# Intro to Machine Learning - Assignment 3+4

Maxmillan Ries s3118134, Cristian Rosiu s3742377

October 2020

## Contents

# 1 Outliers detection and Density-based clustering

## 1.1 DBSCAN Implementation

a In the implementation of the algorithm we used the pseudo-code provided in the slides.

Because of the abstract nature of the pseudo-code, we had to make some design choices regarding the implementation of the $region\_query()$ function. For better performance, we pre-compute the distances of each point to all other points by using the distance matrix and sort indices based on the distances.

```
1  self.distances = distance_matrix(self.points, self.points)
```

```
1  # Compute neighbourhood of a core point
2  def region_query(self, i):
3      neighbours = set()
4      # Check core point's distances to all points
5      for j in range(len(self.points)):
6          if self.distances[i][j] <= self.eps:
7              # If found point distances is in the radius
8              # of the core point, add it to neighbourhood
9              neighbours.add(j)
10     return neighbours
```

It is also worth mentioning that, in order to check if a point is visited, the distance of a point or if a point it's already a member of another cluster, we decided to create lists for every property of a point. (e.g. The list *status_list* is initialized with status 'New' for each element in the points array. This means that each point in the original list is directly mapped using the index to the respective property list).

```
1  # Number of samples / points
2  self.n_pts = len(points)
3  # List of labels
4  self.labels = [0] * self.n_pts
5  # Property Lists
6  self.status_list = [Status.New] * self.n_pts
7  self.member_list = [False] * self.n_pts
8  self.distances = distance_matrix(self.points, self.points)
```

Full code can be found in the appendix.

## 1.2  Evaluating DBSCAN on toy-sets

a  For this part we have chosen to first create 2 arrays of values. One for the values that $\epsilon$ will take and the other array for the values of MinPts variable:

```
1  eps_values = [0.15, 0.2, 0.25, 0.3, 0.32, 0.35, 0.4, 0.45,
   ↪   0.5]
2  min_pts_values = [3, 4, 5, 6, 7, 8, 9, 10]
```

We then proceeded to compute and plot the labels of each combination $(eps, min\_pts)$ in order to better visualize the resulted clusters.

What we observed first is that, changing the value of MinPts had less impact for bigger values on the overall data compared to when the value

of $\epsilon$ changed. However, it is worth mentioning that for bigger values of MinPts, we got better results. Therefore, our observation is that, it doesn't make sense to choose a threshold value which is less than the dimension of your data (e.g. with a value of 1, there will be no outliers detected because each point will be a core point.).
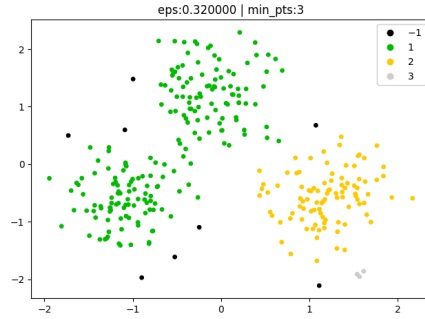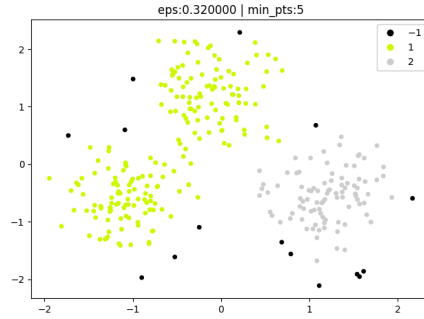


Figure 1: Clustering using min_pts = 3



Figure 2: Clustering using min_pts = 5

On the other hand, as we already stated, the value of $\epsilon$ seems to have the most impact on data. We have observed that for a lower value, more clusters appear whereas if you keep increasing the value, the data tends to converge towards one larger cluster.
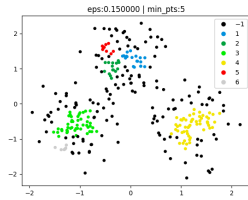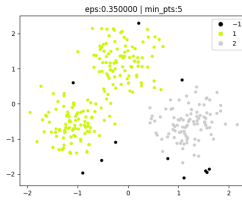


