

Activity 1: *DBSCAN*

- Give a high-level explanation of the *DBSCAN* clustering algorithm.
 - Is it a supervised or unsupervised technique?
 - What does a *cluster* mean in *DBSCAN*?
 - Please also consider what a *cluster* means in *k-means*?
 - What does *noise* mean in *DBSCAN*?
 - Does this algorithm require a specification of the number of clusters in advance (like *k-means*)?
 - Does *DBSCAN* expect a specific shape of clusters?
 - If yes, what is it?
 - If not, what kind of cluster can *DBSCAN* discover?
 - What shapes of clusters *k-means* can discover?
 - What are the roles of the following two parameters in *DBSCAN*:
 - What is *Eps*?
 - What is *MinPts*?
 - How can a *k distance graph* help to guide the selection of *Eps*?
 - What are the advantages and disadvantages of the *DBSCAN* approach to clustering?

Activity 2: *DBSCAN* in *sklearn*

The following is the *DBSCAN* classifier model in *sklearn*:

`sklearn.cluster.DBSCAN(eps, min_samples, metric)`

Please discuss:

- What does *eps* mean?
- What does *min_samples* mean?
- What does *metric* mean?

Please revisit page 24 of Week 9's Lecture slides.

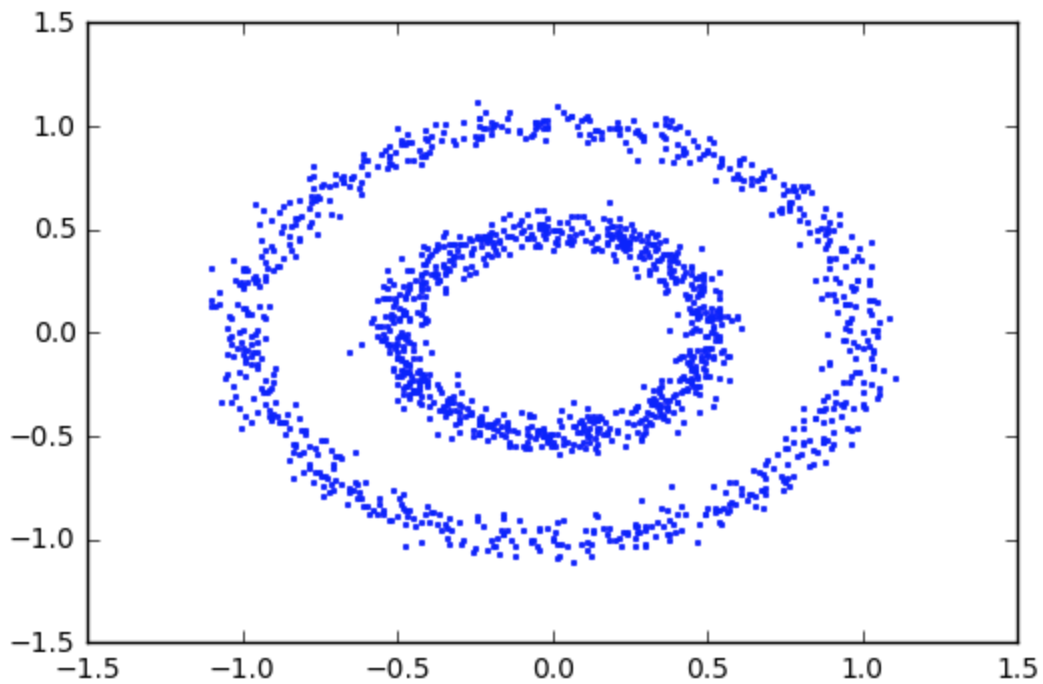
This week, let's explore the DBSCAN on several datasets with different shapes of clusters.

```
In [1]: import numpy as np
In [2]: import matplotlib.pyplot as plt
In [3]: from sklearn import cluster, datasets

In [4]: n_samples = 1500
In [5]: circles = datasets.make_circles(n_samples=n_samples, factor=.5,
                                         noise=.05)
In [6]: moons = datasets.make_moons(n_samples=n_samples, noise=.05)
In [7]: blobs = datasets.make_blobs(n_samples=n_samples, random_state=8)
```

```
Out [8]: (array([[ -0.21009776,   0.46181291],
                [  0.67629075,  -0.68385711],
                [  0.87536459,   0.5651477 ],
                ...,
                [-0.0331735 , -0.46811284],
                [-0.96187452,   0.27081349],
                [-0.21301649, -0.51512252]]), array([1, 0, 0, ..., 1, 0, 1]))
```

