

Generative Models

Lecture 5

Variational Autoencoders

The *Variational Autoencoder* introduces randomness in both Encoder and Decoder.

The idea is that the data are samples from some conditional distribution

$$p(\mathbf{x}|\mathbf{z})$$

where \mathbf{z} is some latent or hidden variable

The Encoder generates the distribution of the latent variable \mathbf{z} dependent on the data points

$$q_\phi(\mathbf{z}|\mathbf{x})$$

The Decoder generates the reconstruction distribution

$$p_\theta(\mathbf{x}|\mathbf{z})$$

where

$$\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$$

is a sample from the latent variable distribution.

ϕ and θ denote the parameters of the Encoder and Decoder networks resp.

Here is how it works: we choose parametrized distributions q_θ and p_ϕ and we then optimize the parameters θ and ϕ such that for each \mathbf{x} in the data set, we sample \mathbf{z} from the distribution $q_\phi(\mathbf{z}|\mathbf{x})$ and then maximize the log-likelihood

$$\log p_\theta(\mathbf{x}|\mathbf{z})$$

We view $q_\phi(\mathbf{z}|\mathbf{x})$ as the variational approximation to the posterior $p_\theta(\mathbf{z}|\mathbf{x})$

Then we can compute the KL-divergence as before

$$\begin{aligned}
D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) &= -\mathbb{E}_{q_\phi} \left(\log \left(\frac{p_\theta(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right) \right) \\
&= - \int \log p_\theta(\mathbf{z}|\mathbf{x}) q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} + \int \log q_\phi(\mathbf{z}|\mathbf{x}) q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\
&= - \int \log \frac{p_\theta(\mathbf{x}|\mathbf{z}) p_{prior}(\mathbf{z})}{p_\theta(\mathbf{x})} q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} + \int \log q_\phi(\mathbf{z}|\mathbf{x}) q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\
&= -\mathbb{E}_{q_\phi} (\log p_\theta(\mathbf{x}|\mathbf{z})) - \mathbb{E}_{q_\phi} (p_{prior}(\mathbf{z})) + \log p_\theta(\mathbf{x}) + \mathbb{E}_{q_\theta} (\log q_\phi(\mathbf{z}|\mathbf{x}))
\end{aligned}$$

where $p_{prior}(\mathbf{z})$ is some suitable prior distribution of the latent variable \mathbf{z}

Rearranging terms we get

$$\begin{aligned}\log p_\theta(\mathbf{x}) &= \mathbb{E}_{q_\phi}(\log p_\theta(\mathbf{x}|\mathbf{z})) + \mathbb{E}_{q_\phi}(p_{prior}(\mathbf{z})) - \mathbb{E}_{q_\phi}(\log q_\phi(\mathbf{z}|\mathbf{x})) + D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_\phi}(p_\theta(\mathbf{x}|\mathbf{z})) + D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_{prior}(\mathbf{z})) + D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))\end{aligned}$$

Thus to maximize the log likelihood $\log p_\theta(\mathbf{x})$ we have to maximize to *ELBO*,

$$\mathbb{E}_{q_\phi}(p_\theta(\mathbf{x}|\mathbf{z})) + D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_{prior}(\mathbf{z}))$$

The first term is the reconstruction log-likelihood (we average over all latent samples).

The second term can be viewed as a regularization term that limits how far the approximation to the posterior can deviate from the prior (one often takes the posterior as $\mathcal{N}(0, I)$)

To do gradient ascent to maximize the *ELBO* we need to estimate the gradients

$$\begin{aligned}\nabla_{\phi,\theta} \text{ELBO} &= \nabla_{\phi,\theta} \mathbb{E}_{q_\phi} (\log p_\theta(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z})) - \log q_\phi(\mathbf{z}|\mathbf{x})) \\ &= \nabla_{\phi,\theta} \mathbb{E}_{q_\phi} (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}))\end{aligned}$$

There is no problem computing ∇_{θ} since we can differentiate with respect to θ under $\mathbb{E}_{q_{\phi}}$

$$\nabla_{\theta} \mathbb{E}_{q_{\phi}} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})) = \mathbb{E}_{q_{\phi}} (\nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z})))$$

In order to differentiate with respect to ϕ we resort to the 'reparametrization trick'.

This means that we cannot arbitrarily choose the distribution output by the Encoder.

In practice this means that we let the Encoder output vectors $\mu_\phi(\mathbf{x})$ and $\sigma_\phi(\mathbf{x})$ and let the latent variable

$$\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})I)$$

i.e. q_ϕ is a multi-variate Gaussian with diagonal covariance matrix.

Then we can change variables

$$\mathbf{z} = \mu_\phi(\mathbf{x}) + \sigma_\phi(\mathbf{x}) \underline{\varepsilon}$$

where $\underline{\varepsilon} \sim \mathcal{N}(0, I)$ so

$$\mathbb{E}_{q_\phi}(\log q_\phi(\mathbf{z}|\mathbf{x})) = \mathbb{E}_{\mathcal{N}(0,I)}(\log q_\phi(\mu_\phi(\mathbf{x}) + \sigma_\phi(\mathbf{x}) \underline{\varepsilon}))$$

So we can estimate gradients by Monte Carlo and do gradient ascent to update the parameters θ and ϕ

