

# Autoencoders

## Introduction:

Autoencoders are a type of artificial neural network designed for unsupervised learning. They consist of an encoder network that maps the input data to a lower-dimensional representation (encoding), and a decoder network that reconstructs the original input from the encoding. The goal of autoencoders is to learn a compact representation of the input data.

## Intuition:

The intuition behind autoencoders is to force the model to learn a compressed representation of the input data by minimizing the reconstruction error. The encoder learns to capture essential features, and the decoder reconstructs the input from this reduced representation. Autoencoders are used for various tasks such as dimensionality reduction, data denoising, and feature learning.

## Algorithm:

### 1. Encoder:

- The encoder network takes the input data and maps it to a lower-dimensional representation (encoding).

### 2. Decoder:

- The decoder network takes the encoding and reconstructs the original input.

### 3. Loss Function:

- The loss function measures the difference between the input and the reconstructed output. Common loss functions include mean squared error or binary cross-entropy.

### 4. Training:

- The model is trained to minimize the reconstruction error by adjusting the weights in both the encoder and decoder.

## Implementation in Python (MNIST Data):

Here's a simple example of implementing an autoencoder on the MNIST dataset using PyTorch.

### Import libraries

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 from torch.utils.data import DataLoader
5 from torchvision import datasets, transforms
6 import matplotlib.pyplot as plt
7 import numpy as np
```

### Define the model

```
1 # Define the autoencoder model
2 class Autoencoder(nn.Module):
3     def __init__(self, encoding_dim):
4         super(Autoencoder, self).__init__()
```

```

5         self.encoder = nn.Sequential(
6             nn.Linear(28 * 28, 512),
7             nn.ReLU(),
8             nn.Linear(512, 256),
9             nn.ReLU(),
10            nn.Linear(256, encoding_dim),
11        )
12        self.decoder = nn.Sequential(
13            nn.Linear(encoding_dim, 256),
14            nn.ReLU(),
15            nn.Linear(256, 512),
16            nn.ReLU(),
17            nn.Linear(512, 28 * 28),
18            nn.Sigmoid(), # Sigmoid activation to ensure outputs are in [0,
19        1]
20    )
21
22    def forward(self, x):
23        x = x.view(x.size(0), -1) # Flatten the input
24        encoded = self.encoder(x)
25        decoded = self.decoder(encoded)
26        return decoded

```

This code defines a simple autoencoder model using the PyTorch framework. An autoencoder is a neural network architecture designed for unsupervised learning of efficient data codings, typically used for dimensionality reduction or feature learning.

Here's a breakdown of the `Autoencoder` class:

### 1. Initialization (`__init__`) Method:

- `encoding_dim`: The dimensionality of the latent space or encoding. This parameter determines the size of the compressed representation of the input data.
- `encoder`: The encoder is defined as a sequential module consisting of linear layers with ReLU activation functions. It takes the flattened input (assumed to be images of size 28x28 pixels) and progressively reduces the dimensionality until reaching the specified encoding dimension.
- `decoder`: The decoder is also defined as a sequential module. It takes the encoded representation and reconstructs the original input size. It mirrors the structure of the encoder but in reverse order, using ReLU activations for intermediate layers and a Sigmoid activation in the final layer to ensure the output values are in the range [0, 1].

### 2. Forward (`forward`) Method:

- `x`: The input data, assumed to be images, is passed through the encoder to obtain the encoded representation (`encoded`). This representation is then passed through the decoder to reconstruct the original input (`decoded`).
- The input is flattened (`view(x.size(0), -1)`) before being processed by the encoder. This is a common practice when working with image data.
- The reconstructed output is returned by the forward method.

The purpose of training such an autoencoder is to learn a compressed representation of the input data in the encoding space. During training, the model aims to minimize the difference between the input and the reconstructed output. This process encourages the autoencoder to capture meaningful features in the data and generate a compact representation in the encoding space.

```

1 from torchsummary import summary
2
3 # Load the trained model
4 complex_autoencoder = Autoencoder(encoding_dim=2)
5 complex_autoencoder.load_state_dict(torch.load('autoencoder_model.pth'))
6
7 # Move the model to the CPU
8 complex_autoencoder.to('cuda:0')
9
10 # Print the summary of the encoder
11 summary(complex_autoencoder.encoder, (128, 1, 784)) # Assuming MNIST images
    (1 channel, 28x28)
12
13 # Print the summary of the decoder
14 summary(complex_autoencoder.decoder, (2,)) # Assuming encoding_dim is 2

```

