

## Project 2: Autoencoders

---

### 1 INTRODUCTION

An autoencoder is simply a neural network which is trained to reproduce an approximation of its input at its output. Typically, the input is mapped to a reduced dimension representation in the hidden layer, referred to as the "bottleneck" layer, and then the "bottleneck" layer representation is used to calculate the reconstruction of the input sample. The *encoder* of the autoencoder is the portion of the neural network that maps the input to the hidden layer representation, and the *decoder* maps the hidden layer representation to an approximate reconstruction of the input at the output. Figure 1.1 below shows the architecture of an autoencoder. Initially, an autoencoder may not seem useful since the objective is to obtain an input reconstruction, which will be slightly different than the input. However, once an autoencoder is trained, its bottleneck layer is of more interest than its reconstruction. Of course, during training, the reconstruction should have a minimum error w.r.t. the input, but, by training a neural network in this way, the bottleneck layer is forced to learn a reduced dimensional representation of the input. In other words, a trained autoencoder can be used for dimensionality reduction. Using an autoencoder for dimensionality reduction is much different than other methods (such as PCA) because the network is learning a non-linear mapping (due to the activation functions) of the input to a reduced dimensional space. The new representation may capture aspects of the input that linear dimensionality reduction algorithms may fail to represent.

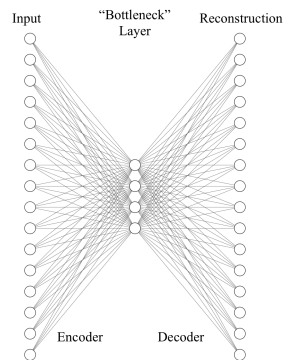


Figure 1.1: The architecture of an autoencoder.

In this project, we will implement various types of autoencoder neural networks using the MNIST Database of Hand Written Digits. This is a dataset of 70,000 images (each a  $28 \times 28$  grayscale image of a hand written digit as shown below in Figure 1.2), which we will divide into a 60000/10000 train/test set. The objective of the autoencoders will be to reconstruct each image at the output with the smallest possible error.

For this project, the MNIST data, which can be loaded using Keras, is done for you in the available code. You are expected to fill in the Jupyter Notebook provided to you as well as generate all deliverables, as outlined in the following sections, in the Notebook. **Submit your final**



Figure 1.2: Example images from the MNIST dataset.

code as a single .ipynb file. Ensure that the deliverables are clearly shown for each part to earn full credit. Name your submission *lastname\_firstname\_autoencoders.ipynb*.

## 2 EXAMPLE AUTOENCODERS USING KERAS IN PYTHON

Below is an example of how to create and train an autoencoder consisting of fully connected layers using Keras in Python with an input dimension of 100, a hidden layer dimension of 25, and an output dimension of 100 (the output will always be the same dimension as the input since our goal is to reconstruct the input at the output). Beginning on line 70, see how to implement an autoencoder consisting of convolutional layers as well.

