

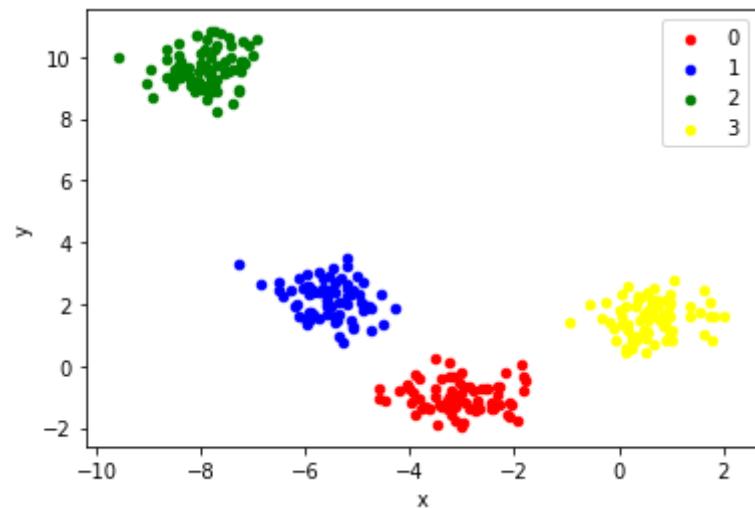
LAB 4 - CLUSTERING

Assignment-1

```
In [1]: import pandas as pd
import numpy as np
from sklearn.datasets.samples_generator import make_blobs
from pandas import DataFrame
from sklearn.cluster import KMeans
from sklearn.metrics.cluster import contingency_matrix
from sklearn.metrics import mean_squared_error
from scipy.cluster.hierarchy import dendrogram, linkage
from sklearn.cluster import DBSCAN
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings("ignore")
```

```
In [9]: X, y = make_blobs(n_samples=300, centers=4, n_features=2, cluster_std=0.6)
```

```
In [10]: df = DataFrame(dict(x=X[:,0], y=X[:,1], label=y))
colors = {0:'red', 1:'blue', 2:'green', 3:'yellow'}
fig, ax = plt.subplots()
grouped = df.groupby('label')
for key, group in grouped:
    group.plot(ax=ax, kind='scatter', x='x', y='y', label=key, color=colors[key])
plt.show()
```



```
In [11]: k = [1,2,3,4,5,6,7,8,9,10]
SSE_iter=[]
for i in k:
    kmeans = KMeans(n_clusters = i, random_state = None)
    y_pred = kmeans.fit_predict(df)
    centroids = kmeans.cluster_centers_
    print('contingency matrix:')
    print(contingency_matrix(y,y_pred))
    print('mean squared error:')
    SSE = mean_squared_error(y,y_pred)
    print(SSE)
    SSE_iter.append(SSE)
    plt.scatter(df['x'], df['y'], c= kmeans.labels_.astype(float), s=50, alpha=0.5)
    plt.scatter(centroids[:, 0], centroids[:, 1], c='red', s=50)
    plt.title(i)
    plt.show()

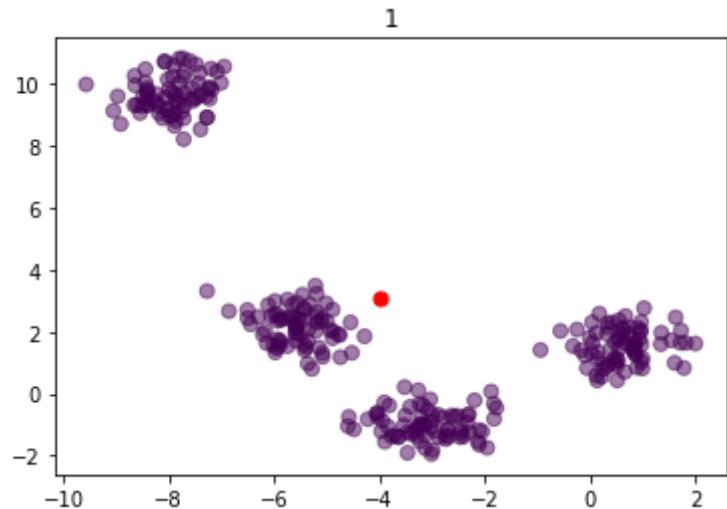
plt.plot(k,SSE_iter)
plt.show()
```

contingency matrix:

```
[[75]
 [75]
 [75]
 [75]]
```

mean squared error:

3.5

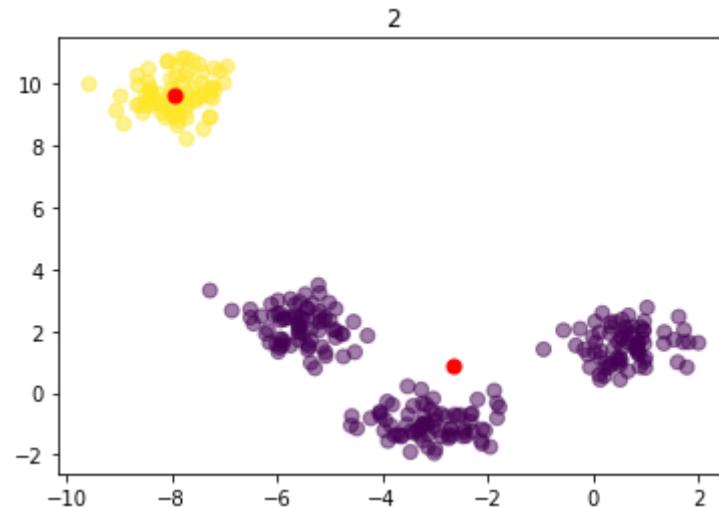


contingency matrix:

```
[[75  0]
 [75  0]
 [ 0 75]
 [75  0]]
```

mean squared error:

2.75

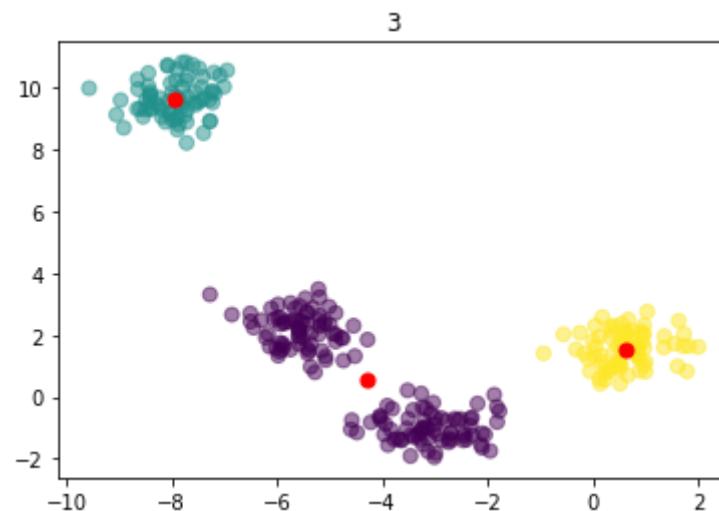


contigency matrix:

```
[[75  0  0]
 [75  0  0]
 [ 0 75  0]
 [ 0  0 75]]
```

mean squared error:

0.75

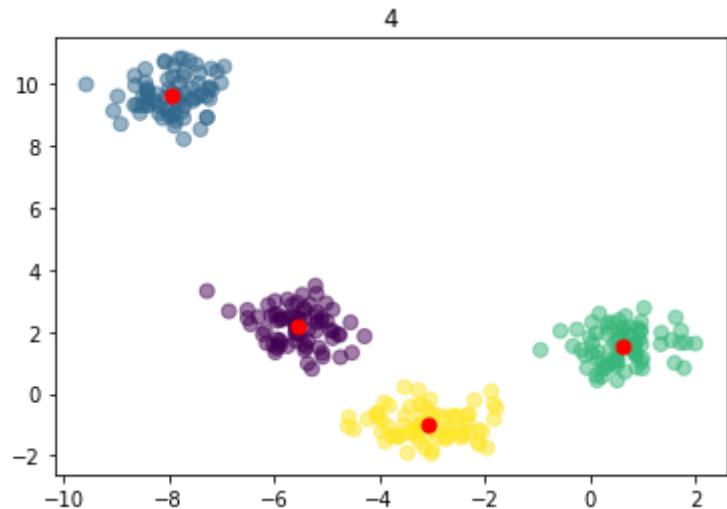


contigency matrix:

```
[[ 0  0  0 75]
 [75  0  0  0]
 [ 0 75  0  0]
 [ 0  0 75  0]]
```

mean squared error:

3.0

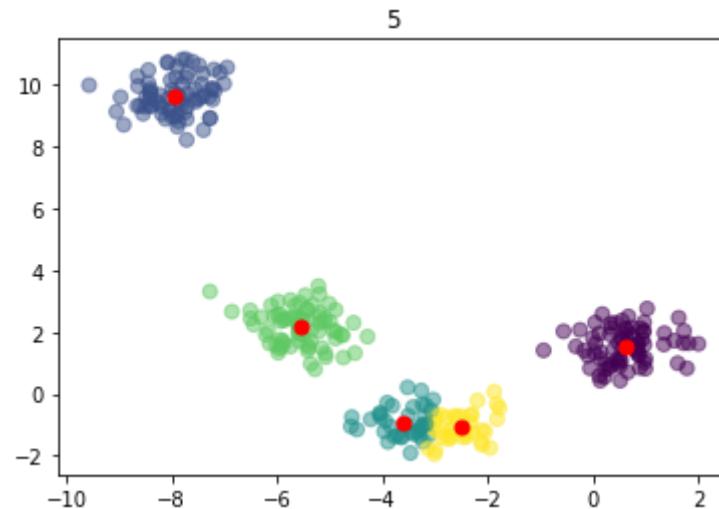


contigency matrix:

```
[[ 0  0 39  0 36]
 [ 0  0  0 75  0]
 [ 0 75  0  0  0]
 [75  0  0  0  0]]
```

mean squared error:

5.94

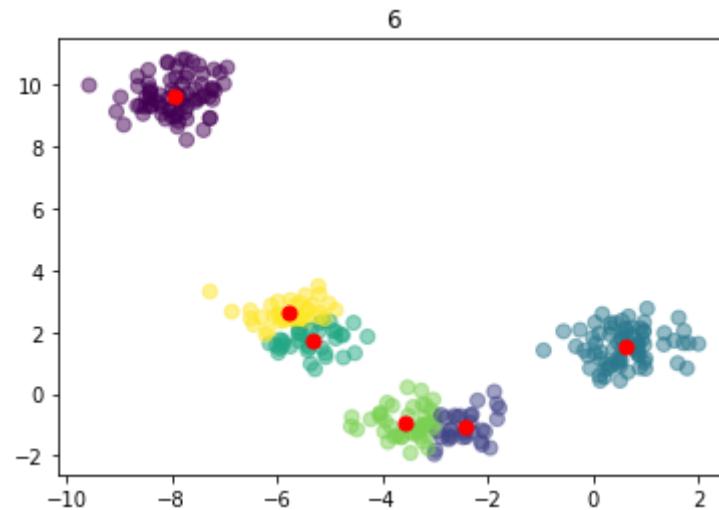


contigency matrix:

```
[[ 0 32 0 0 43 0]
 [ 0 0 0 38 0 37]
 [75 0 0 0 0 0]
 [ 0 0 75 0 0 0]]
```

mean squared error:

6.13

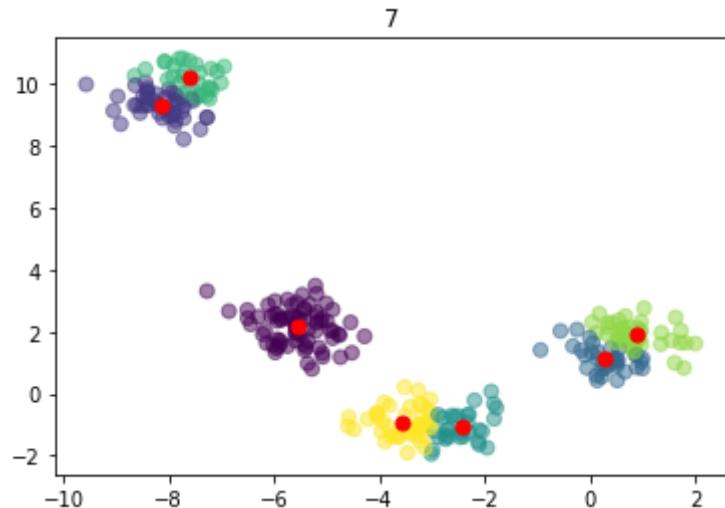


contigency matrix:

```
[[ 0  0  0 32  0  0 43]
 [75  0  0  0  0  0 0]
 [ 0 46  0  0 29  0 0]
 [ 0  0 34  0  0 41 0]]
```

mean squared error:

7.57

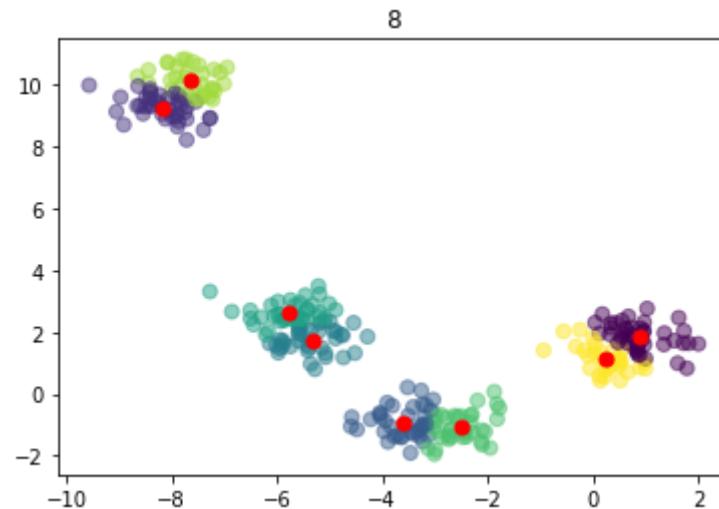


contigency matrix:

```
[[ 0  0 39  0  0 36  0  0]
 [ 0  0  0 37 38  0  0 0]
 [ 0 43  0  0  0  0 32 0]
 [44  0  0  0  0  0  0 31]]
```

mean squared error:

9.976666666666667

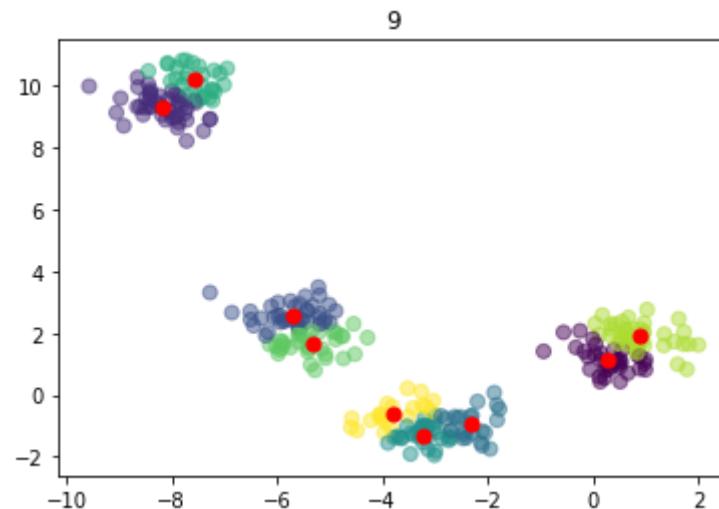


contigency matrix:

```
[[ 0  0  0 25 30  0  0  0 20]
 [ 0  0 42  0  0  0 33  0  0]
 [ 0 47  0  0  0 28  0  0  0]
 [34  0  0  0  0  0  0 41  0]]
```

mean squared error:

13.71

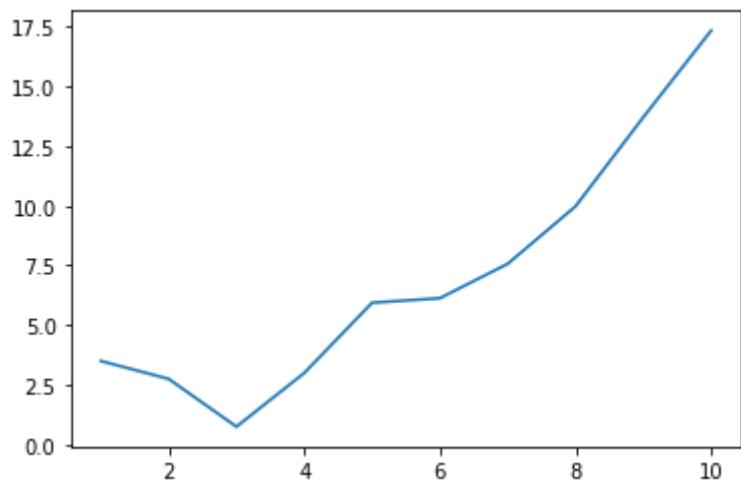
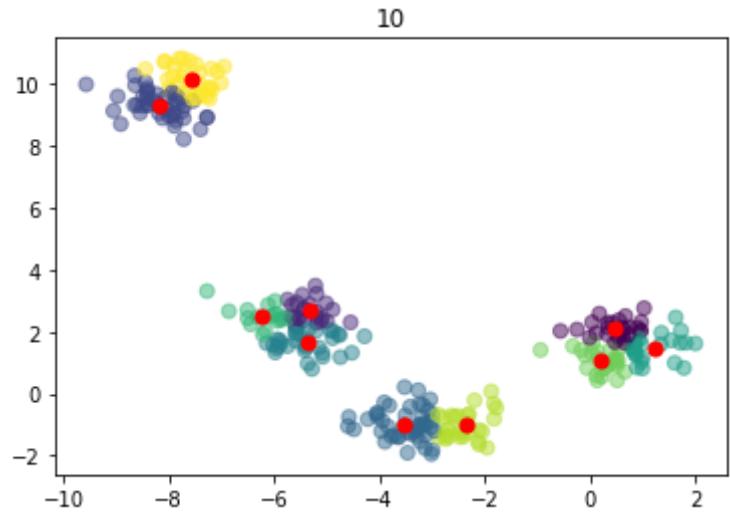


contigency matrix:

```
[[ 0  0  0 46  0  0  0  0 29  0]
 [ 0 24  0  0 34  0 17  0  0  0]
 [ 0  0 45  0  0  0  0  0  0 30]
 [26  0  0  0  0 24  0 25  0  0]]
```

mean squared error:

17.336666666666666

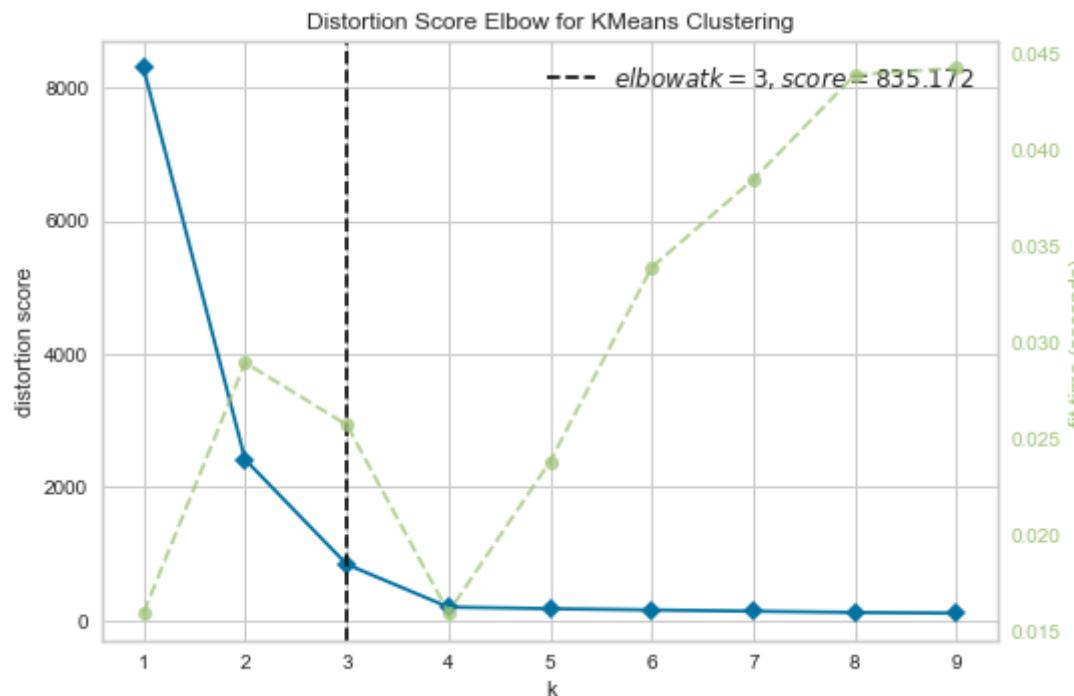


The location of a bend (knee) in the plot is generally considered as an indicator of the appropriate number of clusters., The plot does indicated the natural number of clusters is 4

```
In [17]: from yellowbrick.cluster import KElbowVisualizer
```

```
# Instantiate the clustering model and visualizer
model = KMeans()
visualizer = KElbowVisualizer(model, k=(1,10))

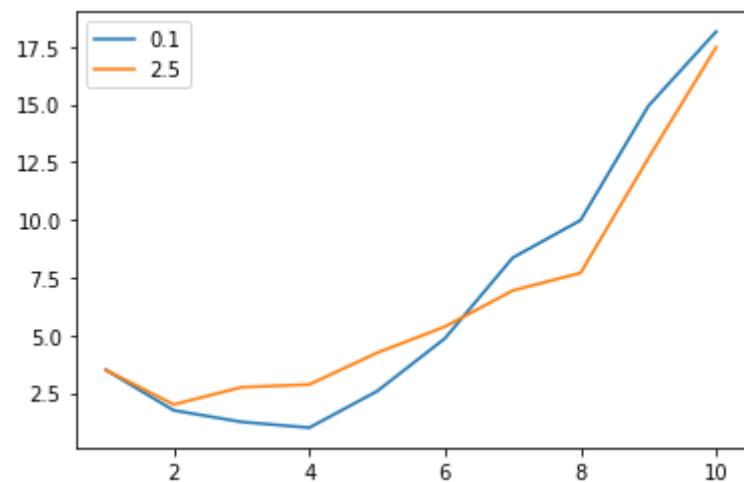
visualizer.fit(df)
visualizer.show()
```



```
Out[17]: <matplotlib.axes._subplots.AxesSubplot at 0x1982f3f44e0>
```

It is also shown by elbow method

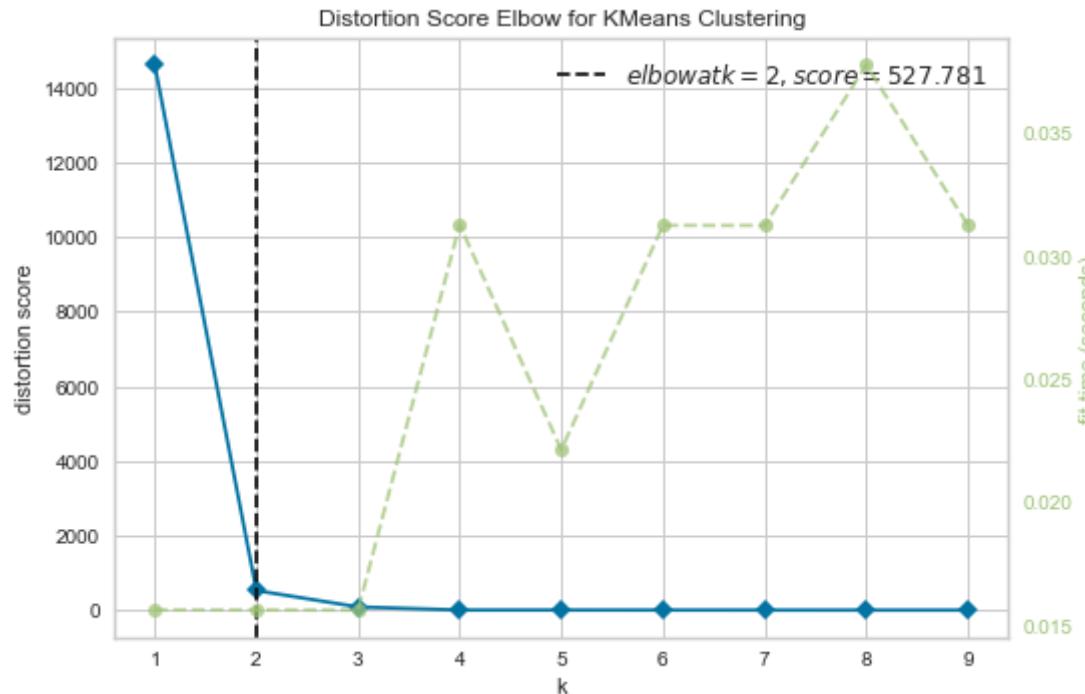
```
In [6]: # d.).  
cluster_std=[0.1,2.5]  
k = [1,2,3,4,5,6,7,8,9,10]  
  
SSE_iter=[]  
for i in cluster_std:  
    X, y = make_blobs(n_samples=300, centers=4, n_features=2, cluster_std=i)  
    df = DataFrame(dict(x=X[:,0], y=X[:,1], label=y))  
    for i in k:  
        kmeans = KMeans(n_clusters = i, random_state = None)  
        y_pred = kmeans.fit_predict(df)  
        SSE = mean_squared_error(y,y_pred)  
        SSE_iter.append(SSE)  
  
plt.plot(k,SSE_iter[0:10],label='0.1')  
plt.plot(k,SSE_iter[10:20],label='2.5')  
plt.legend(loc='upper left')  
  
plt.show()
```

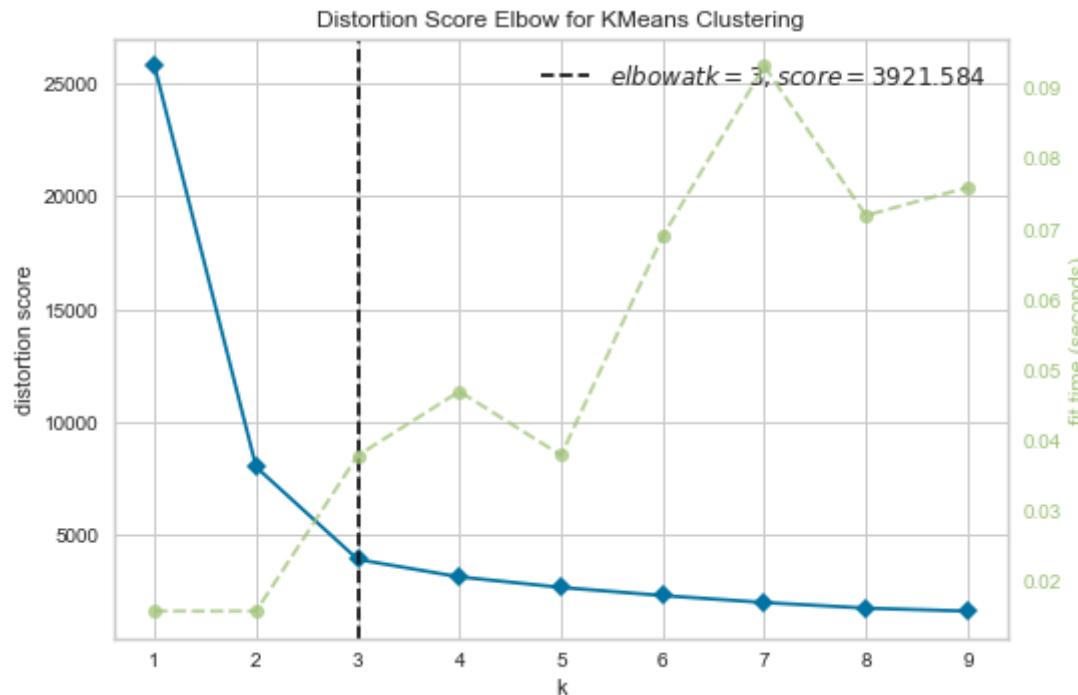


```
In [31]: from yellowbrick.cluster import KElbowVisualizer

# Instantiate the clustering model and visualizer
cluster_std=[0.1,2.5]

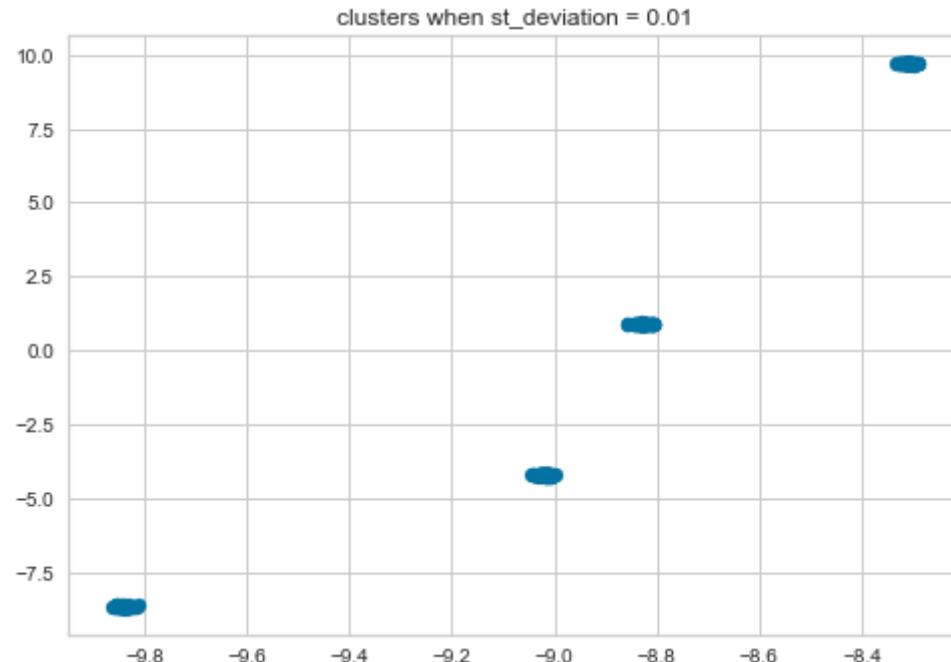
for i in cluster_std:
    A,b = make_blobs(n_samples=300, centers=4, n_features=2,cluster_std=i)
    model = KMeans()
    visualizer = KElbowVisualizer(model, k=(1,10))
    d = DataFrame(dict(x1=A[:,0], y1=A[:,1], label=b))
    visualizer.fit(d)
    visualizer.show()
```





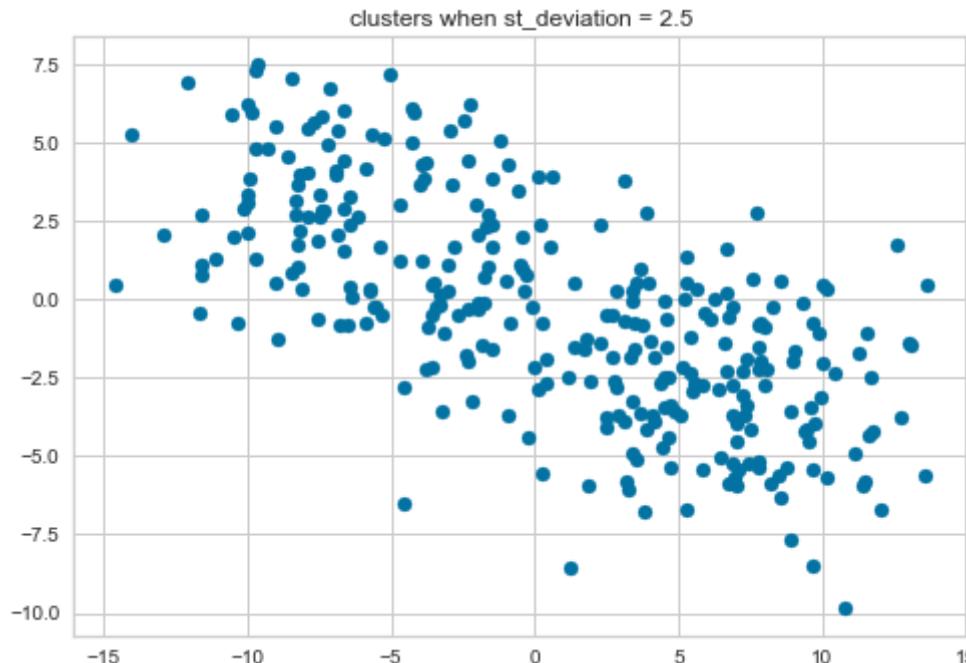
```
In [29]: A,b = make_blobs(n_samples=300, centers=4, n_features=2,cluster_std=0.01)
plt.scatter(A[:,0],A[:,1])
plt.title('clusters when st_deviation = 0.01')
```

```
Out[29]: Text(0.5, 1.0, 'clusters when st_deviation = 0.01')
```



```
In [30]: A,b = make_blobs(n_samples=300, centers=4, n_features=2,cluster_std=2.5)
plt.scatter(A[:,0],A[:,1])
plt.title('clusters when st_deviation = 2.5')
```

```
Out[30]: Text(0.5, 1.0, 'clusters when st_deviation = 2.5')
```



As we can see, all methods, that were used above indicate that natural number of clusters is 4 when standard deviation is 0.6 or 0.1, but when standard deviation is increased (for example to 2.5), we can observe no clear distinction between datapoints and natural number of clusters may be different than 4.