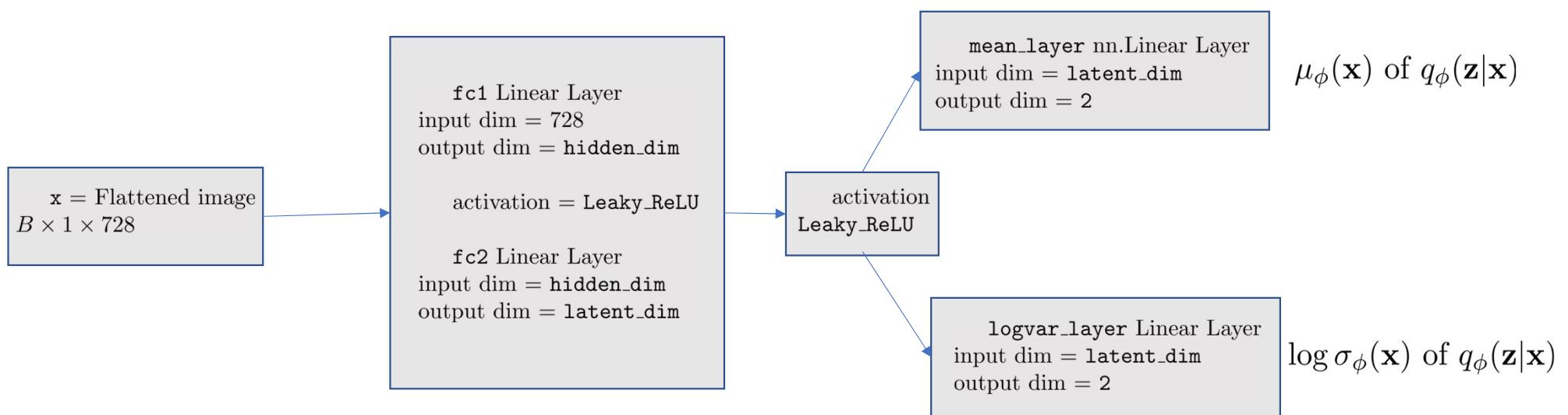
This is done pretty much automatically by the pyro class **SVI**

We shall go through two VAE notebooks, that only uses Pytorch

First the imports

```
1 import torch
2 import numpy as np
3 import torch.nn as nn
4 import torch.nn.functional as F
5 from torch.distributions.bernoulli import Bernoulli
6 from torch.optim import Adam
7 from torchvision.datasets import MNIST
8 from torch.utils.data import DataLoader
9 import torchvision.transforms as transforms
10 from mpl_toolkits.axes_grid1 import ImageGrid
11 from torchvision.utils import save_image, make_grid
12
```

Encoder Layer



```
1 class Encoder(nn.Module):
2
3     def __init__(self,input_dim,hidden_dim,latent_dim):
4         super().__init__()
5
6         self.fc1 = nn.Linear(input_dim,hidden_dim)
7         self.fc2 = nn.Linear(hidden_dim,latent_dim)
8
9         self.mean_layer = nn.Linear(latent_dim,2)
10        self.logvar_layer = nn.Linear(latent_dim,2)
11
12    def forward(self,x):
13
14        x = self.fc1(x)
15        x = F.leaky_relu(x,0.2)
16        x = self.fc2(x)
17        x = F.leaky_relu(x,0.2)
18
19        mean = self.mean_layer(x)
20        logvar = self.logvar_layer(x)
21
22        return mean, logvar
23
```

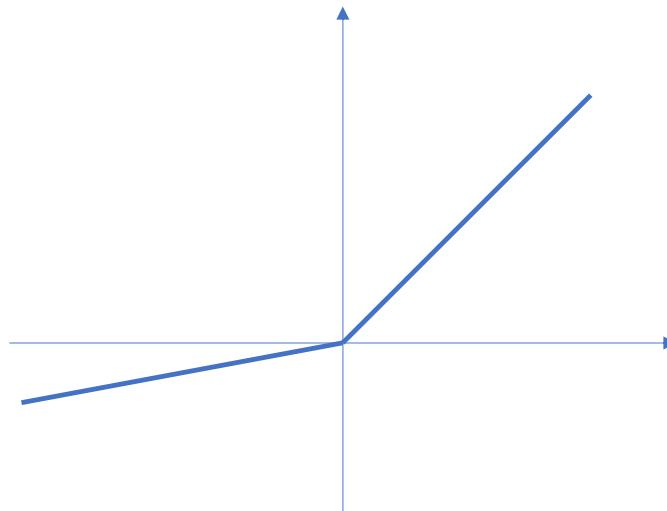
We sample from the latent distribution $\mathcal{N}(\mu_\phi(\mathbf{x}), \exp \log \sigma_\phi(\mathbf{x}))$ using the reparametrization trick i.e. we sample ε from $\mathcal{N}(0, I)$ and let

$$\mathbf{z} = \mu_\phi(\mathbf{x}) + \sqrt{\sigma_\phi(\mathbf{x})} \varepsilon$$

```
o
9  def reparameterize(self,mean,var):
0      epsilon = torch.randn_like(var)
1      z = mean + torch.sqrt(var) * epsilon
2      return z
3
```

The Decoder takes the sample \mathbf{z} and applies several linear layers with **Leaky_ReLU** activations (the leaky_ relu function with parameter $\alpha > 0$ is given by

$$\text{leaky_relu}(x, \alpha) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{if } x < 0 \end{cases}$$

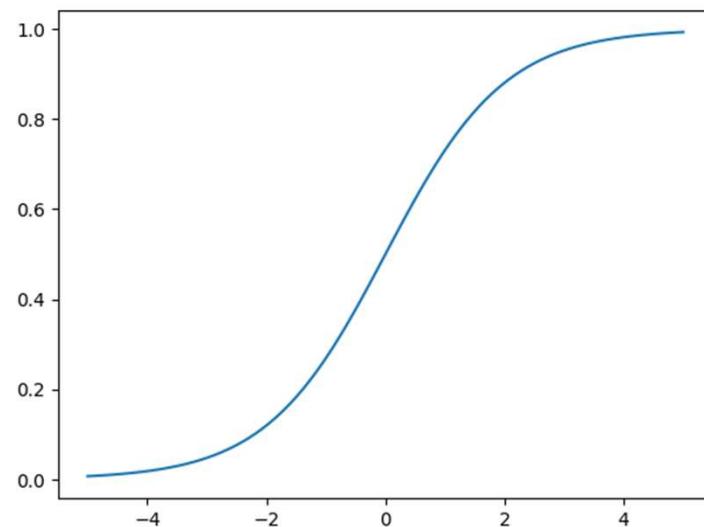


```
1 class Decoder(nn.Module):
2
3     def __init__(self,latent_dim,hidden_dim,input_dim):
4         super().__init__()
5
6         self.fc1 = nn.Linear(2,latent_dim)
7         self.fc2 = nn.Linear(latent_dim,hidden_dim)
8         self.fc3 = nn.Linear(hidden_dim,input_dim)    Last output has dimension 784
9
10    def forward(self,z):
11
12        z = self.fc1(z)
13        z = F.leaky_relu(z,0.2)
14        z = self.fc2(z)
15        z = F.leaky_relu(z,0.2)
16        z = self.fc3(z)
17
18        z = F.sigmoid(z).to(torch.float)
19
20        return z
21
```

The final linear layer outputs a 784 dimensional vector and we apply the sigmoid function σ

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

to each coordinate to get a 784 dimensional vector with coordinates between 0 and 1



The VAE itself is basically just combining the Encoder and Decoder.