```

1 -----
2          Layer (type)          Output Shape          Param #
3 -----
4          Linear-1              [-1, 128, 1, 512]      401,920
5          ReLU-2                [-1, 128, 1, 512]         0
6          Linear-3              [-1, 128, 1, 256]      131,328
7          ReLU-4                [-1, 128, 1, 256]         0
8          Linear-5              [-1, 128, 1, 2]         514
9 -----
10 Total params: 533,762
11 Trainable params: 533,762
12 Non-trainable params: 0
13 -----
14 Input size (MB): 0.38
15 Forward/backward pass size (MB): 1.50
16 Params size (MB): 2.04
17 Estimated Total Size (MB): 3.92
18 -----
19 -----
20          Layer (type)          Output Shape          Param #
21 -----
22          Linear-1              [-1, 256]              768
23          ReLU-2                [-1, 256]                0
24          Linear-3              [-1, 512]             131,584
25          ReLU-4                [-1, 512]                0
26          Linear-5              [-1, 784]             402,192
27          Sigmoid-6             [-1, 784]                0
28 -----
29 Total params: 534,544
30 Trainable params: 534,544
31 Non-trainable params: 0
32 -----
33 Input size (MB): 0.00
34 Forward/backward pass size (MB): 0.02
35 Params size (MB): 2.04
36 Estimated Total Size (MB): 2.06
37 -----

```

The `Autoencoder` class in the provided code has three layers in both the encoder and the decoder. Let's break down the layers:

#### Encoder:

1. Linear layer with input size  $28 * 28$  and output size 512.
2. ReLU activation function.
3. Linear layer with input size 512 and output size 256.
4. ReLU activation function.
5. Linear layer with input size 256 and output size equal to the specified `encoding_dim`.

#### Decoder:

1. Linear layer with input size equal to the specified `encoding_dim` and output size 256.
2. ReLU activation function.
3. Linear layer with input size 256 and output size 512.
4. ReLU activation function.
5. Linear layer with input size 512 and output size  $28 * 28$ .
6. Sigmoid activation function.

Therefore, both the encoder and the decoder consist of five layers. This architecture is a common choice for a basic autoencoder. The encoder progressively reduces the dimensionality of the input, while the decoder reconstructs the original input size. The ReLU activation functions introduce non-linearity, and the Sigmoid activation in the final layer of the decoder ensures that the reconstructed values are in the range  $[0, 1]$ .

#### Training function

```

1  # Function to train the autoencoder
2  def train_autoencoder(model, dataloader, criterion, optimizer,
   num_epochs=10):
3      for epoch in range(num_epochs):
4          for data in dataloader:
5              inputs, _ = data
6
7              optimizer.zero_grad()
8              outputs = model(inputs)
9              loss = criterion(outputs, inputs.view(inputs.size(0), -1))
10             loss.backward()
11             optimizer.step()
12
13         print(f'Epoch {epoch+1}/{num_epochs}, Loss: {loss.item()}')
```

#### Train the model

```

1  # Set random seed for reproducibility
2  torch.manual_seed(42)
3
4  # Define hyperparameters
5  encoding_dim = 2
6  batch_size = 128
7  learning_rate = 0.001
8  num_epochs = 20
9
10 # Load MNIST dataset
```

```

11 transform = transforms.Compose([transforms.ToTensor()])
12 train_dataset = datasets.MNIST(root='./data', train=True, download=True,
    transform=transform)
13 train_loader = DataLoader(dataset=train_dataset, batch_size=batch_size,
    shuffle=True)
14
15 # Initialize the autoencoder model, criterion, and optimizer
16 autoencoder_model = Autoencoder(encoding_dim=encoding_dim)
17 criterion = nn.MSELoss() # Mean Squared Error Loss
18 optimizer = optim.Adam(autoencoder_model.parameters(), lr=learning_rate)
19
20
21 # Train the autoencoder
22 train_autoencoder(autoencoder_model, train_loader, criterion, optimizer,
    num_epochs=num_epochs)

```

### Save the model

```

1 torch.save(autoencoder_model.state_dict(), 'autoencoder_model.pth')
2 print("Trained model saved.")

```

### Visualize embeddings