In [38]: #e.).

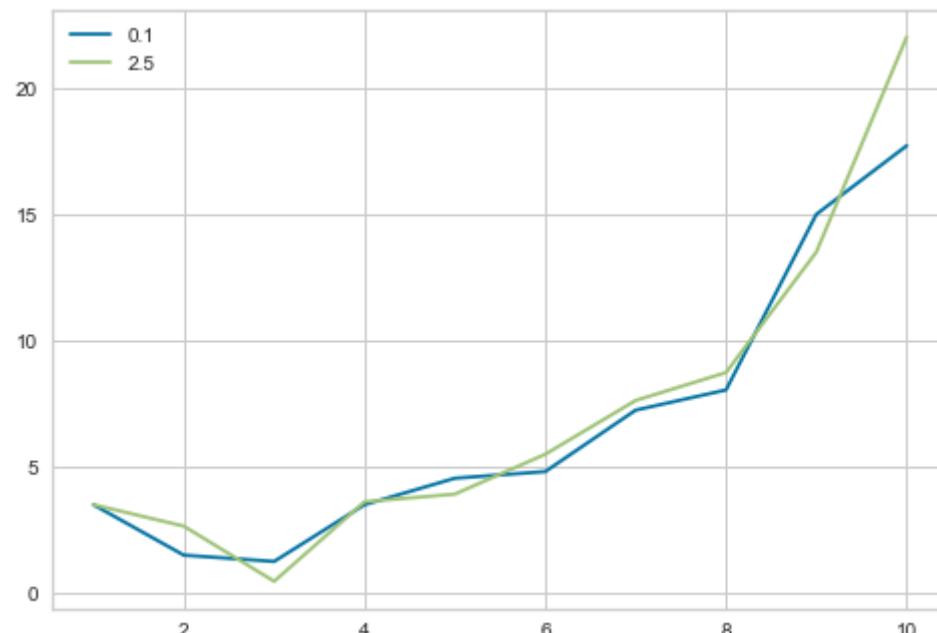
```
cluster_std=[0.1,2.5]
k = [1,2,3,4,5,6,7,8,9,10]

SSE_iter=[]
for i in cluster_std:
    X, y = make_blobs(n_samples=300, centers=4, n_features=2,cluster_std=i)
    df = DataFrame(dict(x=X[:,0], y=X[:,1], label=y))
    for i in k:
        kmeans = KMeans(n_clusters = i, random_state = 40)
        y_pred = kmeans.fit_predict(df)

        SSE = mean_squared_error(y,y_pred)
        SSE_iter.append(SSE)

plt.plot(k,SSE_iter[0:10],label='0.1')
plt.plot(k,SSE_iter[10:20],label='2.5')
plt.legend(loc='upper left')

plt.show()
```



K-means is not a deterministic algorithm, if random_State is set as none, different runs will start at different points, Seeding the pseudo-random number generator ensures that this randomness will always be the same for identical seeds.

In my opinion, when the results vary highly from run to run, this indicates that the data just does not cluster well with k-means at all. If the data is really suited for k-means clustering, the results will be rather stable! If they vary, the clusters may not have the same size, or may be not well separated; and other algorithms may yield better results.

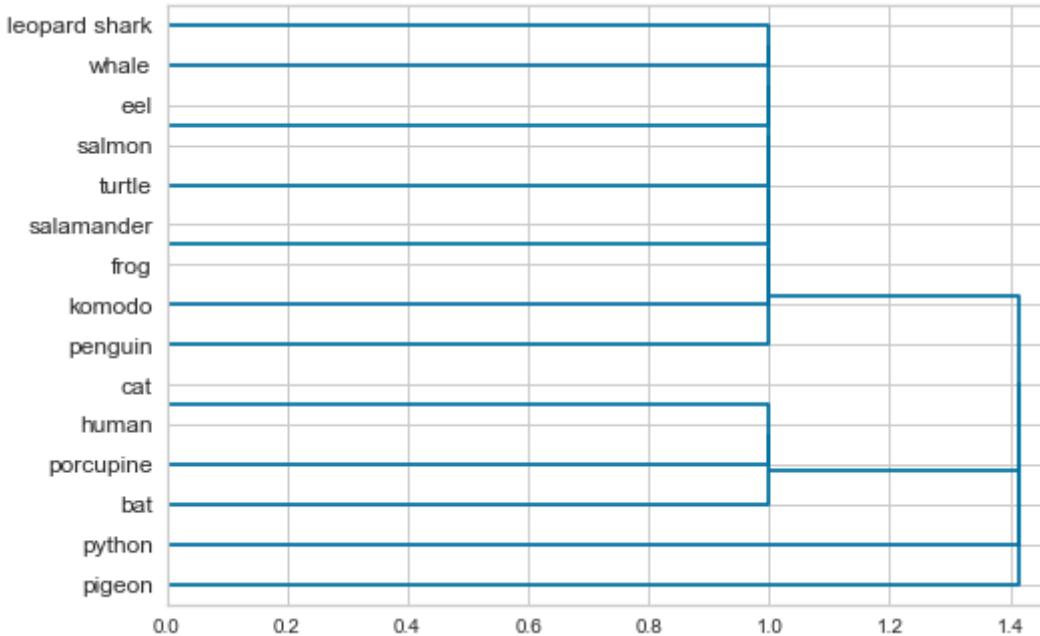
Assignment 2

```
In [51]: X = pd.read_csv('vertebrate.csv',header='infer')  
X.shape
```

```
Out[51]: (15, 8)
```

Single link

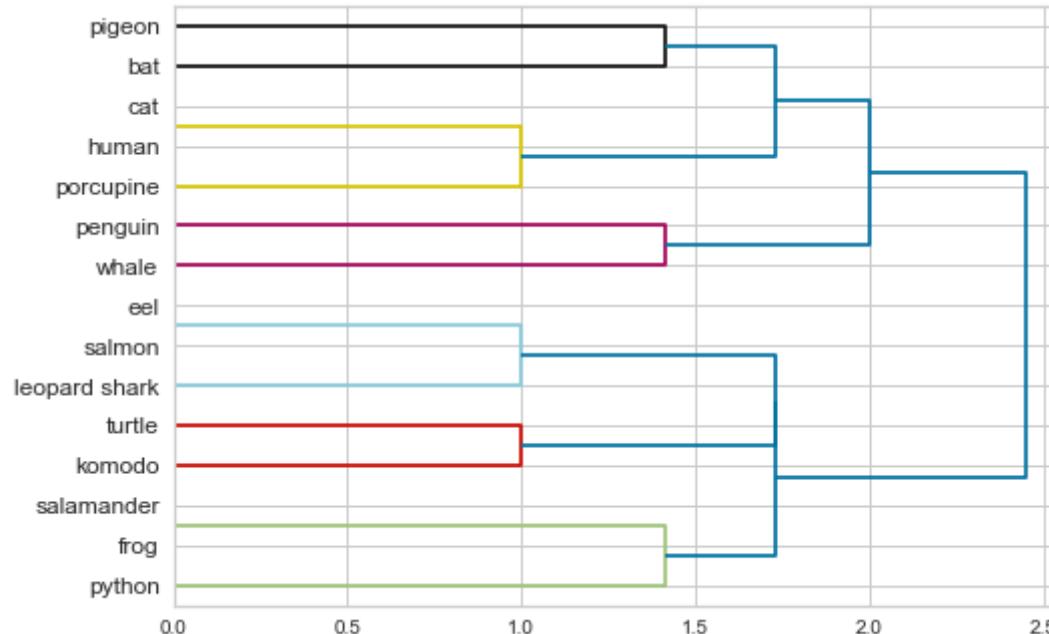
```
In [52]: names = X['Name']
y = X['Class']
X = X.drop(['Name', 'Class'], axis=1)
Z = linkage(X.as_matrix(), 'single')
dn = dendrogram(Z, labels=names.tolist(), orientation='right')
```



Max link

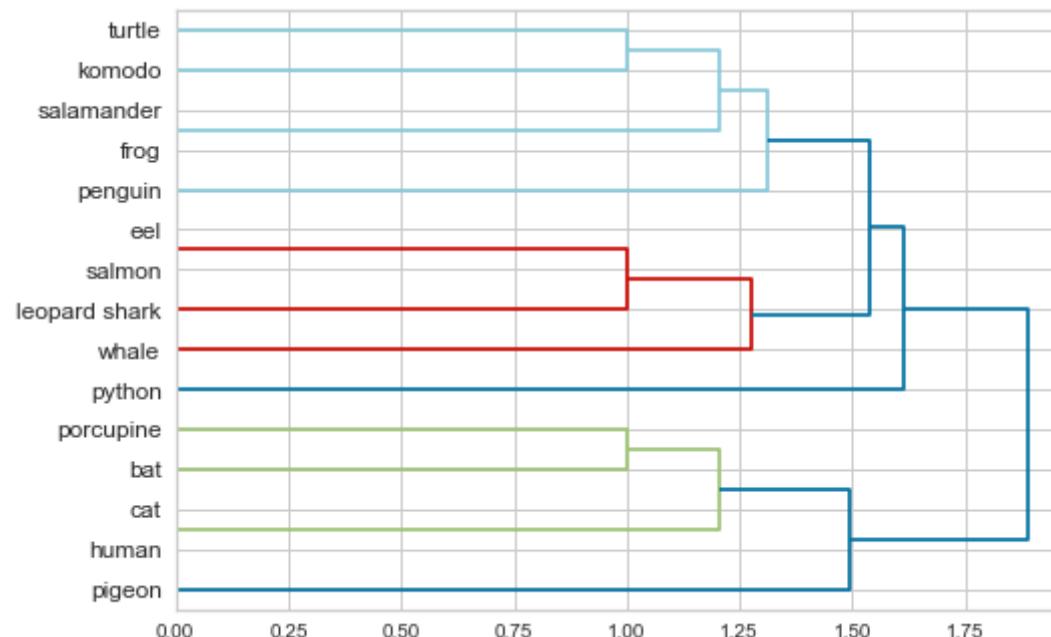
```
In [65]: Z = linkage(X.as_matrix(), 'complete')
dn = dendrogram(Z,labels=names.tolist(),orientation='right')
y.values
```

```
Out[65]: array(['mammals', 'reptiles', 'fishes', 'mammals', 'amphibians',
   'reptiles', 'mammals', 'birds', 'mammals', 'fishes', 'reptiles',
   'birds', 'mammals', 'fishes', 'amphibians'], dtype=object)
```



average link

```
In [35]: Z = linkage(X.as_matrix(), 'average')
dn = dendrogram(Z,labels=names.tolist(),orientation='right')
```

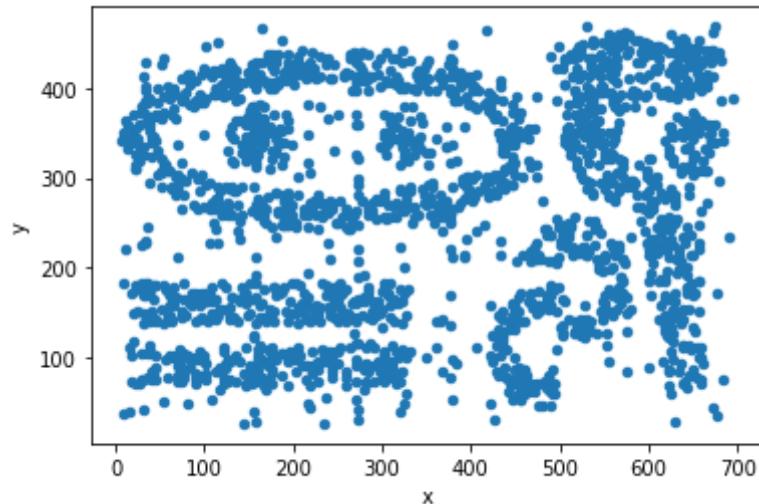


In my opinion the most natural clustering for our data is clustering with Max link. It makes clusters including animals from mostly the same animal kingdom. For example in one cluster we have whale, penguin, percupine, human, cat, bat and pigeon, which are all mammals except from pigeon.

Assignment 3

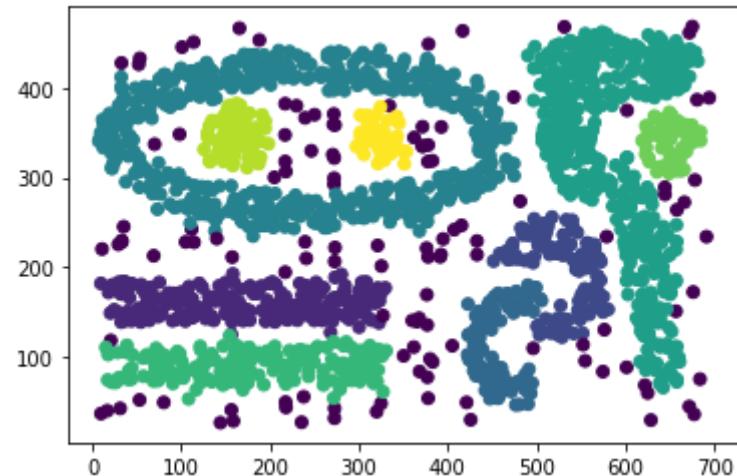
```
In [2]: df = pd.read_csv('chameleon.csv')
# df.head()
df.plot.scatter(x='x',y='y')
```

```
Out[2]: <matplotlib.axes._subplots.AxesSubplot at 0x24a8daeccc0>
```

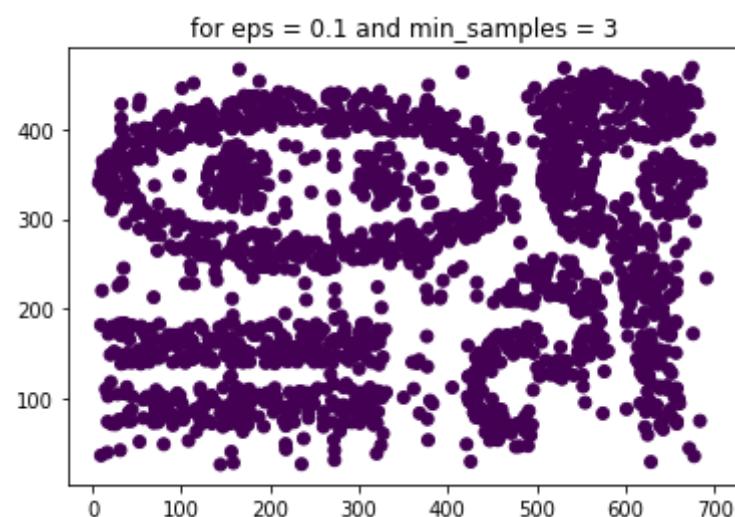
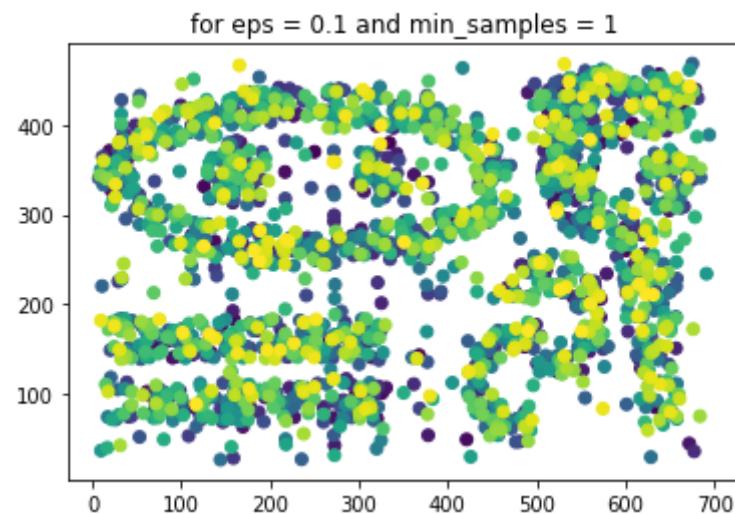


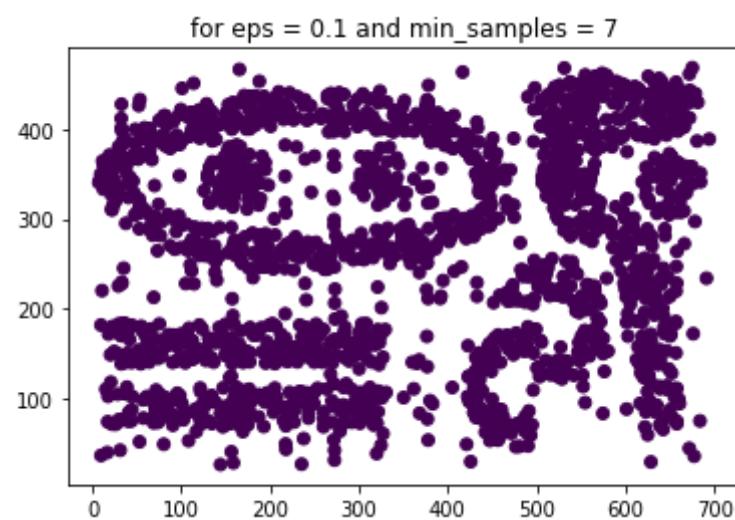
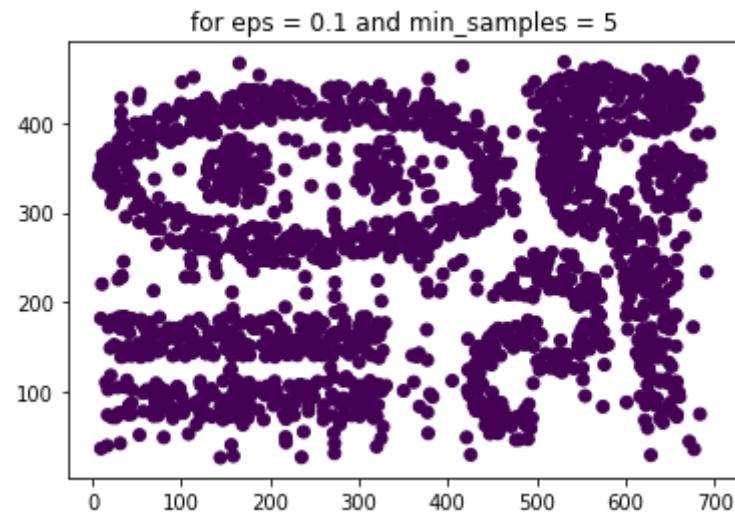
```
In [3]: db = DBSCAN(eps=15.5, min_samples=5).fit_predict(df)
plt.scatter(df['x'],df['y'],c=db)
```

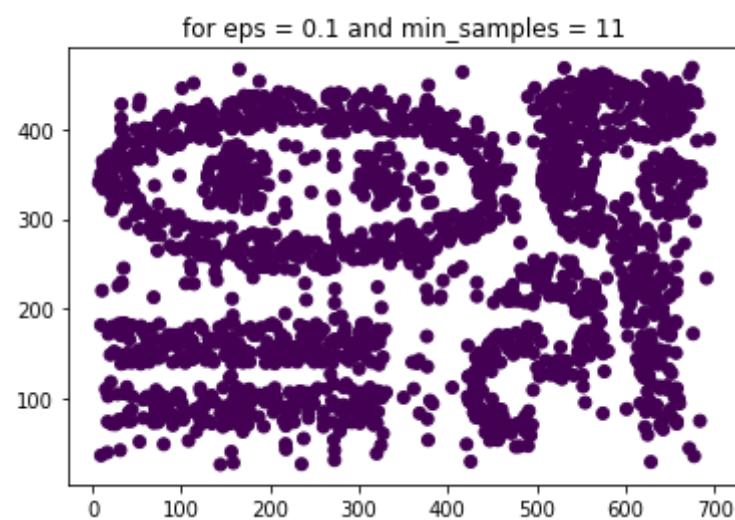
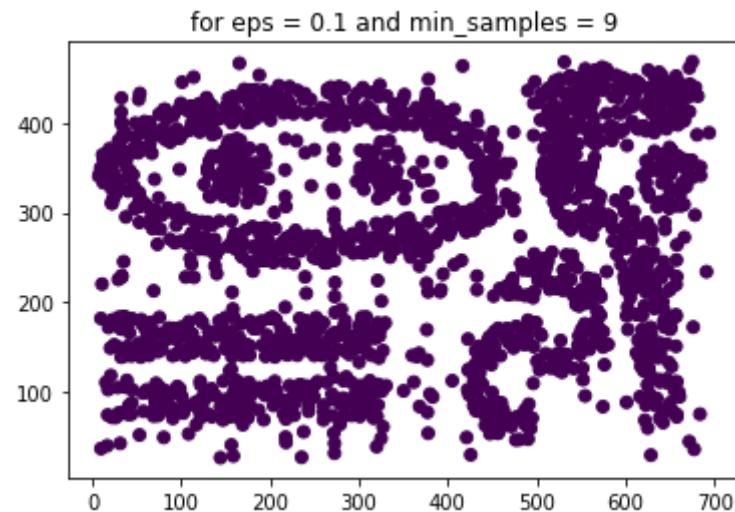
```
Out[3]: <matplotlib.collections.PathCollection at 0x24a8df08a20>
```

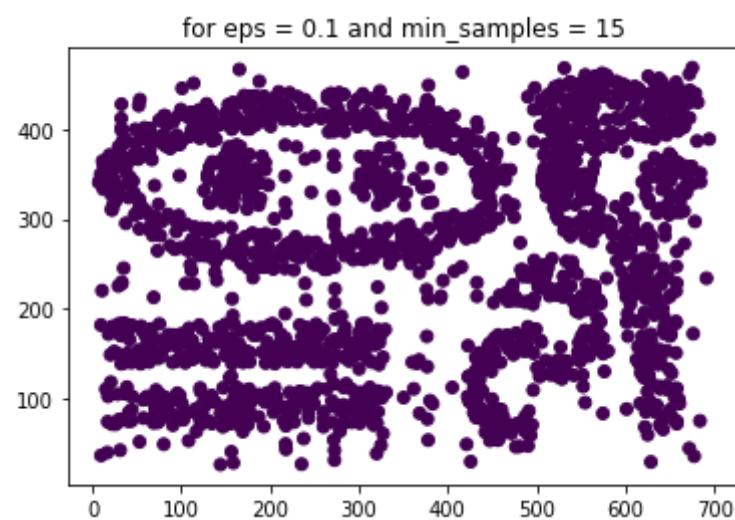
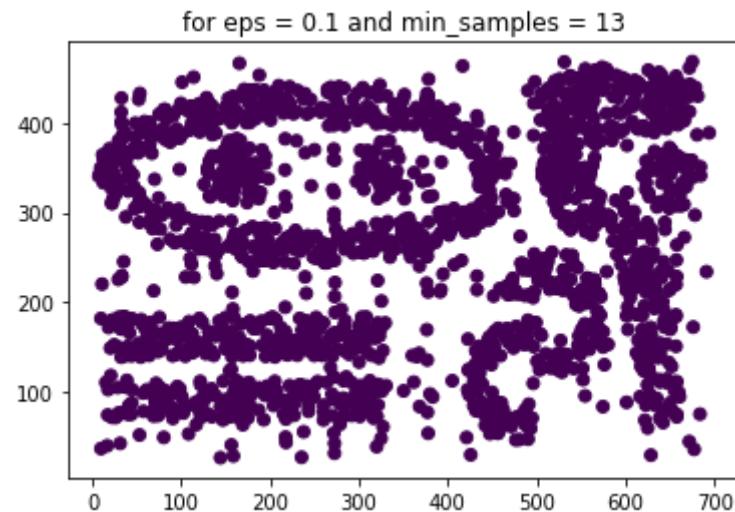


```
In [6]: eps = np.arange(0.1,20,0.5)
min_samples = np.arange(1,20,2)
for i in eps:
    for k in min_samples:
        db = DBSCAN(eps=i, min_samples=k).fit_predict(df)
        plt.scatter(df['x'],df['y'],c=db)
        plt.title('for eps = {} and min_samples = {}'.format(i,k))
        plt.show()
```

