqrt{2} \\ -\sqrt{2} \\ 1 \\ \sqrt{2} \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ -\sqrt{2} \\ \sqrt{2} \\ 1 \\ -\sqrt{2} \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ \sqrt{2} \\ -\sqrt{2} \\ 1 \\ -\sqrt{2} \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ \sqrt{2} \\ \sqrt{2} \\ 1 \\ \sqrt{2} \\ 1 \end{bmatrix} \right] = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ -\frac{\sqrt{2}}{2} \\ 0 \end{bmatrix}$$

$$b^* = \frac{\max\{-1, -1\} + \min\{1, 1\}}{2} = 0$$

- c) The decision boundary are given at $w^{*\top} \phi(\vec{x}) - b^* = 0$. From b), this is when $-x_1 x_2 = 0$; equivalently, $x_1 = 0$ or $x_2 = 0$

The support vectors are when $w^{*\top} \phi(\vec{x}) - b^* = \pm 1$. From b), this is when $x_2 = -\frac{1}{x_1}$ or $x_2 = \frac{1}{x_1}$. All of the vectors are support vectors.



- d) Plot 1 corresponds to when $\gamma = 5$ and Plot 2 corresponds to when $\gamma = 0.05$. In Plot 1 less support vectors are used compared to Plot 2. Large γ penalizes when points are misclassified

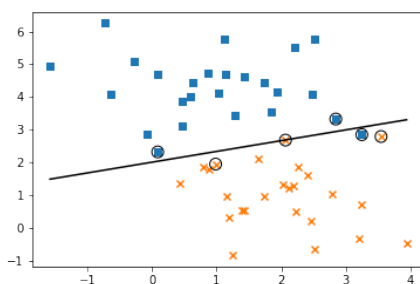


Figure 13: Plot 1

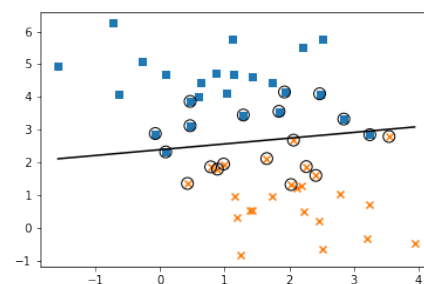


Figure 14: Plot 2

and makes the margin small since SVM would want the least number of misclassification mistakes. Small γ does the opposite so there are more support vectors in Plot 2 in comparison to Plot 1.