



Published in DataSeries

Open in app ↗

Sign up

Sign In



Search Medium



Jun 28, 2021 · 6 min read · ✨ · 🎧 Listen

Save



Convolutional Autoencoder in Pytorch on MNIST dataset

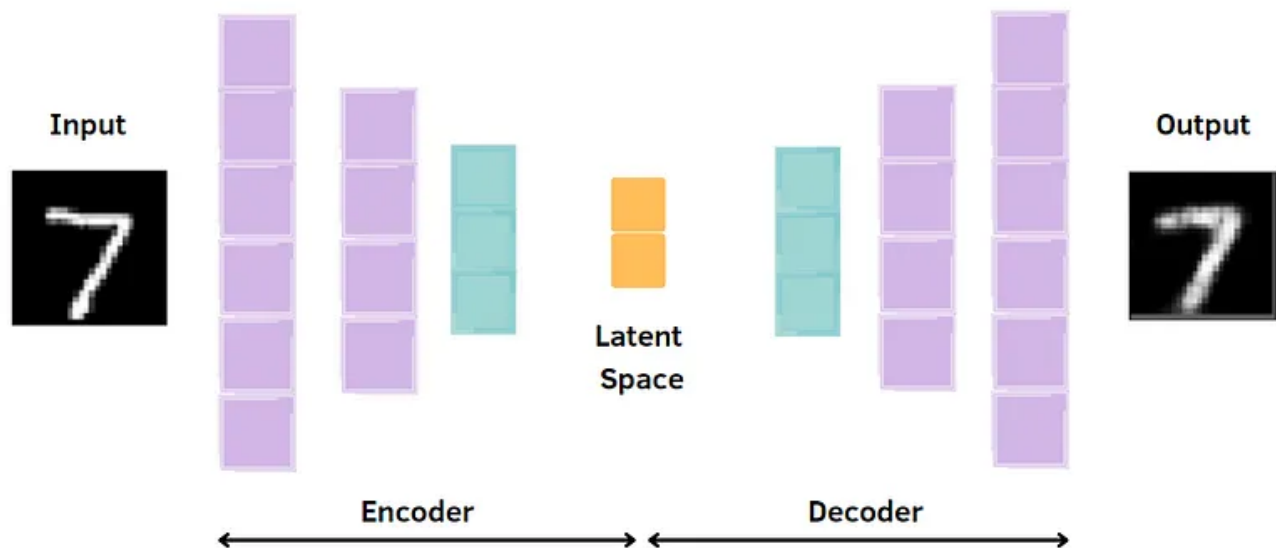


Illustration by Author

The post is the seventh in a series of guides to build deep learning models with Pytorch. Below, there is the full series:

1. [Pytorch Tutorial for Beginners](#)
2. [Manipulating Pytorch Datasets](#)
3. [Understand Tensor Dimensions in DL models](#)
4. [CNN & Feature visualizations](#)
5. [Hyperparameter tuning with Optuna](#)
6. [K Fold Cross Validation](#)
7. [Convolutional Autoencoder \(this post\)](#)
8. [Denoising Autoencoder](#)
9. [Variational Autoencoder](#)



196



6

The goal of the series is to make Pytorch more intuitive and accessible as possible through examples of implementations. There are many tutorials on the Internet to use Pytorch to build many types of challenging models, but it can also be confusing at the same time because there are always slight differences when you pass from a tutorial to another. In this series, I want to start from the simplest topics to the more advanced ones.

Autoencoder

The autoencoder is an unsupervised deep learning algorithm that learns encoded representations of the input data and then reconstructs the same input as output. It consists of two networks, Encoder and Decoder. The Encoder compresses the high-dimensional input into a low-dimensional latent code, called also latent code or encoded space, to extract the most relevant information from it, while the Decoder decompresses the encoded data and recreates the original input.

The goal of this architecture is to maximize the information when encoding and minimize the reconstruction error. But what is the reconstruction error? Its name is also reconstruction loss and is usually the mean-squared error between the

reconstructed input and the original input when the input is real-valued. In case the input data is categorical, the loss function used is the Cross-Entropy Loss.