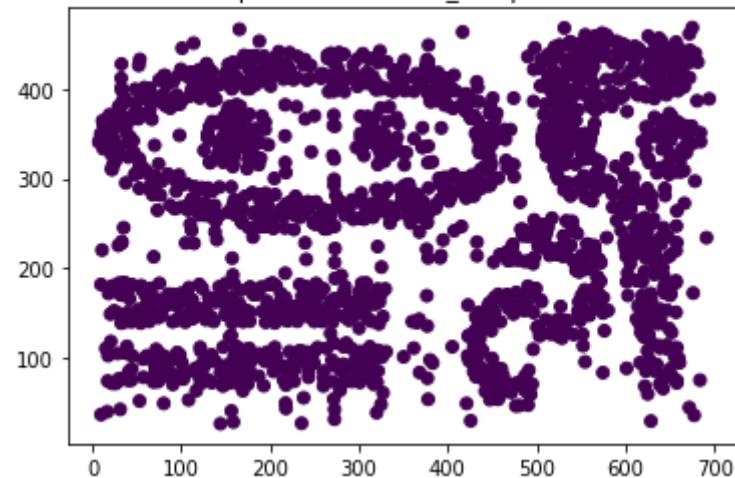




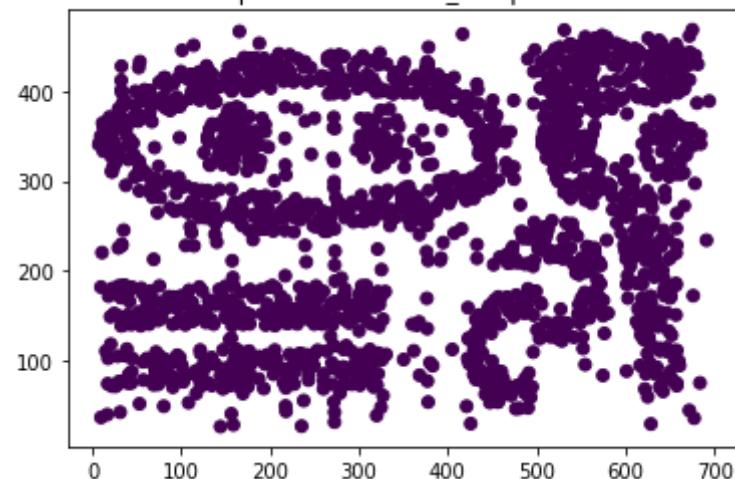


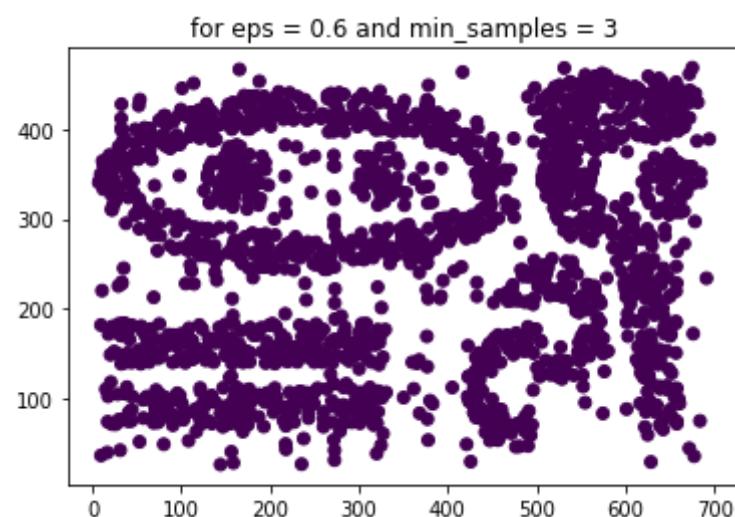
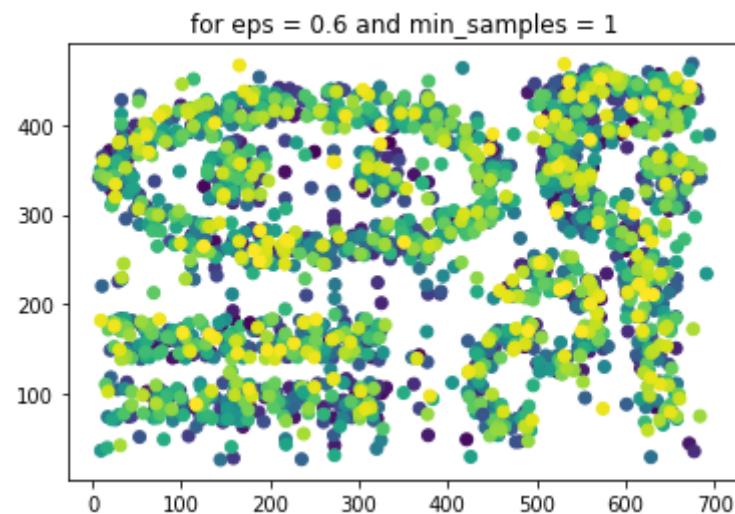


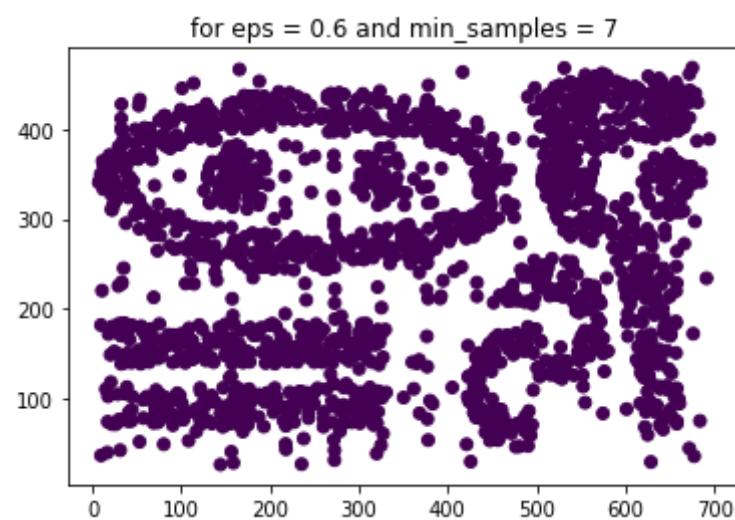
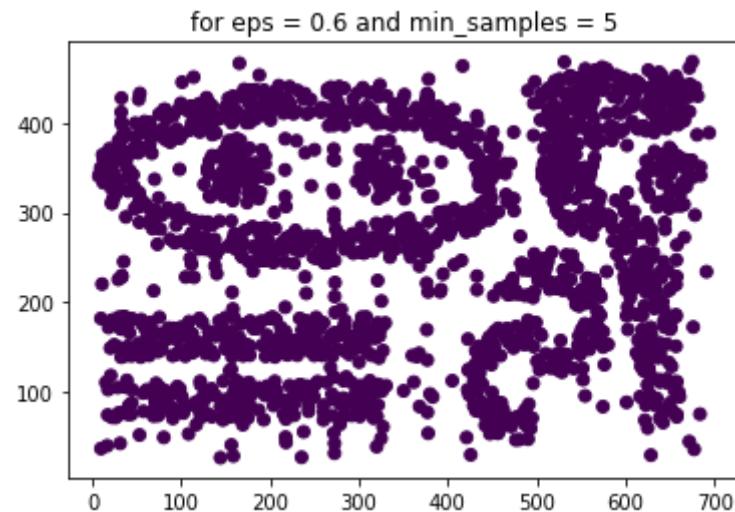
for eps = 0.1 and min_samples = 17

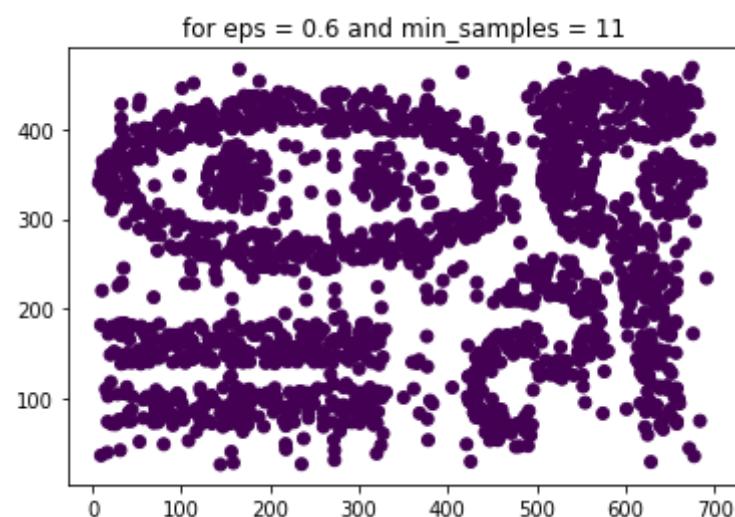
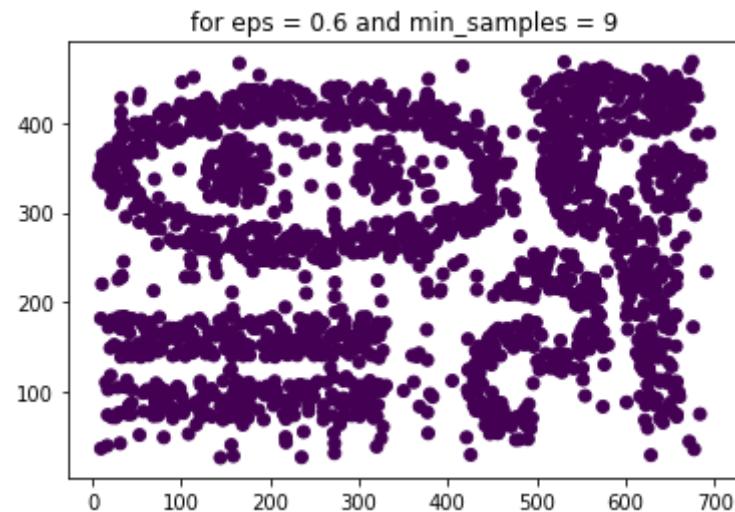


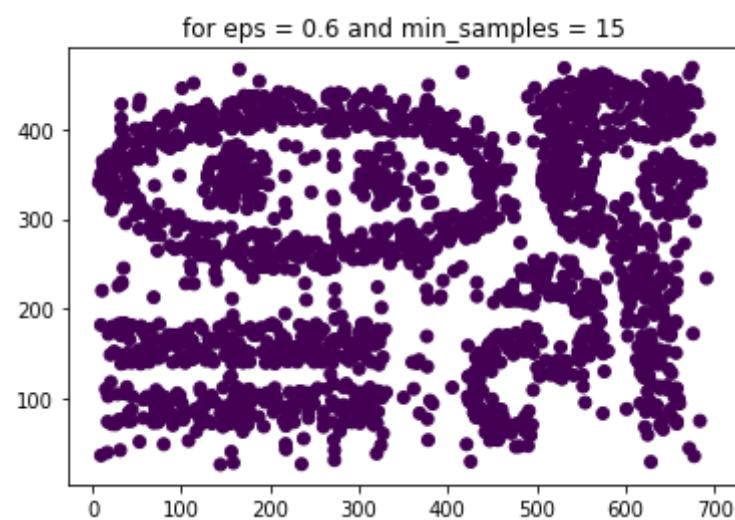
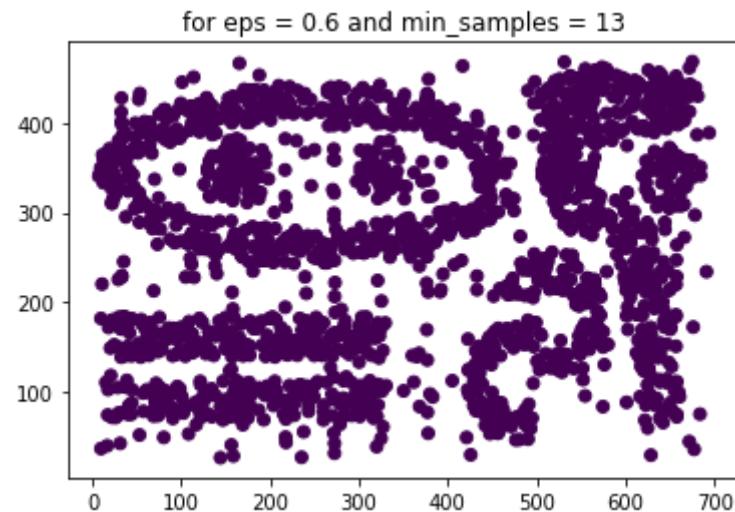
for eps = 0.1 and min_samples = 19

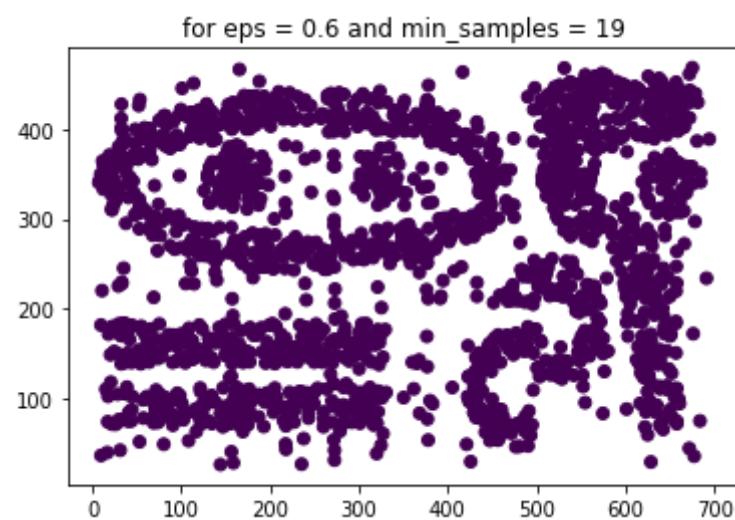
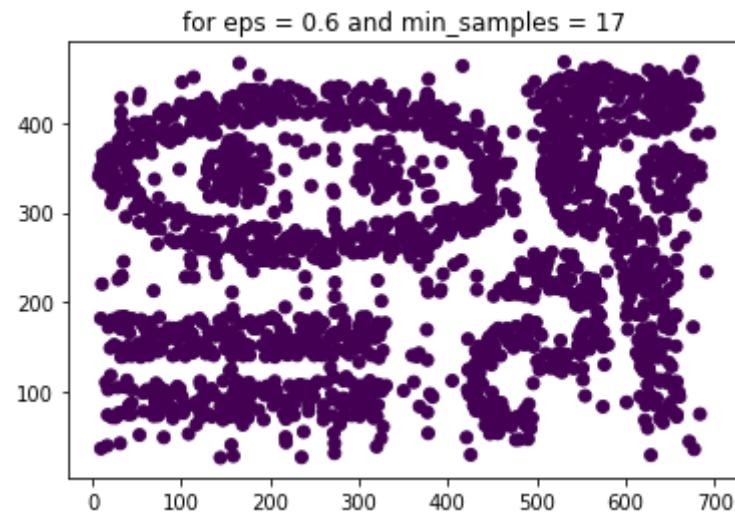


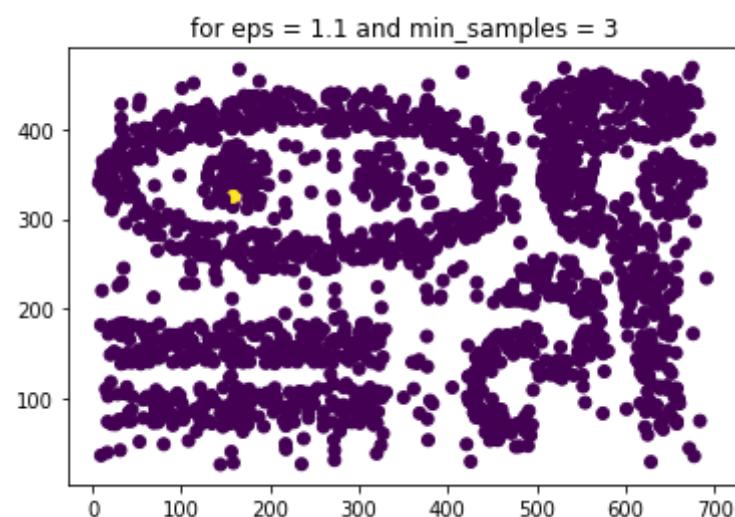
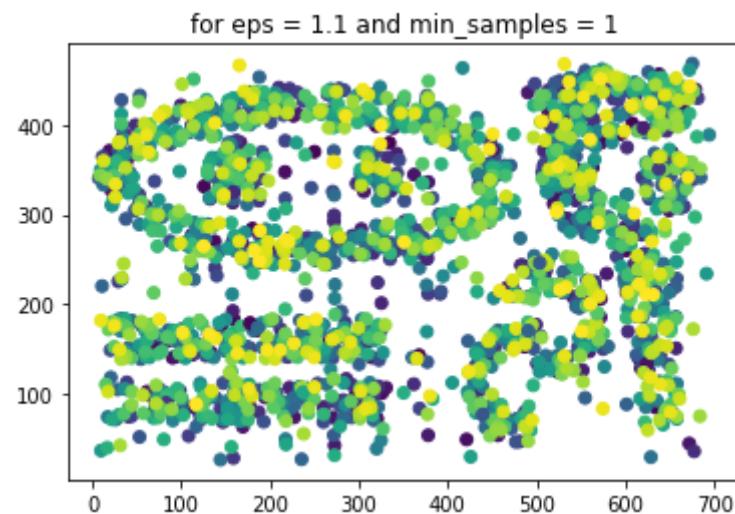


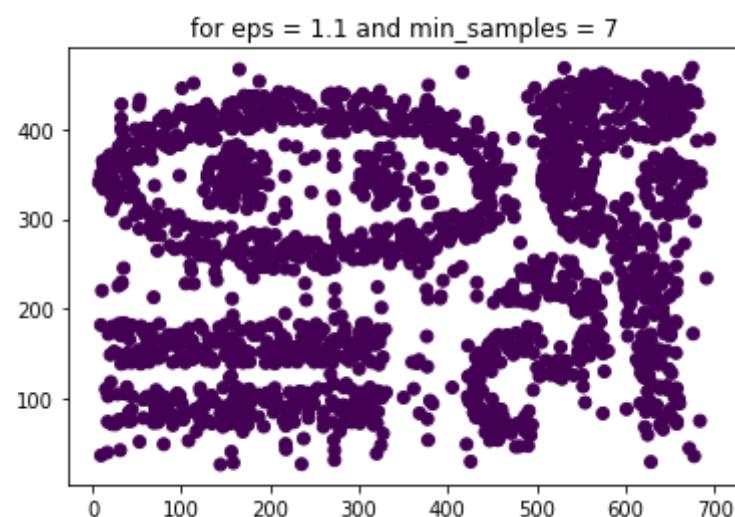
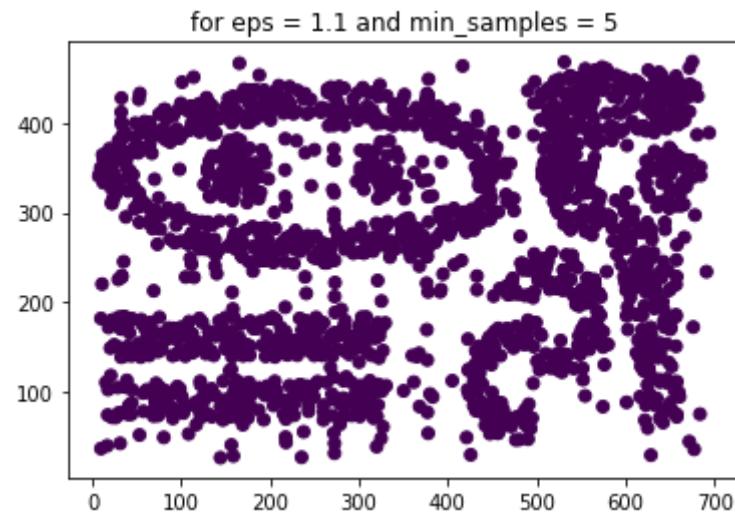


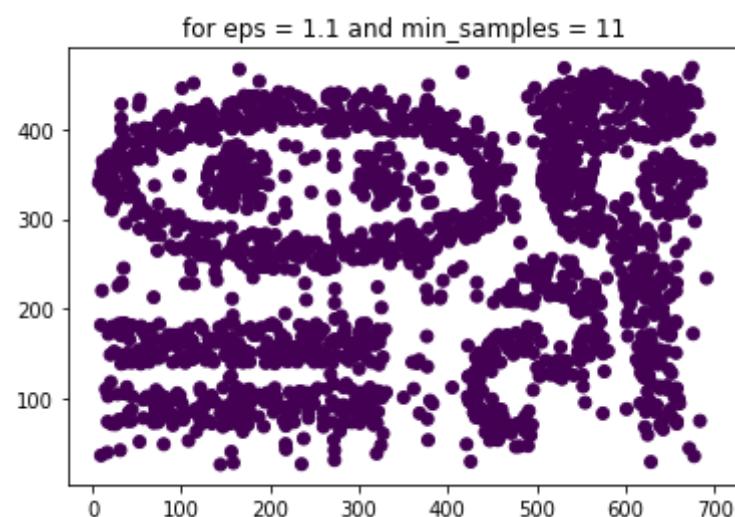
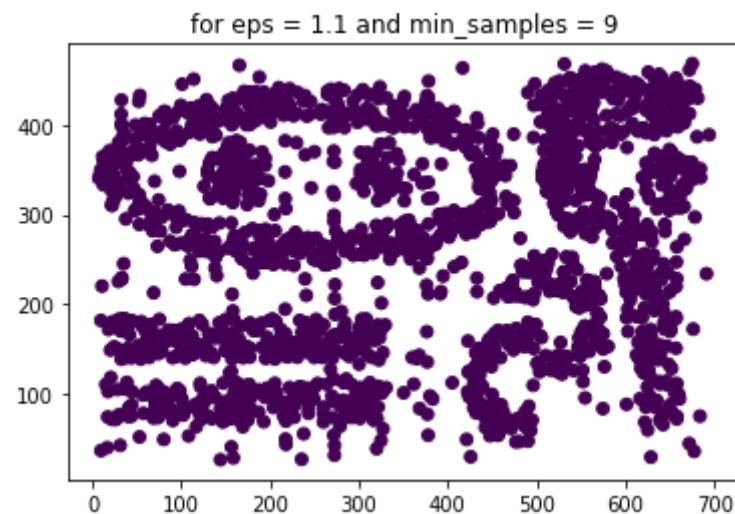


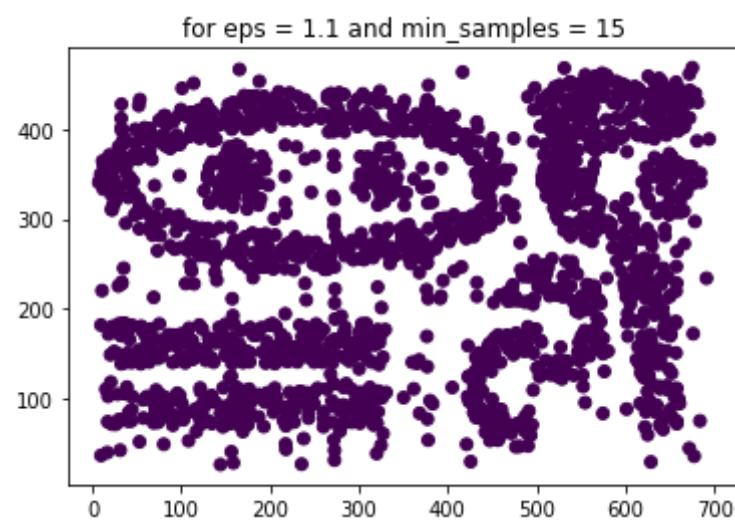
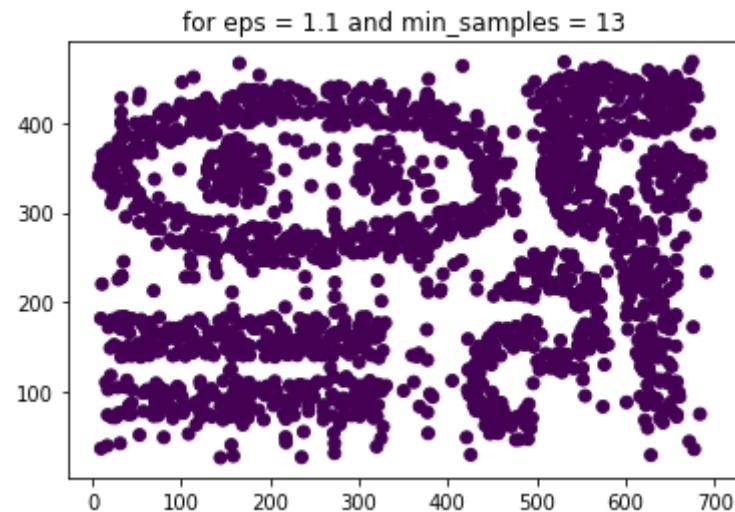


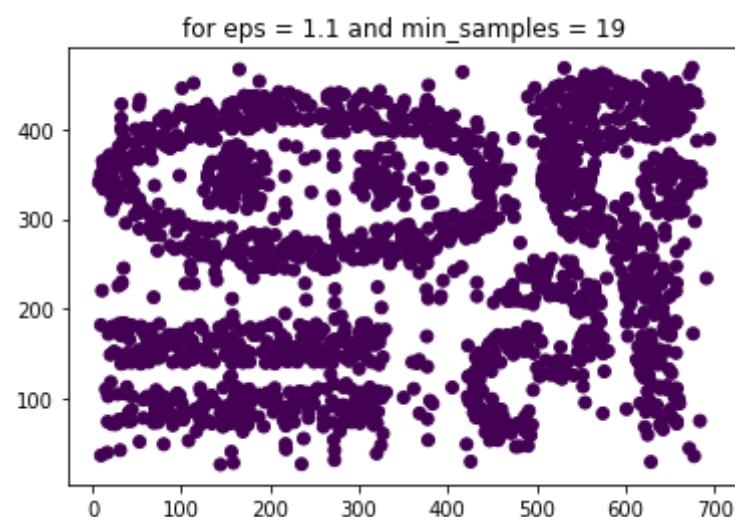
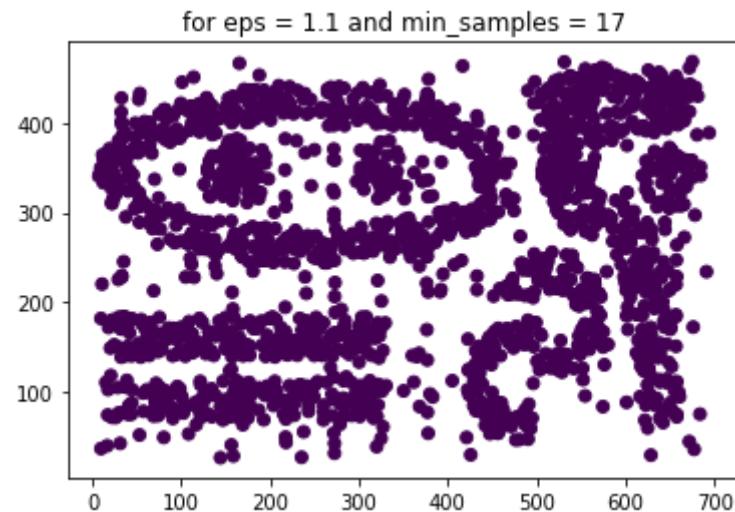


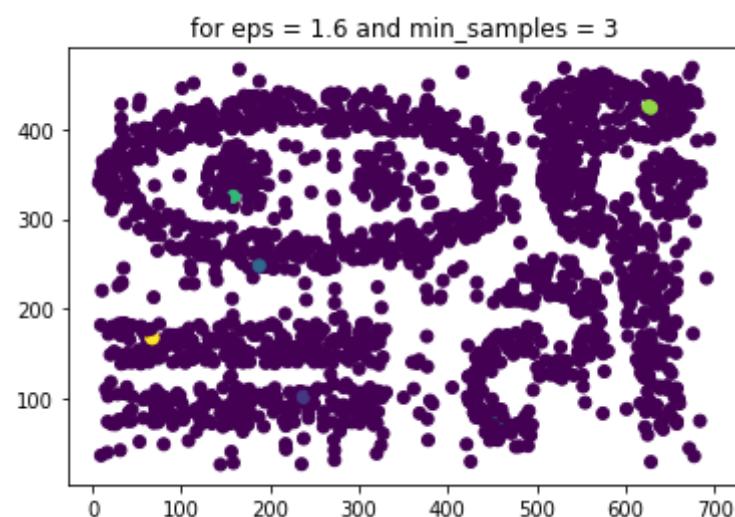
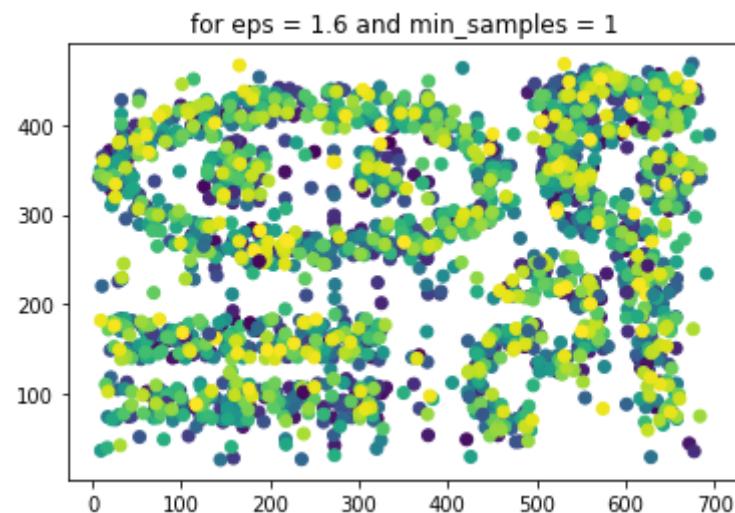