The dataset consists of pairs `img, label`. For this model we are only interested in the `img` part, so we first take out the `img` part and then flatten the image from a $(1, 28, 28)$ tensor to a $1, 784$ tensor

```
class VAE(nn.Module):

    def __init__(self,input_dim=784,hidden_dim=400,latent_dim=200):
        super().__init__()

        self.encoder = Encoder(input_dim,hidden_dim,latent_dim)
        self.decoder = Decoder(latent_dim,hidden_dim,input_dim)

    def reparametrize(self,mean,var):
        epsilon = torch.randn_like(var)
        z = mean + torch.sqrt(var) * epsilon
        return z

    def forward(self,x):

        x = x[0].reshape(-1,1,28*28)

        mean, logvar = self.encoder(x)

        u = self.reparametrize(mean,torch.exp(logvar))

        z = self.decoder(u)

        return z, mean, logvar
```

To train the model we need a loss function. The object of training is to maximize the *ELBO* (or minimize $-\text{ELBO}$)

$$\mathbb{E}_{q_\phi}(\log p_\theta(\mathbf{x}|\mathbf{z})) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_{prior}(\mathbf{z}))$$

The distribution

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \text{Bernoulli}(\mathbf{z})(\mathbf{x})$$

where \mathbf{z} is the output of the decoder

We have

$$\log \text{Bernoulli}(z)(x) = x \log z + (1 - x) \log(1 - z)$$

This is the negative `torch.nn.BCELoss()`(\mathbf{z}, \mathbf{x}) function

```
CLASS torch.nn.BCELoss(weight=None, size_average=None, reduce=None, reduction='mean') [SOURCE]
```

$$\ell(x, y) = L = \{l_1, \dots, l_N\}^\top, \quad l_n = -w_n [y_n \cdot \log x_n + (1 - y_n) \cdot \log(1 - x_n)],$$

$$\log p_{\theta}(\mathbf{x}|\mathbf{z}) = -\text{torch.nn.BCELoss}(\text{reduction} = \text{'sum'})(\mathbf{z}, \mathbf{x})$$

To compute the KL-divergence, we have

$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|mean(\mathbf{x}), \exp logvar(\mathbf{x}))$$

and

$$p_{prior}(\mathbf{z}) = \mathcal{N}(\mathbf{z}|0, I)$$

so we are computing the KL-divergence between two Gaussians

Let $mean(\mathbf{x}) = (\mu_1, \mu_2)$ and $var(\mathbf{x}) = (\sigma_1, \sigma_2)$

$$\begin{aligned} D_{KL}(\mathcal{N}\left((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}\right) || \mathcal{N}(0, I)) \\ = \int \left(\log \mathcal{N}(\mathbf{z}|0, I) - \log \mathcal{N}(\mathbf{z}| \left((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}\right) \right) \mathcal{N}(\mathbf{z}| \left((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}\right) d\mathbf{z} \end{aligned}$$

We have

$$\mathcal{N}(\mathbf{z}|0, I) = \frac{1}{2\pi} \exp(-\frac{1}{2}\mathbf{z}\mathbf{z}^T)$$

and

$$\mathcal{N}(\mathbf{z}|(\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}) = \frac{1}{2\pi \sqrt{\sigma_1 \sigma_2}} \exp\left(-\frac{1}{2}(\mathbf{z} - (\mu_1, \mu_2))^T \begin{pmatrix} \sigma_1^{-1} & 0 \\ 0 & \sigma_2^{-1} \end{pmatrix} (\mathbf{z} - (\mu_1, \mu_2))\right)$$

So

$$\begin{aligned} & \log \mathcal{N}(\mathbf{z}|(\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}) - \log \mathcal{N}(\mathbf{z}|0, I) \\ &= \frac{1}{2} \log \sigma_1 \sigma_2 - \frac{1}{2}(\mathbf{z} - (\mu_1, \mu_2))^T \begin{pmatrix} \sigma_1^{-1} & 0 \\ 0 & \sigma_2^{-1} \end{pmatrix} (\mathbf{z} - (\mu_1, \mu_2))^T + \frac{1}{2}\mathbf{z}\mathbf{z}^T \end{aligned}$$

This expression equals

$$-\frac{1}{2} \log \sigma_1 \sigma_1 - \frac{1}{2} \left(\frac{(z_1 - \mu_1)^2}{\sigma_1} + \frac{(z_2 - \mu_2)^2}{\sigma_2} \right) + \frac{1}{2}(z_1^2 + z_2^2)$$

$$\int \left(-\frac{1}{2} \log \sigma_1 \sigma_2 \mathcal{N}(\mathbf{z} | ((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}) \right) d\mathbf{z} = -\frac{1}{2} \log \sigma_1 \sigma_2$$

$$\int -\frac{1}{2} \left(\frac{(z_1 - \mu_1)^2}{\sigma_1} + \frac{(z_2 - \mu_2)^2}{\sigma_2} \right) \mathcal{N}(\mathbf{z} | ((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix})) d\mathbf{z} = -1$$

Finally

$$\int (z_1 - \mu_1)^2 \mathcal{N}(\mathbf{z} | ((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix})) d\mathbf{z} = \sigma_1$$

and

$$(z_1 - \mu_1)^2 = z_1^2 + \mu_1^2 - 2z_1\mu_1$$

so

$$\begin{aligned} & \int z_1^2 \mathcal{N}(\mathbf{z} | ((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix})) d\mathbf{z} \\ &= \sigma_1^2 - \mu_1^2 + 2\mu_1 \int z_1 \mathcal{N}(\mathbf{z} | ((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix})) d\mathbf{z} \\ &= \sigma_1^2 - \mu_1^2 + 2\mu_1^2 = \sigma_1^2 + \mu_1^2 \end{aligned}$$

