# DSP Assignment 3

## Submitted to: Dr. Mohsen Rashwan

| Name | SN | BN |
|---|---|---|
| Abdelrahman Mostafa Mouse Mohamed Rabah | 3 | 12 |
| Abdallah Adel Ibrahim | 3 | 15 |
| Amr Mohamed | 3 | 30 |

Data Downloading:

```
1  import opendatasets
2
3  opendatasets.download('https://www.kaggle.com/datasets/mohamedgamal07/reduced-mnist')
✓ 24.9s
```

Making sure Data is correct and divided equally:

```
1  classes = os.listdir(train_path)
2  num_classes = len(classes)
3  source_path=[train_path +f'{a}' for a in classes]
4  classes_dir=[f'{a}_dir' for a in classes]
5  for cl_dir,cl_path in zip(classes_dir,source_path):
6      print(cl_dir,': ',len(os.listdir(cl_path)))
✓ 0.1s
```

```
0_dir :  1000
1_dir :  1000
2_dir :  1000
3_dir :  1000
4_dir :  1000
5_dir :  1000
6_dir :  1000
7_dir :  1000
8_dir :  1000
9_dir :  1000
```

```
1  classes = os.listdir(test_path)
2  num_classes = len(classes)
3  source_path=[test_path +f'{a}' for a in classes]
4  classes_dir=[f'{a}_dir' for a in classes]
5  for cl_dir,cl_path in zip(classes_dir,source_path):
6      print(cl_dir,': ',len(os.listdir(cl_path)))
```
✓ 0.2s

```
0_dir :  200
1_dir :  200
2_dir :  200
3_dir :  200
4_dir :  200
5_dir :  200
6_dir :  200
7_dir :  200
8_dir :  200
9_dir :  200
```

Copy data from downloaded folders to arrays:

```
1   def copy(data_path):
2       imgs = []
3       labels = []
4       folders = os.listdir(data_path)
5       for folder in folders:
6           fullpath = os.path.join(data_path, folder)
7           images = glob.glob(os.path.join(fullpath,'*.jpg'))
8           for img in images:
9               labels.append(img.split('/')[-1][0])
10              img = cv2.imread(img)
11              img = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
12              imgs.append(img)
13
14      return np.array(imgs), np.array(labels)
```
✓ 0.2s

Shuffling and Normalizing the data:

```
1  X_train, y_train = copy(train_path)
2  X_train = X_train/255.0
3  X_train, y_train = shuffle(X_train, y_train, random_state=42)
4  X_test, y_test = copy(test_path)
5  X_test = X_test/255.0
6  X_test, y_test = shuffle(X_test, y_test, random_state=42)
7  y_train = y_train.astype(int)
8  y_test = y_test.astype(int)
```
✓  1m 18.2s