Listing 1: Autoencoders in Keras

```
1 #ECE 595 Machine Learning II
2 #Project 1: Autoencoders – Example Code
3 #Written by: Rajeev Sahay and Aly El Gamal
4
5 from keras.models import Sequential
6 from keras.layers import Dense
7 from keras import backend as K
8 import matplotlib.pyplot as plt
9
10 # Assume data_train is a matrix with shape (num_samples, 100)
11 # Assume data_test is a matrix with shape (num_samples, 100)
12
13 # Create a fully connected autoencoder architecture
14 def autoencoder():
15     model = Sequential()
16
17     # Encoder
18     model.add(Dense(25,
19                     activation="choose_activation",
20                     use_bias=True,
21                     kernel_initializer="uniform",
22                     input_dim=100))
23
24     # Decoder
25     model.add(Dense(100,
26                     activation="choose_activation",
27                     kernel_initializer="uniform"))
28
29     return model
30
31 # Create an autoencoder graph
32 ae_model = autoencoder()
33
34 # Compile the graph using a particular loss function and optimizer algorithm
35 ae_model.compile(loss="choose_loss_fcn",
36                 optimizer="choose_optimizer_alg")
37 # Training parameters
38 num_epochs = 100
39 batch_size = 1024
40
41 # Train the model using the inputs as the target outputs
42 ae_model_history = ae_model.fit(data_train, data_train,
```

```

43         validation_data=(data_test, data_test),
44         epochs=num_epochs,
45         batch_size=batch_size,
46         shuffle=True)
47
48 # Generate predicted reconstructions using the trained model
49 reconstructions = ae_model.predict(data_test)
50
51 # Obtain the hidden layer representation of the input data:
52 get_hl = K.function([ae_model.layers[0].input], [ae_model.layers[0].output])
53 data_test_hidden_layer_rep = get_hl([data_test])[0]
54
55 # Plot training history of loss versus epochs
56 plt.plot(ae_model_history.history['loss'])
57 plt.plot(ae_model_history.history['val_loss'])
58
59 # Plot first three data samples in data_test as grayscale image
60 n = 3 # how many digits we will display
61 plt.figure(figsize=(20, 4))
62 for i in range(n):
63     ax = plt.subplot(1, n, i + 1)
64     plt.imshow(data_test[i].reshape(10, 10))
65     plt.gray()
66     ax.get_xaxis().set_visible(False)
67     ax.get_yaxis().set_visible(False)
68
69
70 # One can create a convolutional autoencoder according to the
71 # following architecture, which can be trained identically to
72 # the fully connected model above
73
74 # Assume data_train is a matrix with shape (num_samples, 10, 10, 1) here
75 # Assume data_test is a matrix with shape (num_samples, 10, 10, 1) here
76
77 # Use UpSampling 2D for upsampling layer
78 from keras.layers import Conv2D, MaxPooling2D, UpSampling2D
79
80 def convolutional_autoencoder():
81     model = Sequential()
82
83     #Encoder
84     model.add(Conv2D(64,
85                     kernel_size=4,
86                     activation="choose_activation",
87                     padding="same",
88                     input_shape=(10, 10, 1)))
89     model.add(MaxPooling2D(pool_size=(2, 2),
90                           padding="same"))
91
92     #Decoder

```

```

93     model.add(Conv2D(64,
94                     kernel_size=4,
95                     activation="choose_activation",
96                     padding='same'))
97     model.add(UpSampling2D((2, 2)))
98
99     #Reconstruction
100    model.add(Conv2D(1,
101                    (3, 3),
102                    activation='choose_activation',
103                    padding='same'))
104
105    return model

```

### 3 PART 1: FULLY CONNECTED AUTOENCODER

Implement a fully connected autoencoder according to the following specifications:

- Input layer: 784 perceptrons (reshape each image so that it is a flattened 784-dimensional vector)
- First hidden layer: 400 perceptrons
- Second hidden layer: 200 perceptrons
- Third hidden layer ("bottleneck" layer): 100 perceptrons
- Fourth hidden layer: 200 perceptrons
- Fifth hidden layer: 400 perceptrons
- Output layer (reconstruction): 784 perceptrons

You will have to make decisions about which activation functions to use at each layer, whether or not to use a bias for each layer, which loss function to use for training, which optimizer to use for updating the parameter values during training, and the batch size. **Train your model using 150 epochs.**

Deliverables for Part 1:

1. We want to reconstruct output values representing the intensity of the pixels, which are between 0 and 1. Therefore, the choice of the output layer activation function is important. [Open ended reasoning question]
  - a) Choose 'softmax' or 'sigmoid'. Reason why one is preferred over the other.
2. Show a plot of the model loss versus the epochs for both the training and validation data (the validation data here can be the testing data). Include a legend, appropriate axes labels, and a title.
  - a) Mean Squared Error: Report training and validation loss graph vs epoch.
  - b) Binary Cross Entropy: Report training and validation loss graph vs epoch.
  - c) BONUS [reasoning question]: Compare pairs of the output layer activation and loss function and determine which perform better together.
    - Activation functions: Linear , Sigmoid
    - Loss functions: Binary Cross Entropy, Mean Square error.
3. Show images of the first ten samples in the testing set.
4. Show images of the "bottleneck" layer representation of the first ten images in the testing set (reshape the images to be 10 x 10)
5. Show images of the reconstructed samples of the first ten images in the testing set on both loss functions. Observe the output images and comment on which loss function is better and why?
6. If we were to reconstruct pixels values in the range of [0 to 255] directly at the output of the last layer, should there be an activation function in the last layer? If yes, which activation function and why?, If No, reason why?