Implementation in Pytorch

The following steps will be shown:

1. Import libraries and MNIST dataset
2. Define Convolutional Autoencoder
3. Initialize Loss function and Optimizer
4. Train model and evaluate model
5. Generate new samples from the latent code
6. Visualize the latent space with t-SNE

1. Import libraries and MNIST dataset

We can import the dataset using the library torchvision. We download the training and the test datasets and we transform the image datasets into Tensor. We don't need to normalize the images because the datasets contain colored images. After we divide the training dataset into training and validation sets. The `random_split` provides a random partition for these two sets. The `DataLoader` is used to create data loaders for the training, validation, and test sets, which are split into mini-batches. The `batchsize` is the number of samples used in one iteration during the training of the model.

```
1  import matplotlib.pyplot as plt # plotting library
2  import numpy as np # this module is useful to work with numerical arrays
3  import pandas as pd
4  import random
5  import torch
6  import torchvision
7  from torchvision import transforms
8  from torch.utils.data import DataLoader, random_split
9  from torch import nn
10 import torch.nn.functional as F
11 import torch.optim as optim
12
13 data_dir = 'dataset'
14
15 train_dataset = torchvision.datasets.MNIST(data_dir, train=True, download=True)
16 test_dataset = torchvision.datasets.MNIST(data_dir, train=False, download=True)
17
18 train_transform = transforms.Compose([
19     transforms.ToTensor(),
20 ])
21
22 test_transform = transforms.Compose([
23     transforms.ToTensor(),
24 ])
25
26 train_dataset.transform = train_transform
27 test_dataset.transform = test_transform
28
29 m=len(train_dataset)
30
31 train_data, val_data = random_split(train_dataset, [int(m-m*0.2), int(m*0.2)])
32 batch_size=256
33
34 train_loader = torch.utils.data.DataLoader(train_data, batch_size=batch_size)
35 valid_loader = torch.utils.data.DataLoader(val_data, batch_size=batch_size)
36 test_loader = torch.utils.data.DataLoader(test_dataset, batch_size=batch_size, shuffle=True)
```

.py hosted with ❤ by GitHub

[view raw](#)

2. Define Convolutional Autoencoder

Here, we define the Autoencoder with Convolutional layers. It will be composed of two classes: one for the encoder and one for the decoder. The encoder will contain three

convolutional layers and two fully connected layers. Some batch norm layers are added as regularizers. The decoder will have the same architecture but in inverse order.

```

1  class Encoder(nn.Module):
2
3      def __init__(self, encoded_space_dim,fc2_input_dim):
4          super().__init__()
5
6          ### Convolutional section
7          self.encoder_cnn = nn.Sequential(
8              nn.Conv2d(1, 8, 3, stride=2, padding=1),
9              nn.ReLU(True),
10             nn.Conv2d(8, 16, 3, stride=2, padding=1),
11             nn.BatchNorm2d(16),
12             nn.ReLU(True),
13             nn.Conv2d(16, 32, 3, stride=2, padding=0),
14             nn.ReLU(True)
15         )
16
17         ### Flatten layer
18         self.flatten = nn.Flatten(start_dim=1)
19     ### Linear section
20     self.encoder_lin = nn.Sequential(
21         nn.Linear(3 * 3 * 32, 128),
22         nn.ReLU(True),
23         nn.Linear(128, encoded_space_dim)
24     )
25
26     def forward(self, x):
27         x = self.encoder_cnn(x)
28         x = self.flatten(x)
29         x = self.encoder_lin(x)
30         return x
31
32     class Decoder(nn.Module):
33
34         def __init__(self, encoded_space_dim,fc2_input_dim):
35             super().__init__()
36             self.decoder_lin = nn.Sequential(
37                 nn.Linear(encoded_space_dim, 128),
38                 nn.ReLU(True),
39                 nn.Linear(128, 3 * 3 * 32),
40                 nn.ReLU(True)
41             )
42
43             self.unflatten = nn.Unflatten(dim=1,
44                 unflattened_size=(32, 3, 3))
45
46             self.decoder_cnn = nn.Sequential(