Out [10]:



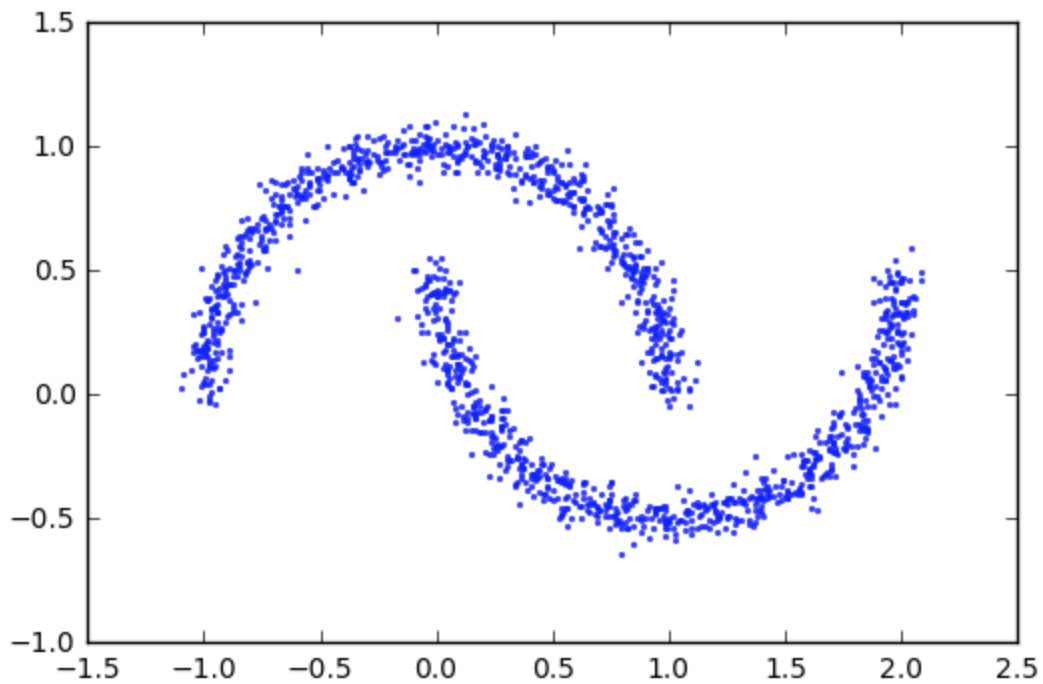
```
In [11]: moons
```

```
Out [11]: (array([[ 0.04454316,  0.0845715 ],
 [ 1.58566568, -0.26181606],
 [ 0.04062815,  0.50069292],
 ...,
 [ 0.73273691,  0.51027295],
 [ 0.57784829, -0.34139774],
 [-0.95526833, -0.03978616]]), array([1, 1, 1, ..., 0, 1, 0]))
```

```
In [12]: plt.scatter(moons[0][:, 0], moons[0][:, 1], alpha = 0.8, s= 5.0, lw=
0)
```

```
In [13]: plt.show()
```

```
Out [13]:
```