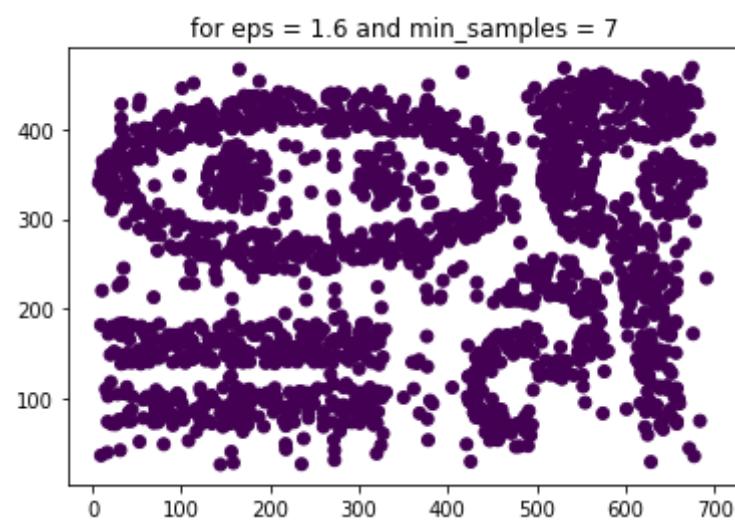
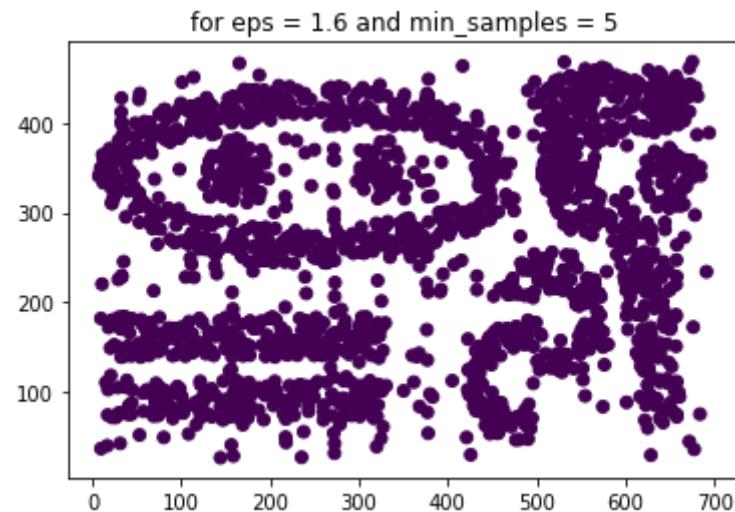


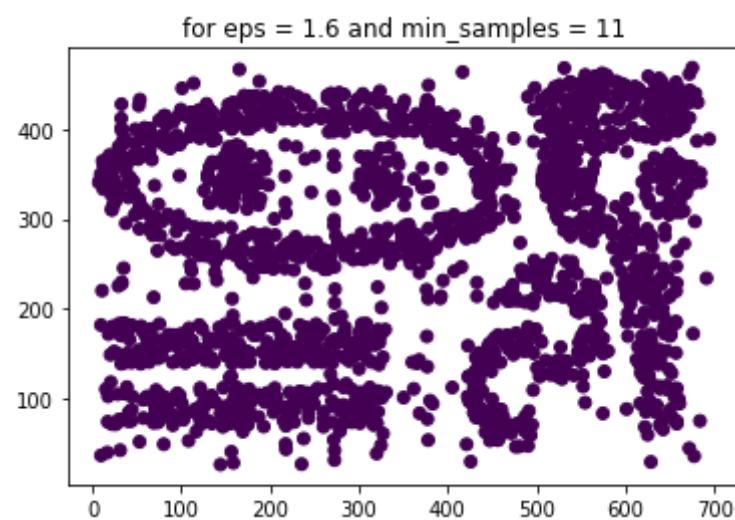
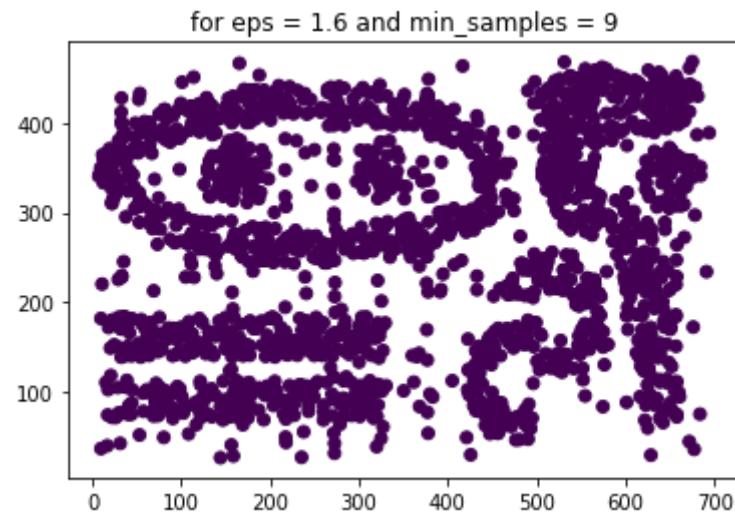


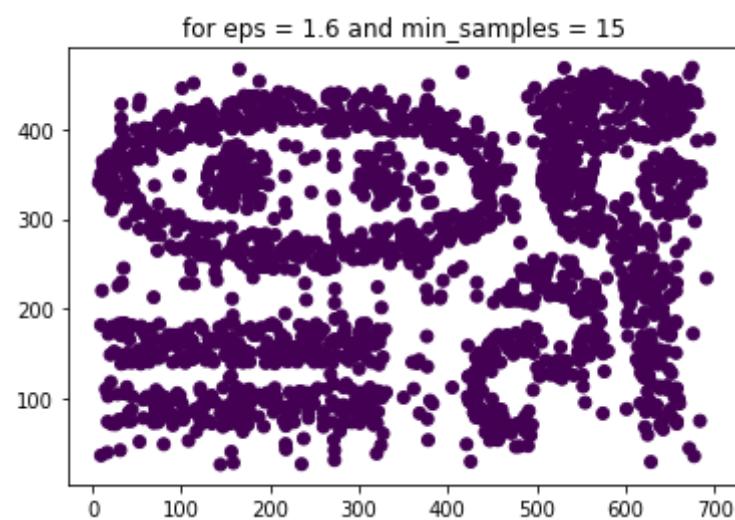
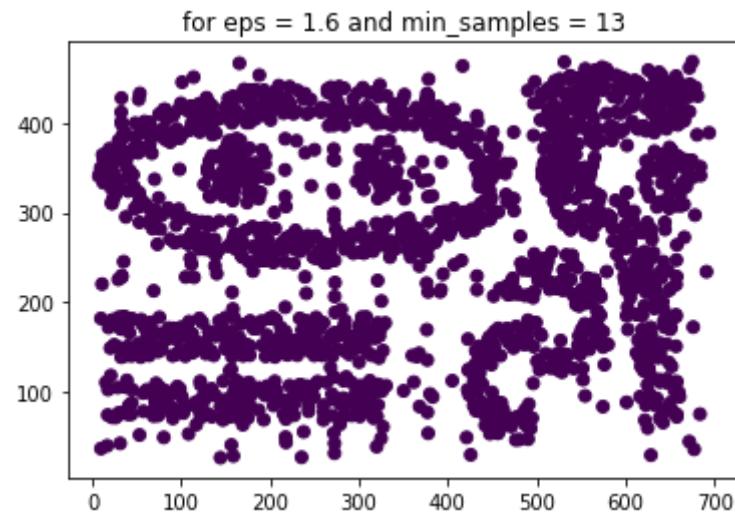


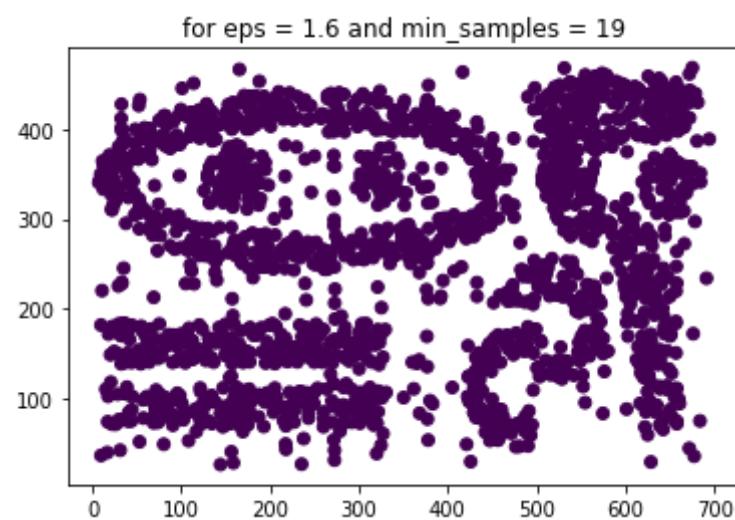
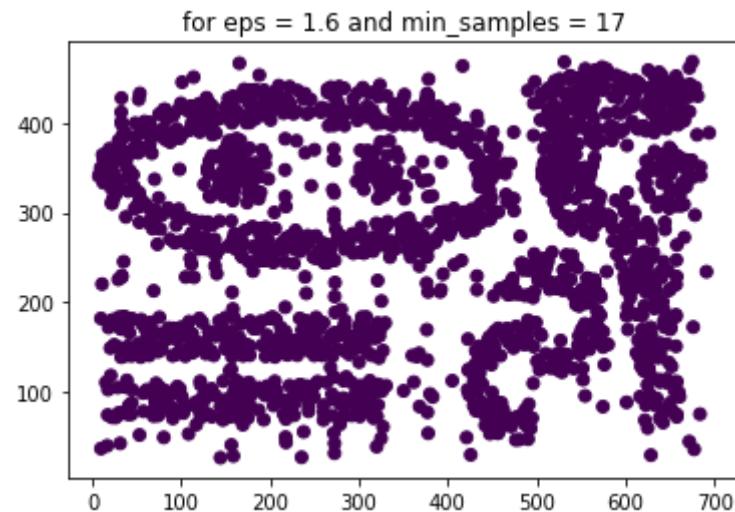


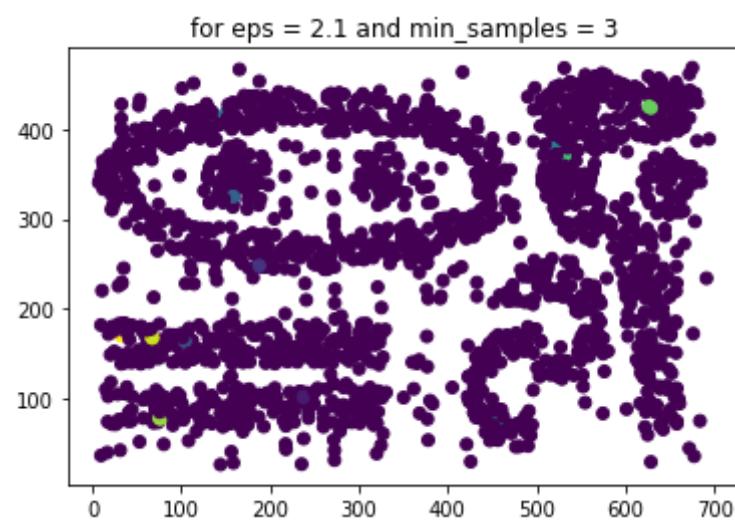
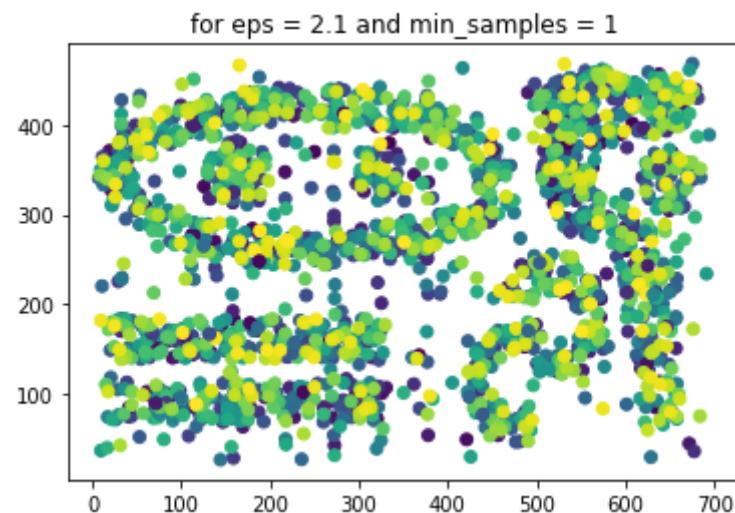


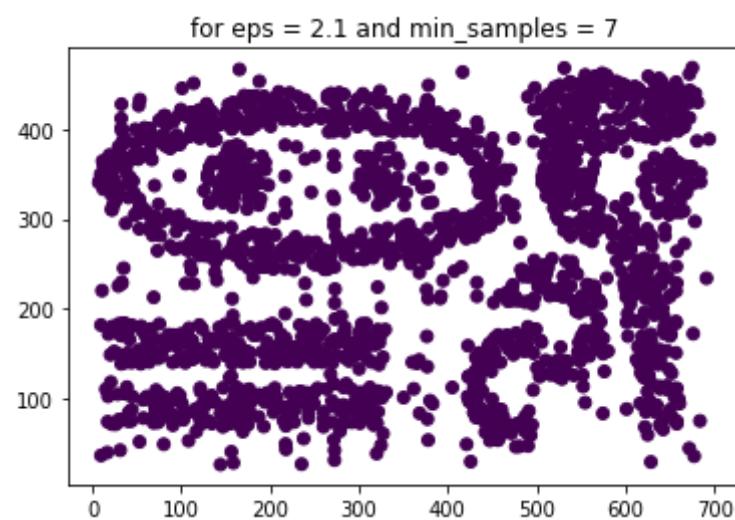
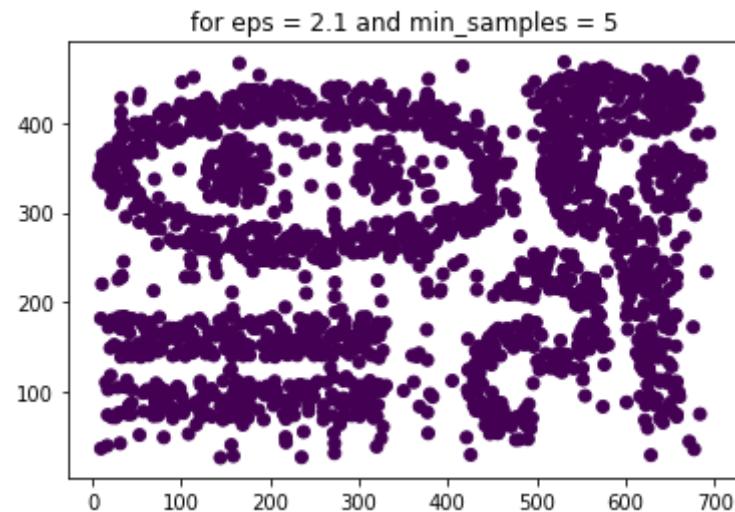


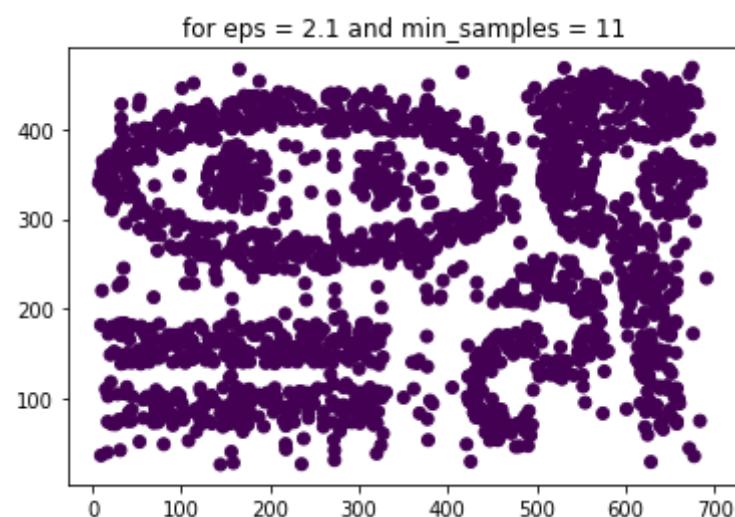
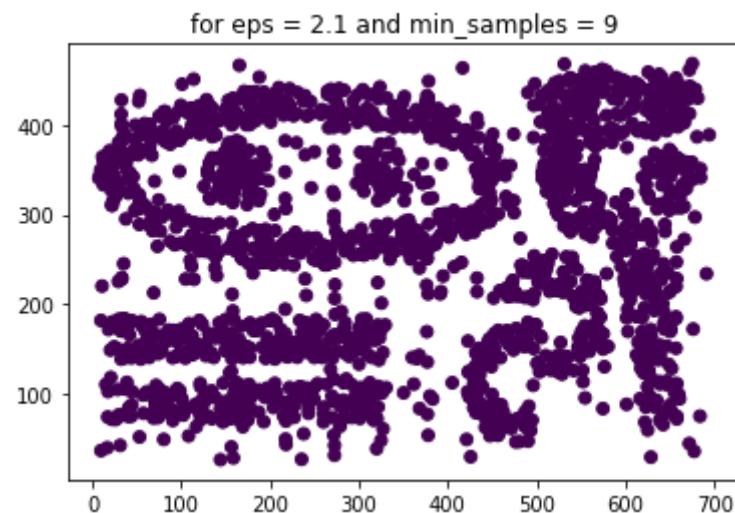


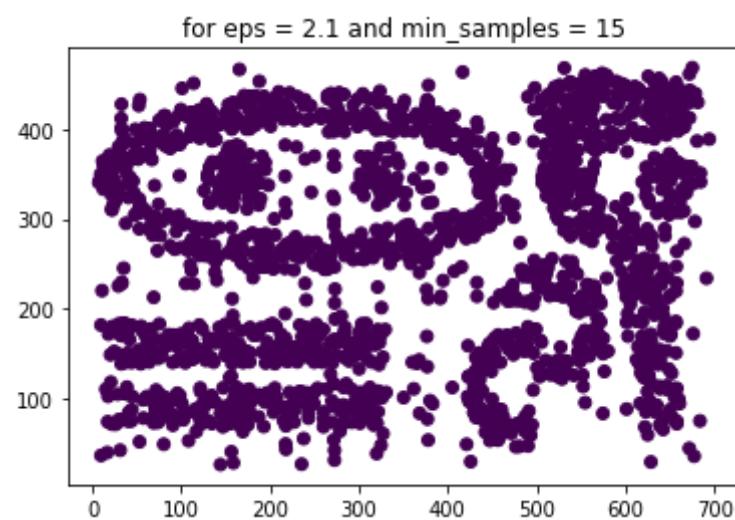
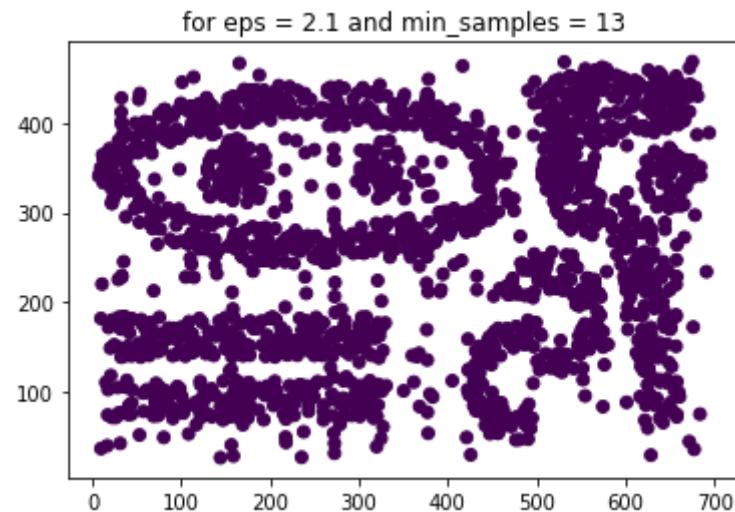


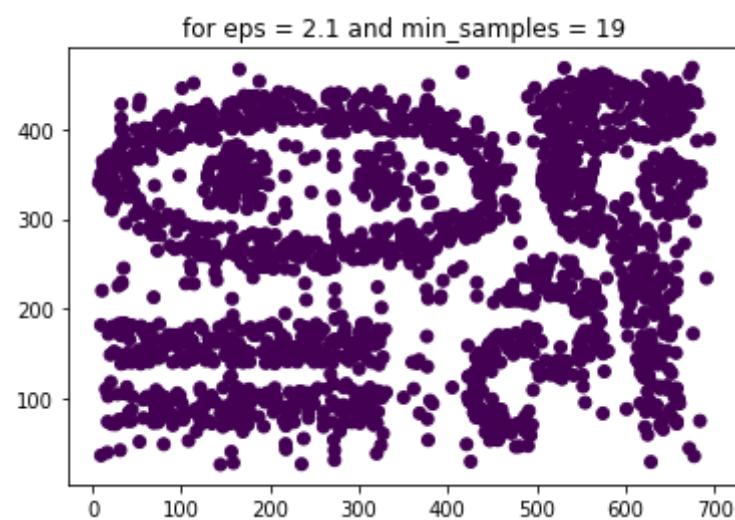
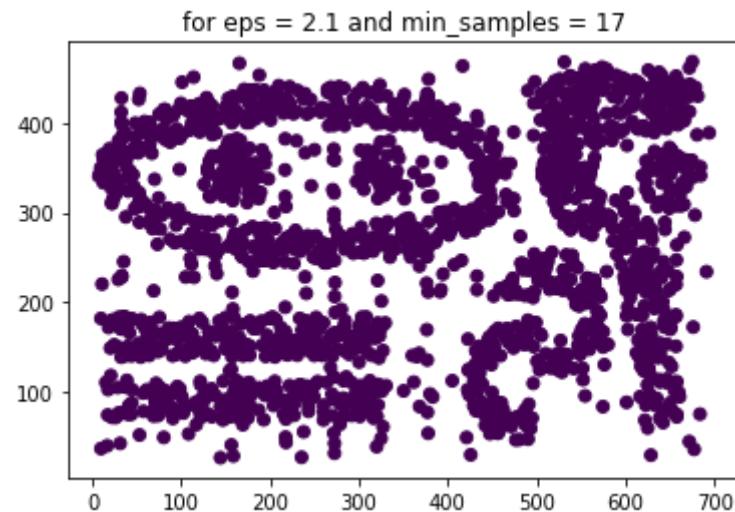


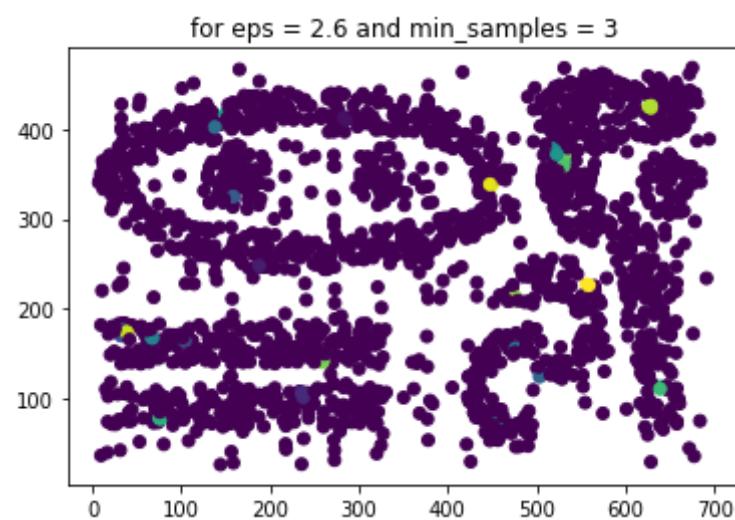
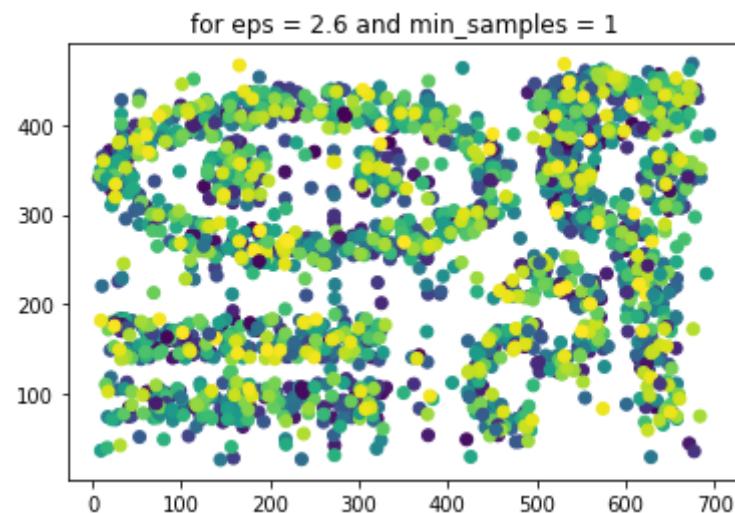


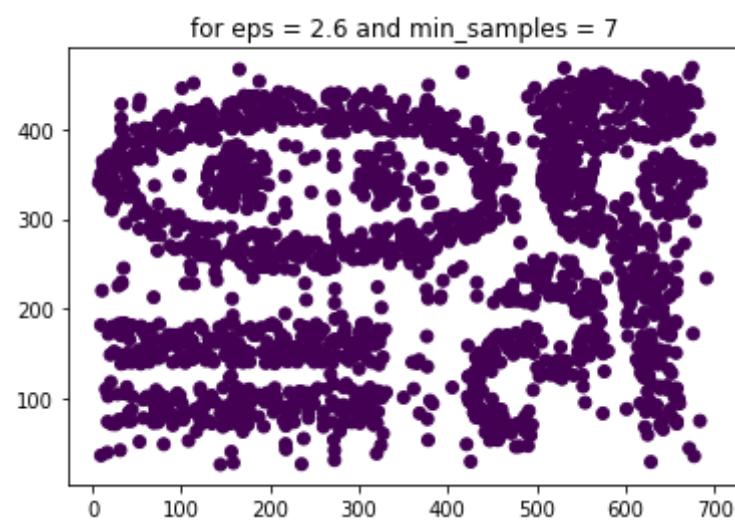
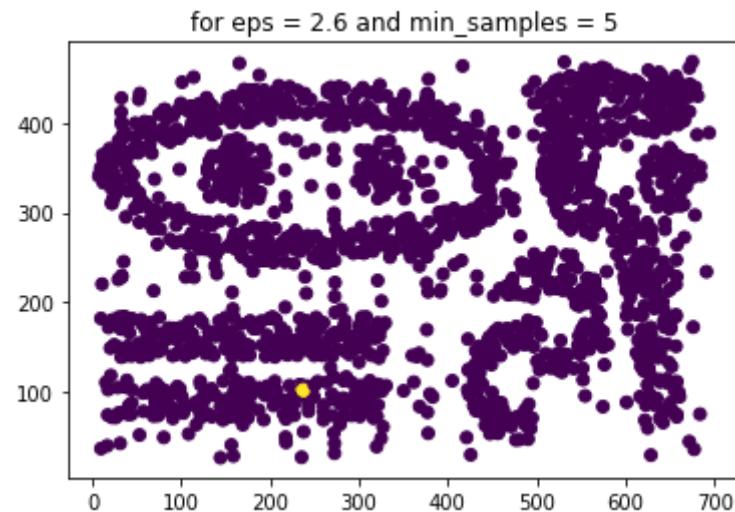


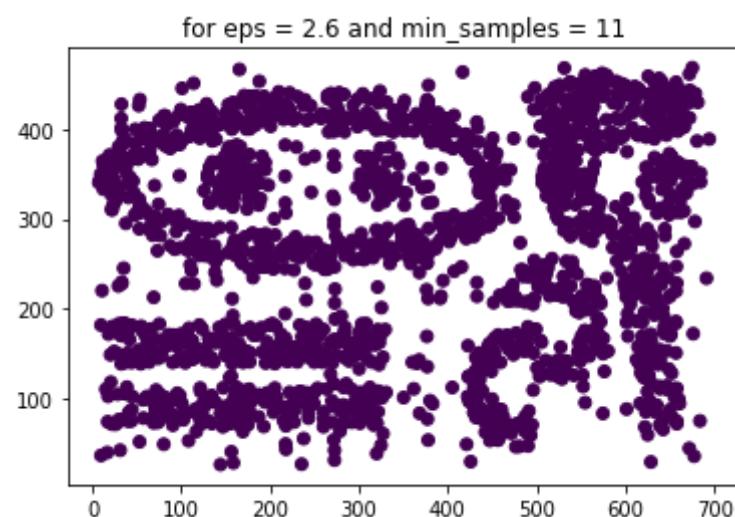
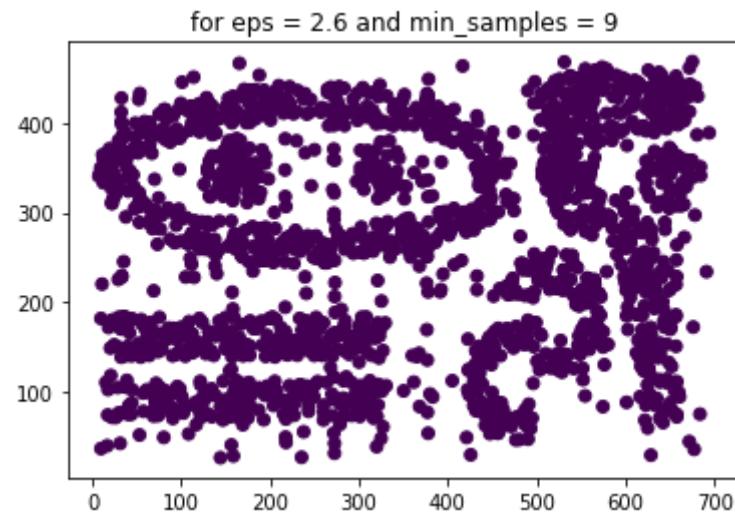