```

```

45         self.decoder_conv = nn.Sequential(
46             nn.ConvTranspose2d(32, 16, 3,
47                 stride=2, output_padding=0),
48             nn.BatchNorm2d(16),
49             nn.ReLU(True),
50             nn.ConvTranspose2d(16, 8, 3, stride=2,
51                 padding=1, output_padding=1),
52             nn.BatchNorm2d(8),
53             nn.ReLU(True),
54             nn.ConvTranspose2d(8, 1, 3, stride=2,
55                 padding=1, output_padding=1)
56         )
57
58     def forward(self, x):
59         # Encode
60         x = self.encoder_conv(x)
61         x = self.encoder_fc(x)
62         # Decode
63         x = self.decoder_fc(x)
64         x = self.decoder_conv(x)
65         return x
66
67
68 ##### Define the loss function
69 loss_fn = torch.nn.MSELoss()
70
71
72 ##### Define an optimizer (both for the encoder and the decoder!)
73 lr= 0.001
74
75
76 ##### Set the random seed for reproducible results
77 torch.manual_seed(0)
78
79
80 ##### Initialize the two networks
81 d = 4
82
83
84 #model = Autoencoder(encoded_space_dim=encoded_space_dim)
85 encoder = Encoder(encoded_space_dim=d,fc2_input_dim=128)
86 decoder = Decoder(encoded_space_dim=d,fc2_input_dim=128)
87 params_to_optimize = [
88     {'params': encoder.parameters()},
89     {'params': decoder.parameters()}
90 ]
91
92
93 optim = torch.optim.Adam(params_to_optimize, lr=lr, weight_decay=1e-05)
94
95
96 # Check if the GPU is available
97 device = torch.device("cuda") if torch.cuda.is_available() else torch.device("cpu")
98 print(f'Selected device: {device}')
99
100
101 # Move both the encoder and the decoder to the selected device
102 encoder.to(device)
103 decoder.to(device)

```

.py hosted with ❤ by GitHub

[view raw](#)

4. Train and evaluate model

We define a function to train the AE model. First, we pass the input images to the encoder. Later, the encoded data is passed to the decoder and then we compute the reconstruction loss with `loss_fn(x_hat, x)`. After we clear the gradient to not accumulate other values, we perform backpropagation and at the end, we compute the gradient by calling `opt.step()`.

```
1  ### Training function
2  def train_epoch(encoder, decoder, device, dataloader, loss_fn, optimizer):
3      # Set train mode for both the encoder and the decoder
4      encoder.train()
5      decoder.train()
6      train_loss = []
7      # Iterate the dataloader (we do not need the label values, this is unsupervised learning)
8      for image_batch, _ in dataloader: # with "_" we just ignore the labels (the second element of
9          # Move tensor to the proper device
10         image_batch = image_batch.to(device)
11         # Encode data
12         encoded_data = encoder(image_batch)
13         # Decode data
14         decoded_data = decoder(encoded_data)
15         # Evaluate loss
16         loss = loss_fn(decoded_data, image_batch)
17         # Backward pass
18         optimizer.zero_grad()
19         loss.backward()
20         optimizer.step()
21         # Print batch loss
22         print('\t partial train loss (single batch): %f' % (loss.data))
23         train_loss.append(loss.detach().cpu().numpy())
24
25     return np.mean(train_loss)
```

.py hosted with ❤ by GitHub

[view raw](#)

Once the training function is created, we define a function to evaluate the performance of the model. As before, we pass the image to the encoder. The encoded image is passed to the decoder. Then, we store all the image batches and the reconstruction into two different lists, that will be used to calculate the test loss.


```
1  ### Testing function
2  def test_epoch(encoder, decoder, device, dataloader, loss_fn):
3      # Set evaluation mode for encoder and decoder
4      encoder.eval()
5      decoder.eval()
6      with torch.no_grad(): # No need to track the gradients
7          # Define the lists to store the outputs for each batch
8          conc_out = []
9          conc_label = []
10         for image_batch, _ in dataloader:
11             # Move tensor to the proper device
12             image_batch = image_batch.to(device)
13             # Encode data
14             encoded_data = encoder(image_batch)
15             # Decode data
16             decoded_data = decoder(encoded_data)
17             # Append the network output and the original image to the lists
18             conc_out.append(decoded_data.cpu())
19             conc_label.append(image_batch.cpu())
20         # Create a single tensor with all the values in the lists
21         conc_out = torch.cat(conc_out)
22         conc_label = torch.cat(conc_label)
23         # Evaluate global loss
24         val_loss = loss_fn(conc_out, conc_label)
25         return val_loss.data
```