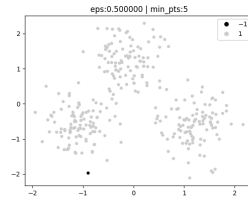Figure 3: Epsion = 0.15 MinPts = 5



Figure 4: Epsion = 0.25 MinPts = 5



Figure 5: Epsion = 0.5 MinPts = 5

**Which set of parameters allow you to cluster the data into 1, 2, and 3 clusters?**
We have found the following set of parameters, which we have written underneath the matching graphs:
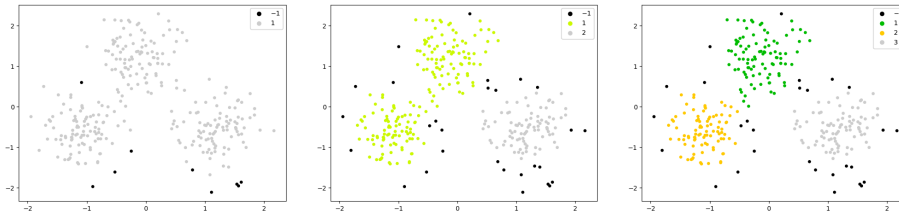
Figure 6: Epsion $= 0.351$  Figure 7: Epsion $= 0.256$  Figure 8: Epsion $= 0.255$
MinPts $= 4$                 MinPts $= 4$                 MinPts $= 4$

b **Compute and evaluate cluster's Silhouette score**

|                   | 1 Cluster | 2 Clusters | 3 Clusters |
|-------------------|-----------|------------|------------|
| Silhouette Score  | 0.1960    | 0.3585     | 0.5316     |

What we can observe from the silhouette score of each cluster number is that, the most overlapping occurs when the algorithm finds a single clusters.

As an obvious observation based upon both logic and graph, is that the best silhouette score occurs when 3 clusters are found by the algorithm.

## 1.3   Outlier Detection On Subset MNIST Dataset

a Data Exploration

We started by retrieving the data from the given .mat file. Note that X array represents the actual values and y their lables

```python
# Retrieve Data
X = mnist_data['X']
y = mnist_data['y']
```

Then we proceeded to standardize the data, apply PCA on the result and then plot the newly reduce data:

```python
X = StandardScaler().fit_transform(X)
# Apply PCA on given data
pca = PCA(n_components=2)
X = pca.fit_transform(X)
# Plot the data and color code using the given labels (ground
    truth)
plt.scatter(X[:, 0], X[:, 1], c=y)
plt.show()
```

Figure 9: Plot Of MNIST Data Set After Dimensionality Reduction

As we can observe, using the ground-truth labels, the graph clearly shows the presence of 2 types of data. The yellow data points, as also stated in the data set folder, are the outliers and are label with the number 0. The magenta data points are labeled with number 1 and they represent inliers.

b  Clustering - Outlier Detection

Firstly, we decided to plot the sorted distance of every point to it's $k^{th}$ nearest neighbour (k = 4):



Figure 10: Graph Which Determines Optimal Epsilon

Here we can observe that the point of maximum curvature is around the interval:

$$0.2, 0.4$$

This implies that the optimal radius of our neighbourhood is in that interval. Everything above this value will be just too big to fit the graph.
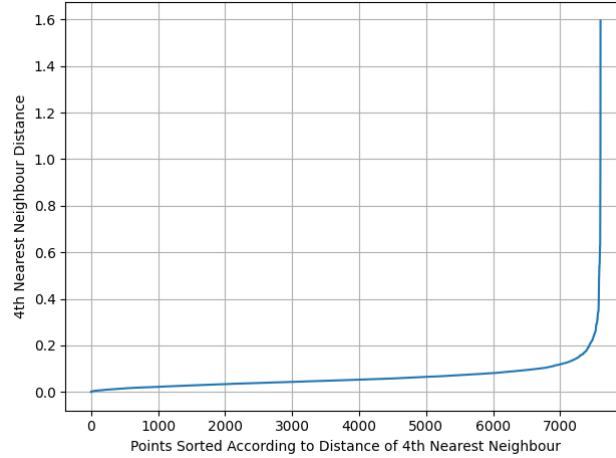
So with a value of 0.35 for $\epsilon$ and MinPts $= 20$ we have got the following graph (-1 represents outliers)