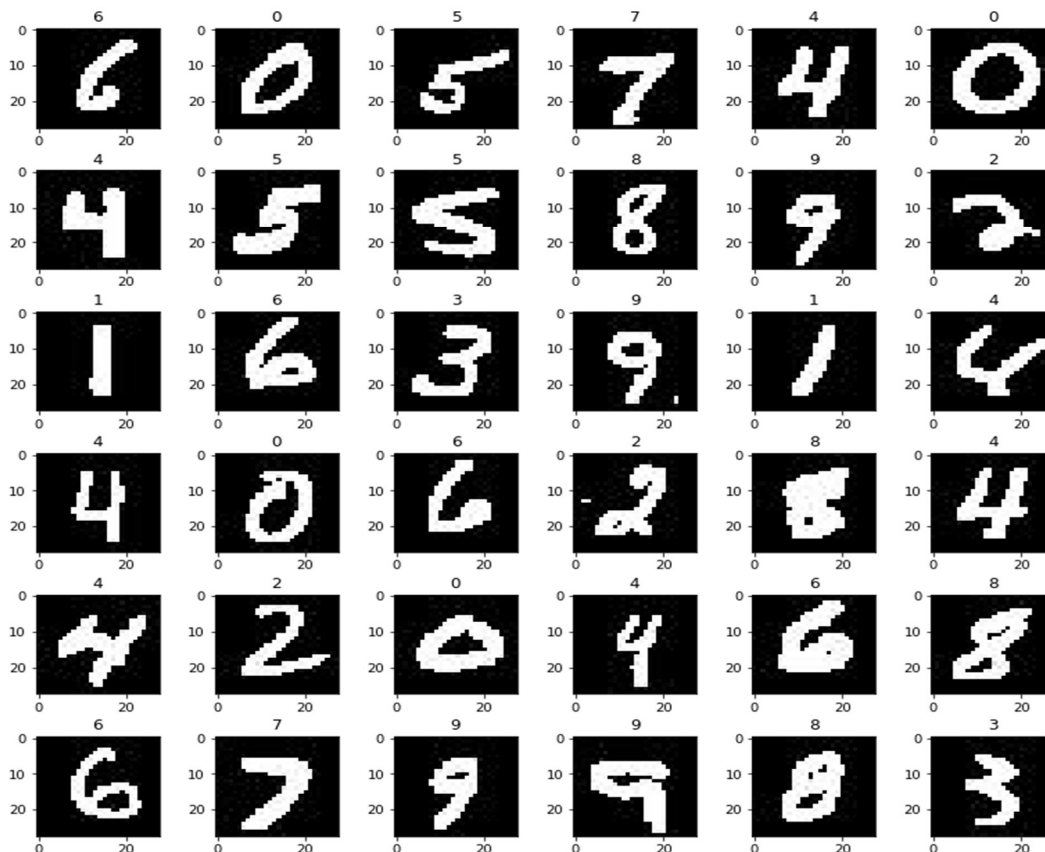Plotting few samples:

```
1  fig, ax = plt.subplots(6, 6, figsize = (12, 12))
2  fig.suptitle('First 36 images in MNIST')
3  fig.tight_layout(pad = 0.3, rect = [0, 0, 0.9, 0.9])
4  for x, y in [(i, j) for i in range(6) for j in range(6)]:
5      ax[x, y].imshow(X_train[x + y * 6].reshape((28, 28)), cmap = 'gray')
6      ax[x, y].set_title(y_train[x + y * 6])
```
✓  5.3s



First 36 images in MNIST

Reshaping the data to make it 1D instead of 2D (Gray Channels):

```python
def unroll(x):
    x = x.reshape(x.shape[0], x.shape[1]*x.shape[2])
    return x
```

The DCT and Zig-Zag method:

```python
def dct2(block):
    l = []
    for b in block:
        l.append(dct(dct(b.T, norm='ortho').T, norm='ortho'))
    return l

def small_zigzag(x):
    l = []
    for a in x:
        l.append(np.concatenate([np.diagonal(a[::-1,:], i)[::(2*(i % 2)-1)] for i in range(1-a.shape[0], a.shape[0])])[0:200])
    return np.array(l)
```

DCT Generation:

```python
1  X_tr_dct, t1 = dct2(X_train)
2  X_tr_dct_zz, t2 = small_zigzag(X_tr_dct)
3  X_train
4  print('Generated Train DCT with Zig-Zag in {} seconds, shape of train{}'.format(t1+t2, X_tr_dct_zz.shape))
5  X_te_dct, t1 = dct2(X_test)
6  X_te_dct_zz, t2 = small_zigzag(X_te_dct)
7  print('Generated Test DCT with Zig-Zag in {} seconds, shape of test{}'.format((t1+t2), X_te_dct_zz.shape))
✓ 2.8s

Generated Train DCT with Zig-Zag in 2.259641170501709 seconds, shape of train(10000, 200)
Generated Test DCT with Zig-Zag in 0.4289977550506592 seconds, shape of test(2000, 200)
```

PCA:

```python
pca_n = 0.9
def pca(x_train, x_test, n):
    start = time.time()
    pca = PCA(n_components=n, random_state=42)
    x_train = pca.fit_transform(x_train)
    x_test = pca.transform(x_test)
    end = time.time()
    print(f'time it took PCA:{end - start}, shape of X train: {x_train.shape}, x_test = {x_test.shape}' )
    return x_train, x_test
```

pca_n is variable which indicates to keep 90% of variance

PCA Generation:

```
1  X_tr_pca = unroll(X_train)
2  X_te_pca = unroll(X_test)
3  X_tr_pca, X_te_pca = pca(X_tr_pca, X_te_pca, pca_n)
✓ 1.1s
```

```
time it took PCA:1.0498216152191162, shape of X train: (10000, 167), x_test = (2000, 167)
```