.py hosted with ❤ by GitHub

[view raw](#)

We also would like to see the reconstructed images during each epoch of the training. The goal is to understand how the autoencoder is learning from the input images.

```

1  def plot_ae_outputs(encoder,decoder,n=10):
2      plt.figure(figsize=(16,4.5))
3      targets = test_dataset.targets.numpy()
4      t_idx = {i:np.where(targets==i)[0][0] for i in range(n)}
5      for i in range(n):
6          ax = plt.subplot(2,n,i+1)
7          img = test_dataset[t_idx[i]][0].unsqueeze(0).to(device)
8          encoder.eval()
9          decoder.eval()
10         with torch.no_grad():
11             rec_img = decoder(encoder(img))
12             plt.imshow(img.cpu().squeeze().numpy(), cmap='gist_gray')
13             ax.get_xaxis().set_visible(False)
14             ax.get_yaxis().set_visible(False)
15             if i == n//2:
16                 ax.set_title('Original images')
17             ax = plt.subplot(2, n, i + 1 + n)
18             plt.imshow(rec_img.cpu().squeeze().numpy(), cmap='gist_gray')
19             ax.get_xaxis().set_visible(False)
20             ax.get_yaxis().set_visible(False)
21             if i == n//2:
22                 ax.set_title('Reconstructed images')
23         plt.show()

```

.py hosted with ❤ by GitHub

[view raw](#)

Let's break the test code into little pieces:

- `test_dataset[i][0].unsqueeze(0)` is used to extract the *i*th image from the test dataset and then it will be increased by 1 dimension on the 0 axis. This step is needed to pass the image to the autoencoder.
- `decoder(encoder(img))` is used to obtain the reconstructed image
- `plt.imshow(img.cpu().squeeze().numpy())` is used to plot the original image. `squeeze()` removes the dimension added before and is essential to visualize the image. `numpy()` transforms a tensor into a n-dimensional array, which is the only type of object accepted by the function `plt.imshow`. `numpy()` returns a copy of the tensor object into CPU memory.

Now we can finally begin to train the model on the training set and evaluate it on the validation set.

```

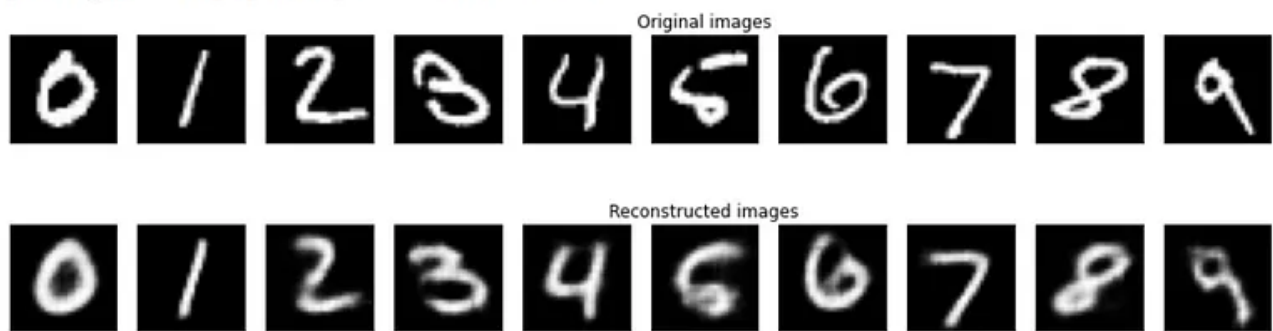
1  num_epochs = 30
2  diz_loss = {'train_loss':[], 'val_loss':[]}
3  for epoch in range(num_epochs):
4      train_loss = train_epoch(encoder, decoder, device,
5                              train_loader, loss_fn, optim)
6      val_loss = test_epoch(encoder, decoder, device, test_loader, loss_fn)
7      print('\n EPOCH {}/{} \t train loss {} \t val loss {}'.format(epoch + 1, num_epochs, train_loss,
8                                                                    val_loss))
9      diz_loss['train_loss'].append(train_loss)
10     diz_loss['val_loss'].append(val_loss)
11     plot_ae_outputs(encoder, decoder, n=10)

