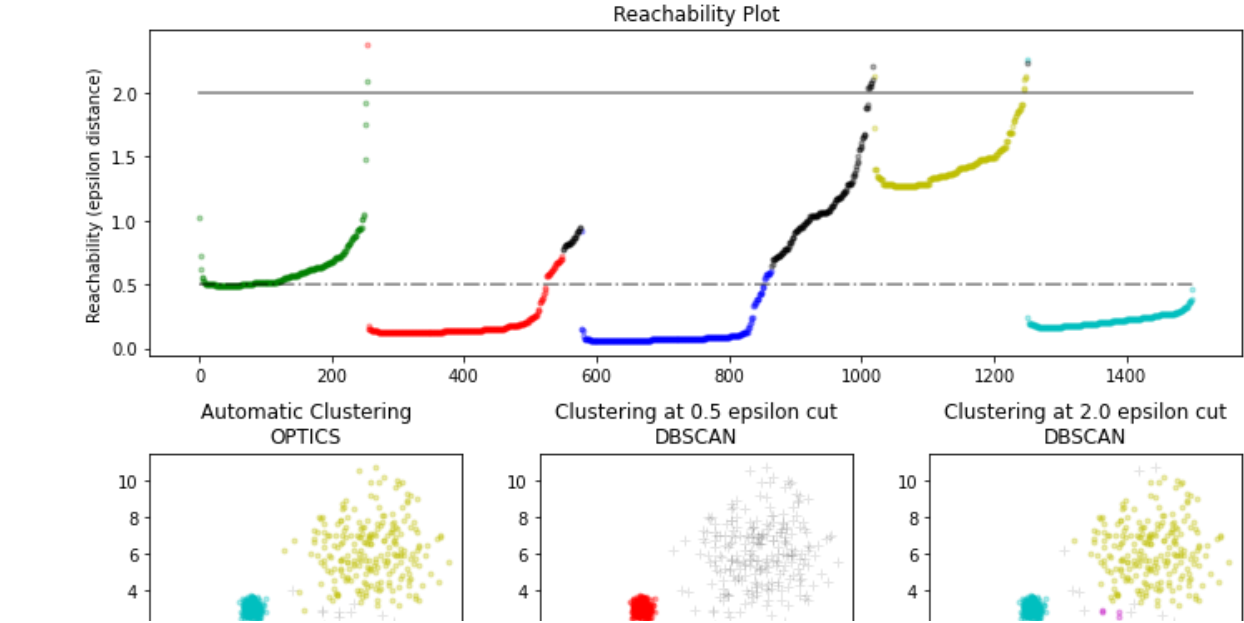Ordering Points To Identify Clustering Structure (OPTICS) is a clustering algorithm that is an improvement of the DBSCAN algorithm. OPTICS can find clusters of varying density as well, which DBSCAN was not able to do due to fixed "eps".