Encoder:

```
1   import tensorflow as tf
2   from tensorflow import keras
3
4   tf.random.set_seed(42)
5   np.random.seed(42)
6
7   stacked_encoder = keras.models.Sequential([
8       keras.layers.Flatten(input_shape=[28, 28]),
9       keras.layers.Dense(100, activation="selu"),
10      keras.layers.Dense(30, activation="selu"),
11  ])
12  stacked_decoder = keras.models.Sequential([
13      keras.layers.Dense(100, activation="selu", input_shape=[30]),
14      keras.layers.Dense(28 * 28, activation="sigmoid"),
15      keras.layers.Reshape([28, 28])
16  ])
17
18  stacked_ae = keras.models.Sequential([stacked_encoder, stacked_decoder])
19  stacked_ae.compile(loss="binary_crossentropy",
20                     optimizer='adam')
21  history = stacked_ae.fit(X_train, X_train, epochs=15,
22                           validation_data=(X_test, X_test))
✓ 28.3s
```

Results:

```
Epoch 14/15
313/313 [==============================] - 2s 6ms/step - loss: 0.1161 - val_loss: 0.1195
Epoch 15/15
313/313 [==============================] - 2s 6ms/step - loss: 0.1154 - val_loss: 0.1184
```

Generation:

```
1  X_train_encoded = stacked_ae.predict(X_train)
2  X_train_encoded_unrolled = unroll(X_train_encoded)
3  X_test_encoded = stacked_ae.predict(X_test)
4  X_test_encoded_unrolled = unroll(X_test_encoded)
```
✓ 0.7s

```
1  print(X_train_encoded_unrolled.shape, X_test_encoded_unrolled.shape)
```
✓ 0.2s

```
(10000, 784) (2000, 784)
```

Plotting some generated images:

```
1  def plot_image(image):
2      plt.imshow(image, cmap="binary")
3      plt.axis("off")
4
5  def show_reconstructions(model, images=X_test, n_images=5):
6      reconstructions = model.predict(images[:n_images])
7      fig = plt.figure(figsize=(n_images * 1.5, 3))
8      for image_index in range(n_images):
9          plt.subplot(2, n_images, 1 + image_index)
10         plot_image(images[image_index])
11         plt.subplot(2, n_images, 1 + n_images + image_index)
12         plot_image(reconstructions[image_index])
13
14  show_reconstructions(stacked_ae)
```
✓ 1.4s

KMeans:

```python
clusters = [1, 4, 16, 32]
def kmeans(x_train, n):
    start = time.time()
    KMS = KMeans(n_clusters=n, random_state=42)
    KMS.fit(x_train)
    end = time.time()
    return KMS, end - start
def retrieve_info_kmeans(cluster_labels, y_train, model):
    reference_labels = {}
    for i in range(len(np.unique(model.labels_))):
        index = np.where(model.labels_ == i,1,0)
        num = np.bincount(y_train[index==1]).argmax()
        reference_labels[i] = num
    return reference_labels
```

Gaussian Mixtures:

```python
mixtures = [1, 2, 4]
def gmm(x_train, n):
    start = time.time()
    GMM = GaussianMixture(n_components=n, random_state=42)
    GMM.fit(x_train)
    end = time.time()
    return GMM, end - start

def retrieve_info_gmm(x, y_train, model):
    labels = model.predict(x)
    reference_labels = {}
    for i in range(len(np.unique(labels))):
        index = np.where(labels == i,1,0)
        num = np.bincount(y_train[index==1]).argmax()
        reference_labels[i] = num
    return reference_labels, labels
```

