# DATAWAREHOUSING AND DATA MINING MODEL LABORATORY EXAMINATION

NAME: KARTHIKAPRIYAA V.S

CLASS: CSE –B

ROLL NO: 18BCS047

SUBJECT CODE :U18CSI6203L

## EXPERIMENT

1**.** Download a suitable dataset for clustering from any Repository. List the attributes and its type in a word Doc.

2. Apply DBSCAN algorithm to make clusters to identify dense region using the dataset.

3. Upload in your github account. Provide the link for access.

## DATA SET DESCRIPTION
## DATASET:

 Make_blobs from sklearn.datasets -(importing dataset )

## ATTRIBUTE DESCRIPTION:

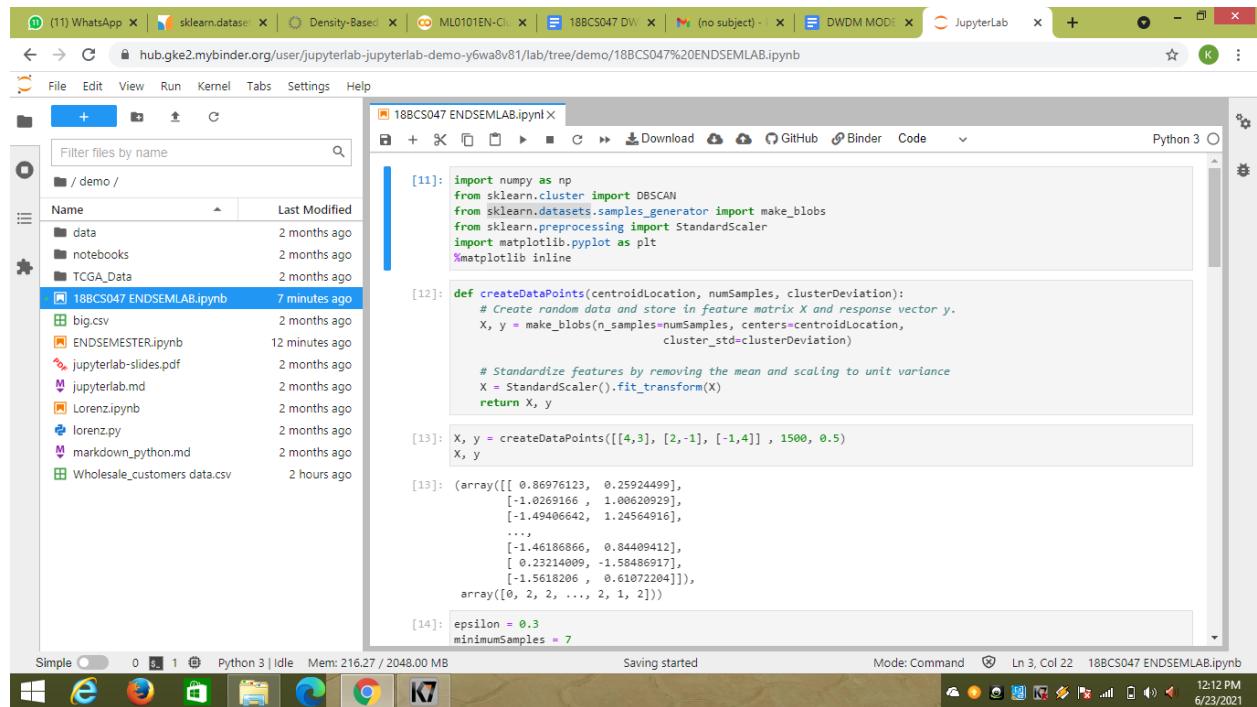*Xndarray of shape (n_samples, n_features)*

The generated samples.

*yndarray of shape (n_samples,)*

The integer labels for cluster membership of each sample.

**centers*ndarray of shape (n_centers, n_features)***

The centers of each cluster. Only returned if `return_centers=True`.

## SCREENSHOT

First screenshot code cells:

```
[14]: epsilon = 0.3
      minimumSamples = 7
      db = DBSCAN(eps=epsilon, min_samples=minimumSamples).fit(X)
      labels = db.labels_
      labels
```

```
[14]: array([0, 1, 1, ..., 1, 2, 1])
```

```
[15]: # Firts, create an array of booleans using the labels from db.
      core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
      core_samples_mask[db.core_sample_indices_] = True
      core_samples_mask
```

```
[15]: array([ True,  True,  True, ...,  True,  True,  True])
```

```
[16]: # Number of clusters in labels, ignoring noise if present.
      n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
      n_clusters_
```

```
[16]: 3
```

```
[17]: # Remove repetition in labels by turning it into a set.
      unique_labels = set(labels)
      unique_labels
```

```
[17]: {-1, 0, 1, 2}
```

```
[18]: # Create colors for the clusters.
      colors = plt.cm.Spectral(np.linspace(0, 1, len(unique_labels)))
```