+ Code    + Text

Example 1

```python
# Authors: Shane Grigsby <refuge@rocktalus.com>
#          Adrin Jalali <adrin.jalali@gmail.com>
# License: BSD 3 clause


from sklearn.cluster import OPTICS, cluster_optics_dbscan
import matplotlib.gridspec as gridspec
import matplotlib.pyplot as plt
import numpy as np

# Generate sample data

np.random.seed(0)
n_points_per_cluster = 250

C1 = [-5, -2] + .8 * np.random.randn(n_points_per_cluster, 2)
C2 = [4, -1] + .1 * np.random.randn(n_points_per_cluster, 2)
C3 = [1, -2] + .2 * np.random.randn(n_points_per_cluster, 2)
C4 = [-2, 3] + .3 * np.random.randn(n_points_per_cluster, 2)
C5 = [3, -2] + 1.6 * np.random.randn(n_points_per_cluster, 2)
C6 = [5, 6] + 2 * np.random.randn(n_points_per_cluster, 2)
X = np.vstack((C1, C2, C3, C4, C5, C6))

clust = OPTICS(min_samples=50, xi=.05, min_cluster_size=.05)

# Run the fit
clust.fit(X)

labels_050 = cluster_optics_dbscan(reachability=clust.reachability_,
                                   core_distances=clust.core_distances_,
                                   ordering=clust.ordering_, eps=0.5)
labels_200 = cluster_optics_dbscan(reachability=clust.reachability_,
                                   core_distances=clust.core_distances_,
                                   ordering=clust.ordering_, eps=2)

space = np.arange(len(X))
reachability = clust.reachability_[clust.ordering_]
labels = clust.labels_[clust.ordering_]

plt.figure(figsize=(10, 7))
G = gridspec.GridSpec(2, 3)
ax1 = plt.subplot(G[0, :])
ax2 = plt.subplot(G[1, 0])
ax3 = plt.subplot(G[1, 1])
ax4 = plt.subplot(G[1, 2])

# Reachability plot
colors = ['g.', 'r.', 'b.', 'y.', 'c.']
for klass, color in zip(range(0, 5), colors):
    Xk = space[labels == klass]
    Rk = reachability[labels == klass]
    ax1.plot(Xk, Rk, color, alpha=0.3)
ax1.plot(space[labels == -1], reachability[labels == -1], 'k.', alpha=0.3)
ax1.plot(space, np.full_like(space, 2., dtype=float), 'k-', alpha=0.5)
ax1.plot(space, np.full_like(space, 0.5, dtype=float), 'k-.', alpha=0.5)
ax1.set_ylabel('Reachability (epsilon distance)')
ax1.set_title('Reachability Plot')

# OPTICS
colors = ['g.', 'r.', 'b.', 'y.', 'c.']
for klass, color in zip(range(0, 5), colors):
    Xk = X[clust.labels_ == klass]
    ax2.plot(Xk[:, 0], Xk[:, 1], color, alpha=0.3)
ax2.plot(X[clust.labels_ == -1, 0], X[clust.labels_ == -1, 1], 'k+', alpha=0.1)
ax2.set_title('Automatic Clustering\nOPTICS')

# DBSCAN at 0.5
colors = ['g', 'greenyellow', 'olive', 'r', 'b', 'c']
for klass, color in zip(range(0, 6), colors):
    Xk = X[labels_050 == klass]
    ax3.plot(Xk[:, 0], Xk[:, 1], color, alpha=0.3, marker='.')
ax3.plot(X[labels_050 == -1, 0], X[labels_050 == -1, 1], 'k+', alpha=0.1)
ax3.set_title('Clustering at 0.5 epsilon cut\nDBSCAN')

# DBSCAN at 2.
colors = ['g.', 'm.', 'y.', 'c.']
for klass, color in zip(range(0, 4), colors):
    Xk = X[labels_200 == klass]
    ax4.plot(Xk[:, 0], Xk[:, 1], color, alpha=0.3)
ax4.plot(X[labels_200 == -1, 0], X[labels_200 == -1, 1], 'k+', alpha=0.1)
ax4.set_title('Clustering at 2.0 epsilon cut\nDBSCAN')

plt.tight_layout()
plt.show()
```

## Reachability Plot



Example 2

```
1   from sklearn.cluster import OPTICS
2   import numpy as np
3   X = np.array([[1, 2], [2, 5], [3, 6],
4                 [8, 7], [8, 8], [7, 3]])
5   clustering = OPTICS(min_samples=2).fit(X)
6   clustering.labels_

    array([0, 0, 0, 1, 1, 1])
```