Figure 11: Epsilon $= 0.35$  MinPts $= 5$

c Evaluation
For the computation of accuracy and F1 scores we decided to use the confusion matrix in the following way:

```python
from sklearn.metrics import confusion_matrix

def evaluate(self, y_true):
    # Get the values of the confusion matrix
    tn, fp, fn, tp = confusion_matrix(y_true,
        self.labels).ravel()
    # Compute Accuracy score.
    accuracy = float((tp + tn)/(tp + tn + fp + fn))
    # Precision and recall will help calculate the F1 score.
    precision = float(tp/(tp + fp))
    recall = float(tp/(tp + fn))
    F1 = 2*precision*recall/(precision+recall)
    # Print the results.
    print('Accuracy: %f' % accuracy)
```

```
14        print('F1 Score: %f' % F1)
```

Here is a table with the final scores:

| Accuracy Score | 0.157832 |
|---|---|
| F1 Score | 0.149781 |

As DBSCAN clusters according to the density of the data, it is not a suitable choice for the given data, which consists of handwritten labelled numbers. As such, we can see in the table above, that the accuracy and fscore are lower than expected.

A more suitable algorithm for this kind of clustering, we think, is one that focuses on detecting clusters, rather than detecting outliers based on the density of the data.

Therefore, DBSCAN is not suitable for the given data set (handwritten numbers from 0 to 9).

# 2 Vector Quantization

## 2.1 Winner-Takes-All unsupervised competitive learning (VQ)

For this assignment we decided to put every implemented algorithm in a single package and import them separately when testing their functionalities. Therefore, the VQ and DBSCAN algorithms can be found in the AssignmentAlgorithms.py, and can be imported as follows:

```
1 from AssignmentAlgorithms import VQ, DBSCAN
```

Further more, the way we implemented Winter-Takes-All algorithms is quite simple and pretty intuitive. It's important to mention that our implementation is heavily inspired by the slides and instruction that we have got.

First of all, the most important function of our class is the following one:

```
1 # X is the desired dataset.
2 def fit(self, X):
3     # Randomly initialize K prototypes with actual points from the
      ↪   dataset.
4     self.init_prototypes(X)
5     # Plot first epoch where the prototypes are randomly
      ↪   positioned.
6     self.plot_epoch(X, 0)
7     for epoch in range(self.epochs):
8         # Change data order in every epoch.
9         randomized_data = permutation(X)
```

```
10          sum_squared = 0.0
11          for point in randomized_data:
12              # Get the closest prototype to current point.
13              winner_index, distance = self.evaluate_winner(point)
14              # Calculate Squared error.
15              for i in range(len(point)):
16                  error = point[i] - self.prototypes[winner_index][i]
17                  sum_squared += error**2
18              # Move prototype towards point using the learning_rate.
19              self.update_prototype(winner_index, point)
20          self.squared_errors.append(sum_squared)
21      # Plot final result.
22      self.plot_epoch(X, 0)
```

It updates each prototype, plots the result of each epoch and computes the squared error distance. Moreover, another somewhat important function which we thought it would be important to talk about is the one that finds the winner prototype:

```
1   # Determines the winner prototype for a specific point
2   def evaluate_winner(self, data_point):
3       winner = (0, sys.maxsize)
4       # Go over the list of prototypes and find the closest using
        ↪ euclidean distance
5       for index, prototype in enumerate(self.prototypes):
6           distance = np.linalg.norm(prototype - data_point)
7           # Find prototype with the minumim distance.
8           if distance < winner[1]:
9               winner = (index, distance)
10      # Return the winner index
11      return winner[0], winner[1]
```

The update winner prototype function follows the pseudocode described closely, iterating over each prototype and updating it if its the winner prototype, rpeating the process for each data point.
An important detail is that we chose to calculate the Hvq directly as part of our running algorithm, rather than create a separate function.
After implementing this, we then tested the functionality on the given data set.

## 2.2 Simple VQ Data Set

a As stated in the assignment, we proceeded to test our algorithm on the given data set by using different values for K and 0.1 for the learning rate. We would like to precise that the given data set did not contain any names for each component of the vectors, and as such as have kept them blank in our graphs. First, we found the following results using K=2:
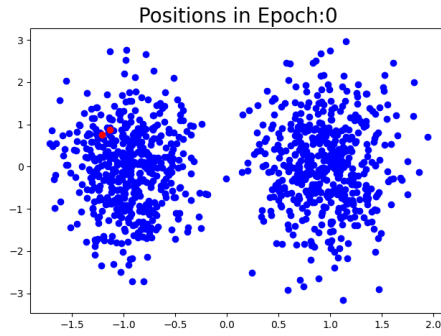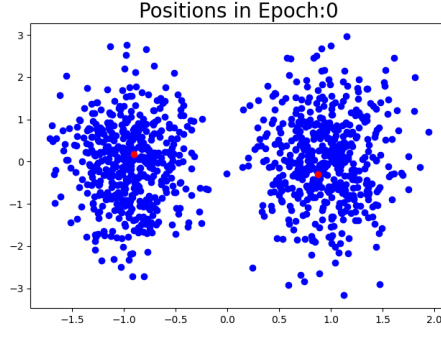
Figure 12: Initial guess using K = 2