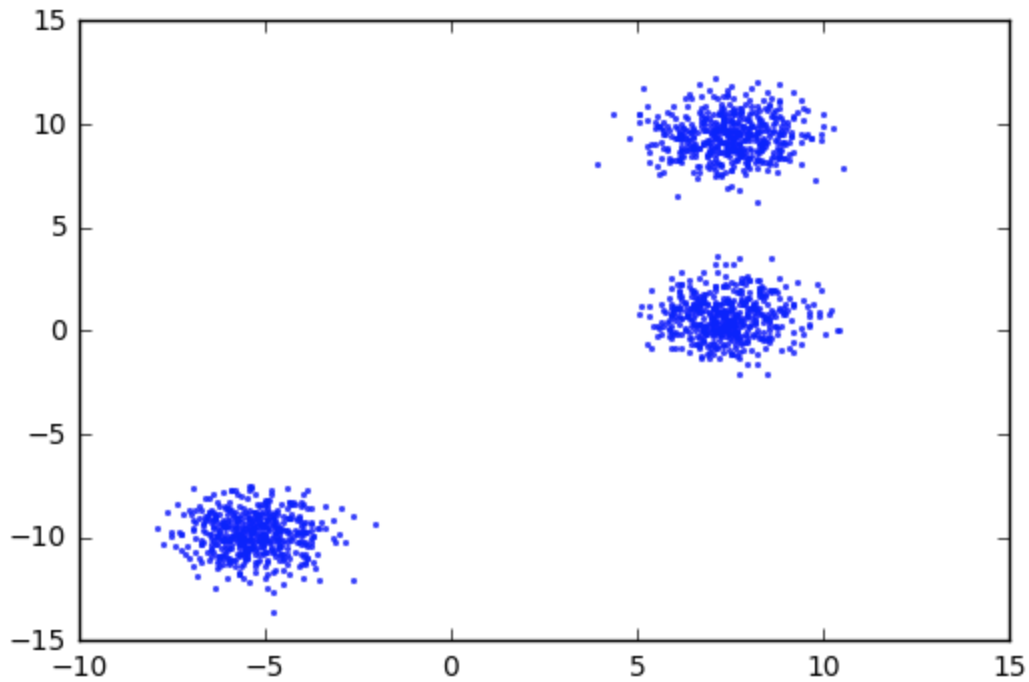
In [14]: blobs

```
Out [14]: (array([[ 5.86749807,  8.17715188],
 [ 5.61369982,  9.93295527],
 [ 7.22508428, 10.44886194],
 ...,
 [ 7.73674097, 10.82855388],
 [-4.61701094, -9.64855983],
 [-3.48640175, -9.25766922]]), array([0, 0, 0, ..., 0, 2, 2]))
```

```
In [15]: plt.scatter(blobs[0][:, 0], blobs[0][:, 1], alpha = 0.8, s= 5.0, lw=
0)
```

```
In [16]: plt.show()
```

```
Out [16]:
```



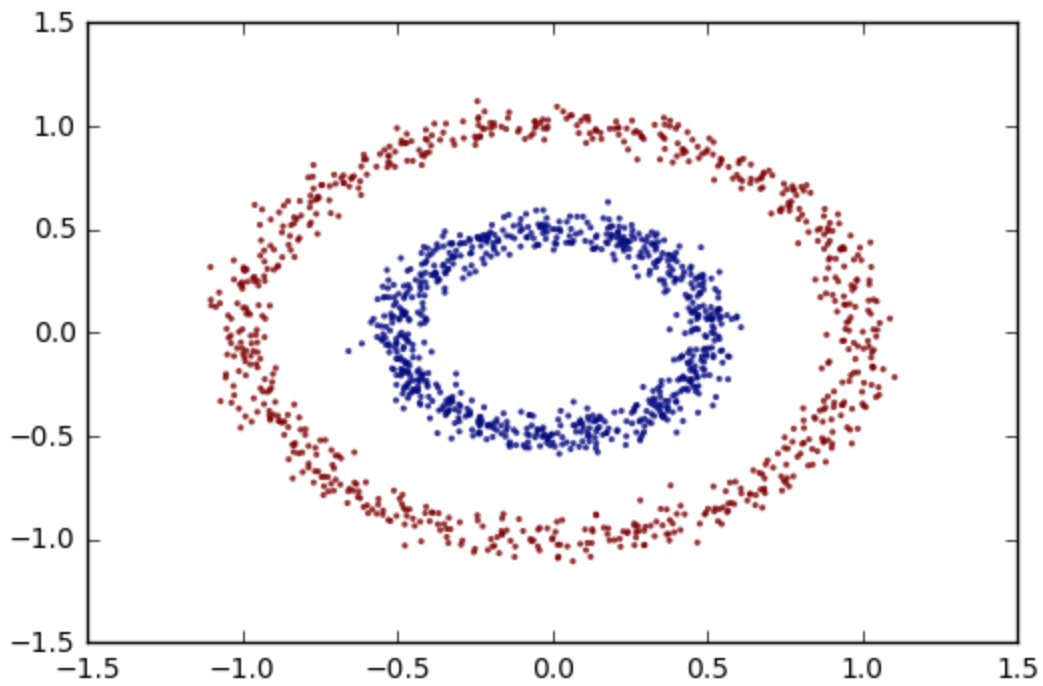
Practical exercise 2: *DBSCAN* Clustering

Next, we would like to build the *DBSCAN* model for each dataset.