Example 3

```
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   from matplotlib import gridspec
5   from sklearn.cluster import OPTICS, cluster_optics_dbscan
6   from sklearn.preprocessing import normalize, StandardScaler
```

```
1   X = pd.read_csv('/content/Mall_Customers.csv')
2
3   # Dropping irrelevant columns
4   drop_features = ['CustomerID', 'Gender']
5   X = X.drop(drop_features, axis = 1)
6
7   # Handling the missing values if any
8   X.fillna(method ='ffill', inplace = True)
9
10  X.head()
```

|   | Age | Annual Income (k$) | Spending Score (1-100) |
|---|-----|--------------------|------------------------|
| 0 | 19  | 15                 | 39                     |
| 1 | 21  | 15                 | 81                     |
| 2 | 20  | 16                 | 6                      |
| 3 | 23  | 16                 | 77                     |
| 4 | 31  | 17                 | 40                     |

```
1   # Scaling the data to bring all the attributes to a comparable level
2   scaler = StandardScaler()
3   X_scaled = scaler.fit_transform(X)
4
5   # Normalizing the data so that the data
6   # approximately follows a Gaussian distribution
7   X_normalized = normalize(X_scaled)
8
9   # Converting the numpy array into a pandas DataFrame
10  X_normalized = pd.DataFrame(X_normalized)
11
12  # Renaming the columns
13  X_normalized.columns = X.columns
14
15  X_normalized.head()
```