```

.py hosted with ❤ by GitHub

[view raw](#)

EPOCH 30/30 train loss 0.025 val loss 0.026



It's possible to notice that the autoencoder is able to reconstruct well the images after 30 epochs, even if there are some imperfections. But since this model is really simple, it performed very well. Now the model is trained and we want to do a final evaluation of the test set:

```

1  test_epoch(encoder, decoder, device, test_loader, loss_fn).item()

```

.py hosted with ❤ by GitHub

[view raw](#)

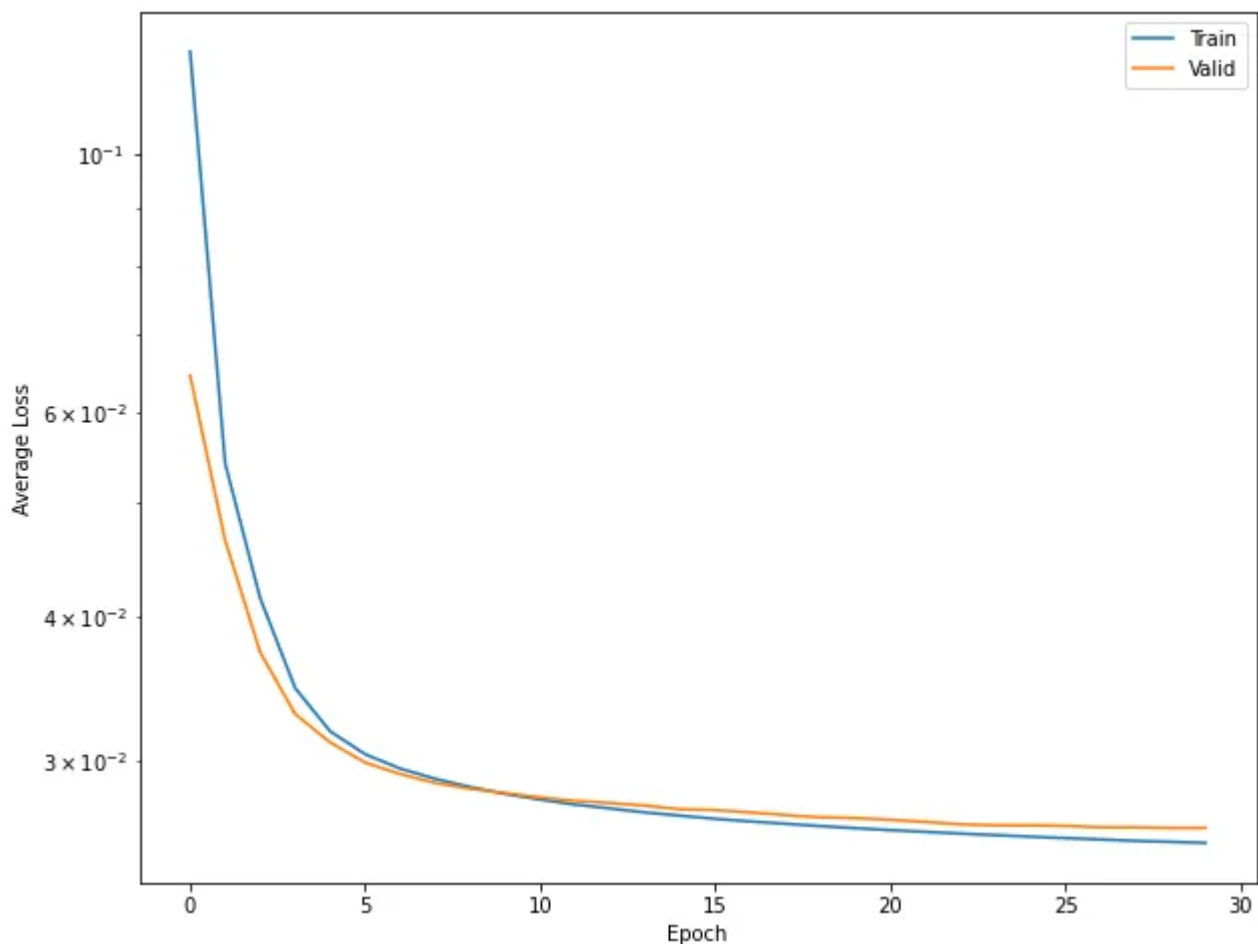
0.026236288249492645

We can also observe how the reconstruction losses decrease over the epochs:

```
1 # Plot losses
2 plt.figure(figsize=(10,8))
3 plt.semilogy(diz_loss['train_loss'], label='Train')
4 plt.semilogy(diz_loss['val_loss'], label='Valid')
5 plt.xlabel('Epoch')
6 plt.ylabel('Average Loss')
7 #plt.grid()
8 plt.legend()
9 #plt.title('loss')
10 plt.show()
```

.py hosted with ❤ by GitHub

[view raw](#)



5. Generate new samples from the random latent code

To generate new images from the latent code, we sample randomly from a normal distribution with the mean and standard deviation of the encoded data. These samples will be passed to the decoder, which will create the reconstructed images.

```
1  def show_image(img):
2      npimg = img.numpy()
3      plt.imshow(np.transpose(npimg, (1, 2, 0)))
4
5  encoder.eval()
6  decoder.eval()
7
8  with torch.no_grad():
9      # calculate mean and std of latent code, generated taking in test images as inputs
10     images, labels = iter(test_loader).next()
11     images = images.to(device)
12     latent = encoder(images)
13     latent = latent.cpu()
14
15     mean = latent.mean(dim=0)
16     print(mean)