Support Vector Machines Classifier:

```python
kinds = ['linear', 'rbf']
def svm(x_train, x_test, kind):
    start = time.time()
    SVM = SVC(kernel=kind, random_state=42)
    SVM.fit(x_train, y_train)
    y_hat = SVM.predict(x_train)
    y_pred = SVM.predict(x_test)
    end = time.time()
```

| Classifier | Features | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | DCT | | | | PCA | | | | Autoencoder | | | |
| | Train | | Test | | Train | | Test | | Train | | Test | |
| | Accuracy | Time | Accuracy | Time | Accuracy | Time | Accuracy | Time | Accuracy | Time | Accuracy | Time |
| Kmeans | | | | | | | | | | | | |
| Kmeans-1 | 10.0% | 2.777 | 10.0% | 2.556 | 10.0% | 2.703 | 10.0% | 2.511 | 10.0% | 3.644 | 10.0% | 2.736 |
| Kmeans-4 | 36.52% | 1.630 | 38.3% | 0.467 | 36.51% | 1.450 | 38.35% | 0.500 | 36.01% | 7.261 | 38.25% | 1.947 |
| Kmeans-16 | 72.76% | 5.002 | 77.149% | 1.132 | 73.34% | 5.123 | 76.1% | 1.011 | 73.78% | 28.83 | 77.60% | 5.47 |
| Kmeans-32 | 80.63% | 9.391 | 86.1% | 1.652 | 82.06% | 9.507 | 86.55% | 2.328 | 80.88% | 47.386 | 85.1% | 9.656 |
| GMM | | | | | | | | | | | | |
| GMM-1 | 10.0% | 0.974 | 10.0% | 0.705 | 10.0% | 1.07 | 10.0% | 0.449 | 10.0% | 4.706 | 10.0% | 1.2409 |
| GMM-2 | 19.16% | 9.217 | 19.8% | 1.021 | 18.98% | 24.626 | 19.90% | 0.889 | 19.58% | 88.09 | 19.90% | 4.558 |
| GMM-4 | 30.83% | 41.42 | 38.45% | 1.423 | 35.14% | 67.23 | 38.55% | 1.448 | 29.18% | 233.52 | 37.7% | 6.1389 |
| SVM | Accuracy Train | Accuracy Test | Total Time | | Accuracy Train | Accuracy Test | Total Time | | Accuracy Train | Accuracy Test | Total Time | |
| SVM-Linear | 98.07% | 94.35% | 9.264 | | 98.8% | 93.8% | 5.716 | | 99.039% | 94.69% | 22.00 | |
| SVM-rbf | 97.66% | 97.35% | 17.31 | | 99.19% | 97.7% | 21.12 | | 97.09% | 96.35% | 49.86 | |

Time is in seconds. Processing speed depends on CPU/GPU speed. Autoencoder takes so much time since it just generates new noisy images with all features.

an RBF (radial basis function) kernel is used when the boundaries are hypothesized to be curve-shaped.

RBF kernel uses two main parameters, gamma and C that are related to:

1. the decision region (how spread the region is), and
2. the penalty for misclassifying a data point

respectively.

So, in contrast, it helps the SVM to become nonlinear rather than linear. RBF kernel function is similar to Normal distribution.

- For Kmeans-32 best model (PCA) confusion matrix:

Train:

Seaborn Confusion Matrix with labels

| Actual \ Predicted | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 921 | 0 | 2 | 16 | 2 | 15 | 28 | 0 | 15 | 1 |
| 1 | 0 | 987 | 3 | 2 | 0 | 2 | 1 | 1 | 3 | 1 |
| 2 | 9 | 23 | 893 | 11 | 10 | 3 | 10 | 10 | 22 | 9 |
| 3 | 9 | 11 | 7 | 908 | 4 | 14 | 1 | 4 | 27 | 15 |
| 4 | 2 | 16 | 3 | 0 | 635 | 0 | 19 | 17 | 0 | 308 |
| 5 | 19 | 8 | 2 | 197 | 17 | 689 | 34 | 5 | 4 | 25 |
| 6 | 8 | 6 | 0 | 0 | 1 | 7 | 978 | 0 | 0 | 0 |
| 7 | 0 | 35 | 4 | 2 | 27 | 3 | 0 | 863 | 1 | 65 |
| 8 | 13 | 26 | 6 | 136 | 14 | 30 | 5 | 17 | 743 | 10 |
| 9 | 6 | 14 | 1 | 17 | 231 | 3 | 0 | 131 | 8 | 589 |