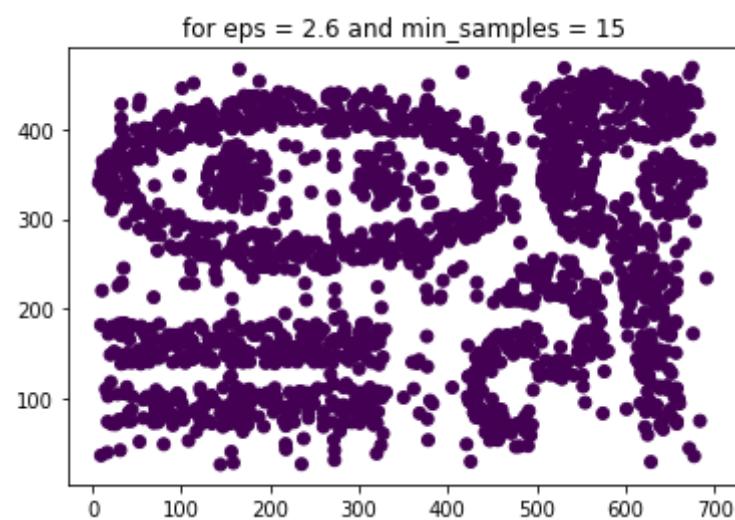
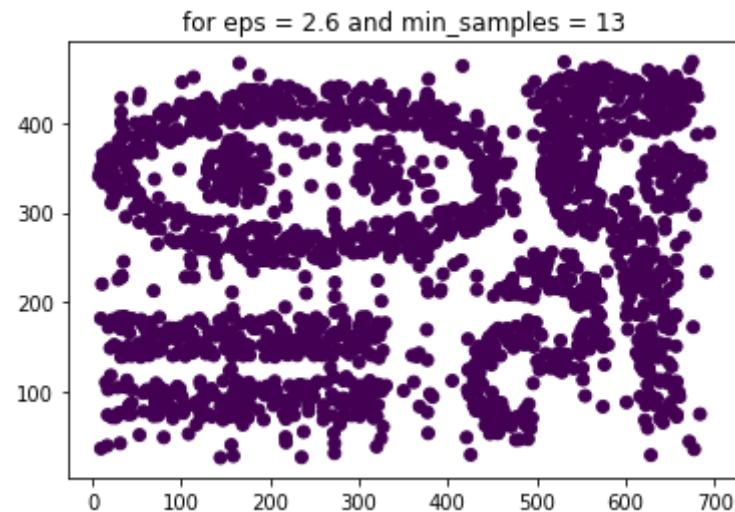


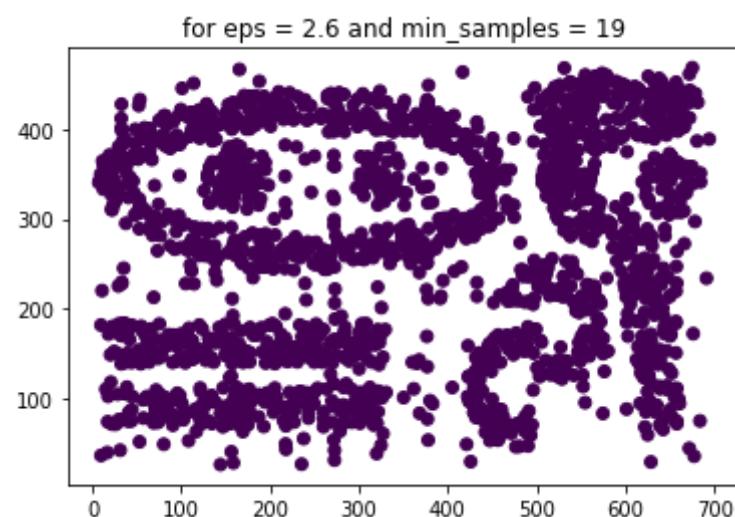
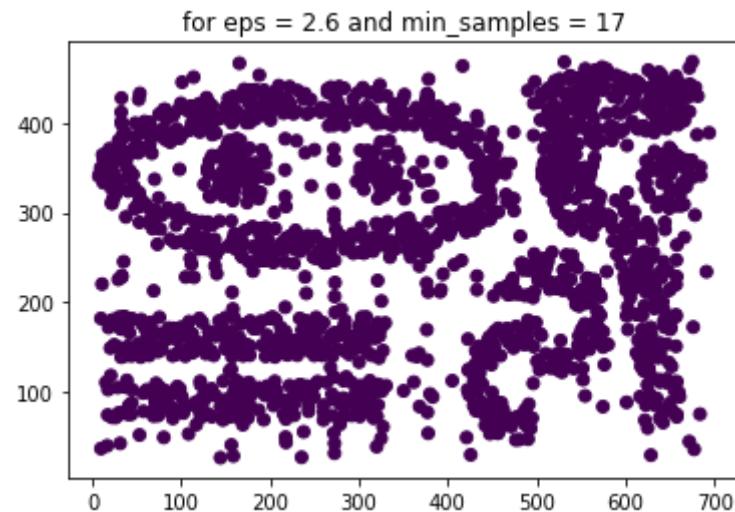


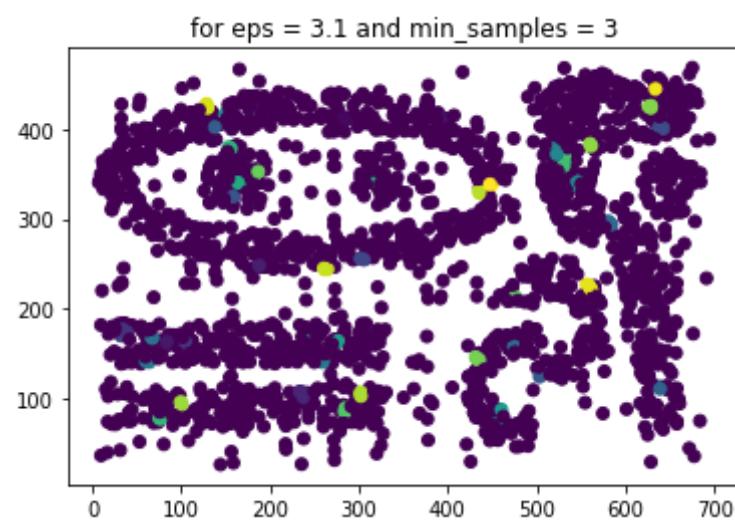
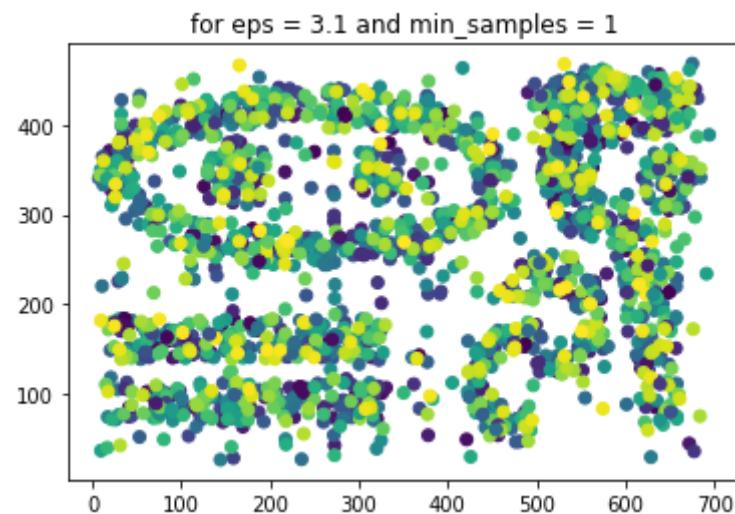


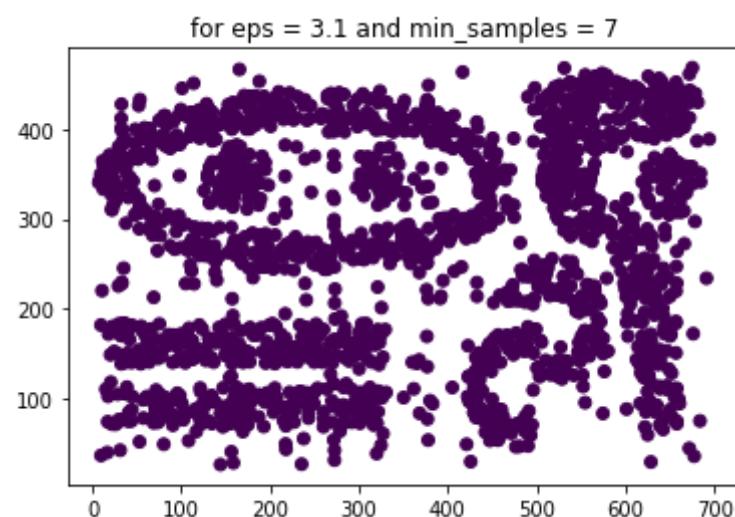
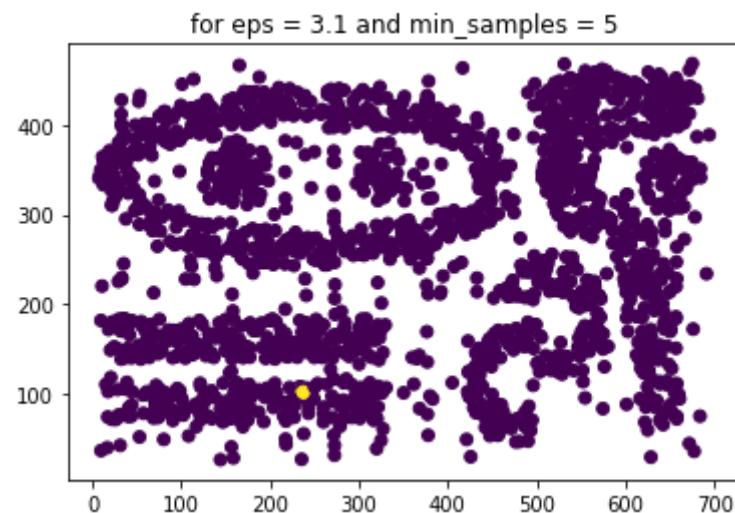


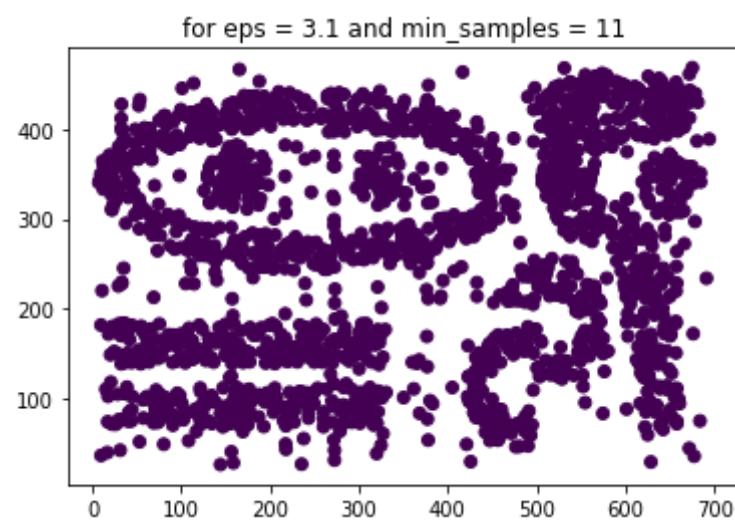
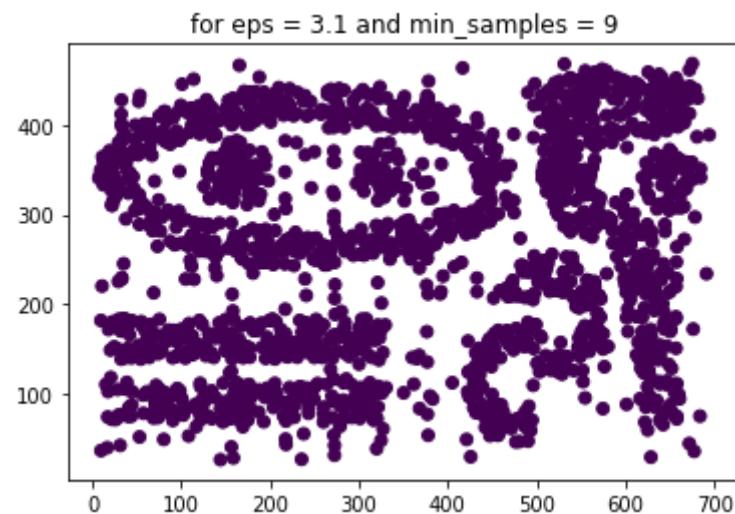


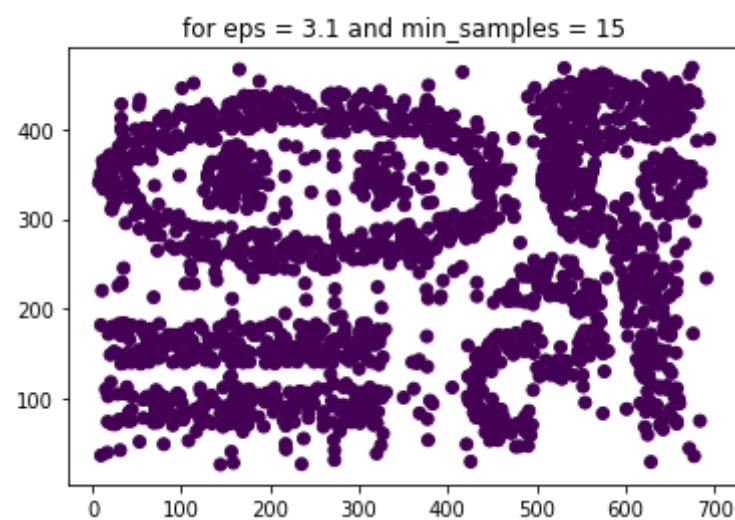
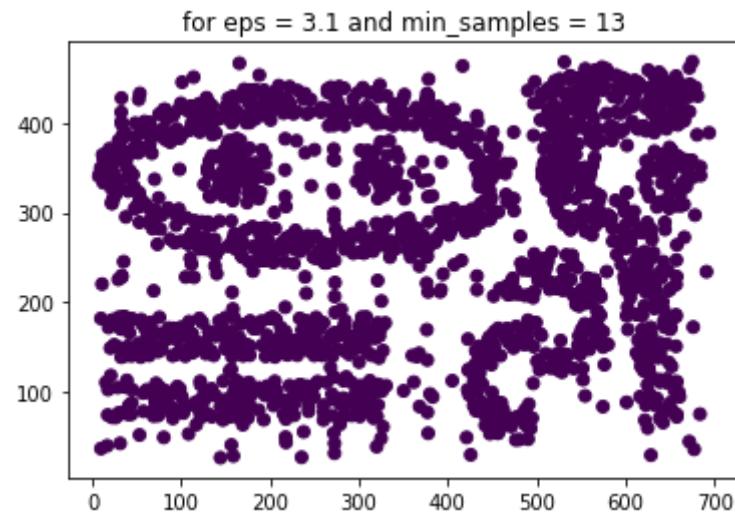


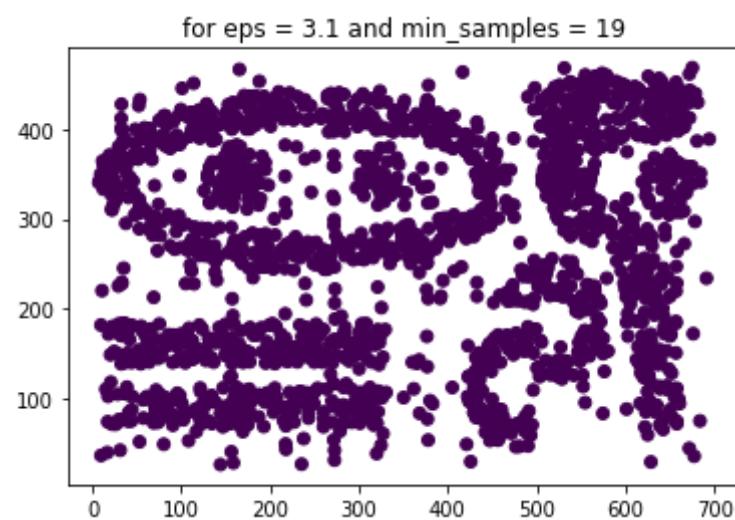
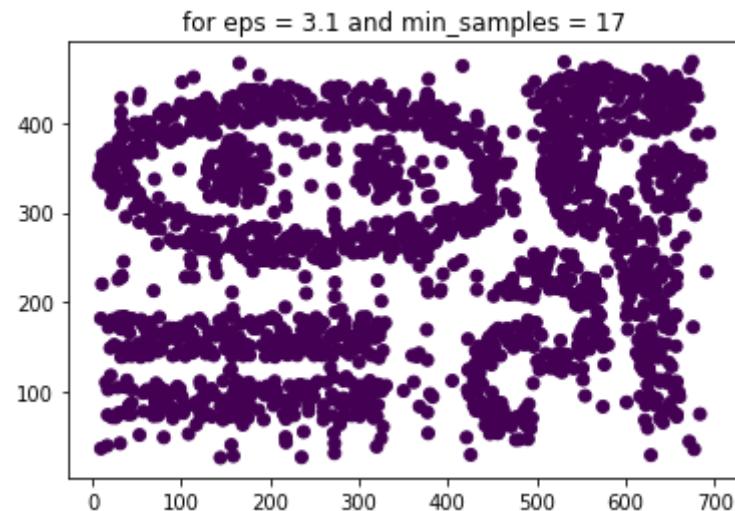


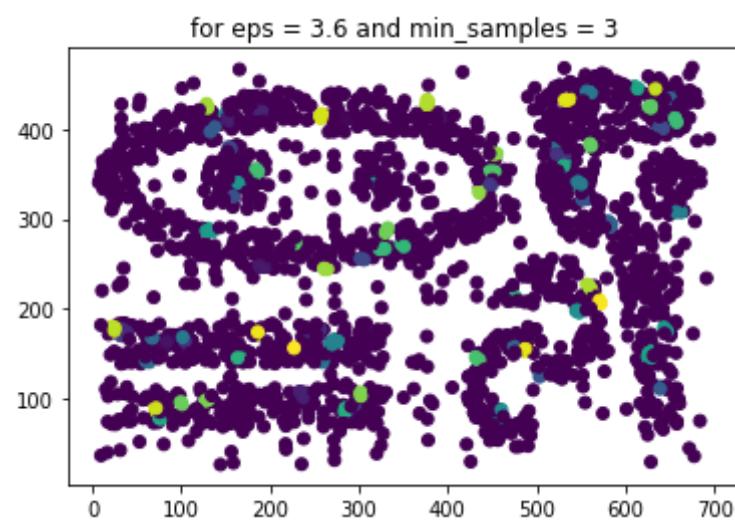
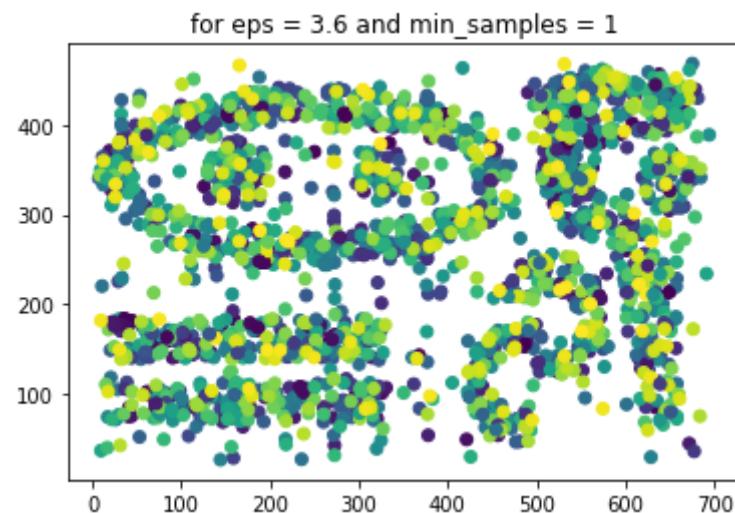


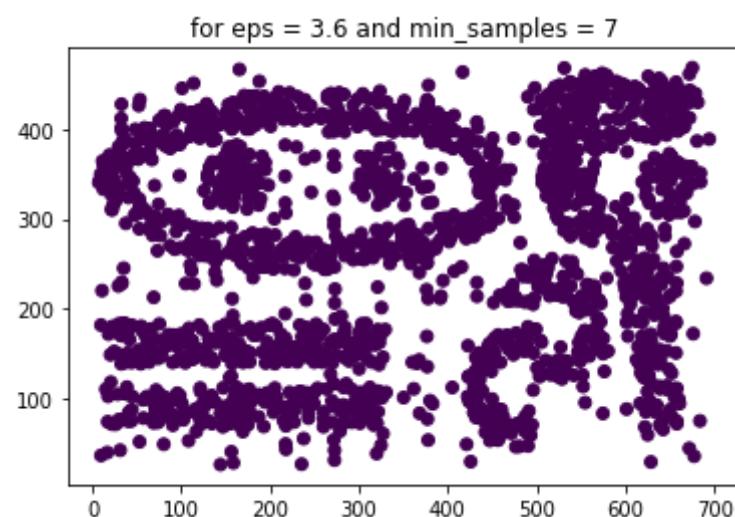
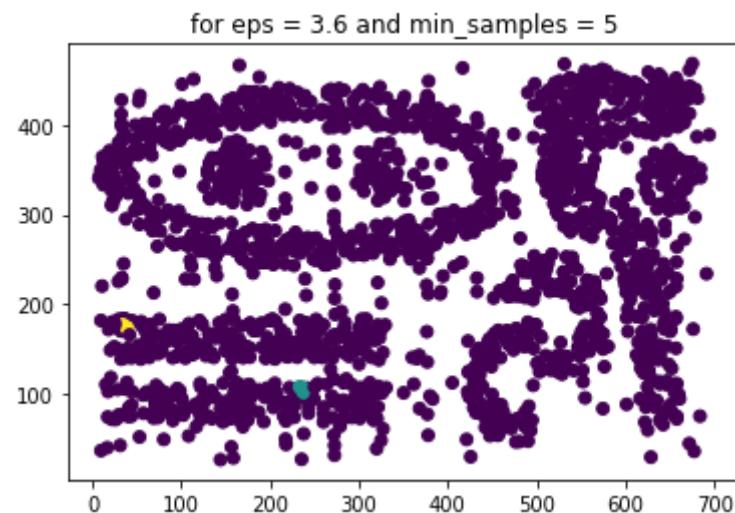


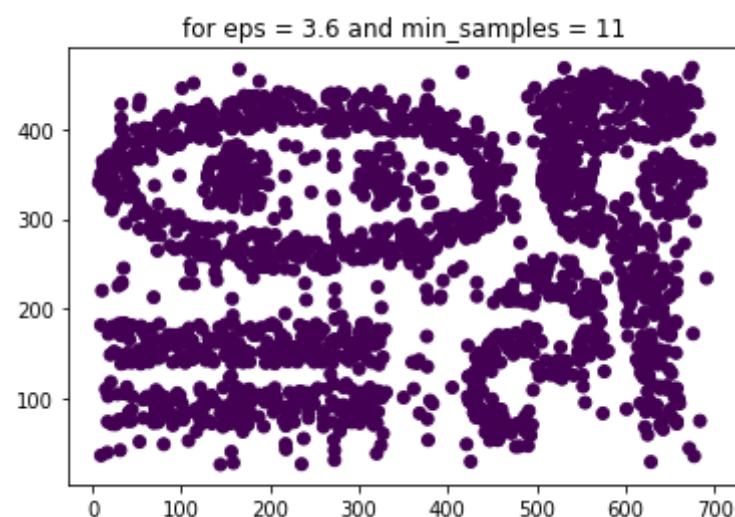
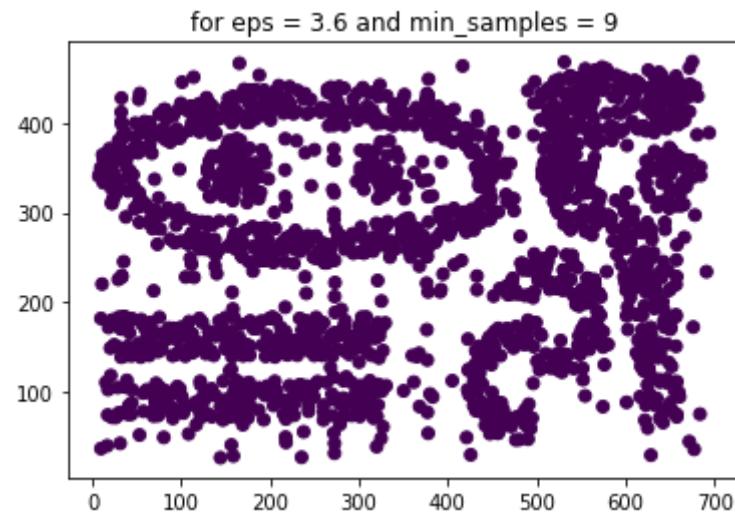


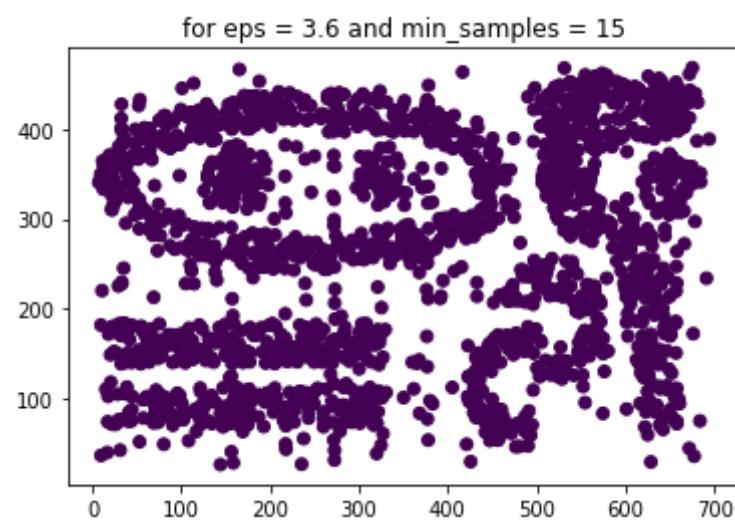
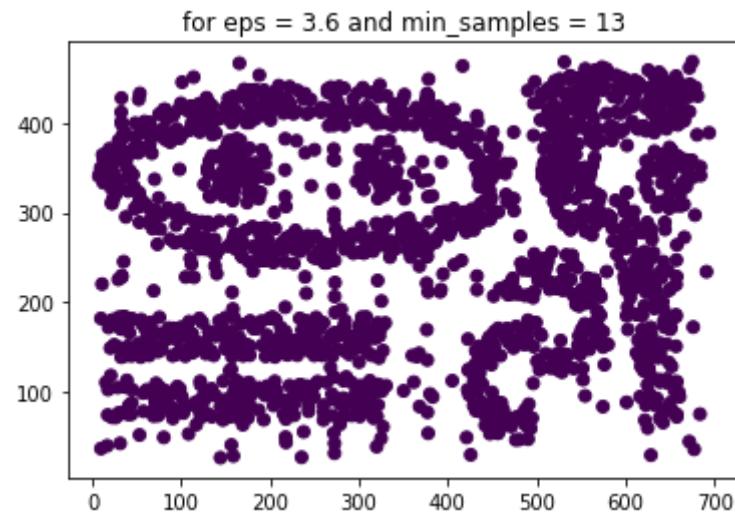