17     std = (latent - mean).pow(2).mean(dim=0).sqrt()
18     print(std)
19
20     # sample latent vectors from the normal distribution
21     latent = torch.randn(128, d)*std + mean
22
23     # reconstruct images from the random latent vectors
24     latent = latent.to(device)
25     img_recon = decoder(latent)
26     img_recon = img_recon.cpu()
27
28     fig, ax = plt.subplots(figsize=(20, 8.5))
29     show_image(torchvision.utils.make_grid(img_recon[:100],10,5))
30     plt.show()
```

.py hosted with ❤ by GitHub

[view raw](#)



It should notice that this procedure allows the sampling to be in the same region as the latent code but at the same time there are digits that make no sense. This aspect is explained by the fact that the **latent space of the autoencoder is extremely irregular**: close points in the latent space can produce very different and meaningless patterns over visible units. For this reason, the autoencoder doesn't perform well for generative purposes.

6. Visualize the latent space with t-SNE

After we can observe dynamic visualization to see the latent space learned by the autoencoder. First, we create the encoded samples using the test set.

```
1  encoded_samples = []
2  for sample in tqdm(test_dataset):
3      img = sample[0].unsqueeze(0).to(device)
4      label = sample[1]
5      # Encode image
6      encoder.eval()
7      with torch.no_grad():
8          encoded_img = encoder(img)
9      # Append to list
10     encoded_img = encoded_img.flatten().cpu().numpy()
11     encoded_sample = {"Enc. Variable {i}": enc for i, enc in enumerate(encoded_img)}
12     encoded_sample['label'] = label
13     encoded_samples.append(encoded_sample)
14 encoded_samples = pd.DataFrame(encoded_samples)
15 encoded_samples
```