Predicted Values

Test:

Seaborn Confusion Matrix with labels

| Actual \ Predicted | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 189 | 1 | 0 | 3 | 1 | 3 | 2 | 0 | 1 | 0 |
| 1 | 0 | 199 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 2 | 7 | 2 | 170 | 2 | 0 | 2 | 1 | 2 | 14 | 0 |
| 3 | 0 | 1 | 2 | 173 | 0 | 2 | 0 | 0 | 19 | 3 |
| 4 | 0 | 1 | 0 | 0 | 174 | 0 | 3 | 0 | 3 | 19 |
| 5 | 0 | 5 | 0 | 37 | 1 | 140 | 6 | 0 | 9 | 2 |
| 6 | 1 | 0 | 0 | 0 | 0 | 0 | 197 | 0 | 2 | 0 |
| 7 | 0 | 1 | 6 | 0 | 12 | 0 | 0 | 175 | 0 | 6 |
| 8 | 2 | 1 | 0 | 5 | 1 | 9 | 0 | 0 | 181 | 1 |
| 9 | 0 | 0 | 0 | 0 | 54 | 0 | 0 | 13 | 0 | 133 |

Predicted Values

- For GMM-4 best model (PCA) confusion matrix:

Train:

**Seaborn Confusion Matrix with labels**

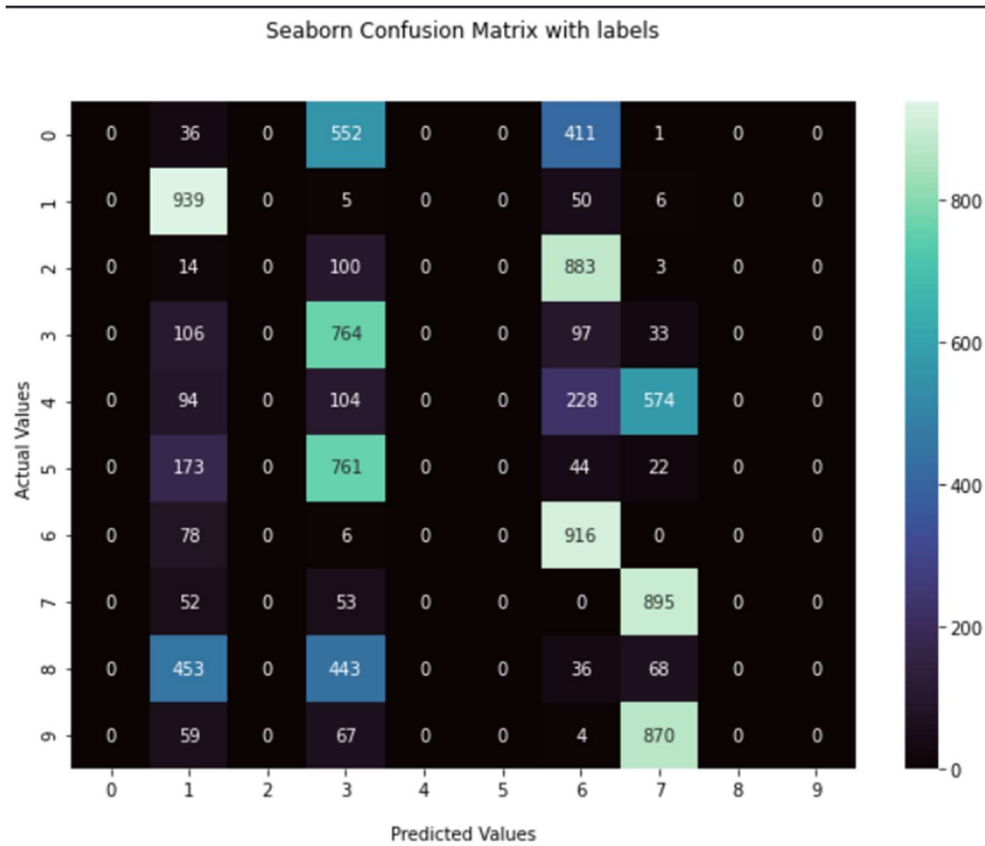| Actual \ Predicted | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 36 | 0 | 552 | 0 | 0 | 411 | 1 | 0 | 0 |
| 1 | 0 | 939 | 0 | 5 | 0 | 0 | 50 | 6 | 0 | 0 |
| 2 | 0 | 14 | 0 | 100 | 0 | 0 | 883 | 3 | 0 | 0 |
| 3 | 0 | 106 | 0 | 764 | 0 | 0 | 97 | 33 | 0 | 0 |
| 4 | 0 | 94 | 0 | 104 | 0 | 0 | 228 | 574 | 0 | 0 |
| 5 | 0 | 173 | 0 | 761 | 0 | 0 | 44 | 22 | 0 | 0 |
| 6 | 0 | 78 | 0 | 6 | 0 | 0 | 916 | 0 | 0 | 0 |
| 7 | 0 | 52 | 0 | 53 | 0 | 0 | 0 | 895 | 0 | 0 |
| 8 | 0 | 453 | 0 | 443 | 0 | 0 | 36 | 68 | 0 | 0 |
| 9 | 0 | 59 | 0 | 67 | 0 | 0 | 4 | 870 | 0 | 0 |