```

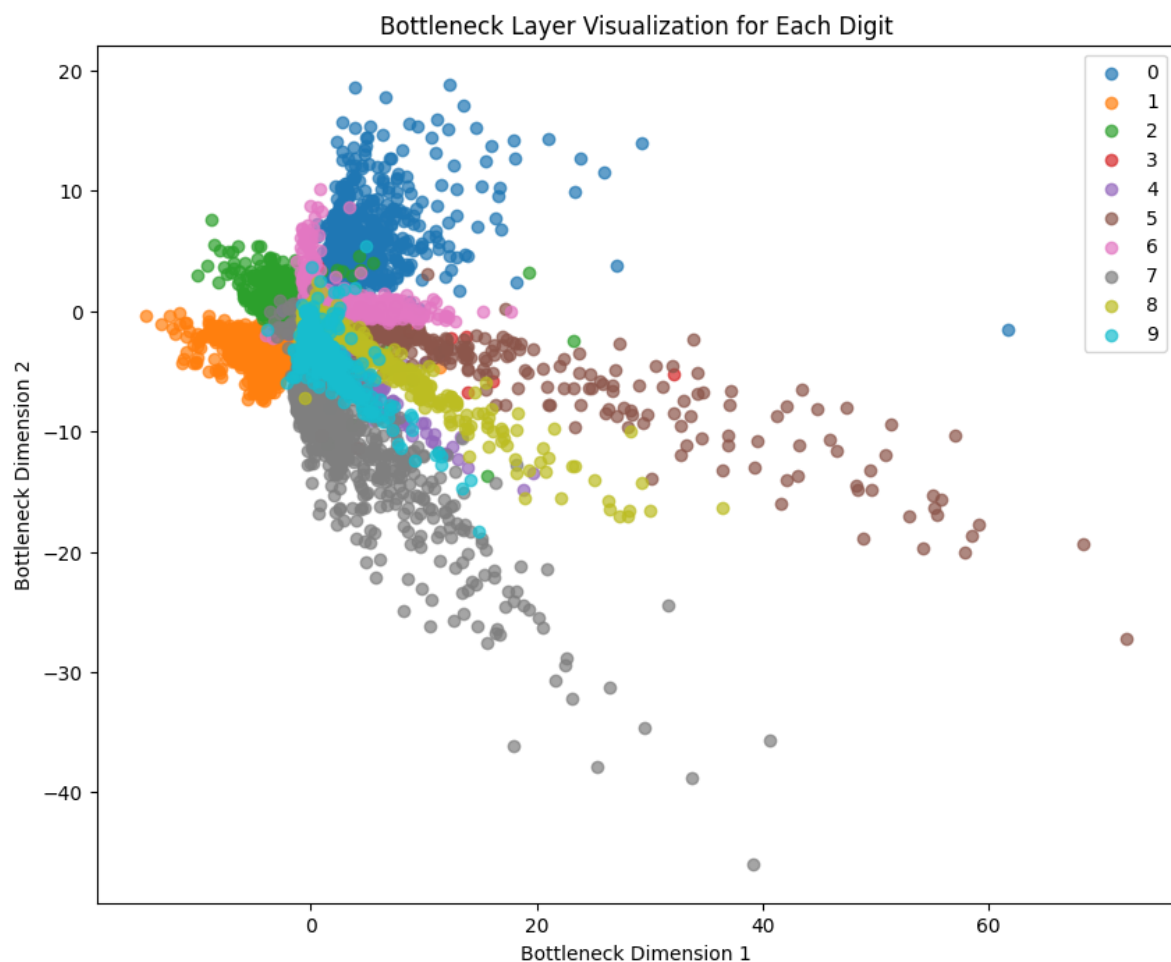
1 # Function to visualize the bottleneck layer for each digit
2 def visualize_bottleneck(encoder, dataloader):
3     encoder.eval()
4     all_embeddings = []
5     all_labels = []
6
7     with torch.no_grad():
8         for data in dataloader:
9             inputs, labels = data
10            embeddings = encoder(inputs.view(inputs.size(0),
-1)).detach().numpy()
11            all_embeddings.append(torch.from_numpy(embeddings)) # Convert
to PyTorch tensor
12            all_labels.append(labels)
13
14            all_embeddings = torch.cat(all_embeddings, dim=0)
15            all_labels = torch.cat(all_labels, dim=0).numpy() # Convert to NumPy
array
16
17            # Plot 2D representations, color-coded by digit label
18            plt.figure(figsize=(10, 8))
19            print(all_labels)
20            for digit in range(10):
21                digit_indices = (all_labels == digit)
22                plt.scatter(all_embeddings[digit_indices, 0],
all_embeddings[digit_indices, 1], label=str(digit), alpha=0.7)
23
24            plt.title('Bottleneck Layer Visualization for Each Digit')
25            plt.xlabel('Bottleneck Dimension 1')
26            plt.ylabel('Bottleneck Dimension 2')
27            plt.legend()

```

```

28 plt.show()
29
30 # Load the trained autoencoder
31 autoencoder_model = Autoencoder(encoding_dim=2) # Assuming you have a
    trained Autoencoder
32 autoencoder_model.load_state_dict(torch.load('autoencoder_model.pth')) #
    Load the saved model
33
34 # Load MNIST dataset
35 transform = transforms.Compose([transforms.ToTensor()])
36 test_dataset = datasets.MNIST(root='./data', train=False, download=True,
    transform=transform)
37 test_loader = DataLoader(dataset=test_dataset, batch_size=len(test_dataset),
    shuffle=False)
38
39 # visualize the bottleneck layer for each digit
40 visualize_bottleneck(autoencoder_model.encoder, test_loader)

```



In this example:

- The autoencoder consists of a 5 fully connected layers for both the encoder and decoder.
- The model is trained to minimize the MSE loss between the input and the reconstructed output.
- After training, the autoencoder is used to encode and decode the test set, and the results are visualized.

Feel free to adjust the architecture, hyperparameters, and visualize more samples based on your needs.