Thus we finally get

$$D_{KL}(\mathcal{N}\left((\mu_1, \mu_2), \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}\right) || \mathcal{N}(0, I)) = \frac{1}{2}(\mu_1^2 + \mu_2^2 + \sigma_1 + \sigma_2 - 2 - \log \sigma_1 \sigma_2)$$

If `z,mean,logvar = vae(x)` we can code this as

```
KL_div = - 0.5 * torch.sum(1+ logvar - mean**2 - torch.exp(logvar))
```

The total loss function i.e. the $-ELBO$ is then

```
1 def loss_function(z,x,mean,logvar):
2     reconstruction_loss = F.binary_cross_entropy(x,z,reduction='sum')
3     KL_div = - 0.5 * torch.sum(1+ logvar - mean**2 - torch.exp(logvar))
4
5     return reconstruction_loss + KL_div
```

We can now run the training loop

We use the Adam optimizer

```
1 optimizer = Adam(vae.parameters(), lr=1e-3)
```

```
22
1 epochs = 50
2 for e in range(epochs):
3     overall_loss = 0.0
4     for i,img in enumerate(train_loader):
5
6         x = img[0].reshape(-1,28*28)
7
8         z,mean,logvar = vae(x)
9
10        optimizer.zero_grad()
11
12        loss = loss_function(x,z,mean,logvar)
13
14        overall_loss += loss.item()
15
16        loss.backward()
17        optimizer.step()
18
19        print("\tEpoch", e + 1, "\tAverage Loss: ", overall_loss/(i*batch_size))
20
21
22
```

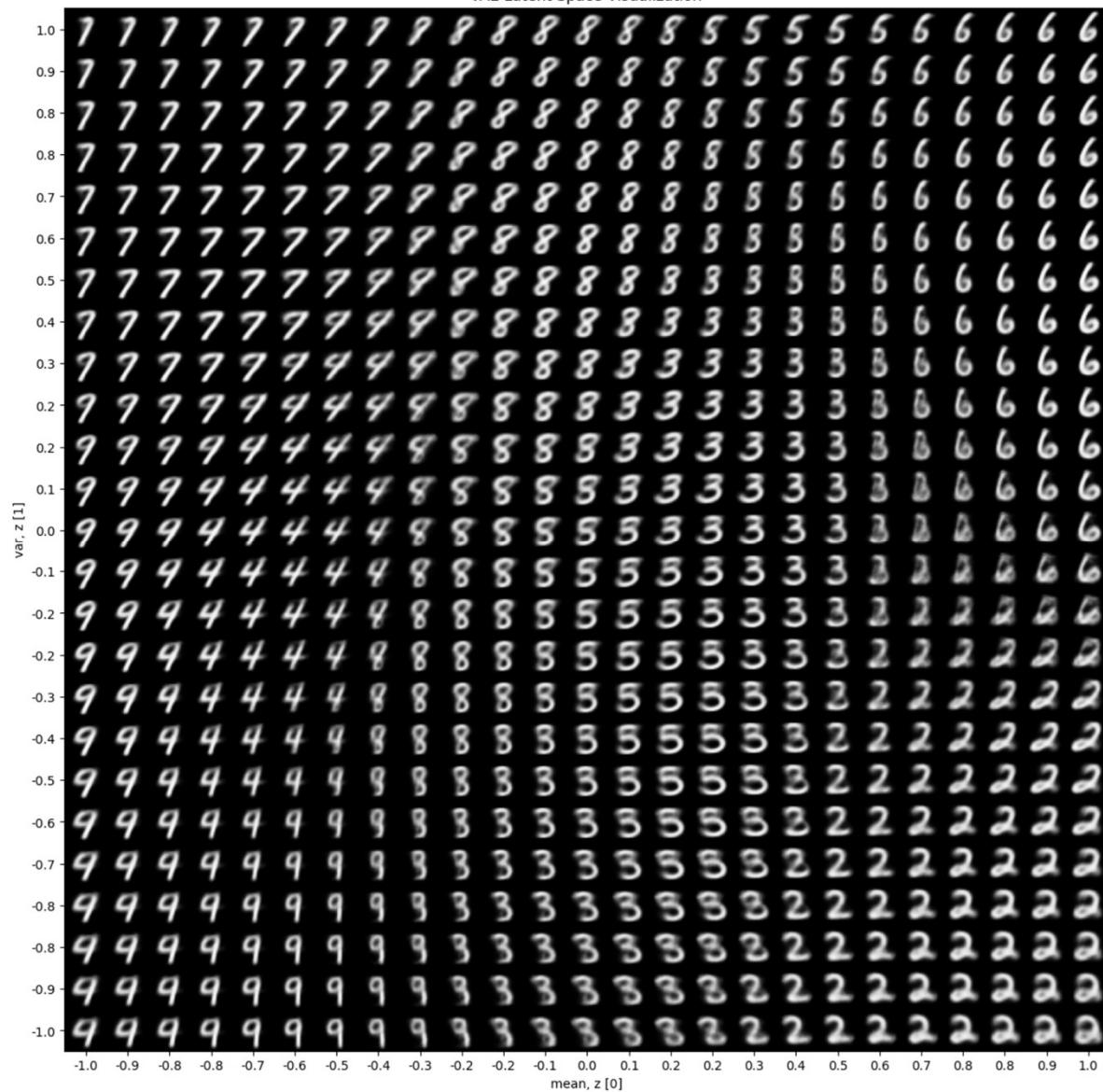
```
Epoch 1      Average Loss: 169.05433931226523
Epoch 2      Average Loss: 160.5248147139764
Epoch 3      Average Loss: 156.35052939795492
Epoch 4      Average Loss: 153.65856944190838
Epoch 5      Average Loss: 151.8359131103923
Epoch 6      Average Loss: 150.75318972375834
Epoch 7      Average Loss: 149.655738600793
Epoch 8      Average Loss: 149.00604903028486
Epoch 9      Average Loss: 148.3741138225167
Epoch 10     Average Loss: 147.7428735718385
Epoch 11     Average Loss: 147.18598841493636
Epoch 12     Average Loss: 146.71981404554467
```

This function moves around in (mean,logvar) latent space and samples a new \mathbf{z} for every (mean,logvar) and then shows the sampled image.