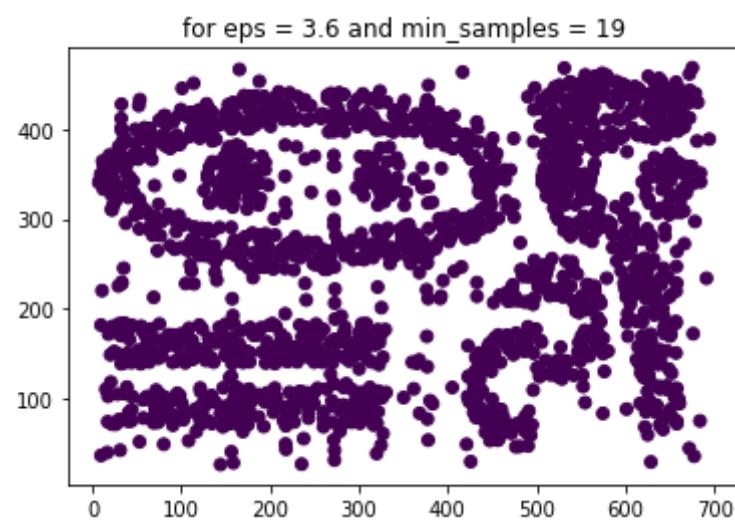
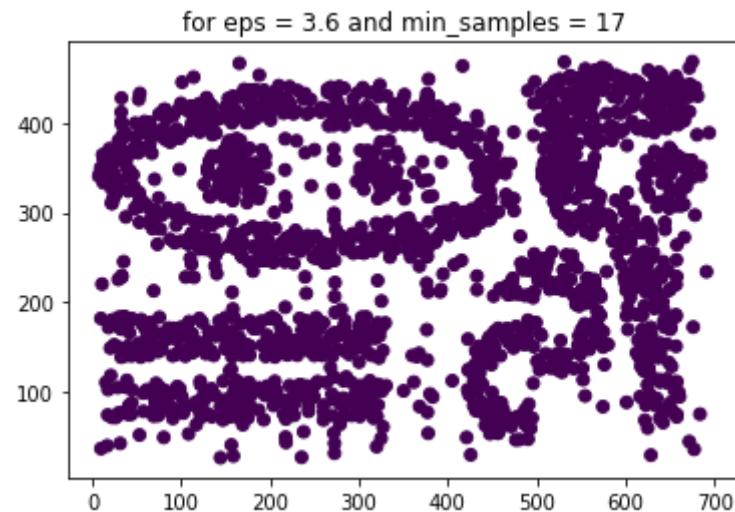


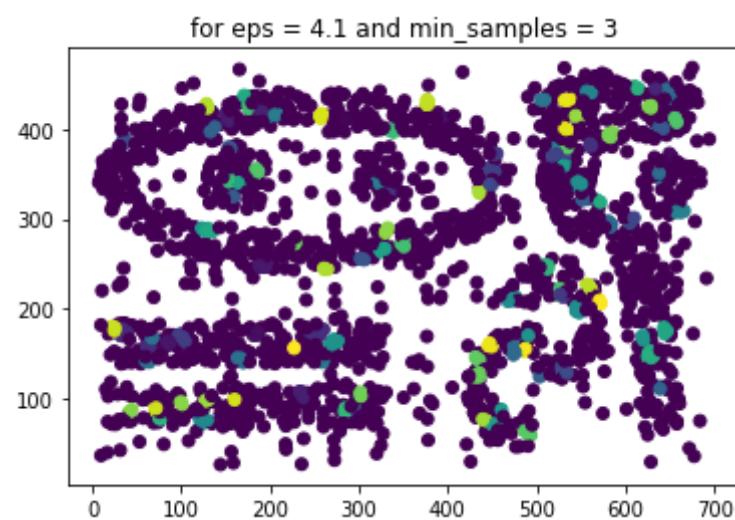
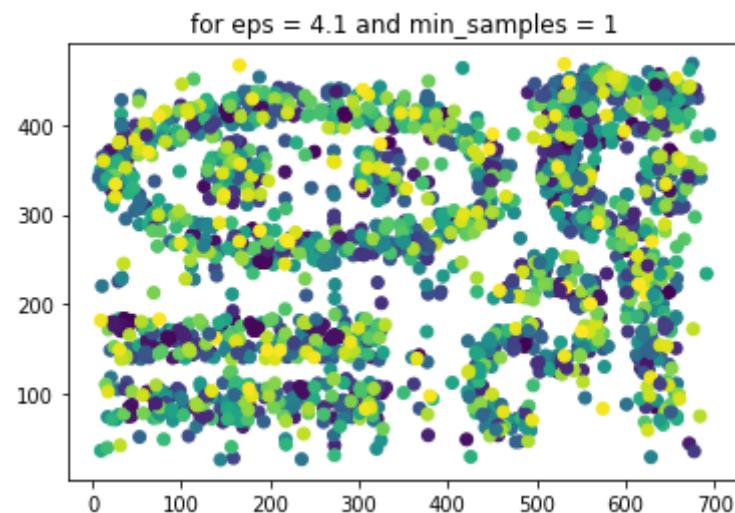


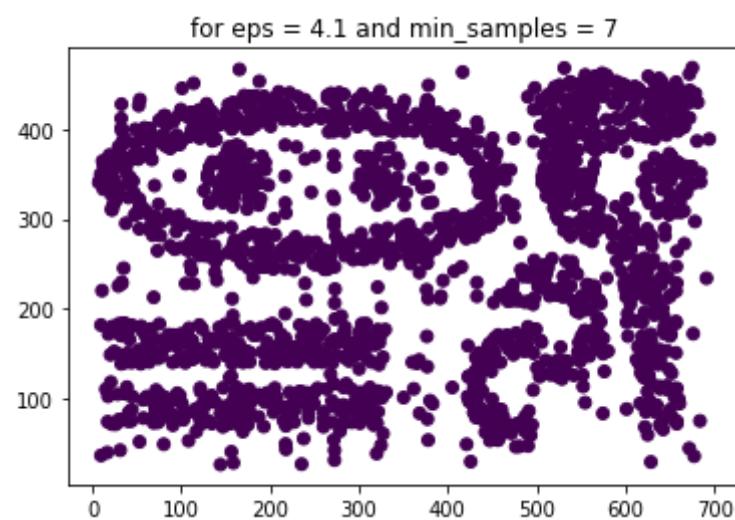
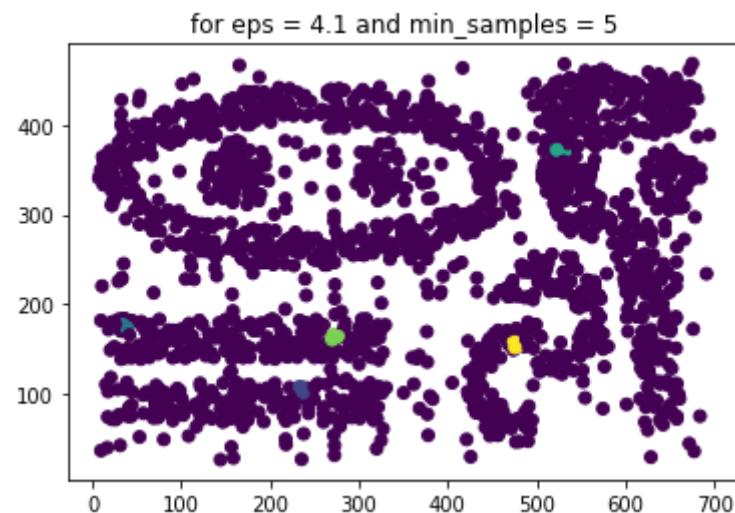


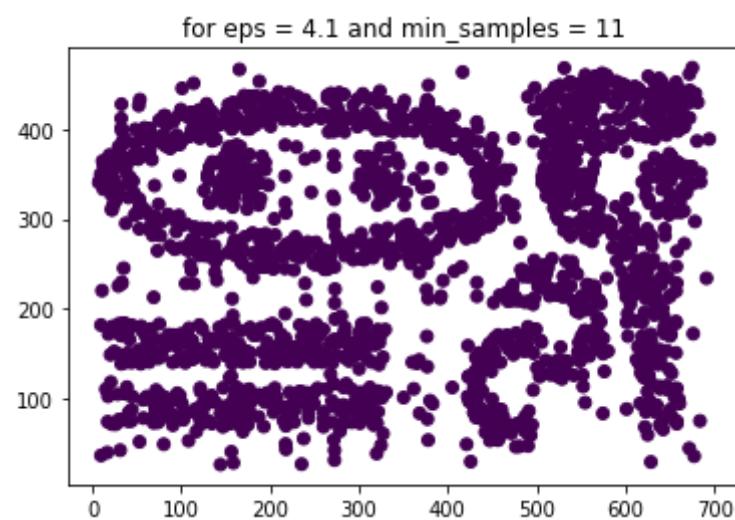
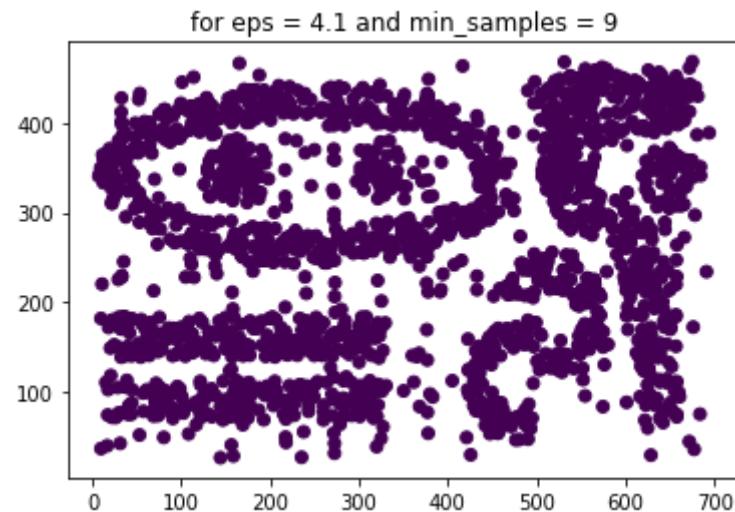


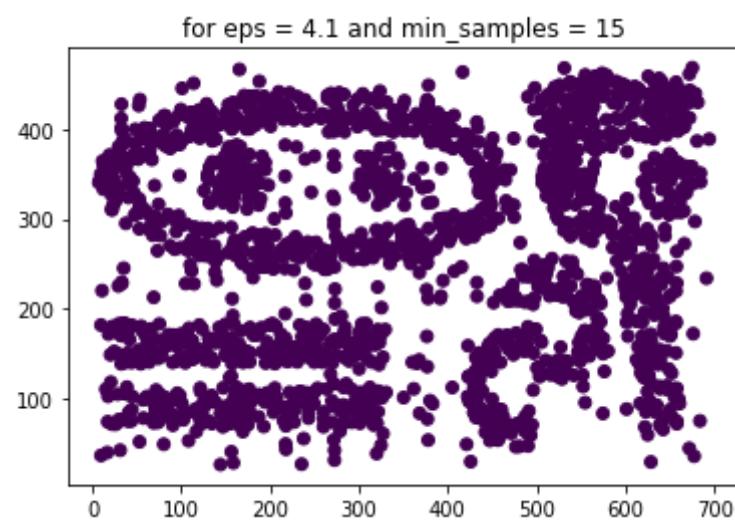
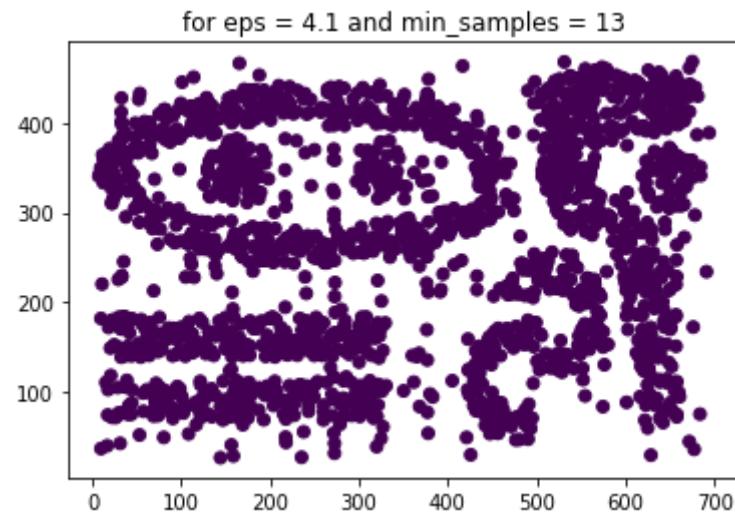




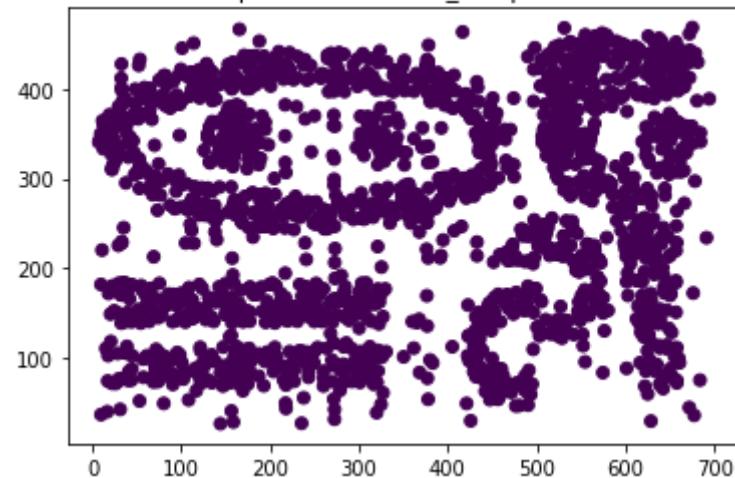




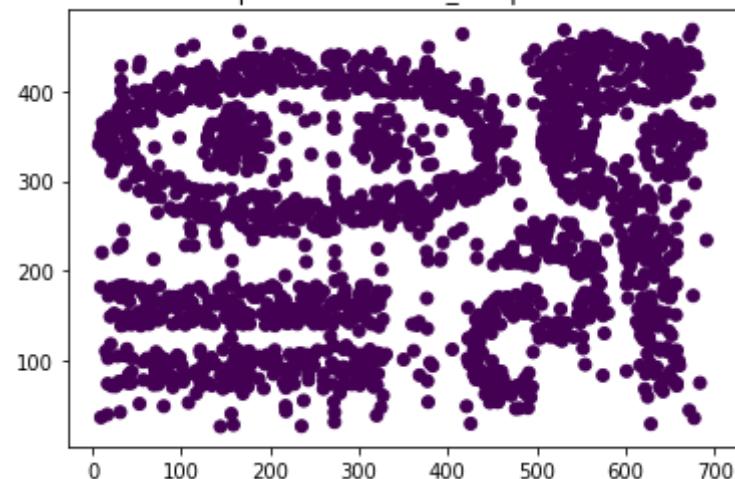


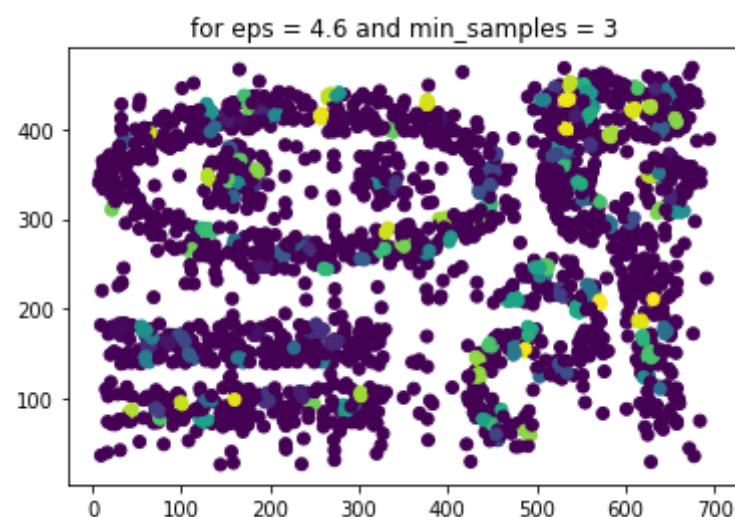
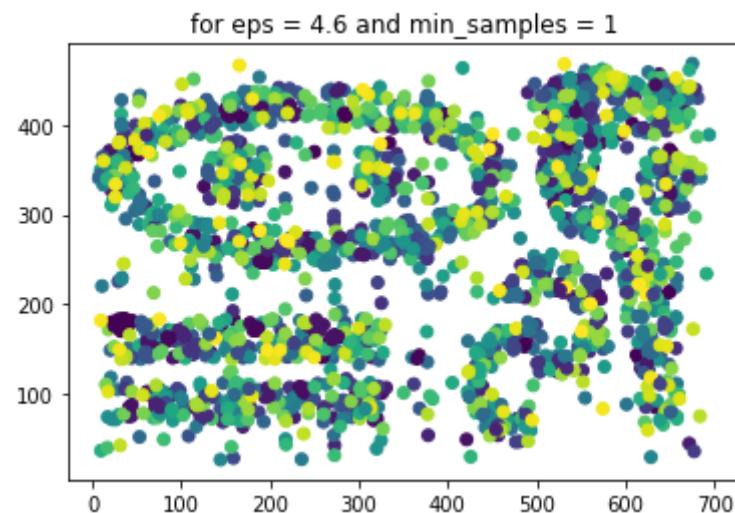


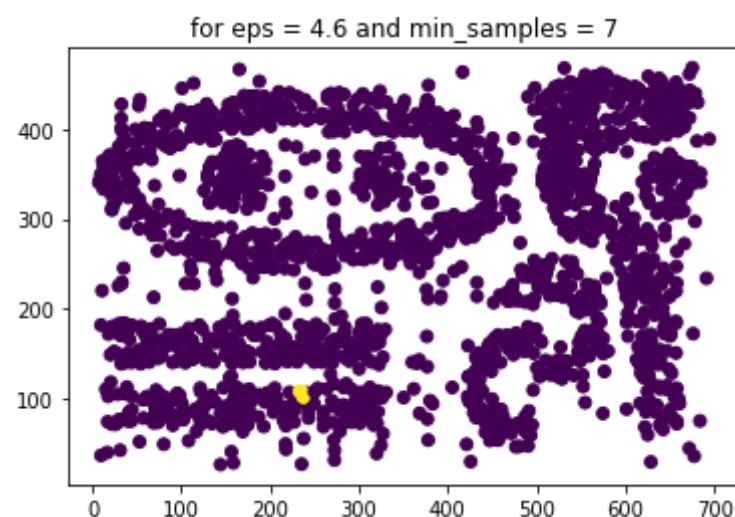
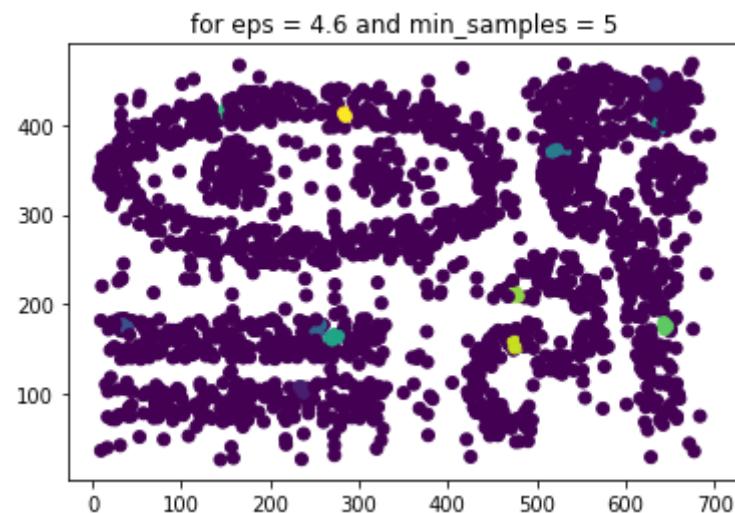
for eps = 4.1 and min_samples = 17

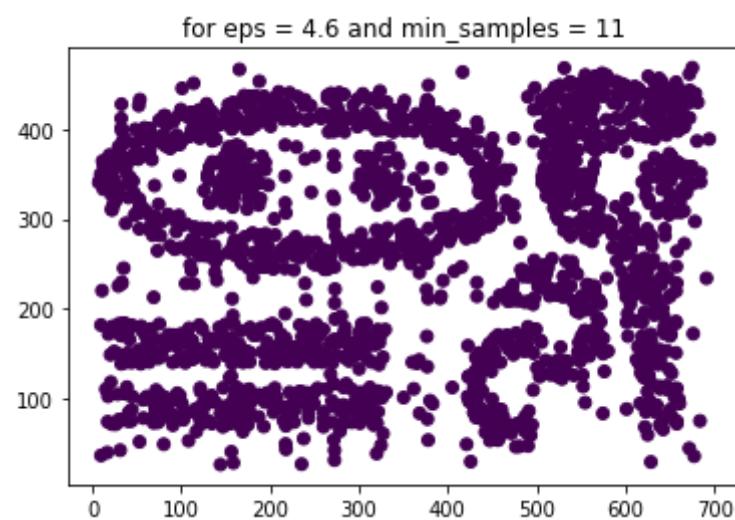
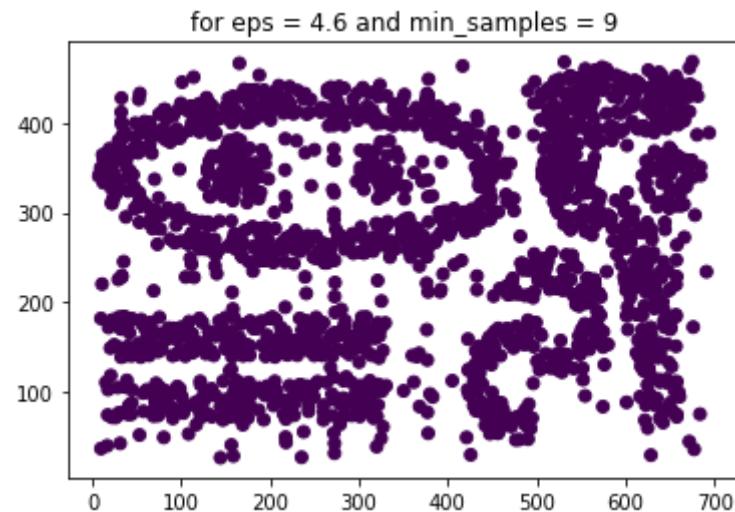


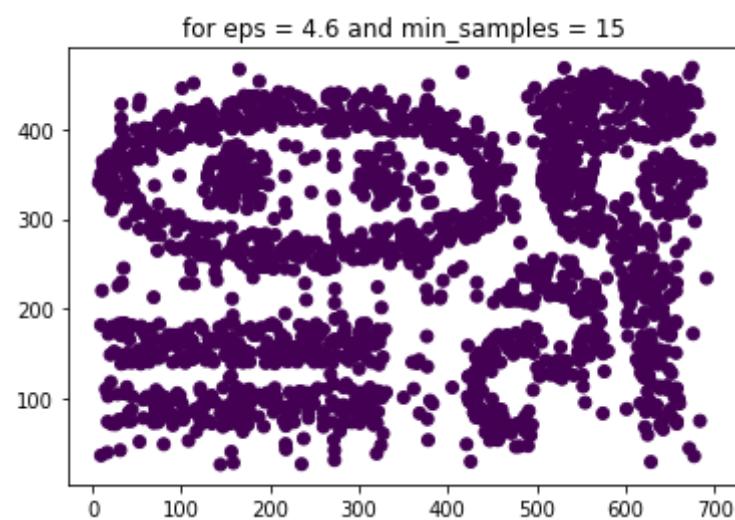
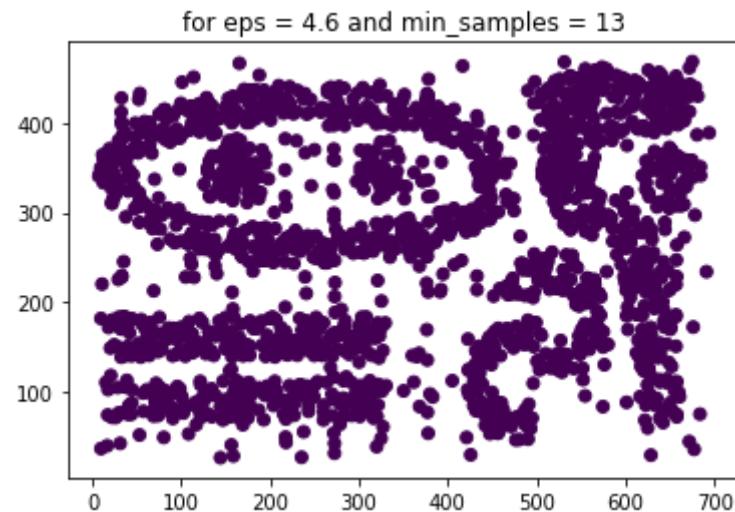
for eps = 4.1 and min_samples = 19

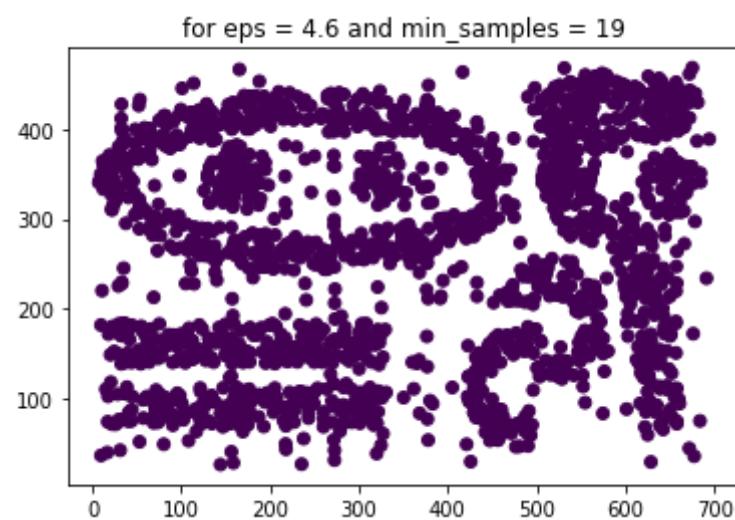
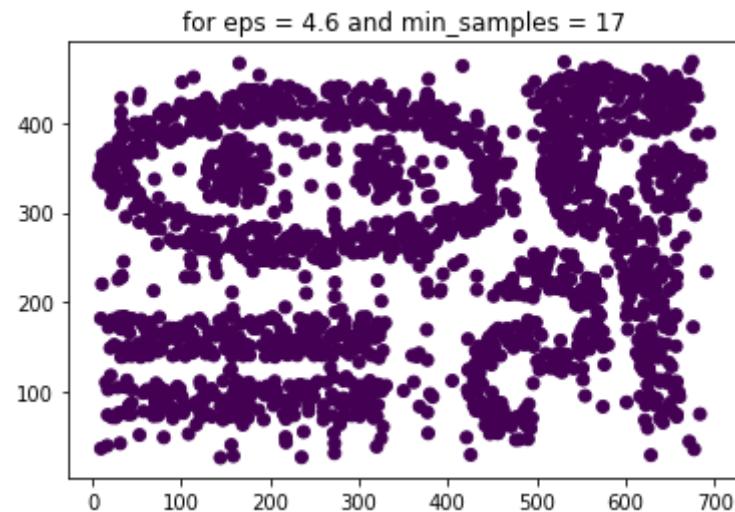


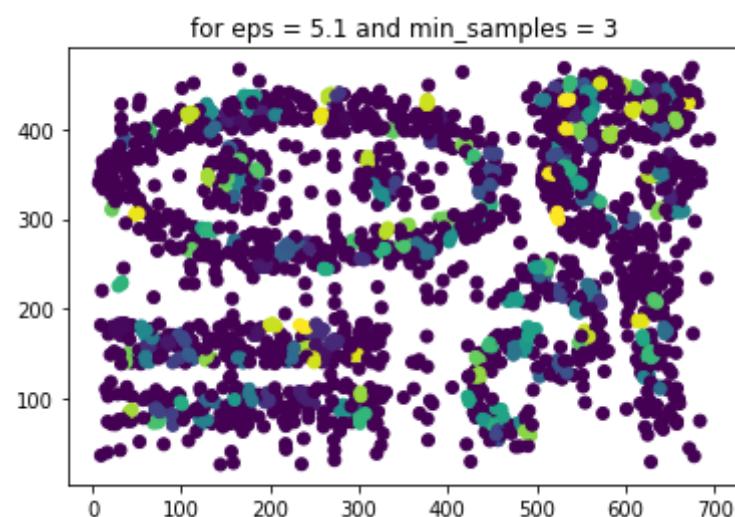
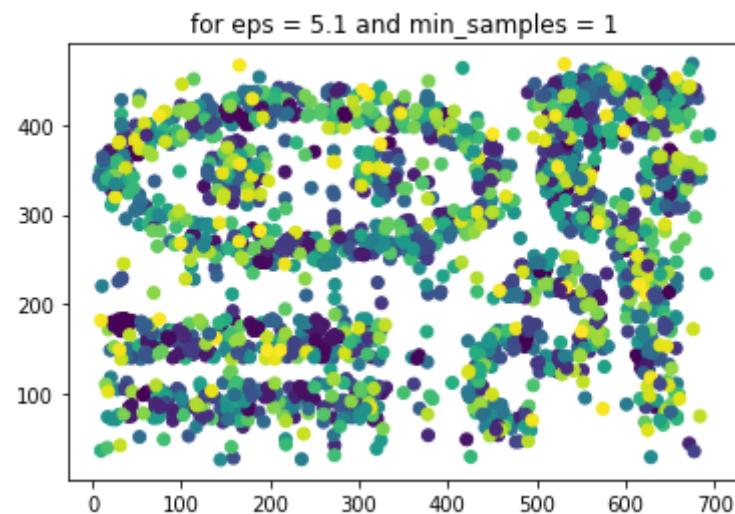