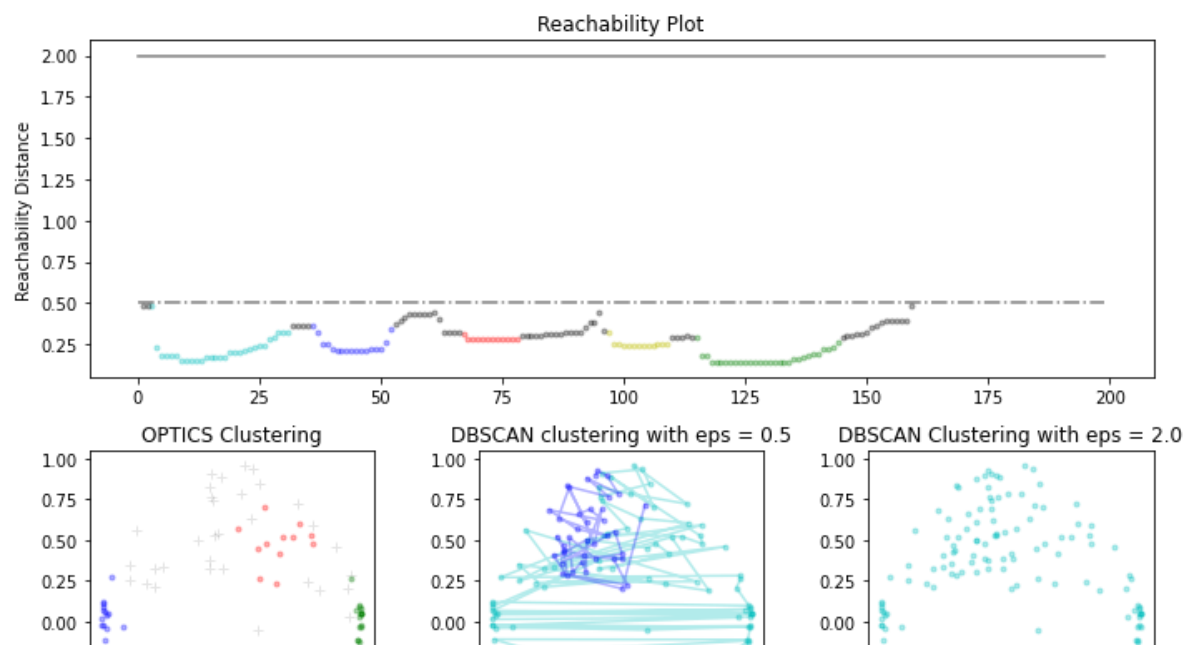|   | Age | Annual Income (k$) | Spending Score (1-100) |
|---|-----|--------------------|------------------------|
| 0 | -0.622173 | -0.759499 | -0.189897 |
| 1 | -0.518894 | -0.704396 | 0.484330 |
| 2 | -0.488556 | -0.614244 | -0.619691 |
| 3 | -0.495541 | -0.740949 | 0.453247 |
| 4 | -0.313049 | -0.923896 | -0.220036 |

```
1   # Building the OPTICS Clustering model
2   optics_model = OPTICS(min_samples = 10, xi = 0.05, min_cluster_size = 0.05)
3
4   # Training the model
5   optics_model.fit(X_normalized)
```

```
OPTICS(algorithm='auto', cluster_method='xi', eps=None, leaf_size=30,
       max_eps=inf, metric='minkowski', metric_params=None,
       min_cluster_size=0.05, min_samples=10, n_jobs=None, p=2,
       predecessor_correction=True, xi=0.05)
```

```
1   # Producing the labels according to the DBSCAN technique with eps = 0.5
2   labels1 = cluster_optics_dbscan(reachability = optics_model.reachability_,
3                                   core_distances = optics_model.core_distances_,
4                                   ordering = optics_model.ordering_, eps = 0.5)
5
6   # Producing the labels according to the DBSCAN technique with eps = 2.0
7   labels2 = cluster_optics_dbscan(reachability = optics_model.reachability_,
8                                   core_distances = optics_model.core_distances_,
9                                   ordering = optics_model.ordering_, eps = 2)
10
11  # Creating a numpy array with numbers at equal spaces till
12  # the specified range
13  space = np.arange(len(X_normalized))
14
15  # Storing the reachability distance of each point
16  reachability = optics_model.reachability_[optics_model.ordering_]
17
18  # Storing the cluster labels of each point
19  labels = optics_model.labels_[optics_model.ordering_]
20
21  print(labels)
```

```
[-1 -1 -1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
  0  0  0  0  0  0  0  0 -1 -1 -1 -1  1  1  1  1  1  1  1  1  1  1  1  1
  1  1  1  1  1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1  2  2  2  2  2
  2  2  2  2  2  2  2 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
 -1  3  3  3  3  3  3  3  3  3  3  3  3  3 -1 -1 -1 -1 -1  4  4  4  4  4
  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4  4
  4 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1  5  5  5  5  5  5  5  5  5