Figure 13: Final guess after all iterations using K = 2

The left picture shows the initial location of the 2 prototypes after randomly assigned positions whereas the right one shows the same prototypes, but in their final position. We can clearly observe that the algorithm managed to find the 2 obvious clusters. However, the two prototypes are not exactly in the center of each cluster, but slightly off. There may well be a combination of causes which result in this, but our current hypothesis is that, due to the randomly selected data, the final few evaluated points consisted of the outer values of the respective clusters, causing the prototypes to be "pulled" a little away from the center.

We are also fairly confident in stating that the learning rate has a strong impact on the visual and mathematical shifting of the winner prototype and may partially be responsible for the lack of a perfect center.

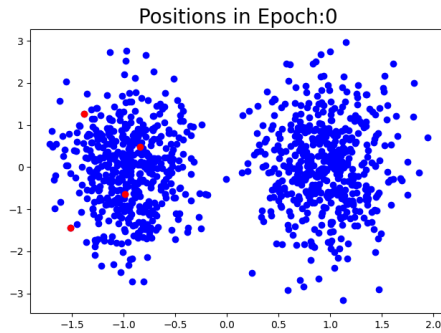On the other hand, we found the following results for K=4 and the same learning rate
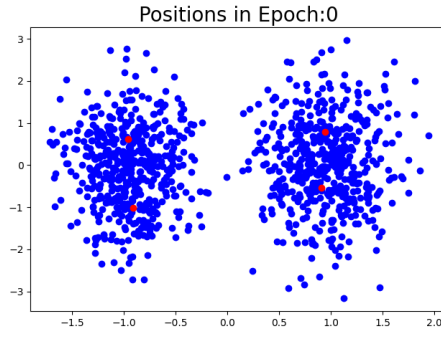


Figure 14: Initial guess using K = 4



Figure 15: Final guess after all iterations using K = 4

Here, we can observe a similar behaviour to the example where we had K = 2. However, the distance between one prototype and the center of a cluster is much bigger now. This is caused by what is described as the "ill-defined" issue with clustering. From a mathematical and logical standpoint, the algorithm has function, of splitting the data both horizontally and vertically once.

Visually for us as observers however, the clustering using K=4 does not seem to have resulted in something tangible and properly done. To evaluate this further, the elbow method will be used later in the assignment, as to find the best number of clusters for this set of data. Our current hypothesis is that the elbow method will result in an optimum number of clusters being equal to 2.

For comparison, we have also tried a smaller value for the learning rate (learning_rate=0.05). We have got the following results:
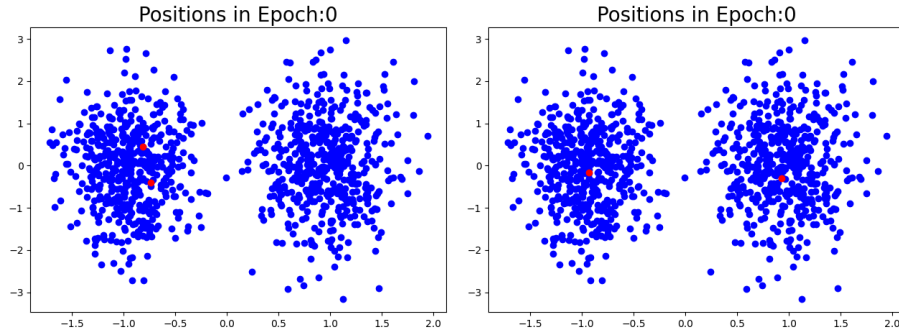


Figure 16: Initial guess using learning_rate = 0.05    Figure 17: Final guess after all iterations using learning_rate = 0.05

It seems that, running the algorithm with a smaller learning rate value, significantly increased the chance of a better guess. This was expected as the distance traveled by a winner prototype is directly proportional with the value of the learning_rate. Formula proof: (n is learning rate, e is current data point):

$$w^{i*} = w^{i*} + n(e^{t}w^{i*})$$

N, or the learning rate, is a ratio that scales the difference in the distance (between the winner and the data point), before adding it to the current position of the winner prototype, effectively determining how far the prototype will travel. A lower value of n, will result in a shorter distance travelled, which can prevent any cluster outliars/values from causing the prototypes to shift drastically away from the center of the cluster.