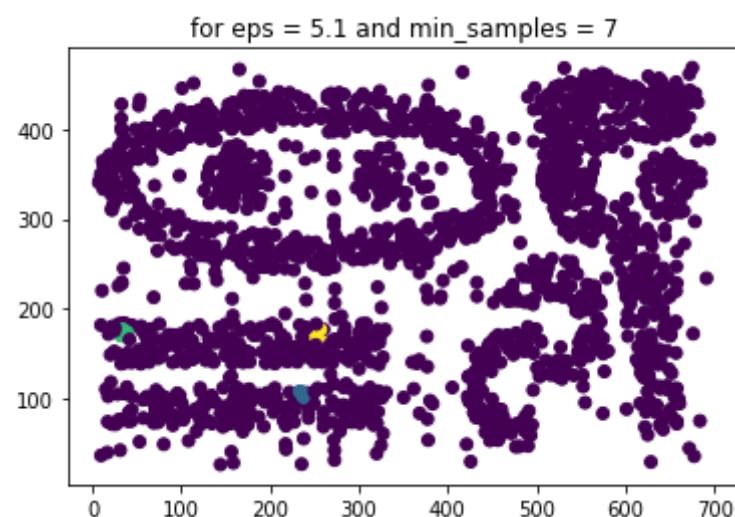
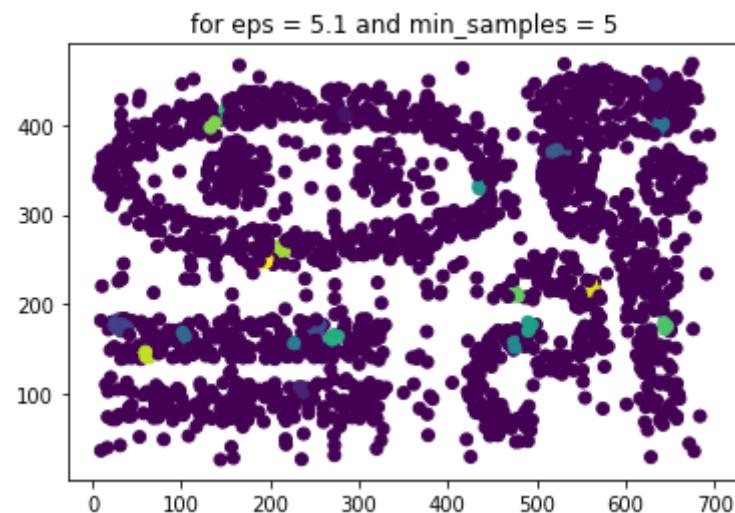


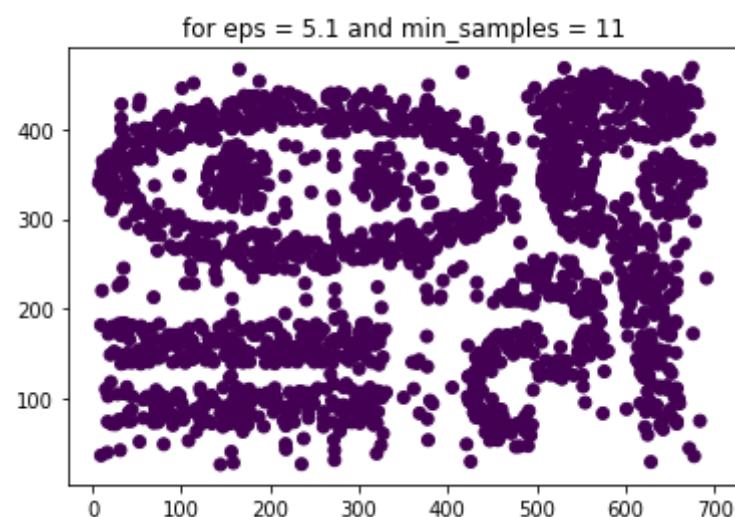
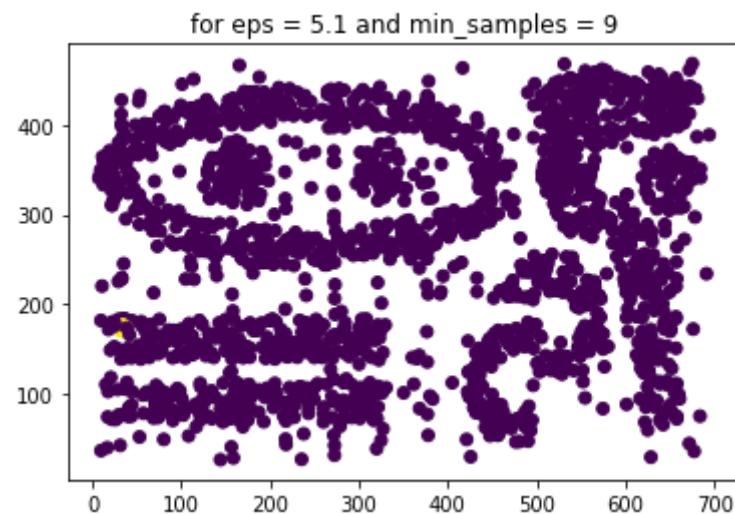


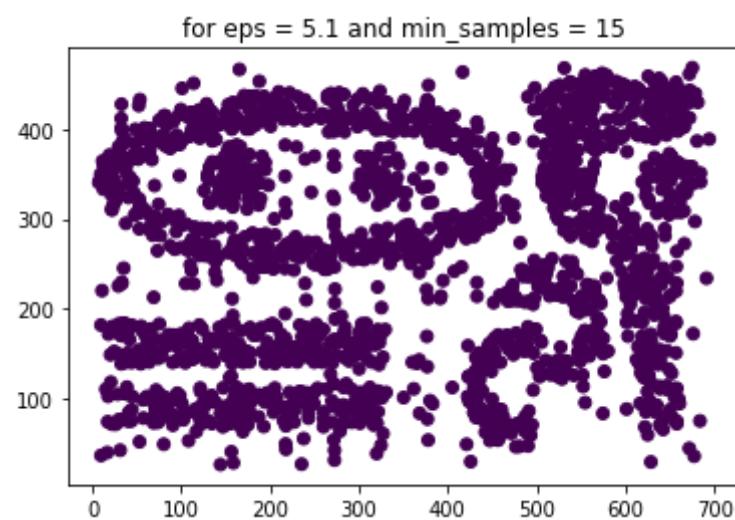
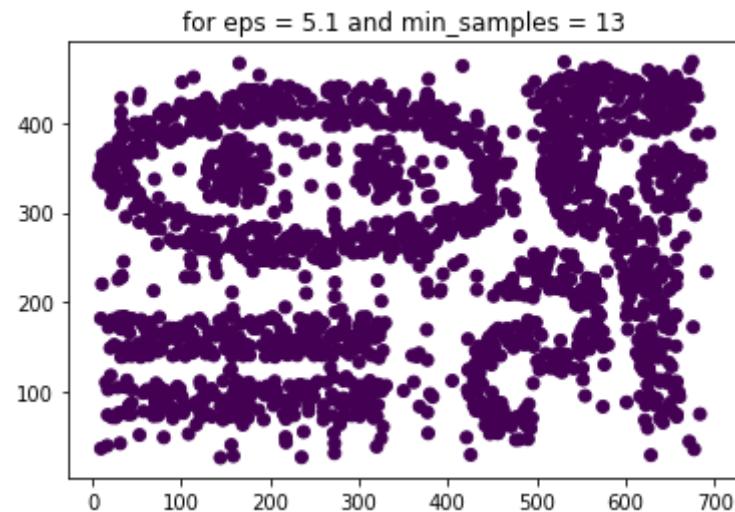


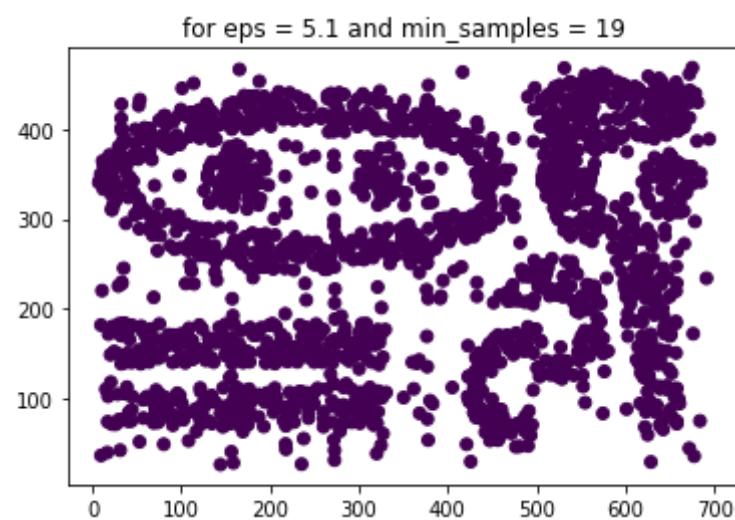
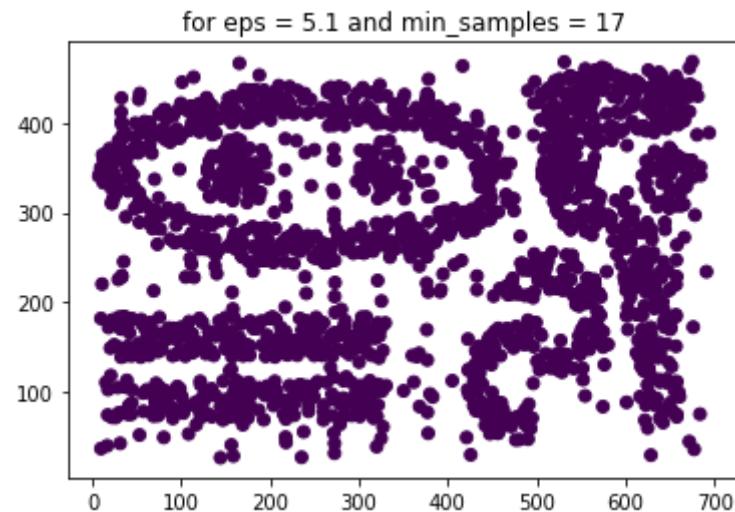


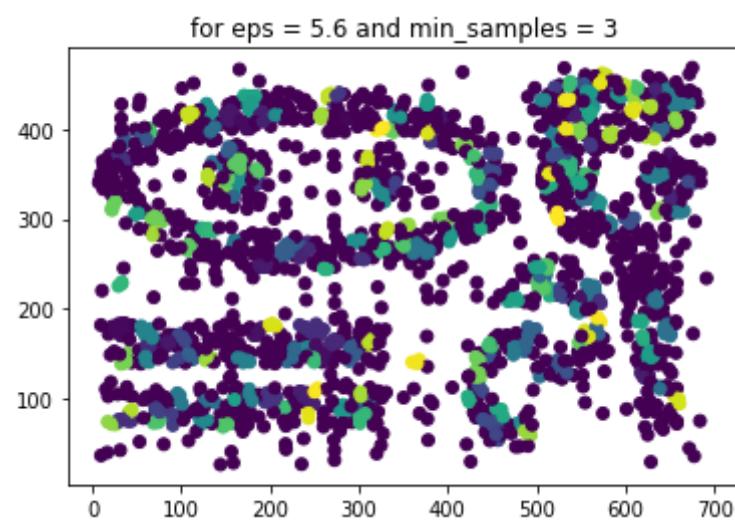
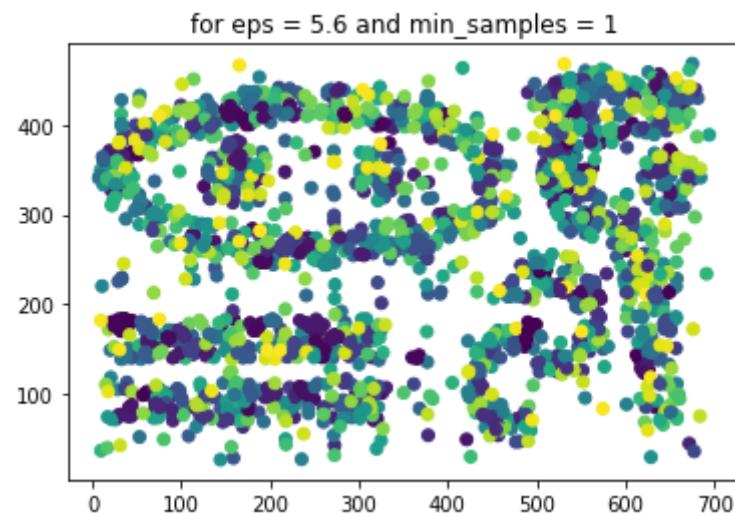


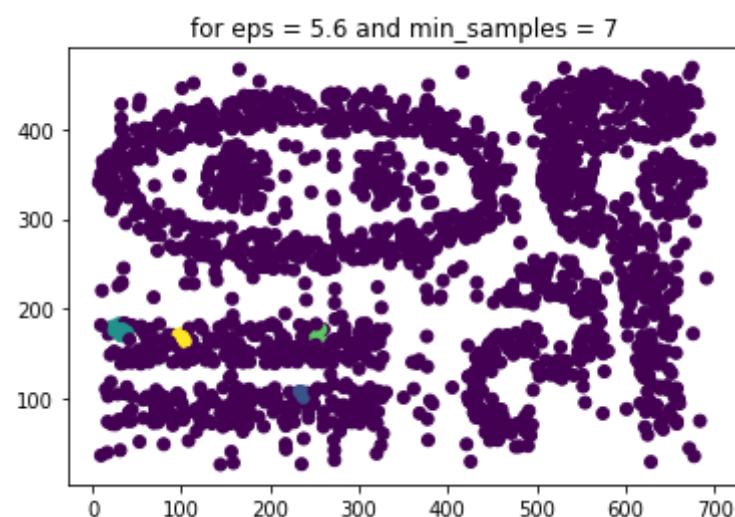
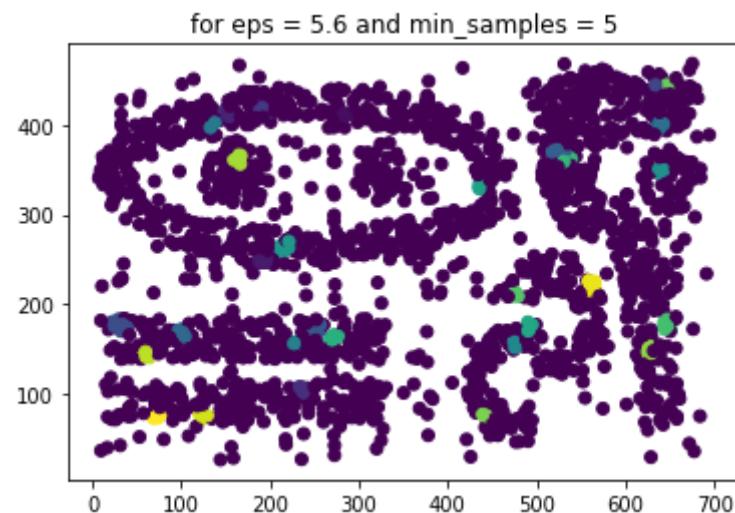


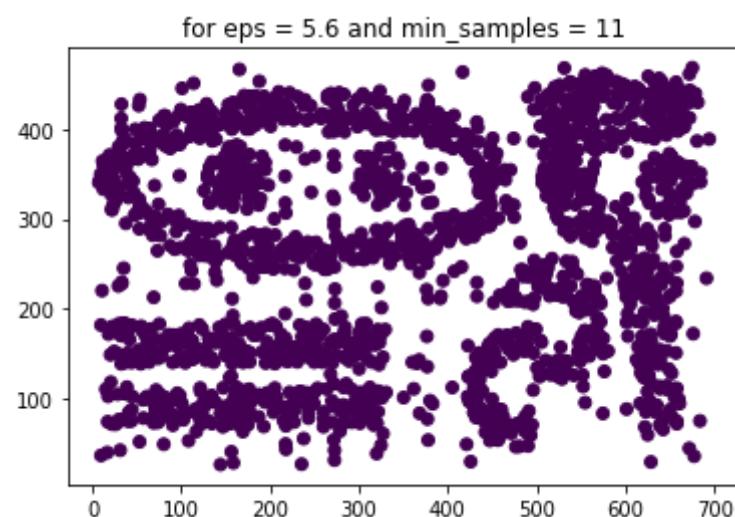
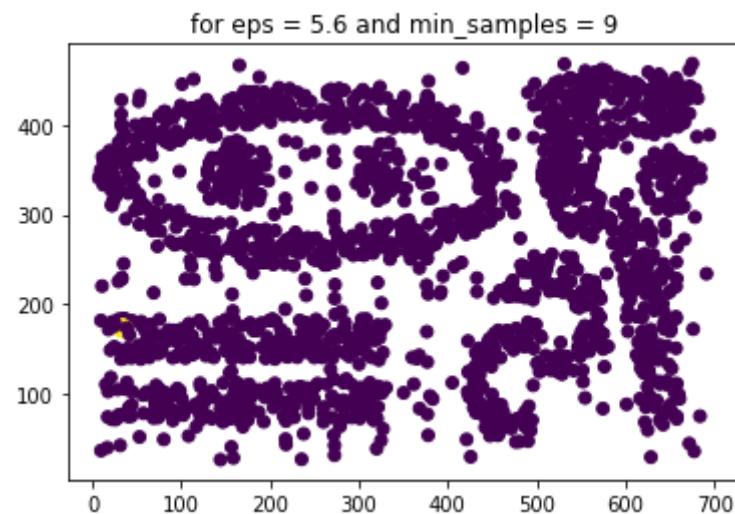


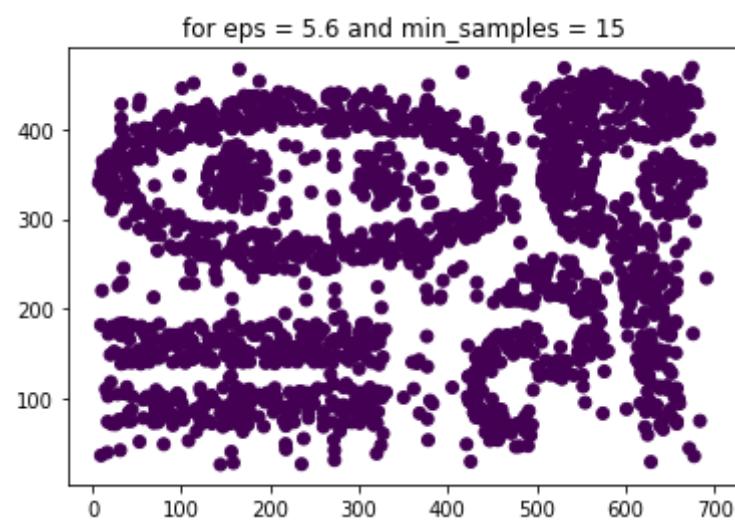
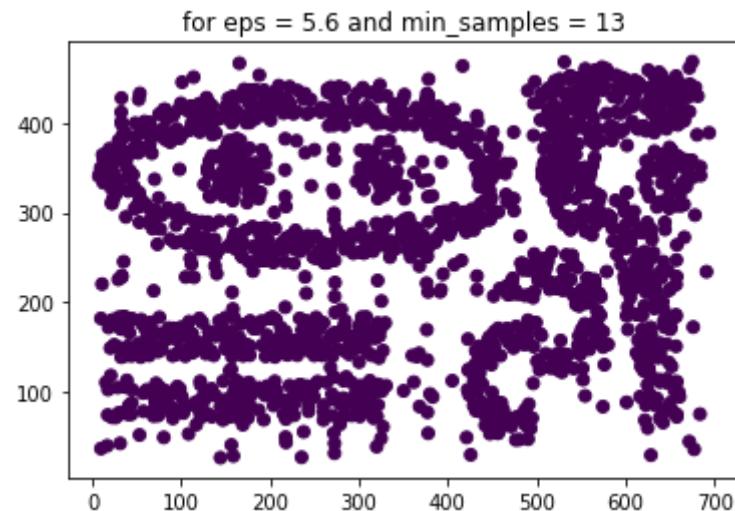


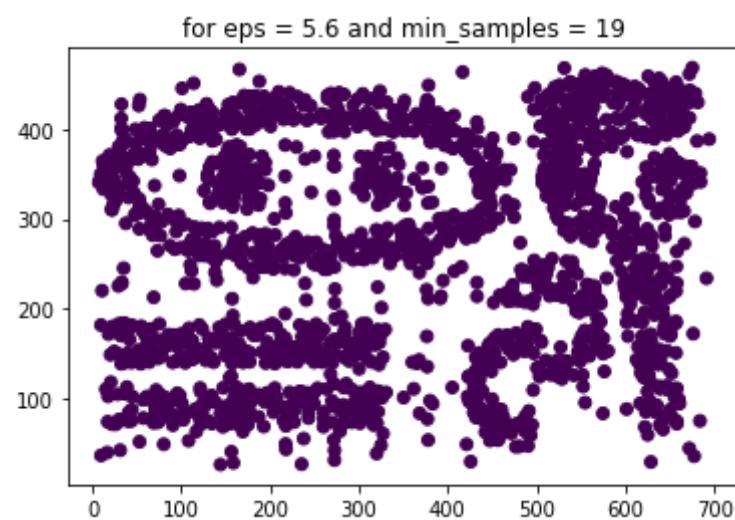
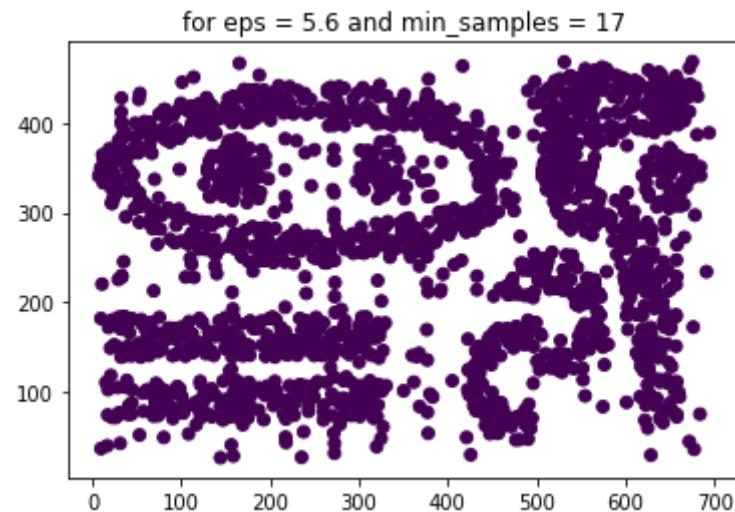


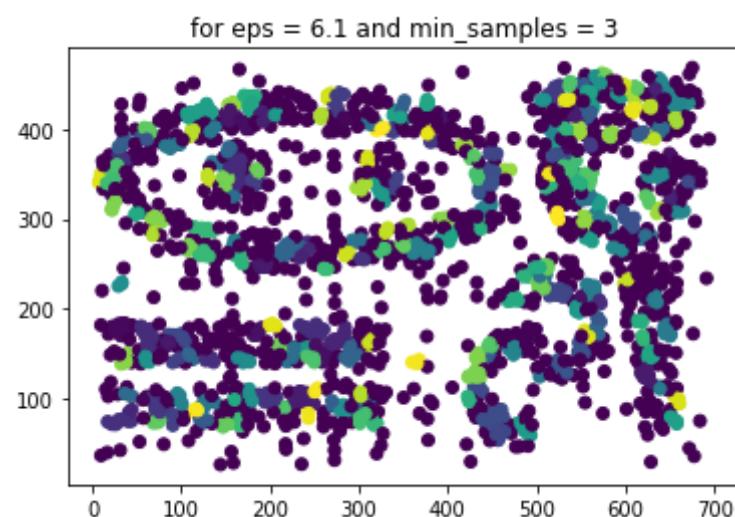
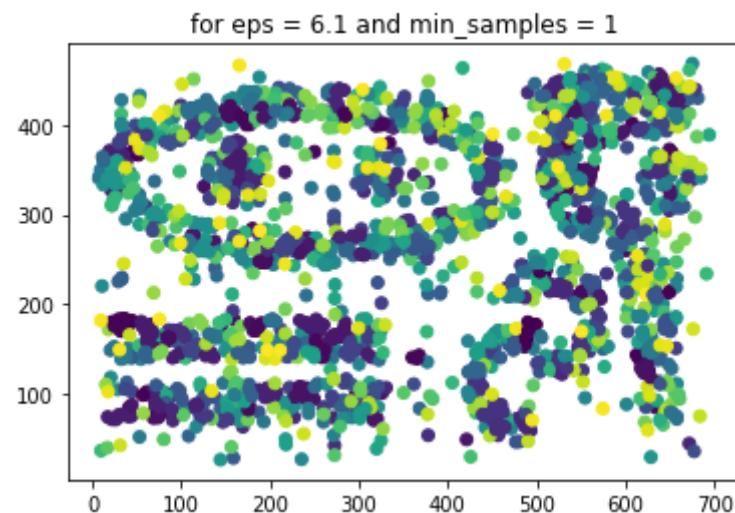


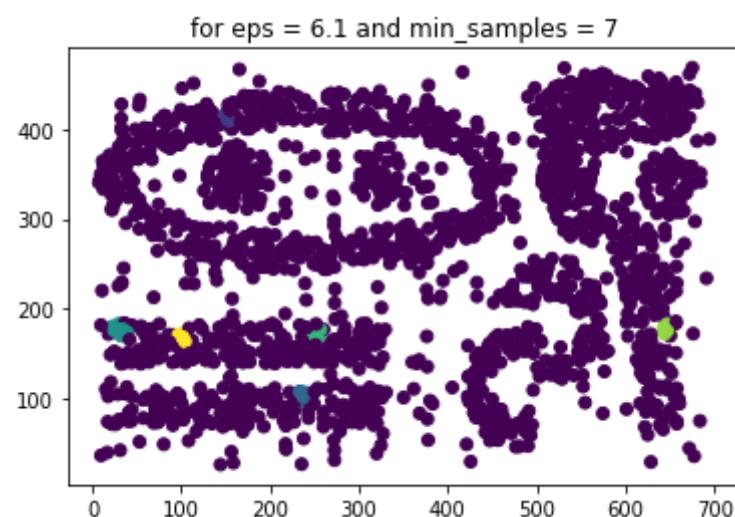
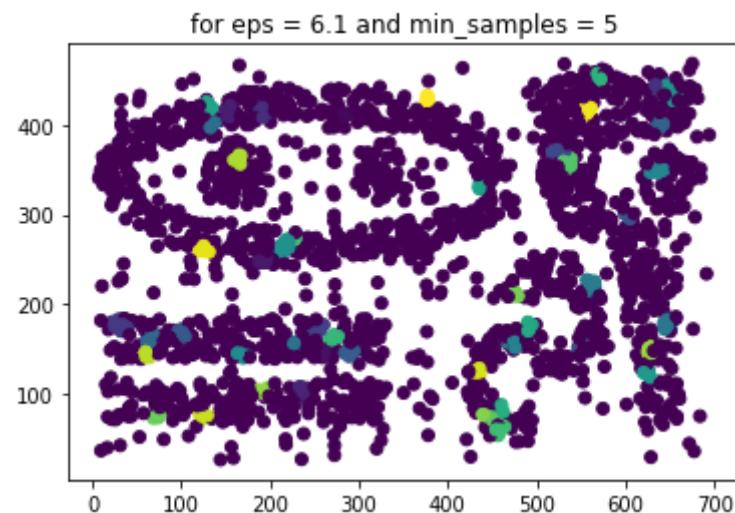


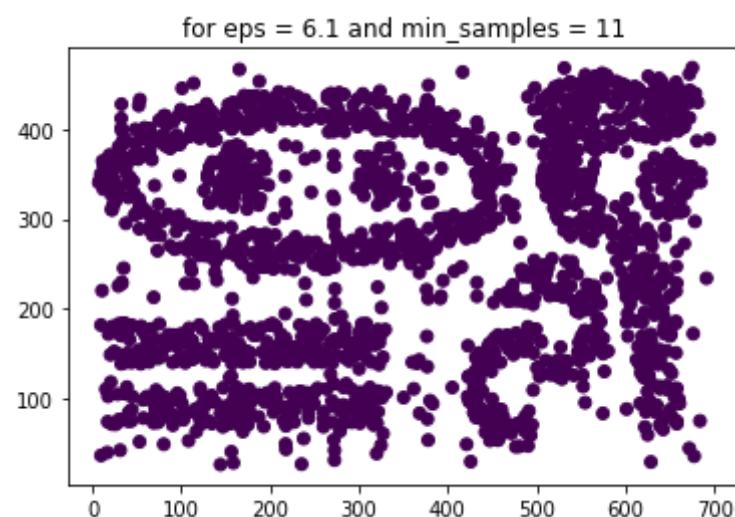
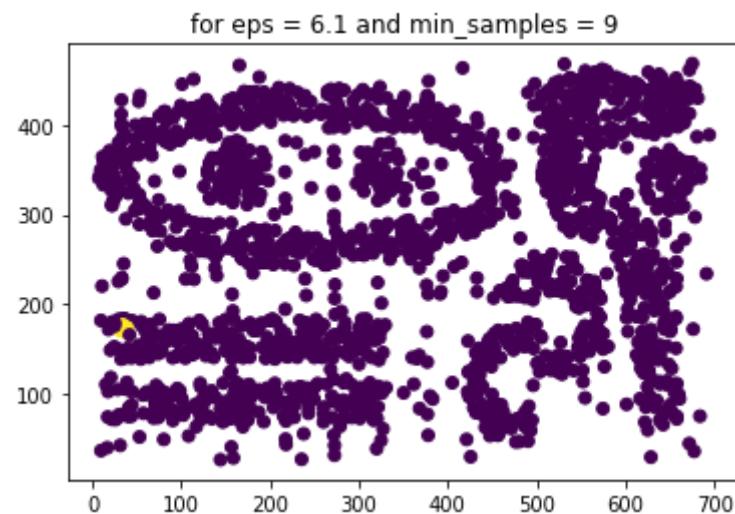