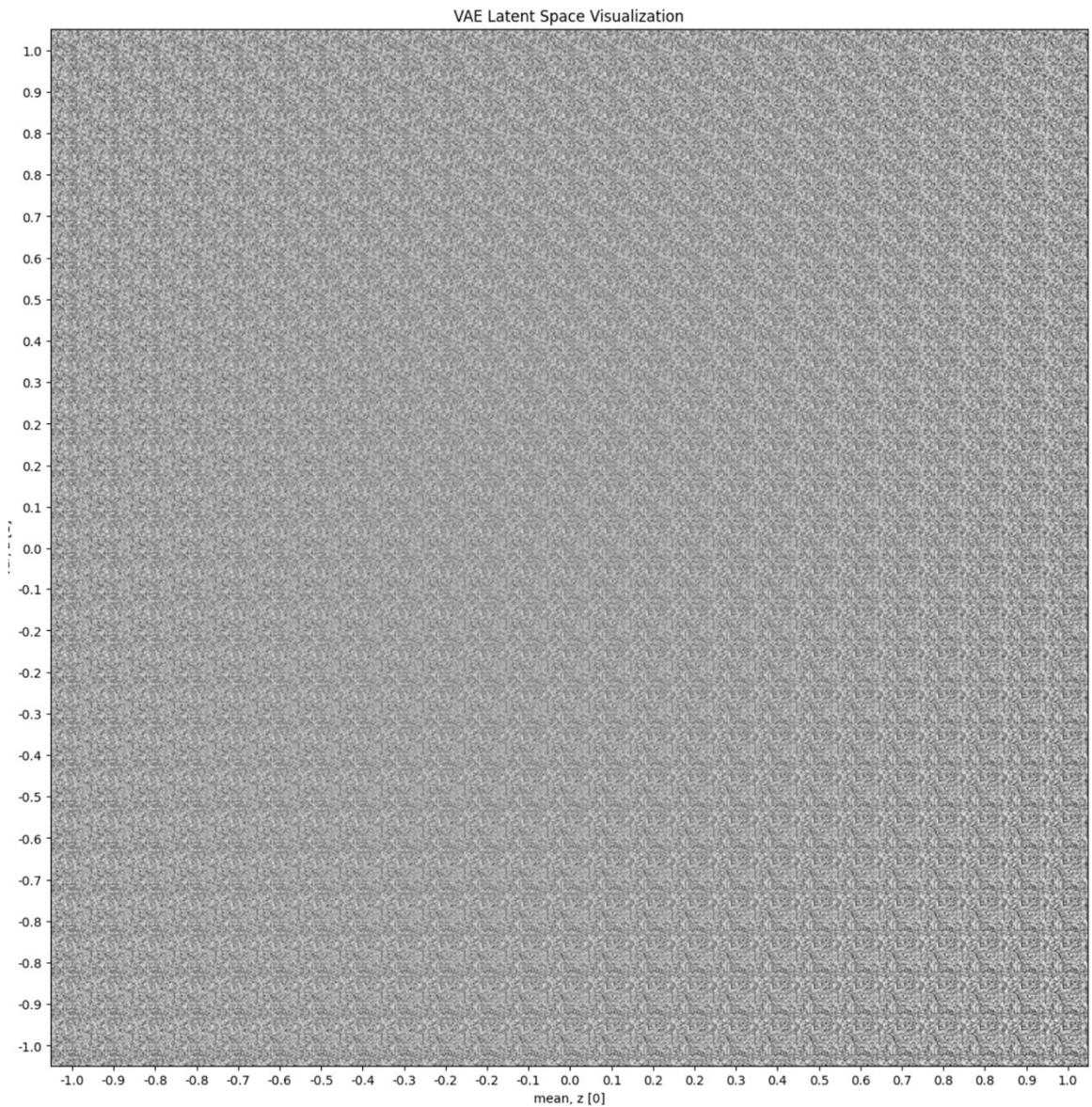
Remark how the images slowly change between different digits

```
1 def plot_latent_space(model, scale=5.0, n=25, digit_size=28, figsize=15)
2     # display a n*n 2D manifold of digits
3     figure = np.zeros((digit_size * n, digit_size * n))
4
5     # construct a grid
6     grid_x = np.linspace(-scale, scale, n)
7     grid_y = np.linspace(-scale, scale, n)[::-1]
8
9     for i, yi in enumerate(grid_y):
10         for j, xi in enumerate(grid_x):
11             z_sample = torch.tensor([[xi, yi]], dtype=torch.float)
12             x_decoded = vae.decoder(z_sample)
13             digit = x_decoded[0].detach().cpu().reshape(digit_size, digit_size)
14             figure[i * digit_size : (i + 1) * digit_size, j * digit_size] = digit
15
16     plt.figure(figsize=(figsize, figsize))
17     plt.title('VAE Latent Space Visualization')
18     start_range = digit_size // 2
19     end_range = n * digit_size + start_range
20     pixel_range = np.arange(start_range, end_range, digit_size)
21     sample_range_x = np.round(grid_x, 1)
22     sample_range_y = np.round(grid_y, 1)
23     plt.xticks(pixel_range, sample_range_x)
24     plt.yticks(pixel_range, sample_range_y)
25     plt.xlabel("mean, z [0]")
26     plt.ylabel("var, z [1]")
27     plt.imshow(figure, cmap="Greys_r")
28     plt.show()
```

VAE Latent Space Visualization



Same function with an
untrained model



Next we are going through another notebook, implementing a VAE using pyro

As usual, first the imports

```
1 import os
2
3 import numpy as np
4 import torch
5 from pyro.contrib.examples.util import MNIST
6 import torch.nn as nn
7 import torchvision.transforms as transforms
8
9 import pyro
10 import pyro.distributions as dist
11 import pyro.contrib.examples.util # patches torchvision
12 from pyro.infer import SVI, Trace_ELBO
13 from pyro.optim import Adam
```