10

For a K value of 2, we have found that a reasonable value of epochs is 100. Beyond this number of epochs, the value of Hvq does not seem to steadily decrease towards a minimum any further.



b The following graphs present the variation of Hvq in each epoch, for K=2 and K=4 for the following values of learning rate: 0.04, 0.1, 0.6:



In the above iterations (K = 2), we can clearly observe the larger variation in error with a slow learning rate compared to the increasing variation of the other graphs. This seems to go against our hypothesis. This will be further discussed below.

11

On the other hand, it seems that if the number of prototypes K is increased, the variation in error is much smaller compared to a smaller K.
With a value of K, the data is "split" according the winner takes all rules K times. As the value of K goes up, the prototypes will be scattered such that the difference between the data surrounding each prototype and each respective prototype will be smaller. As such, due to the smaller euclidean distances obtained as a result, the value of Hvq, which is the sum of these distances, will be smaller.
It is for this reason, that, assuming perfect clustering, if K = number of data points, the value of Hvq will go to 0.

**How does the final value of the cost function change with n?**
The learning rate is a ratio which affects the extent to which the prototype is updated. Our hypothesis, is that a too high or too low learning will lead to insufficient clustering.
To further prove our point we decided to apply the algorithm with extreme values for the learning rate (i.e. 0.0001 and 1) and plot it to see what would the result be.
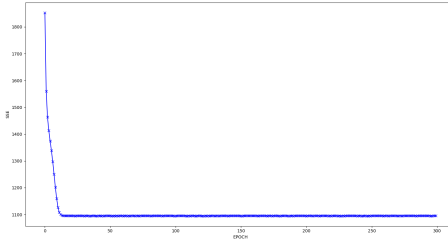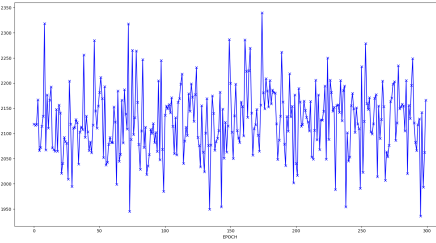
Figure 18: Learning_rate = 0.001



Figure 19: Learning_rate = 1

From the graph with a learning rate of 0.01, we can see exactly what was expected. With a too low learning rate, the equation will simply equal the following:

w^{i*} = w^{i*}

This means that, whatever data points are chosen, the prototypes are simply not updated, or sufficiently updated, and as such, the clustering will equal the random selection was was initially done before the first epoch. This translates to a flat line for the Hvq value as nothing changes.

If the value of the learning rate is too high (in this case 1), the Hvq will continuously spike as the prototypes will continuously be updated to the position of the data point they are compared to. As such, the clustering will not really occur, in the sense that the final prototype values will simply equal the coordinates of the data points processed last as a result of the randomization.

Beyond the extremes, we found that a good learning rate value, n, seemed to be located in the range between 0.05 and 0.1. With this value, or range of values, we found that the clustering seemed optimal, in the sense that the final prototype placements were in the center of the respective clusters and that the travel path of each prototype, which is graphed further down in this report, shows a smooth and continuous shift.

**Repeat the training from at least two different initializations (one of them "clever", one of them "stupid"). Discuss the observed effects.**

Taking a learning rate of 0.05 and a K value of 4, we took what we considered to be a "clever" and a "stupid" start and plotted the trajectories of each prototype.
**Clever:** For the clever start, we positioned the prototypes initially close to their final value and observed that small incremental shifts caused the

13

points to get closer and closer to the center. We found that the Hvq values were significantly lower with the clever start on initialization, due to the euclidean distances between each data point and matching prototype being smaller.

We can see in the graph below that the trajectories spend most of the time shifting around the "ideal" center, something which is caused by the fact that we are using a static learning rate. This will be further addressed below.



Figure 20: Plot Of Prototype Trajectories with Clever Start. Yellow = Initial Position

**Stupid:** The stupid start was to position all 4 prototypes in the right cluster.

When observing their trajectories, we can immediately see that the first update for the 2 prototypes closest to the left consists of both being immediately shifted into the cluster they will remain in.