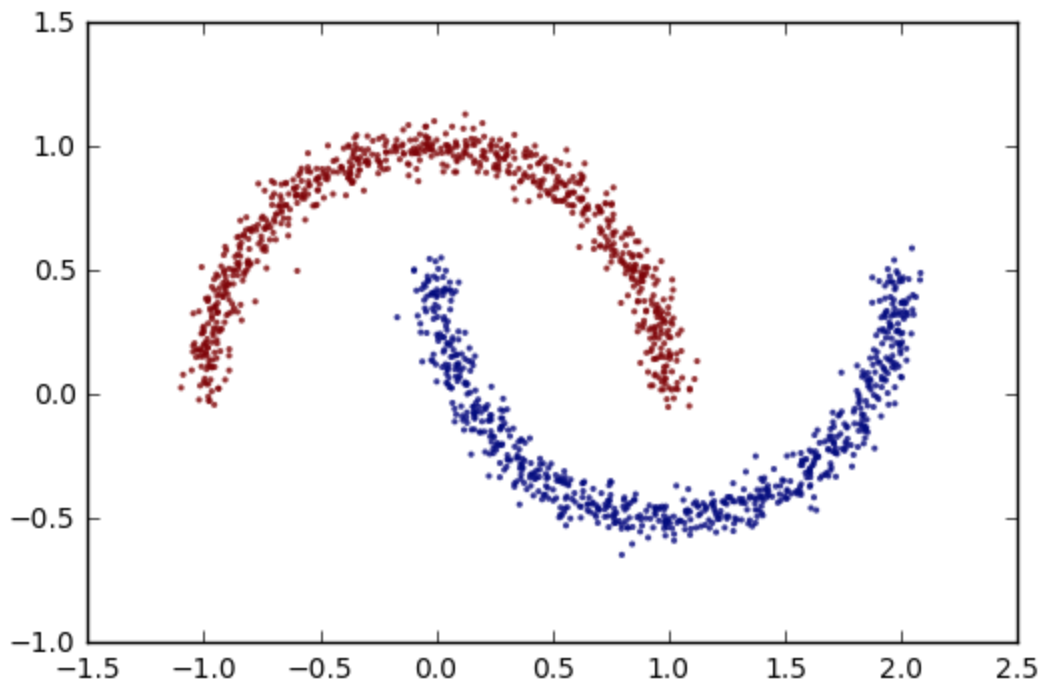
```
In [17]: dbs_1 = cluster.DBSCAN(eps=.2)
In [18]: dbs_fit = dbs_1.fit(circles[0])
```

Please think about why we use `circles[0]` here? Hint: please look at `Out [8]`.