  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5  5
  5  5  5  5  5  5  5  5]
```

```
1   # Defining the framework of the visualization
2   plt.figure(figsize =(10, 7))
3   G = gridspec.GridSpec(2, 3)
4   ax1 = plt.subplot(G[0, :])
5   ax2 = plt.subplot(G[1, 0])
6   ax3 = plt.subplot(G[1, 1])
7   ax4 = plt.subplot(G[1, 2])
8
9   # Plotting the Reachability-Distance Plot
10  colors = ['c.', 'b.', 'r.', 'y.', 'g.']
11  for Class, colour in zip(range(0, 5), colors):
12      Xk = space[labels == Class]
13      Rk = reachability[labels == Class]
14      ax1.plot(Xk, Rk, colour, alpha = 0.3)
15  ax1.plot(space[labels == -1], reachability[labels == -1], 'k.', alpha = 0.3)
16  ax1.plot(space, np.full_like(space, 2., dtype = float), 'k-', alpha = 0.5)
17  ax1.plot(space, np.full_like(space, 0.5, dtype = float), 'k-.', alpha = 0.5)
18  ax1.set_ylabel('Reachability Distance')
19  ax1.set_title('Reachability Plot')
20
21  # Plotting the OPTICS Clustering
22  colors = ['c.', 'b.', 'r.', 'y.', 'g.']
23  for Class, colour in zip(range(0, 5), colors):
24      Xk = X_normalized[optics_model.labels_ == Class]
25      ax2.plot(Xk.iloc[:, 0], Xk.iloc[:, 1], colour, alpha = 0.3)
26
27  ax2.plot(X_normalized.iloc[optics_model.labels_ == -1, 0],
28           X_normalized.iloc[optics_model.labels_ == -1, 1],
29           'k+', alpha = 0.1)
30  ax2.set_title('OPTICS Clustering')
31
32  # Plotting the DBSCAN Clustering with eps = 0.5
33  colors = ['c', 'b', 'r', 'y', 'g', 'greenyellow']
34  for Class, colour in zip(range(0, 6), colors):
35      Xk = X_normalized[labels1 == Class]
36      ax3.plot(Xk.iloc[:, 0], Xk.iloc[:, 1], colour, alpha = 0.3, marker ='.')
37
38  ax3.plot(X_normalized.iloc[labels1 == -1, 0],
39           X_normalized.iloc[labels1 == -1, 1],
40           'k+', alpha = 0.1)
41  ax3.set_title('DBSCAN clustering with eps = 0.5')
42
43  # Plotting the DBSCAN Clustering with eps = 2.0
44  colors = ['c.', 'y.', 'm.', 'g.']
45  for Class, colour in zip(range(0, 4), colors):
46      Xk = X_normalized.iloc[labels2 == Class]
47      ax4.plot(Xk.iloc[:, 0], Xk.iloc[:, 1], colour, alpha = 0.3)
48
49  ax4.plot(X_normalized.iloc[labels2 == -1, 0],
50           X_normalized.iloc[labels2 == -1, 1],
51           'k+', alpha = 0.1)
52  ax4.set_title('DBSCAN Clustering with eps = 2.0')
53
54
55  plt.tight_layout()
56  plt.show()
```
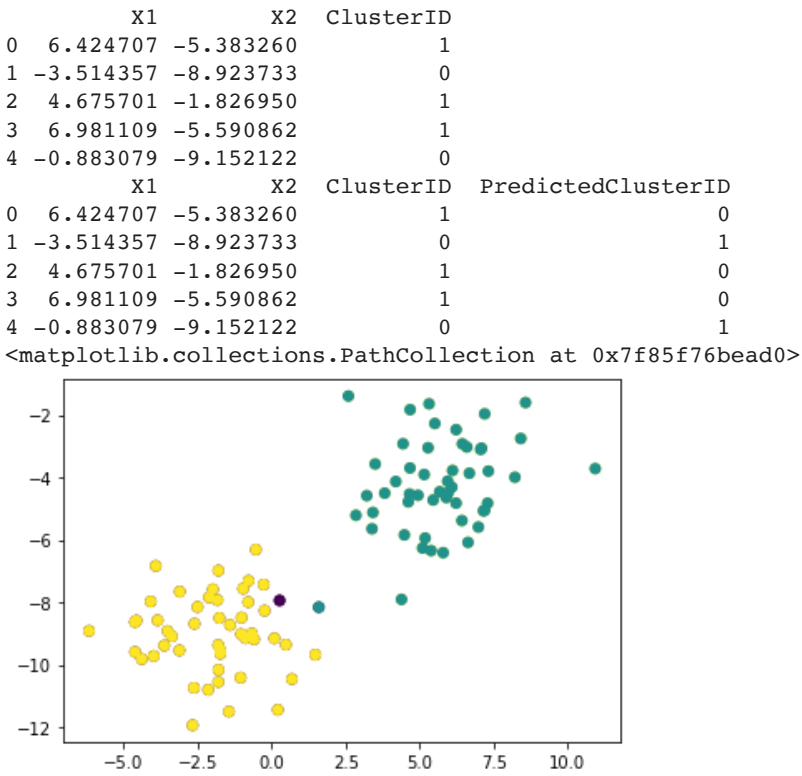
Reachability Plot