.py hosted with ❤ by GitHub

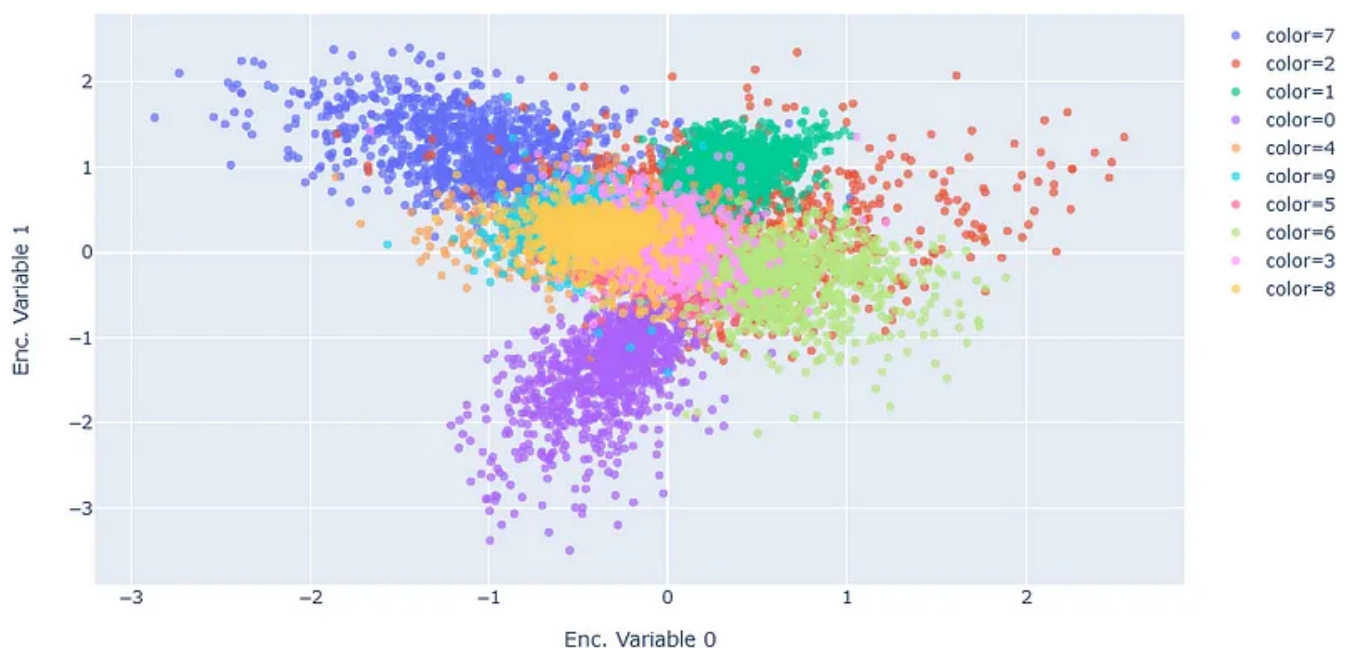
[view raw](#)

Let's plot the latent space representation using plotly express library:

```
1  import plotly.express as px
2
3  px.scatter(encoded_samples, x='Enc. Variable 0', y='Enc. Variable 1',
4             color=encoded_samples.label.astype(str), opacity=0.7)
```

.py hosted with ❤ by GitHub

[view raw](#)