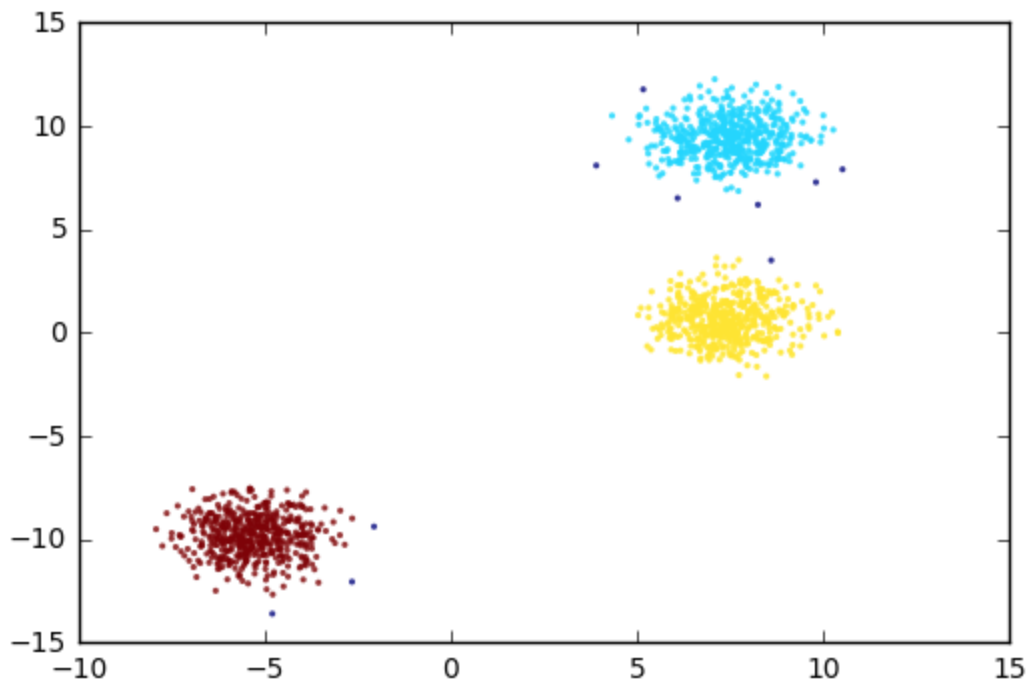
```
In [19]: labels_1 = dbs_fit.labels_
In [20]: plt.scatter(circles[0][:, 0], circles[0][:, 1], c=labels_1, alpha =
0.8, s= 5.0, lw= 0)
In [21]: plt.show()
Out [21]:
```



```
In [22]: dbs_2 = cluster.DBSCAN(eps=.2)
In [23]: dbs_fit = dbs_2.fit(moons[0])
In [24]: labels_2 = dbs_fit.labels_
In [25]: plt.scatter(moons[0][:, 0], moons[0][:, 1], c=labels_2, alpha = 0.8,
s= 5.0, lw= 0)
In [26]: plt.show()
Out [26]:
```



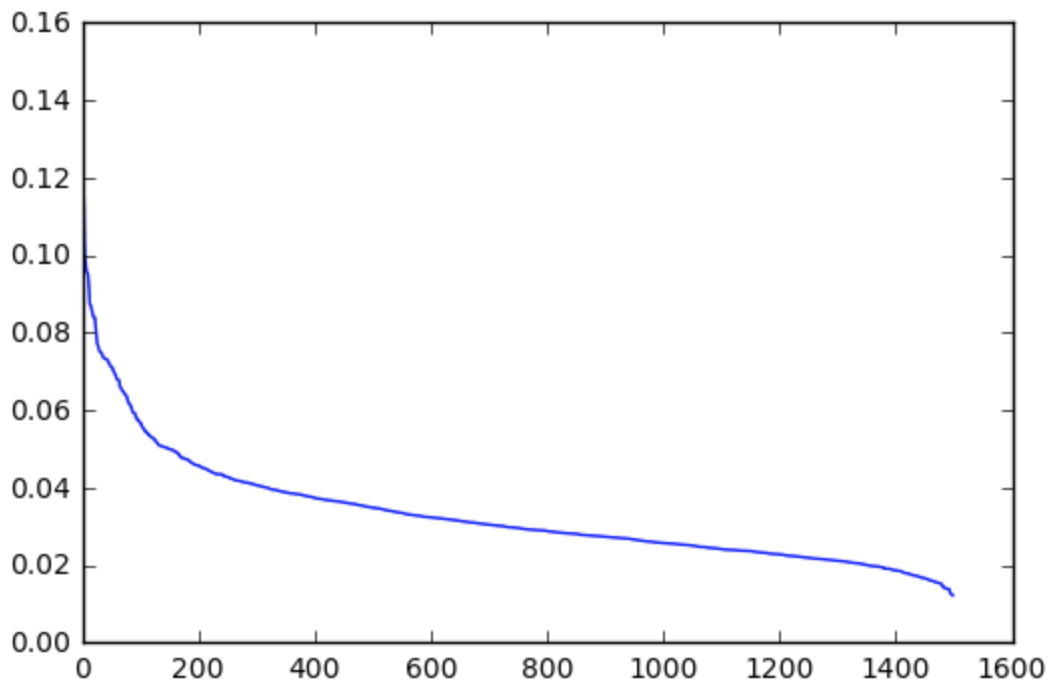
```
In [27]: dbs_3 = cluster.DBSCAN(eps=.8)
In [28]: dbs_fit = dbs_3.fit(blobs[0])
In [29]: labels_3 = dbs_fit.labels_
In [30]: plt.scatter(blobs[0][:, 0], blobs[0][:, 1], c=labels_3, alpha = 0.8,
s= 5.0, lw= 0)
In [31]: plt.show()
Out [31]:
```