Predicted Values

Test:

**Seaborn Confusion Matrix with labels**

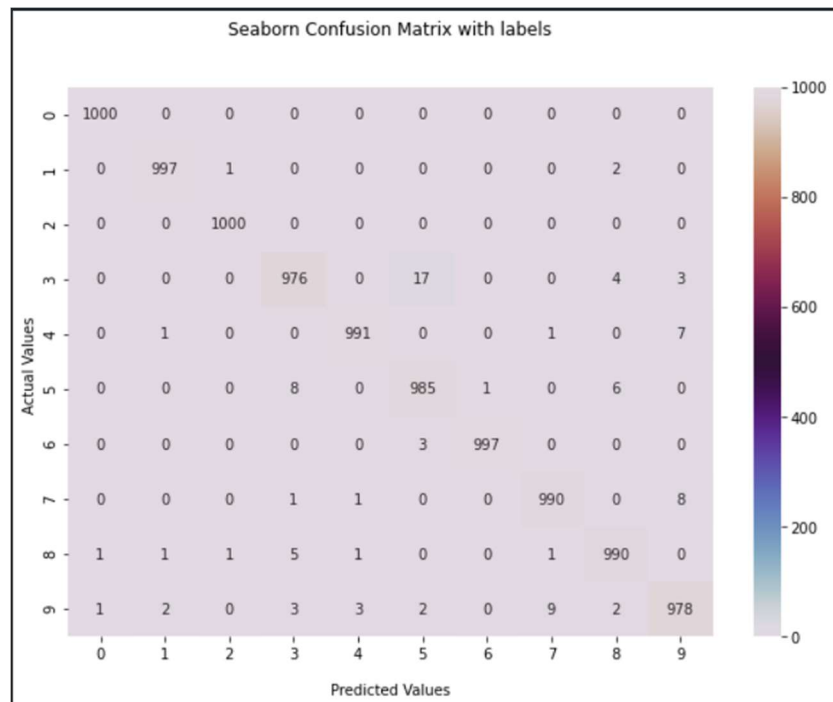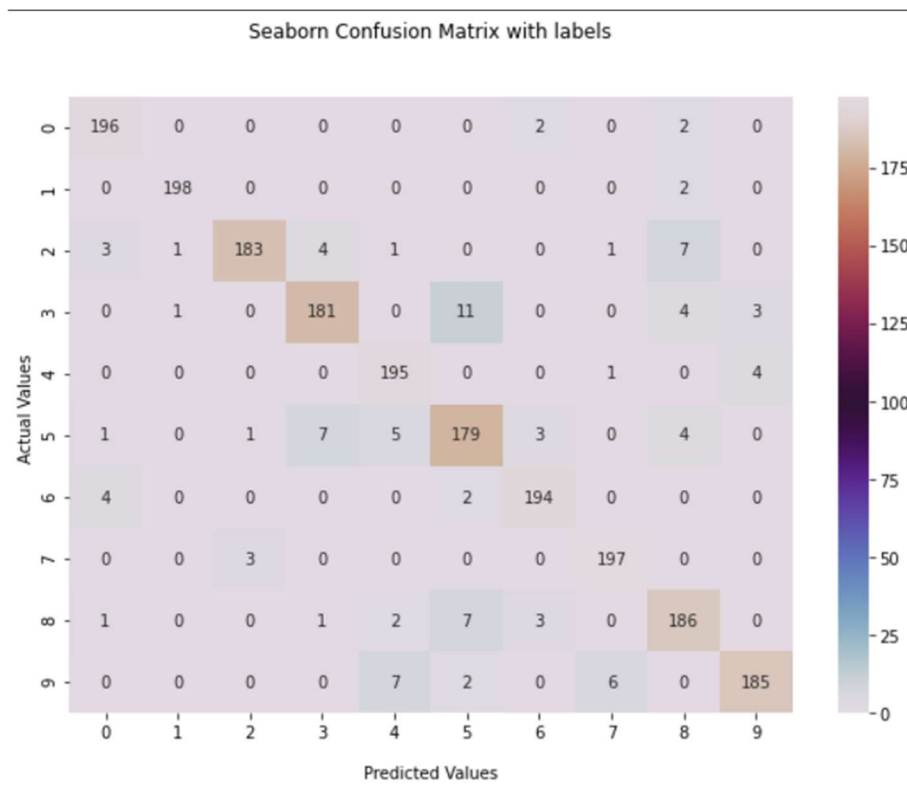| Actual \ Predicted | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 192 | 3 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 198 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 92 | 32 | 0 | 74 | 0 | 0 | 0 | 0 | 0 | 2 |
| 3 | 2 | 8 | 0 | 188 | 0 | 0 | 0 | 0 | 0 | 2 |
| 4 | 8 | 29 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 162 |
| 5 | 12 | 83 | 0 | 91 | 0 | 0 | 0 | 0 | 0 | 14 |
| 6 | 180 | 17 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 1 | 7 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 189 |
| 8 | 3 | 49 | 0 | 126 | 0 | 0 | 0 | 0 | 0 | 22 |
| 9 | 0 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 193 |