The $H_{VQ}$ values we observed were higher in this case as the difference first iteration yielded a big euclidean difference for 2 of the prototypes.

Figure 21: Plot Of Prototype Trajectories with Stupid Start. Yellow = Initial Position

As a conclusion to this question, while a clever start is a very useful starting point, we experimented further with the learning curve and found a better solution which yielded lower Hvq values over time.

This solution consisted of adapting the learning rate at every iteration, decreasing it at each epoch, as at each time unit, the prototypes get closer to their ideal centering and hence do not require as high of a translation. This solution however, was not expected of us, and shall not be discussed any further in this report.

## 2.3 Bonus - Elbow Method

For the bonus question we have chosen to try and use the elbow method in order to find the most suitable K for our data set. Below we plotted $H_{VQ}$ as a function of K.

Figure 22: Elbow Method Plot

From the graph, we can clearly observe that the "elbow" is at K = 4. It became apparent the moment we plotted the $H_{VQ}$ graph, that our initial hypothesis of K = 2 being the optimal number of clusters was incorrect and that K = 4 was the better number of prototypes.

Despite K = 4 being the optimal number however, we can observe that K = 2 is not far off in terms of efficiency and when comparing the $H_{VQ}$ values between each, while there is a difference, it is noticeably smaller than the difference between K = 4 and other value than 2.

# 3 Appendix

## 3.1 Code

### 3.1.1 DBSCAN

```python
# Create simple class used as Enum
class Status:
    def __init__(self):
        pass

    New, Visited, Noise = range(3)


def distance(a, b):
    return np.linalg.norm(a - b)
```

```python
class DBscan():
    def __init__(self, points, eps=0.15, min_pts=2):
        self.points = points
        # Maximum radius of a point in order to be regarded as
        ↪   neighbour
        self.eps = eps
        # Minimum number of neighbours that a point needs to have
        ↪   in order for it's density
        # to be regarded as "high density"
        self.min_pts = min_pts
        # Number of samples / points
        self.n_pts = len(points)
        # List of labels
        self.labels = [0] * self.n_pts
        # Property Lists
        self.status_list = [Status.New] * self.n_pts
        self.member_list = [False] * self.n_pts
        self.distances = distance_matrix(self.points, self.points)

    def fit(self):
        # Initial cluster label
        cluster = 0
        for i in range(self.n_pts):
            if self.status_list[i] == Status.Visited:
                continue
            self.status_list[i] = Status.Visited
            neighbours = self.region_query(i)
            if len(neighbours) < self.min_pts:
                self.status_list[i] = Status.Noise
            else:
                cluster = cluster + 1
                self.add_to_cluster(cluster, i)
                self.expand_cluster(neighbours, cluster)

        # Set noise points in the label list as the value -1
        for i, label in enumerate(self.labels):
            if self.status_list[i] == Status.Noise:
                self.labels[i] = -1

    # Compute neighbourhood of a core point
    def region_query(self, i):
        neighbours = set()
        # Check core point's distances to all points
        for j in range(len(self.points)):
            if self.distances[i][j] <= self.eps:
                # If found point distances is in the radius
                # of the core point, add it to neighbourhood
                neighbours.add(j)
        return neighbours
```

```
60
61     def expand_cluster(self, neighbours, cluster):
62         while neighbours:
63             index = neighbours.pop()
64             if self.status_list[index] != Status.Visited:
65                 # If Noise or Undefined, transform the point into a
                   ↪  border
66                 self.status_list[index] = Status.Visited
67                 extended_neighbours = self.region_query(index)
68                 if len(extended_neighbours) >= self.min_pts:
69                     neighbours.update(extended_neighbours)
70             if not self.member_list[index]:
71                 self.add_to_cluster(cluster, index)
72
73     def add_to_cluster(self, cluster, i):
74         self.labels[i] = cluster
75         self.member_list[i] = True
76
77     def evaluate(self, y_true):
78         # Get the values of the confusion matrix
79         tn, fp, fn, tp = confusion_matrix(y_true,
            ↪  self.labels).ravel()
80         # Compute Accuracy score.
81         accuracy = float((tp + tn)/(tp + tn + fp + fn))
82         # Precision and recall will help calculate the F1 score.
83         precision = float(tp/(tp + fp))
84         recall = float(tp/(tp + fn))
85         # Print the results.
86         print('Accuracy: %f' % accuracy)
87         print('F1 Score: %f' % float((2 * precision *
            ↪  recall)/(precision+recall)))
```