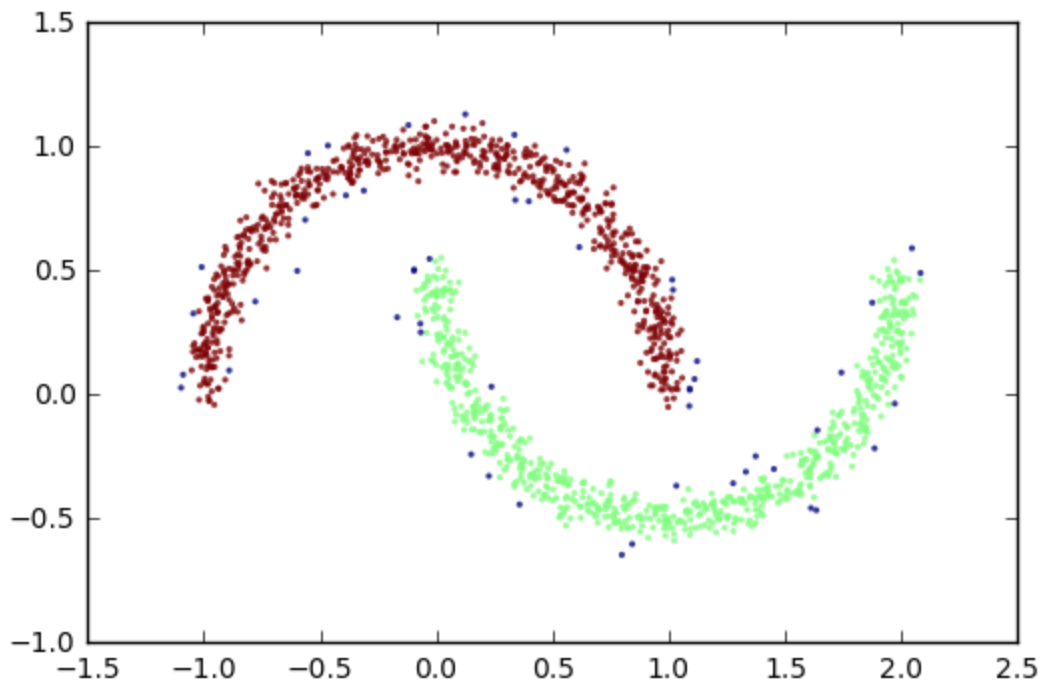
Practical exercise 3: k distance graph

We explore the k -distance graph for the moons dataset:

```
In [32]: from sklearn.neighbors import NearestNeighbors
In [33]: nbrs = NearestNeighbors().fit(moons[0])
In [34]: distances, indices = nbrs.kneighbors(moons[0], 20)
In [35]: kDis = distances[:, 4]
In [36]: kDis.sort()
In [37]: kDis = kDis[range(len(kDis)-1, 0, -1)]
In [38]: plt.plot(range(0, len(kDis)), kDis)
In [39]: plt.show()
Out [39]:
```

```
In [40]: dbs_2 = cluster.DBSCAN(eps=.05)
In [41]: dbs_fit = dbs_2.fit(moons[0])
In [42]: labels_2 = dbs_fit.labels_
In [43]: plt.scatter(moons[0][:, 0], moons[0][:, 1], c=labels_2, alpha = 0.8,
s= 5.0, lw= 0)
In [44]: plt.show()
Out [44]:
```



Please compare `Out [44]` with `Out [26]`, which was generated using the parameter $Eps = 0.2$.

- What are the differences?
- Why could this happen?

All materials copyright RMIT University. Students are welcome to download and print for the purpose of studying for this course.