## 4 PART 2: CONVOLUTIONAL AUTOENCODER

Implement a convolutional autoencoder according to the following specifications:

- Input layer that accepts a  $28 \times 28$  grayscale image
- First hidden layer: A convolutional layer with 16 feature maps and a  $3 \times 3$  filter size
- Second hidden layer: A  $2 \times 2$  maxpooling layer
- Third hidden layer: A convolutional layer with 8 feature maps and a  $3 \times 3$  filter size
- Fourth hidden layer ("bottleneck" layer): A  $2 \times 2$  maxpooling layer
- Fifth hidden layer: A convolutional layer with 8 feature maps and a  $3 \times 3$  filter size
- Sixth hidden layer: A  $2 \times 2$  upsampling layer
- Seventh hidden layer: A convolutional layer with 16 feature maps and a  $3 \times 3$  filter size
- Eighth hidden layer: A  $2 \times 2$  upsampling layer
- Output layer (reconstruction): A convolutional layer with a single feature map and a  $3 \times 3$  filter size

You will have to make decisions about which activation functions to use at each layer, whether or not to use a bias for each layer, which loss function to use for training, which optimizer to use for updating the parameter values during training, and the batch size. **Train your model using 150 epochs.**

Deliverables for Part 2:

1. Show a plot of the model loss versus the epochs for both the training and validation data (the validation data here can be the testing data). Include a legend, appropriate axes labels, and a title.
  - a) Mean Squared Error: Report training and validation loss graph vs. epoch.
  - b) Binary Cross Entropy: Report training and validation loss graph vs. epoch.
2. Show images of the first ten samples in the testing set
3. Show images of the "bottleneck" layer representation of the first ten images in the testing set (reshape the images to be  $28 \times 14$ )
4. Show images of the reconstructed samples of the first ten images in the testing set for both loss functions. Observe the output images and comment on which loss function is better and why?

## 5 PART 3: DENOISING AUTOENCODER

An autoencoder can be used to filter noise (denoise) from data samples. The architecture of a denoising autoencoder is no different than that of a regular autoencoder. In fact, the only difference between a normal autoencoder and a denoising autoencoder is the training data. A denoising autoencoder is trained to filter noise from the input and produce a denoised version of the input as the reconstructed output. This is achieved during training by using the noisy sample as the input and the original noise-free sample as the target reconstruction. In Part 3, you are going to implement a denoising autoencoder that can filter Gaussian noise injected into the MNIST dataset. The noisy data is created by generating 784 values, per sample, from a zero mean unit variance Gaussian distribution, scaling them by 0.25 (to make the additional noise more intense), and adding them to each point in the original image. Each image is then clipped to stay in the valid pixel range. The autoencoder will then be trained to eliminate this added noise.

Begin by selecting an architecture (you may use your architecture from Part 1 or Part 2 or create a new architecture) for your denoising autoencoder. Then, compile and train the model appropriately using the noisy data as inputs and the noise-free data as the target labels. Finally, show ten samples of original images from the testing set, their corresponding noisy counterparts, and their denoised DAE outputs.

### Deliverables for Part 3:

1. Show images of the first ten samples in the training set
2. Show images of the first ten samples in the training set after injecting Gaussian noise
3. Show a plot of the model loss versus the epochs for both the training and validation data (the validation data here can be the testing data) of the DAE. Make sure to include a legend, appropriate axes labels, and a title.
4. Show images of the first ten samples in the testing set
5. Show images of the first ten samples in the testing set after adding noise
6. Show images of the first ten samples in the testing set after “denoising”