The data import and setting up the data loaders is the same as the previous notebook

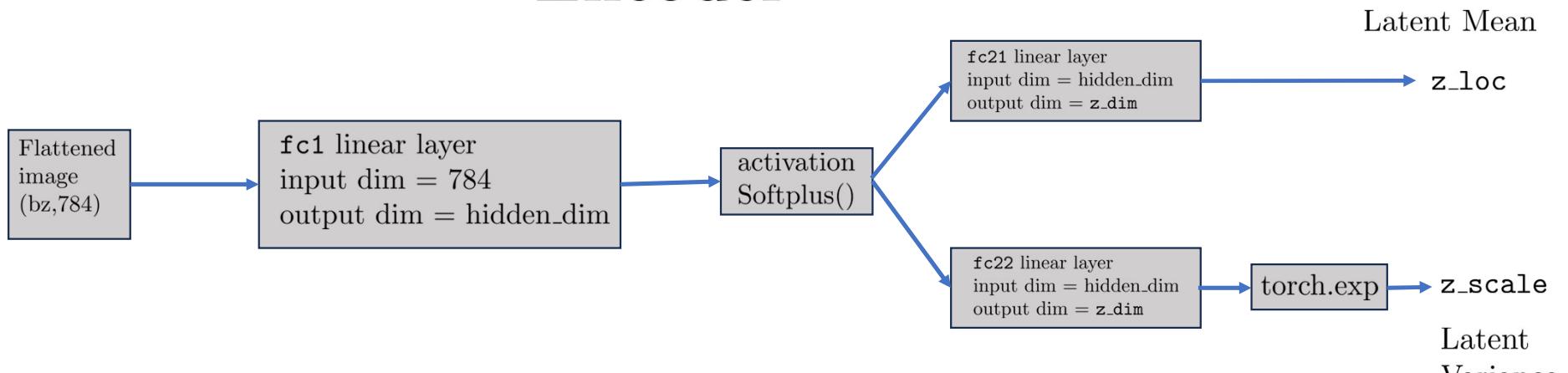
```
1 transform = transforms.Compose([transforms.ToTensor()])
2
3 path = './data'
4
5 train_dataset = MNIST(path,transform=transform,download=True)
6 test_dataset = MNIST(path,transform=transform,download=True)
7
8 batch_size=100
9 train_loader = DataLoader(train_dataset,batch_size=batch_size,shuffle=True)
10 test_loader = DataLoader(test_dataset,batch_size=batch_size,shuffle=False)
```

We write an Encoder and a Decoder

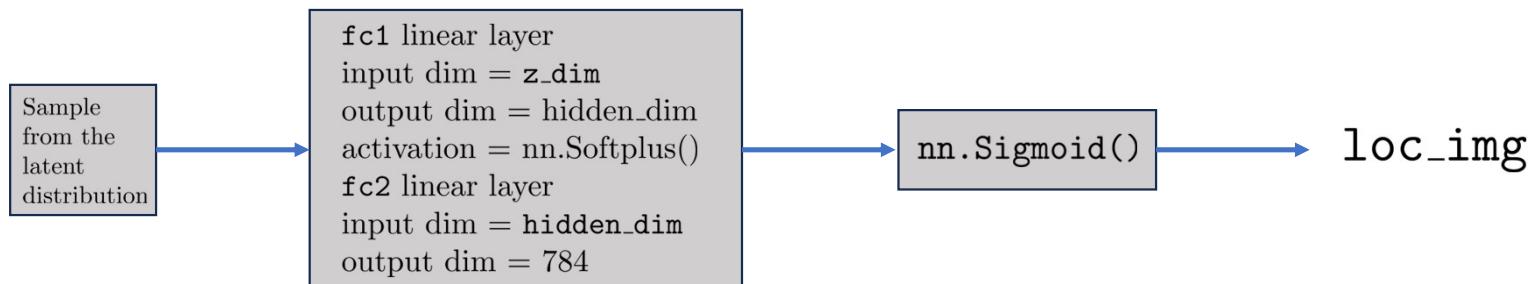
```
1 class Encoder(nn.Module):
2     def __init__(self, z_dim, hidden_dim):
3         super().__init__()
4         # setup the three linear transformations used
5         self.fc1 = nn.Linear(784, hidden_dim)
6         self.fc21 = nn.Linear(hidden_dim, z_dim)
7         self.fc22 = nn.Linear(hidden_dim, z_dim)
8         # setup the non-linearities
9         self.softplus = nn.Softplus()
10
11    def forward(self, x):
12        # define the forward computation on the image x
13        # first shape the mini-batch to have pixels in the rightmost dim
14        x = x.reshape(-1, 784)
15        # then compute the hidden units
16        hidden = self.softplus(self.fc1(x))
17        # then return a mean vector and a (positive) square root covariance
18        # each of size batch_size x z_dim
19        z_loc = self.fc21(hidden)
20        z_scale = torch.exp(self.fc22(hidden))
21        return z_loc, z_scale
```

```
1 class Decoder(nn.Module):
2     def __init__(self, z_dim, hidden_dim):
3         super().__init__()
4         # setup the two linear transformations used
5         self.fc1 = nn.Linear(z_dim, hidden_dim)
6         self.fc2 = nn.Linear(hidden_dim, 784)
7         # setup the non-linearities
8         self.softplus = nn.Softplus()
9         self.sigmoid = nn.Sigmoid()
10
11    def forward(self, z):
12        # define the forward computation on the latent z
13        # first compute the hidden units
14        hidden = self.softplus(self.fc1(z))
15        # return the parameter for the output Bernoulli
16        # each is of size batch_size x 784
17        loc_img = self.sigmoid(self.fc2(hidden))
18        return loc_img
```