### 3.1.2   Toy Dataset Code

```
1   from AssignmentAlgorithms import DBscan
2   import numpy as np
3   from sklearn.metrics import silhouette_score
4   import matplotlib.pyplot as plt
5   from sklearn.preprocessing import StandardScaler
6
7
8   def plot_data(data, labels):
9       fig, ax = plt.subplots()
10
11      scatter = ax.scatter(data[:, 0], data[:, 1], c=labels,
        ↪  cmap='nipy_spectral', s=15)
12      legend = ax.legend(*scatter.legend_elements())
```

```
13        ax.add_artist(legend)
14        plt.show()
15
16
17  if __name__ == "__main__":
18        X = np.load('toy_set.npy')
19        X = StandardScaler().fit_transform(X)
20
21        eps_values = [0.15, 0.2, 0.25, 0.3, 0.32, 0.35, 0.4, 0.45, 0.5]
22        min_pts_values = [3, 4, 5, 6, 7, 8, 9, 10]
23        for epsilon in eps_values:
24            for min_samples in min_pts_values:
25                dbscan = DBscan(X, epsilon, min_samples)
26                dbscan.fit()
27                plot_data(X, dbscan.labels)
28                print(silhouette_score(X, dbscan.labels))
```

### 3.1.3 Minst Dataset Code

```
1   from sklearn.decomposition import PCA
2   from sklearn.preprocessing import StandardScaler
3   import scipy.io as sio
4   import matplotlib.pyplot as plt
5   from AssignmentAlgorithms import DBscan
6   from sklearn.neighbors import NearestNeighbors
7   import numpy as np
8
9
10  def plot_data(data, labels):
11      fig, ax = plt.subplots()
12
13      scatter = ax.scatter(data[:, 0], data[:, 1], c=labels,
        ↪  cmap='nipy_spectral', s=20)
14      legend = ax.legend(*scatter.legend_elements())
15      ax.add_artist(legend)
16      plt.show()
17
18
19  mnist_data = sio.loadmat('mnist.mat')
20
21  # --------- Data Exploration ---------
22  X = mnist_data['X']
23  y = mnist_data['y']
24  # Standardize the data
25  X = StandardScaler().fit_transform(X)
26  # Apply PCA on given data
27  pca = PCA(n_components=2)
28  X = pca.fit_transform(X)
```

19

```
29   # Plot the data and color code using the given labels (ground
     ↪  truth)
30   plt.scatter(X[:, 0], X[:, 1], c=y, s=5)
31   plt.show()
32
33   # --------- Clustering - Outlier Detection ---------
34
35   # Find optimal value for Epsilon and plot the graph
36   neigh = NearestNeighbors(n_neighbors=4)
37   nbrs = neigh.fit(X)
38   distances, indices = nbrs.kneighbors(X)
39   distances = np.sort(distances, axis=0)
40   distances = distances[:, 1]
41   plt.xlabel('Points Sorted According to Distance of 4th Nearest
     ↪  Neighbour')
42   plt.ylabel('4th Nearest Neighbour Distance')
43   plt.grid()
44   plt.plot(distances)
45
46   # Apply algorithm using the optimal value
47   dbs = DBscan(X, eps=0.35, min_pts=20)
48   dbs.fit()
49   y_pred = dbs.labels
50   plot_data(X, y_pred)
51   # Transform -1 labels to 0 in order to mach the ground-truth data
52   for i in range(len(y_pred)):
53       if y_pred[i] == -1:
54           y_pred[i] = 0
55   # Evaluate Scores
56   dbs.evaluate(list(y))
```