Predicted Values

- For SVM, best training accuracy for Linear Kernel + Autoencoder

Train:



Test:

- But highest test accuracy for rbf Kernel + PCA

Train:

Seaborn Confusion Matrix with labels

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **1** | 0 | 1000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **2** | 0 | 0 | 995 | 1 | 1 | 0 | 0 | 0 | 3 | 0 |
| **3** | 1 | 1 | 1 | 977 | 0 | 14 | 0 | 0 | 4 | 2 |
| **4** | 0 | 0 | 0 | 0 | 989 | 0 | 0 | 1 | 0 | 10 |
| **5** | 1 | 1 | 1 | 11 | 3 | 981 | 0 | 0 | 2 | 0 |
| **6** | 0 | 0 | 0 | 0 | 0 | 0 | 1000 | 0 | 0 | 0 |
| **7** | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 988 | 0 | 10 |
| **8** | 1 | 2 | 3 | 9 | 0 | 8 | 0 | 0 | 977 | 0 |
| **9** | 1 | 0 | 0 | 1 | 9 | 0 | 0 | 14 | 2 | 973 |

Actual Values / Predicted Values

Test:

Seaborn Confusion Matrix with labels

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 198 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| **1** | 0 | 199 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| **2** | 2 | 0 | 181 | 3 | 1 | 1 | 2 | 1 | 9 | 0 |
| **3** | 1 | 0 | 3 | 181 | 0 | 5 | 0 | 3 | 5 | 2 |
| **4** | 2 | 0 | 2 | 0 | 189 | 0 | 2 | 2 | 0 | 3 |
| **5** | 5 | 1 | 2 | 2 | 7 | 173 | 3 | 0 | 7 | 0 |
| **6** | 4 | 0 | 0 | 0 | 0 | 3 | 192 | 0 | 1 | 0 |
| **7** | 0 | 0 | 4 | 0 | 1 | 0 | 0 | 191 | 0 | 4 |
| **8** | 1 | 0 | 1 | 3 | 0 | 3 | 4 | 0 | 188 | 0 |
| **9** | 0 | 0 | 0 | 1 | 10 | 0 | 0 | 4 | 1 | 184 |