Visualization of latent space

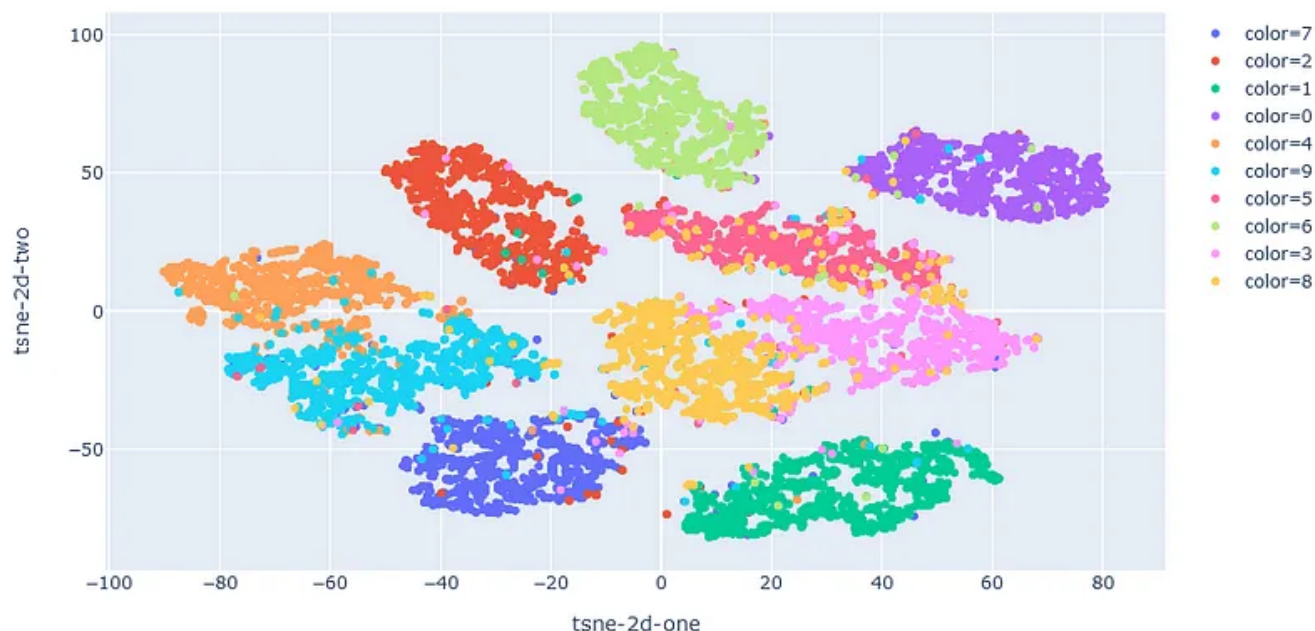
From this plot, we see that similar digits are clustered together. For example “4” overlap with “9” and “5”.

In order to have an easier representation to read, we can apply a dimensionality reduction, called t-SNE, to visualize the latent code in a 2d space. We'll fix the number of components equal to 2 for this reason.

```
1 from sklearn.manifold import TSNE
2
3 tsne = TSNE(n_components=2)
4 tsne_results = tsne.fit_transform(encoded_samples.drop(['label'],axis=1))
5 fig = px.scatter(tsne_results, x=0, y=1,
6                 color=encoded_samples.label.astype(str),
7                 labels={'0': 'tsne-2d-one', '1': 'tsne-2d-two'})
8 fig.show()
```

.py hosted with ❤ by GitHub

[view raw](#)



You can see that it distinguishes clearly one digit from the other. There are some exceptions with points that fall in other categories, but t-SNE still remains an improvement compared to the previous representation.

Final thoughts:

Congratulations! You have learned to implement a Convolutional autoencoder. There aren't many tutorials that talk about autoencoders with convolutional layers with Pytorch, so I wanted to contribute in some way. The autoencoder provides a way to compress images and extract the most important information. There are also many extensions of this model to improve the performance, some of these are the Denoising Autoencoder, the Variational Autoencoder, and the Generative Adversarial Networks. The GitHub code is [here](https://github.com/smartgeometry-ucl/dl4g). Thanks for reading. Have a nice day.

Reference:

[1] <https://github.com/smartgeometry-ucl/dl4g>

Did you like my article? *[Become a member](#) and get unlimited access to new data science posts every day! It's an indirect way of supporting me without any extra cost to you. If you are already a member, [subscribe](#) to get emails whenever I publish new data science and python guides!*

[Machine Learning](#)[Data Science](#)[Convolutional](#)[Autoencoder](#)[Pytorch](#)

Get an email whenever Eugenia Anello publishes.

Discover interesting topics about data science

Your email

[Subscribe](#)

By signing up, you will create a Medium account if you don't already have one. Review our [Privacy Policy](#) for more information about our privacy practices.

[About](#) [Help](#) [Terms](#) [Privacy](#)

Get the Medium app