### 3.1.4   VQ Code

```
1   class VQ:
2       def __init__(self, K, learning_rate, epochs):
3           self.K = K
4           self.learning_rate = learning_rate
5           self.epochs = epochs
6           self.prototypes = []
7           self.squared_errors = []
8           self.prototypes_trajectory = {}
9           self.is_fit = False
10
11      # Initialize the prototypes list with random points from the
        ↪  dataset.
12      def init_prototypes(self, data):
13          for i in range(self.K):
```

```python
14            self.prototypes.append(np.copy(data[random.randint(0,
      ↪  len(data))]))
15            self.prototypes_trajectory[i] = []
16        self.prototypes = np.array(self.prototypes)
17
18    # Used to update a winner prototype.
19    def update_prototype(self, prototype_index, point):
20        self.prototypes[prototype_index] += self.learning_rate *
      ↪  (point - self.prototypes[prototype_index])
21
22    # Plot all prototypes as red points.
23    def plot_epoch(self, X, epoch_number, trajectory=False):
24        if not self.is_fit:
25            raise Exception('Use fit before you plot')
26        if not trajectory:
27            plt.title('Positions in Epoch:%d' % epoch_number,
      ↪  fontsize=20)
28            plt.scatter(X[:, 0], X[:, 1], c='b')
29            plt.scatter(self.prototypes[:, 0], self.prototypes[:,
      ↪  1], c='r')
30            plt.show()
31        else:
32            plt.title('Prototypes trajectory', fontsize=20)
33            if epoch_number == 0:
34                plt.scatter(X[:, 0], X[:, 1], c='b')
35                plt.scatter(self.prototypes[:, 0],
      ↪  self.prototypes[:, 1], c='yellow')
36            else:
37                plt.scatter(self.prototypes[:, 0],
      ↪  self.prototypes[:, 1], c='r')
38
39            for element in self.prototypes_trajectory.values():
40                plt.plot(np.array(element)[:, 0],
      ↪  np.array(element)[:, 1])
41
42    # Plots the error curve.
43    def plot_error(self):
44        if not self.is_fit:
45            raise Exception('Use fit before you plot')
46        epochs = range(self.epochs)
47        plt.figure(figsize=(15, 8))
48        plt.plot(epochs, self.squared_errors, 'bx-')
49        plt.xlabel('EPOCH')
50        plt.ylabel('SSE')
51        plt.show()
52
53    def evaluate_winner(self, data_point):
54        winner = (0, sys.maxsize)
55        for index, prototype in enumerate(self.prototypes):
```

```
56             distance = np.linalg.norm(prototype - data_point)
57             if distance < winner[1]:
58                 winner = (index, distance)
59         return winner[0], winner[1]
60
61     # X is the desired dataset.
62     def fit(self, X, show_trajectory=False, show_plot=False):
63         self.is_fit = True
64         # Randomly initialize K prototypes with actual points from
          ↪  the dataset.
65         self.init_prototypes(X)
66         for epoch in range(self.epochs):
67             # Change data order in every epoch.
68             randomized_data = permutation(X)
69             sum_squared = 0.0
70             for i in range(len(self.prototypes)):
71                 # Keeps track of the history of prototype
                  ↪  positions.
72
                  ↪  self.prototypes_trajectory[i].append(list(self.prototypes[i]))
73             if show_trajectory:
74                 self.plot_epoch(X, epoch, show_trajectory)
75             if show_plot:
76                 self.plot_epoch(X, epoch)
77             for point in randomized_data:
78                 # Get the closest prototype to current point.
79                 winner_index, distance =
                  ↪  self.evaluate_winner(point)
80                 # Calculate Squared error.
81                 for i in range(len(point)):
82                     error = point[i] -
                      ↪  self.prototypes[winner_index][i]
83                     sum_squared += error ** 2
84                 # Move prototype towards point using the
                  ↪  learning_rate.
85                 self.update_prototype(winner_index, point)
86             self.squared_errors.append(sum_squared)
87         if show_trajectory:
88             plt.show()
```

### 3.1.5   VQ Application Code

```
1 from AssignmentAlgorithms import VQ
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import StandardScaler
4 import pandas
5
```

```python
6
7   def plot_elbow(X):
8       SSE = []
9       for k in range(1, 21):
10          vq = VQ(k, 0.1, 100)
11          vq.fit(X)
12          SSE.append(sum(vq.squared_errors))
13      K = range(1, 21)
14      plt.figure(figsize=(15, 8))
15      plt.plot(K, SSE, 'bx-')
16      plt.xlabel('K')
17      plt.ylabel('SSE')
18      plt.show()
19
20
21  if __name__ == '__main__':
22      X = pandas.read_csv('simplevqdata.csv').to_numpy()
23      X = StandardScaler().fit_transform(X)
24
25      # Plot the points and the error graph for various values of K
         ↪   and learning rate
26      K = [2, 4]
27      rates = [0.04, 0.1, 0.7]
28      for k in K:
29          for rate in rates:
30              vq = VQ(k, rate, 100)
31              vq.fit(X)
32              vq.plot_error()
33
34      # Plot the K vs SSE graph in order to find the most optimal K
35      plot_elbow(X)
```