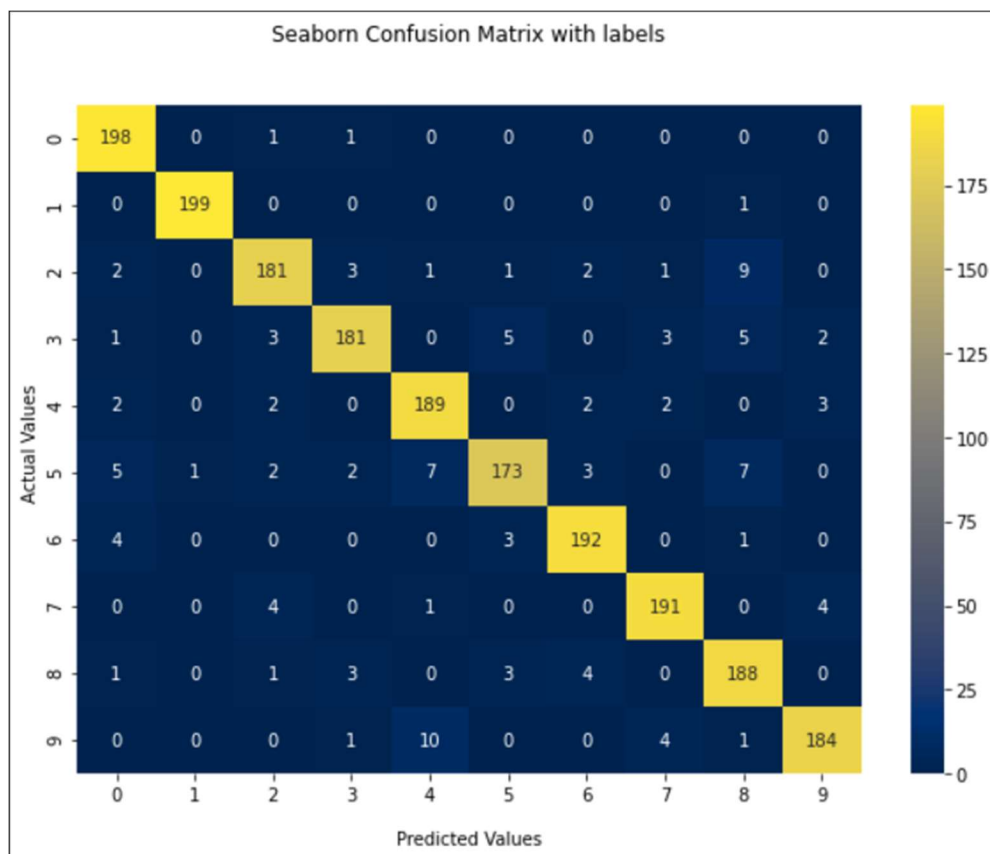Actual Values / Predicted Values

- Conclusion:

- As the number of clusters increases the accuracy increases, but of course it will be more computational expensive.
- As the number of Gaussian means increases the accuracy increases, but of course it will be more computational expensive, same as previous.
- It's noted that KMeans is faster and performed better than Gaussian Mixture even though Gaussian Mixture is a more of generalization of Kmeans.
- Autoencoder takes much time because it's not just a mathematical formula unlike other methods, besides choosing number of epochs, batch size, all of that affect the timing and processing speed. We can increase number of hidden layers or choose Convolutional layers but all of that besides it will require more time with more processing power (The training of the neural network used GPU not CPU), the better the network the worse the output, because the output may copy the input exactly and won't learn many features, so we stick with simple models.
- For this high dimensional data, rbf filter was much better than linear filter and this was noted in the test accuracy but with almost around 2x *more or less* time required for computations, so does it worth this almost 2-3% increase? it depends.