Second screenshot code cells:

```
[17]: # Remove repetition in labels by turning it into a set.
      unique_labels = set(labels)
      unique_labels
```

```
[17]: {-1, 0, 1, 2}
```

```
[18]: # Create colors for the clusters.
      colors = plt.cm.Spectral(np.linspace(0, 1, len(unique_labels)))
      colors
```

```
[18]: array([[0.61960784, 0.00392157, 0.25882353, 1.        ],
             [0.99346405, 0.74771242, 0.43529412, 1.        ],
             [0.74771242, 0.89803922, 0.62745098, 1.        ],
             [0.36862745, 0.30980392, 0.63529412, 1.        ]])
```

```
[21]: # Plot the points with colors
      for k, col in zip(unique_labels, colors):
          if k == -1:
              # Black used for noise.
              col = 'k'

          class_member_mask = (labels == k)

          # Plot the datapoints that are clustered
          xy = X[class_member_mask & core_samples_mask]
          plt.scatter(xy[:, 0], xy[:, 1],s=50, c=col, marker=u'o', alpha=0.5)

          # Plot the outliers
          xy = X[class_member_mask & ~core_samples_mask]
          plt.scatter(xy[:, 0], xy[:, 1],s=50, c=col, marker=u'o', alpha=0.5)
```
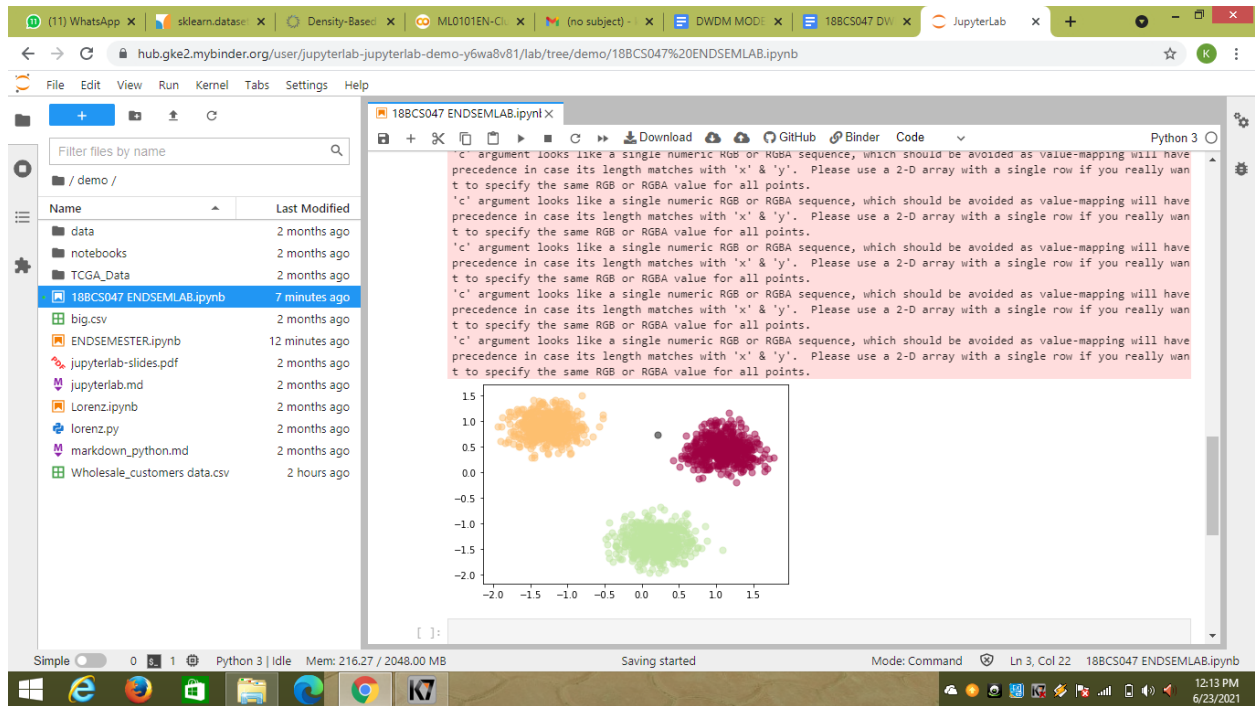
'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have

# OUTPUT

IMPLEMENTATION:

GITHUB LINK:

https://github.com/Karthika1712/END-SEMESTER-DWDM-LAB/upload