Encoder

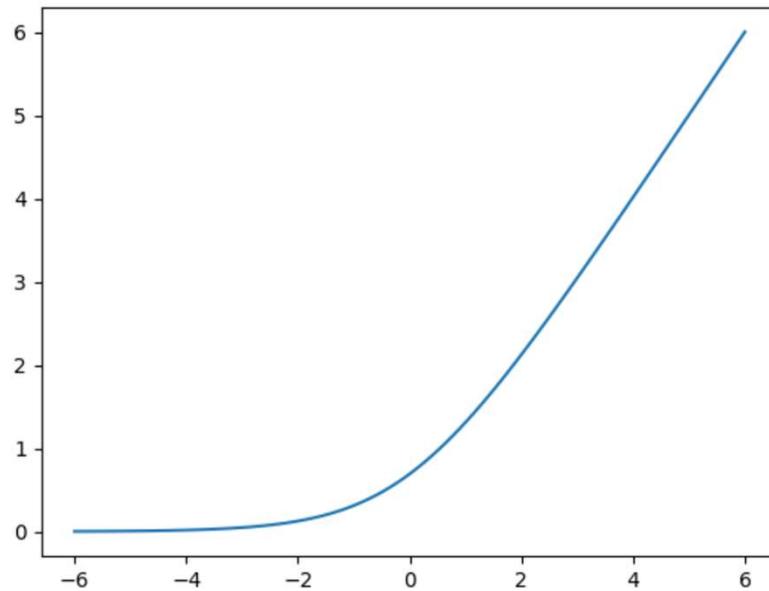


Decoder



Softplus activation function

$$\text{softplus}(x) = \log(1 + \exp(x))$$



The Encoder and Decoder are basically the same as in the Pytorch VAE.

The differences come in the VAE class

Using the `pyro` `SVI` class we need a `model` which implements

$$p_{\theta}(\mathbf{x}|\mathbf{z})p_{prior}(\mathbf{z})$$

The prior is $\mathcal{N}(0, I)$ where I is the `z_dim` \times `z_dim` identity matrix

```

1  class VAE(nn.Module):
2      # by default our latent space is 50-dimensional
3      # and we use 400 hidden units
4      def __init__(self, z_dim=50, hidden_dim=400, use_cuda=False):
5          super().__init__()
6          # create the encoder and decoder networks
7          self.encoder = Encoder(z_dim, hidden_dim)
8          self.decoder = Decoder(z_dim, hidden_dim)
9
10         if use_cuda:

```

We need to register the decoder with pyro

The `z_loc` and `z_scale` will be updated through the svi updates

Sampling from the latent distribution and applying the decoder

The likelihood $p(\mathbf{x}|\mathbf{z})$ is the Bernoulli distribution with probabilities `loc_img`

```

16
17      # define the model  $p(\mathbf{x}/\mathbf{z})p(\mathbf{z})$ 
18      def model(self, x):
19          # register PyTorch module `decoder` with Pyro
20          pyro.module("decoder", self.decoder)
21          with pyro.plate("data", x.shape[0]):
22              # setup hyperparameters for prior  $p(\mathbf{z})$ 
23              z_loc = x.new_zeros(torch.Size((x.shape[0], self.z_dim)))
24              z_scale = x.new_ones(torch.Size((x.shape[0], self.z_dim)))
25              # sample from prior (value will be sampled by guide when computing the ELBO)
26              z = pyro.sample("latent", dist.Normal(z_loc, z_scale).to_event(1))
27              # decode the latent code z
28              loc_img = self.decoder(z)
29              # score against actual images
30              pyro.sample("obs", dist.Bernoulli(loc_img, validate_args=False).to_event(1), obs=x.reshape(-1, 784))
31

```

The `guide` is the Variational approximation to the posterior $p(\mathbf{z}|\mathbf{x})$.

In the VAE this is the latent distribution $q_\phi(\mathbf{z}|\mathbf{x})$ with the parameters mean = `z_loc` and variance = `z_scale`

The `encoder` has to be registered with pyro

The `encoder` outputs `z_loc` and `z_scale`, the mean and variance of the latent distribution

Sampling from the latent distribution

```
32  # define the guide (i.e. variational distribution) q(z/x)
33  def guide(self, x):
34      # register PyTorch module `encoder` with Pyro
35      pyro.module("encoder", self.encoder)
36      with pyro.plate("data", x.shape[0]):
37          # use the encoder to get the parameters used to define q(z/x)
38          z_loc, z_scale = self.encoder(x)
39          # sample the Latent code z
40          pyro.sample("latent", dist.Normal(z_loc, z_scale).to_event(1))
41
```

Both the `model` and the `guide` are methods in the VAE class.

We can now write the training loop

We first instantiate the VAE, set the optimizer and instantiate the SVI

Instantiate the
VAE

```
1 vae = VAE()
```

Setting the optimizer

```
1 optimizer = Adam({"lr": 1.0e-3})
```

Instantiating the SVI

```
1 svi = SVI(vae.model, vae.guide, optimizer, loss=Trace_ELBO())
```

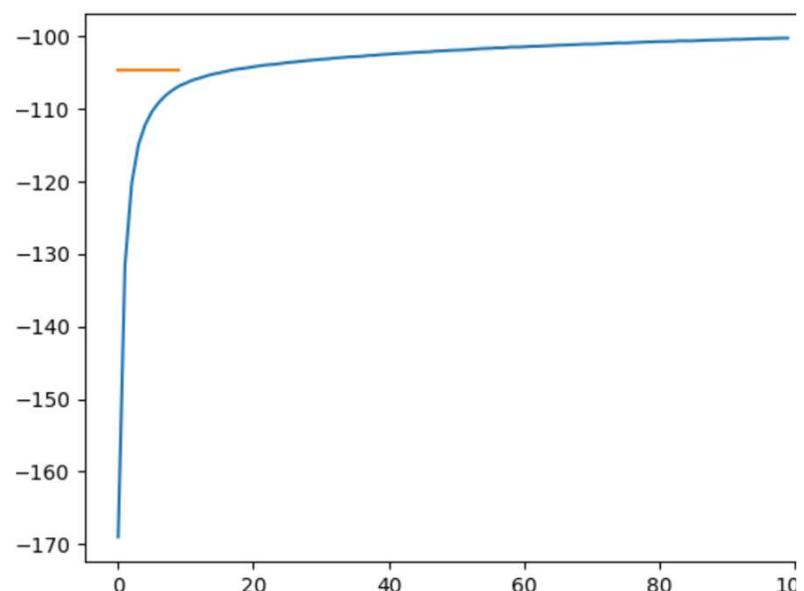
The training loop

The `svi` computes the loss (i.e. the *ELBO*), computes the gradients of the *ELBO* and does the gradient step to update the parameters i.e. the parameters in the `encoder` and the `decoder`. This is why they have to be registered with `pyro`