OPTICS Clustering    DBSCAN clustering with eps = 0.5    DBSCAN Clustering with eps = 2.0

### Example 4

```python
# Sample code to create OPTICS Clustering in Python
# Creating the sample data for clustering
from sklearn.datasets import make_blobs
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd

# create sample data for clustering
SampleData = make_blobs(n_samples=100, n_features=2, centers=2, cluster_std=1.5, random_state=40)

#create np array for data points
X = SampleData[0]
y = SampleData[1]

# Creating a Data Frame to represent the data with labels
ClusterData=pd.DataFrame(list(zip(X[:,0],X[:,1],y)), columns=['X1','X2','ClusterID'])
print(ClusterData.head())

# create scatter plot to visualize the data
%matplotlib inline
plt.scatter(ClusterData['X1'], ClusterData['X2'], c=ClusterData['ClusterID'])

############################################################################
# This function is not present in python version 3.6
# Other option is pyclustering.cluster.optics but its not neat
from sklearn.cluster import OPTICS
op = OPTICS(min_samples=40, xi=0.02, min_cluster_size=0.1)

# Generating cluster id for each row using DBSCAN algorithm
ClusterData['PredictedClusterID']=op.fit_predict(X)
print(ClusterData.head())

# Plotting the predicted clusters
plt.scatter(ClusterData['X1'], ClusterData['X2'], c=ClusterData['PredictedClusterID'])
```
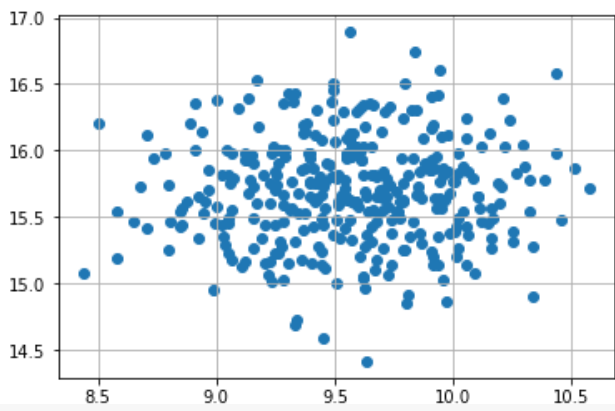
```
         X1        X2  ClusterID
0   6.424707 -5.383260          1
1  -3.514357 -8.923733          0
2   4.675701 -1.826950          1
3   6.981109 -5.590862          1
4  -0.883079 -9.152122          0
         X1        X2  ClusterID  PredictedClusterID
0   6.424707 -5.383260          1                   0
1  -3.514357 -8.923733          0                   1
2   4.675701 -1.826950          1                   0
3   6.981109 -5.590862          1                   0
4  -0.883079 -9.152122          0                   1
<matplotlib.collections.PathCollection at 0x7f85f76bead0>
```



### Example 5 **Anomaly Detection with OPTICS**

```python
from sklearn.cluster import OPTICS
from sklearn.datasets import make_blobs
from numpy import quantile, where, random
import matplotlib.pyplot as plt
random.seed(123)
x, _ = make_blobs(n_samples=350, centers=1, cluster_std=.4, center_box=(20, 5))

plt.scatter(x[:,0], x[:,1])
plt.grid(True)
plt.show()
```

```
1  model = OPTICS().fit(x)
2  print(model)
```

```
OPTICS(algorithm='auto', cluster_method='xi', eps=None, leaf_size=30,
       max_eps=inf, metric='minkowski', metric_params=None,
       min_cluster_size=None, min_samples=5, n_jobs=None, p=2,
       predecessor_correction=True, xi=0.05)
```

We determine the scores of each sample of x data by using core_distance_ property of the model. thresh = quantile(scores, .98) print(thresh)

```
1  scores = model.core_distances_
2  thresh = quantile(scores, .98)
3  print(thresh)
```

```
0.35064484877392416
```

By using threshold value, we'll find the samples with the scores that are equal to or higher than the threshold value.

```
1  index = where(scores >= thresh)
2  values = x[index]
3  print(values)
```

```
[[ 9.45071447 14.58847433]
 [ 8.500387   16.2113985 ]
 [ 9.56481939 16.89136015]
 [ 9.63176979 14.41548797]
 [ 8.43771706 15.07302741]
 [10.33672675 14.89789167]
 [10.43533425 16.58262441]]
```

We visualize the results in a plot by highlighting the anomalies with a color.

```
1  plt.scatter(x[:,0], x[:,1])
2  plt.scatter(values[:,0],values[:,1], color='r')
3  plt.legend(("normal", "anomal"), loc="best", fancybox=True, shadow=True)
4  plt.grid(True)
5  plt.show()
```