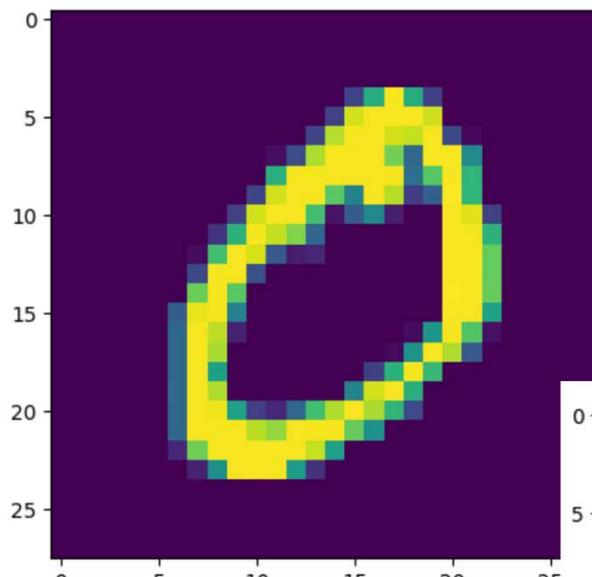
Here we evaluate the loss on the test set using the `svi.evaluate_loss` command

```
1 epochs = 100
2
3 len_train_set = len(train_loader.dataset)
4 len_test_set = len(test_loader.dataset)
5
6 train_elbo = []
7 test_elbo = []
8 # training Loop
9 for epoch in range(epochs):
10
11     epoch_loss = 0.
12
13     # do a training epoch over each mini-batch x returned
14     # by the data loader
15     for x, _ in train_loader:
16
17         pyro.clear_param_store()
18
19         loss = svi.step(x)
20
21         epoch_loss += loss
22
23     total_epoch_loss_train = epoch_loss / len_train_set
24
25     train_elbo.append(-total_epoch_loss_train)
26
27     print("[epoch %03d] average training loss: %.4f" % (epoch, total_epoch_loss_train))
28
29     if epoch % 10 == 0:
30
31         test_loss = 0.
32         # compute the loss over the entire test set
33         for x, _ in test_loader:
34             test_loss += svi.evaluate_loss(x)
35
36
37         average_test_loss = test_loss / len_test_set
38         test_elbo.append(-total_epoch_loss_test)
39         print("[epoch %03d] average test loss: %.4f" % (epoch, total_epoch_loss_test))
```

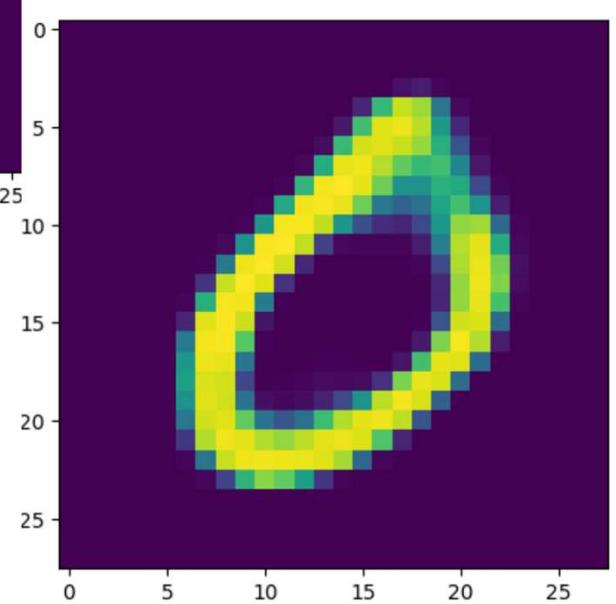
ELBO



Original



Reconstruction



Using a different dataset

```
1 transform = transforms.Compose(  
2     [transforms.Pad(padding=2), transforms.ToTensor()])  
3 )  
4  
5 path = './data'  
6  
7 train_dataset = FashionMNIST(  
8     path, train=True, download=True, transform=transform)  
9 test_dataset = FashionMNIST(  
10    path, train=False, download=True, transform=transform)  
11  
12 batch_size=100  
13 train_loader = DataLoader(train_dataset,batch_size=batch_size,shuffle=True)  
14 test_loader = DataLoader(test_dataset,batch_size=batch_size,shuffle=False)  
15
```

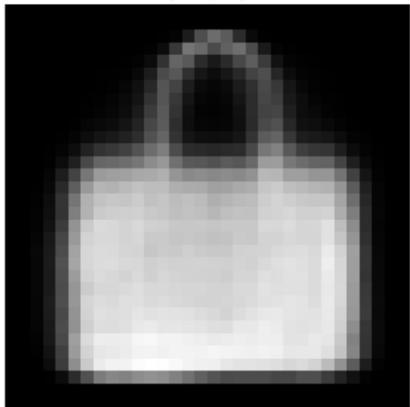
These images are 32×32 pixels

```
1 num_samples = 25  
2 sample_images = [image[0][i,0] for i in range(num_samples)]  
3  
4 fig = plt.figure(figsize=(5,5))  
5 grid = ImageGrid(fig,111,nrows_ncols=(5,5), axes_pad=0.1)  
6  
7 for ax, im in zip(grid,sample_images):  
8     ax.imshow(im,cmap='gray')  
9     ax.axis('off')
```

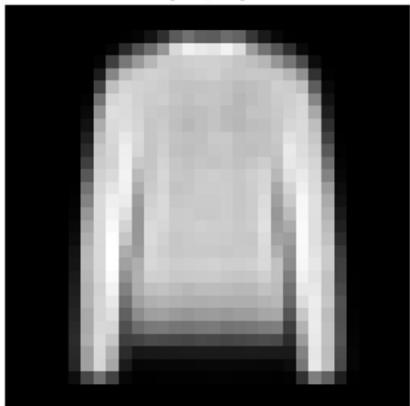


Generated Images

[0.0,1.0]



[1.0,1.0]



VAE Latent Space Visualization