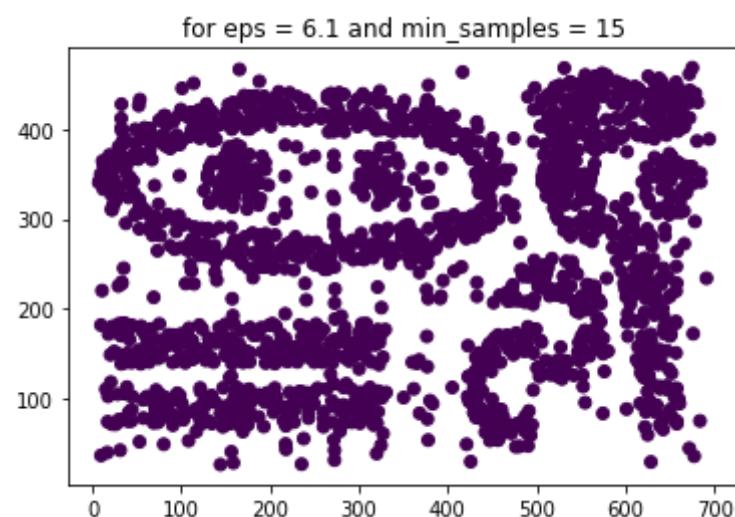
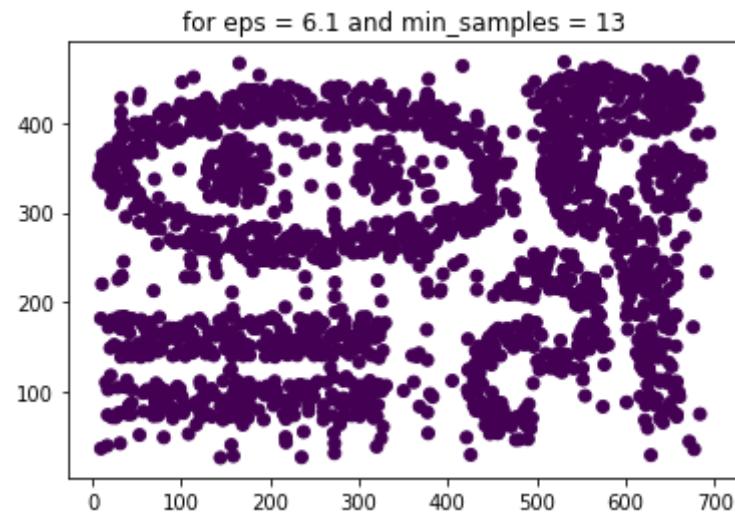


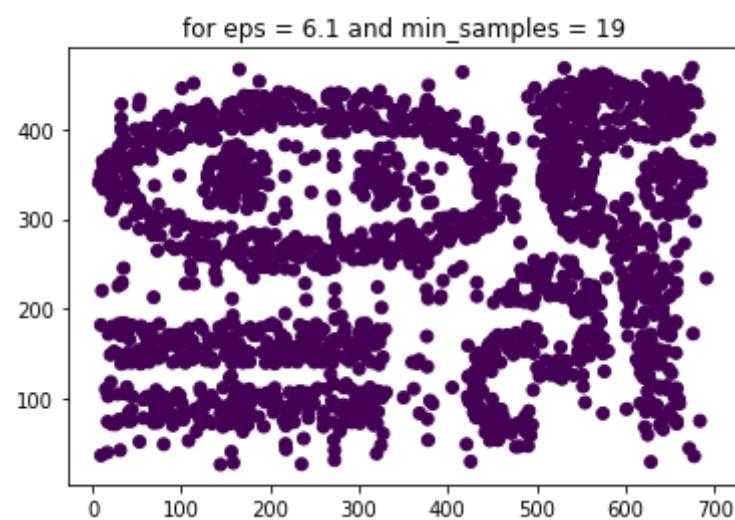
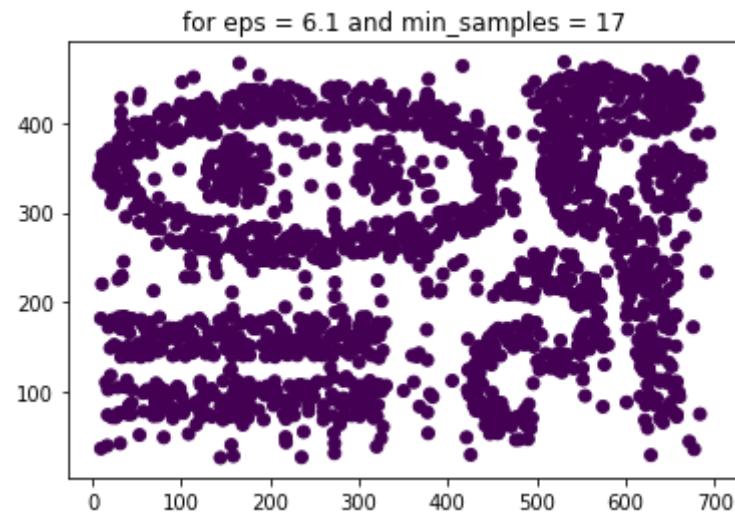


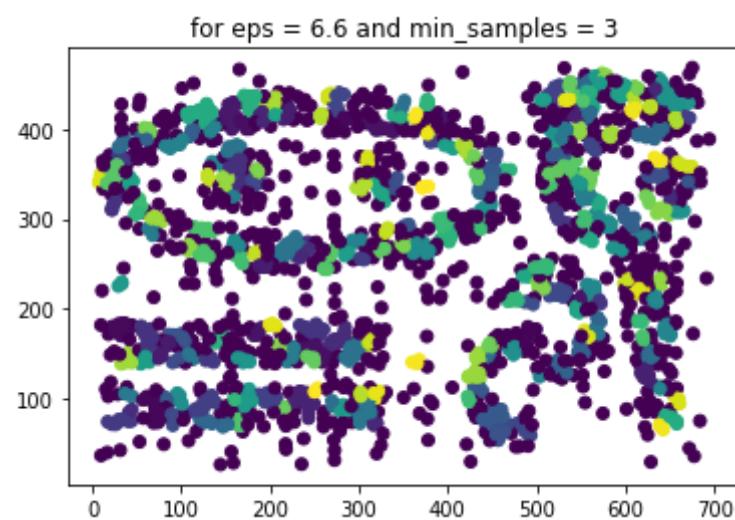
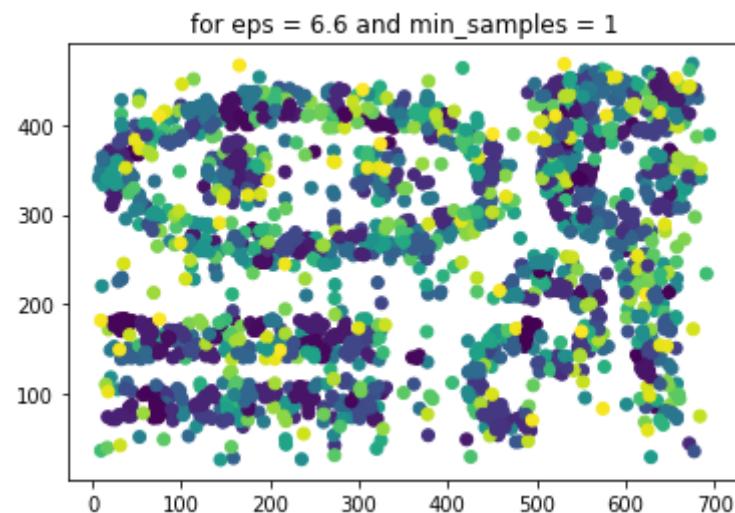


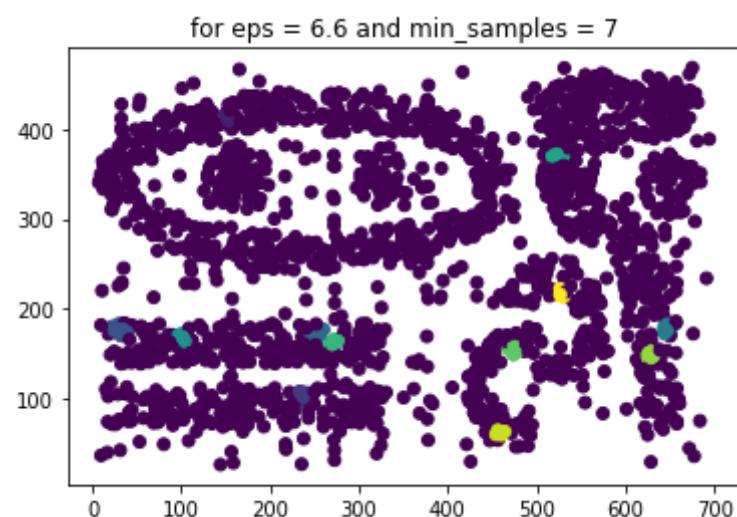
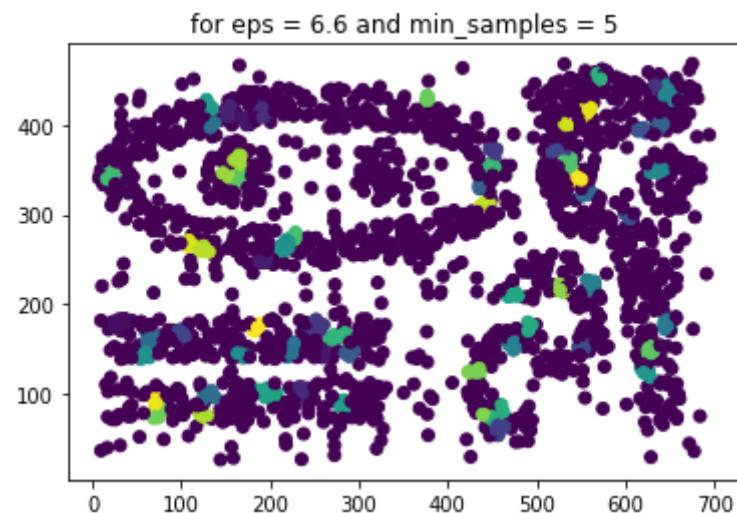


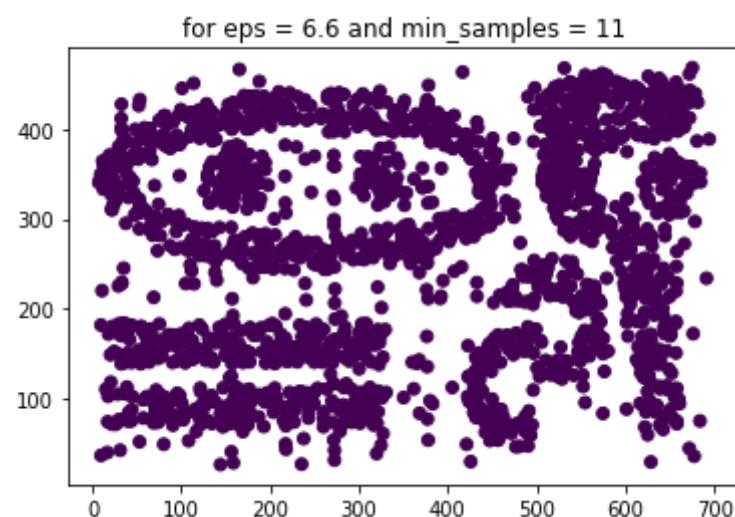
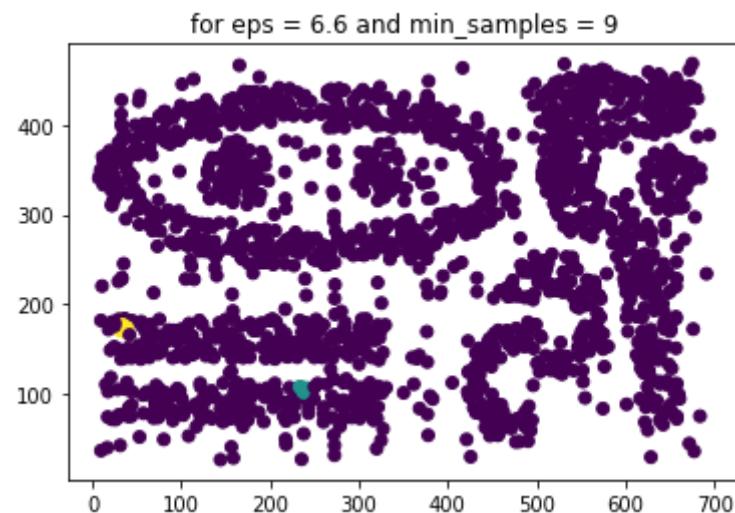


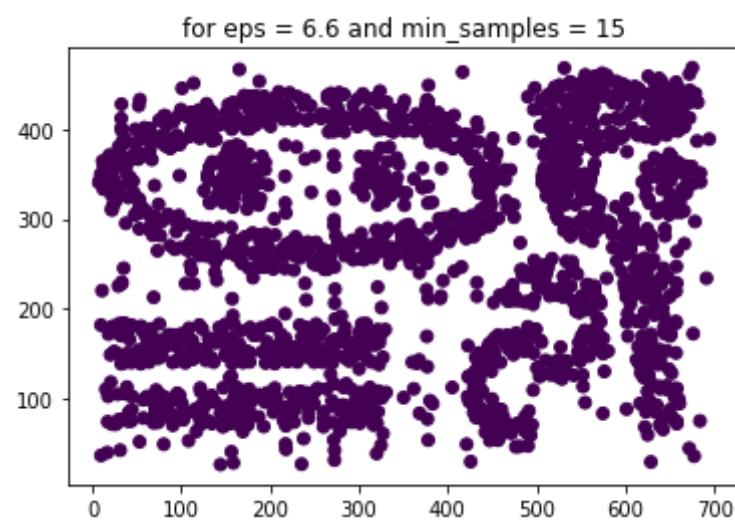
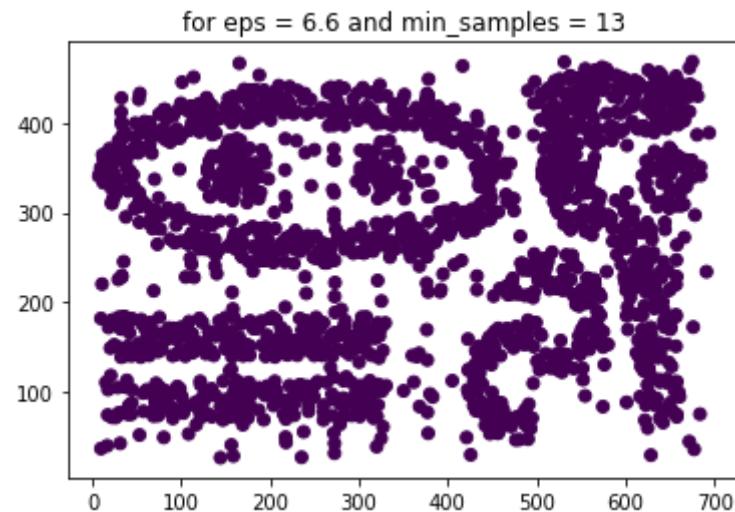


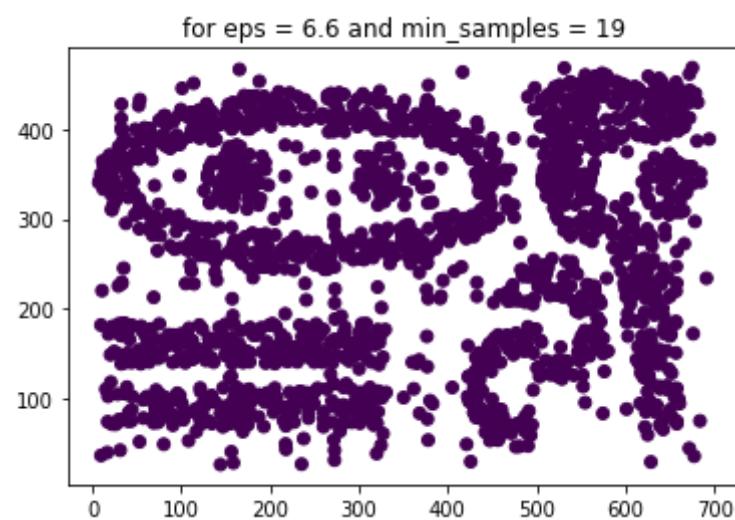
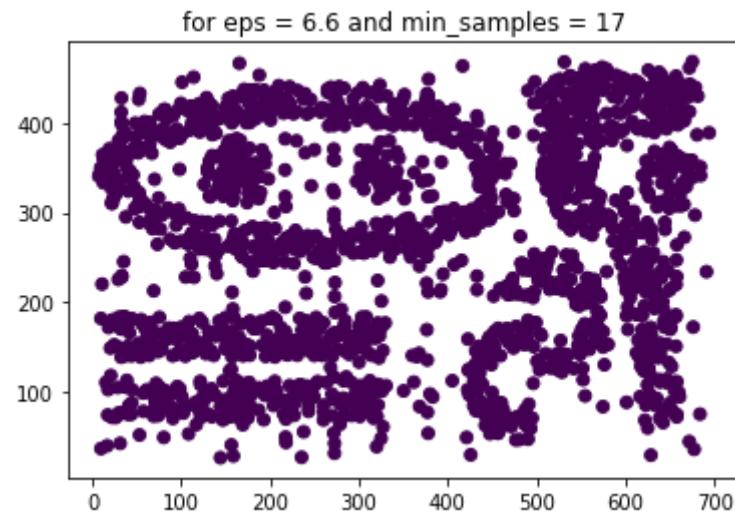




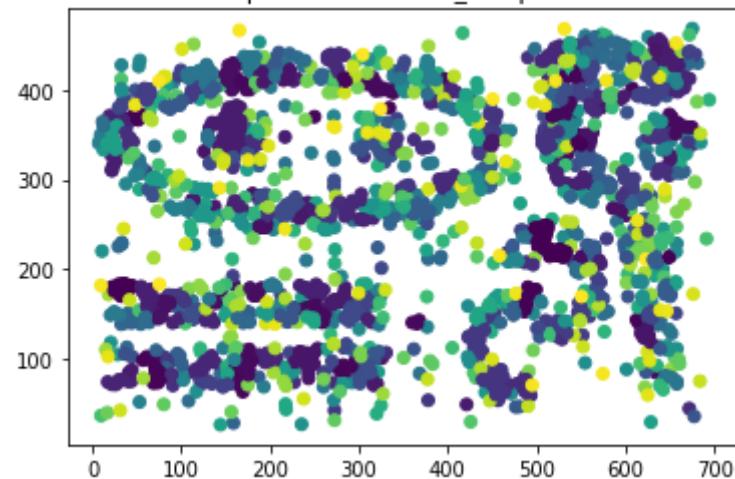




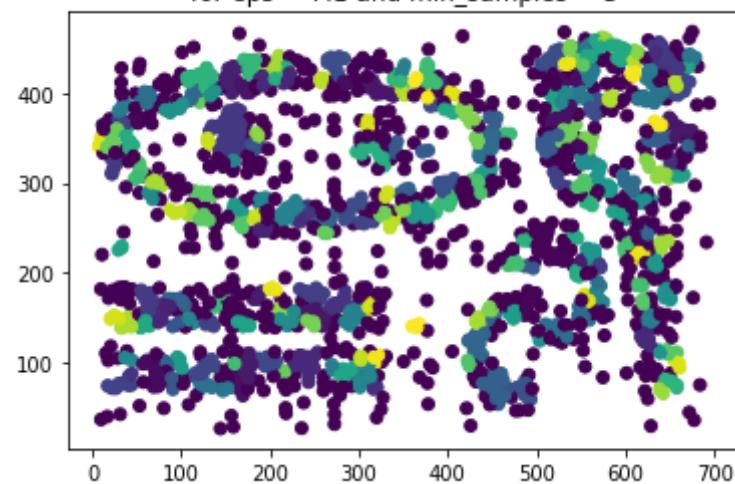


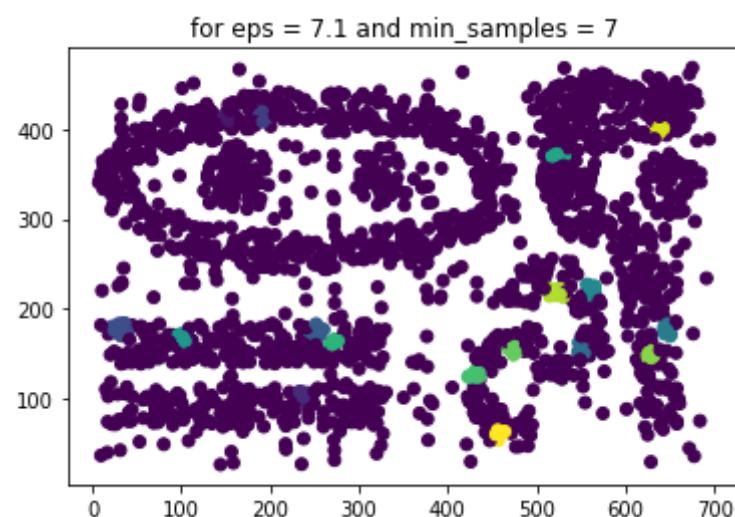
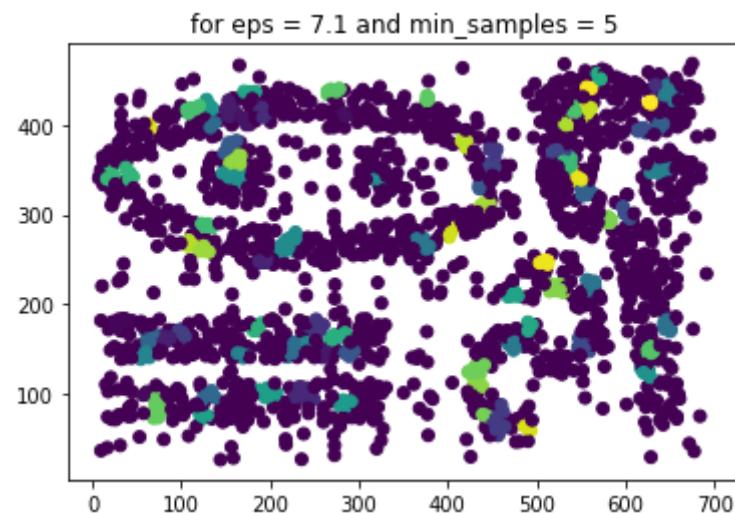


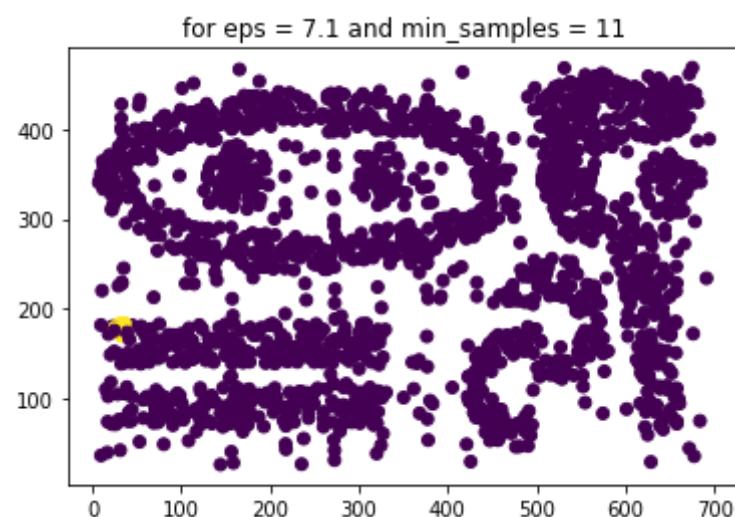
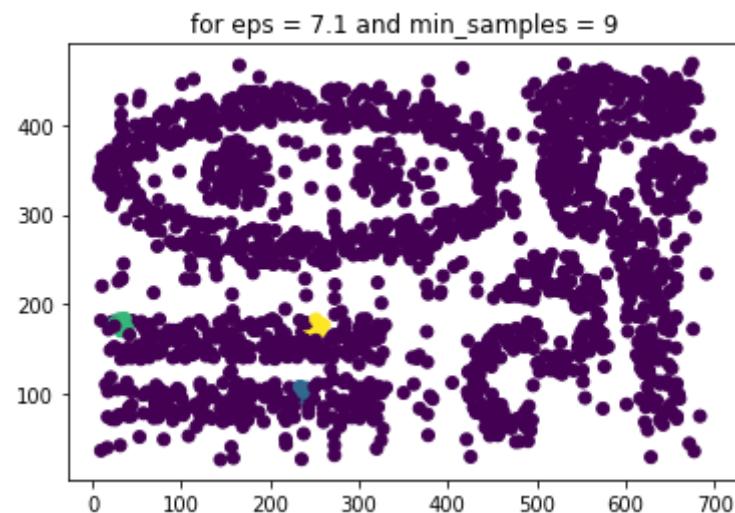
for eps = 7.1 and min_samples = 1

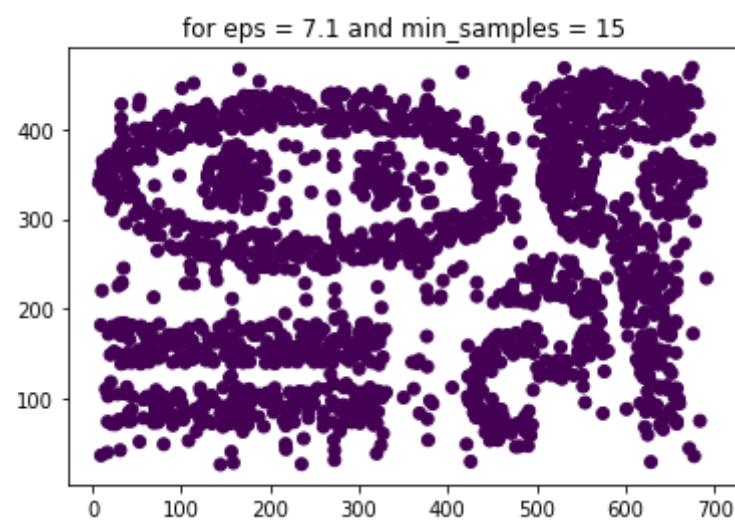
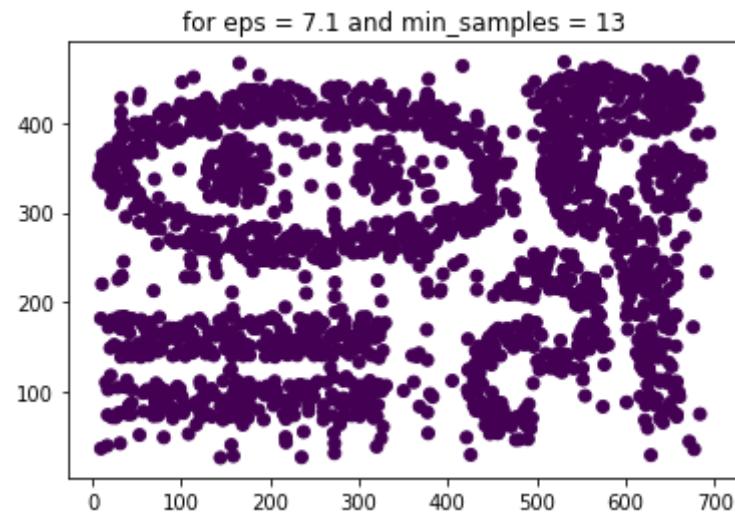


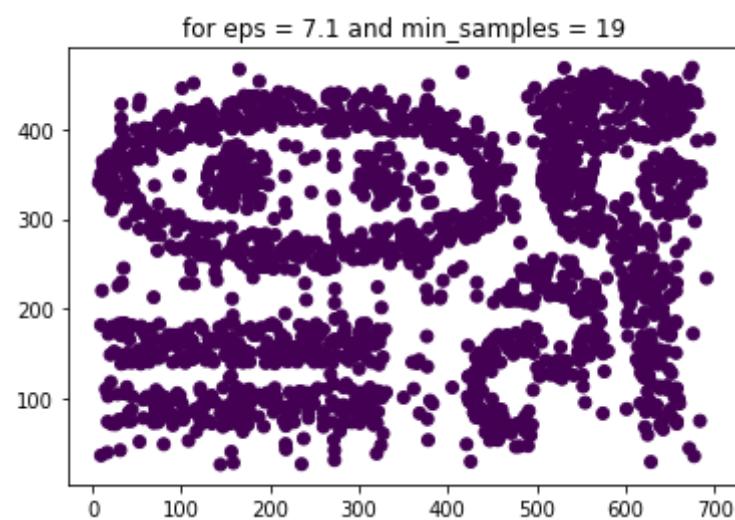
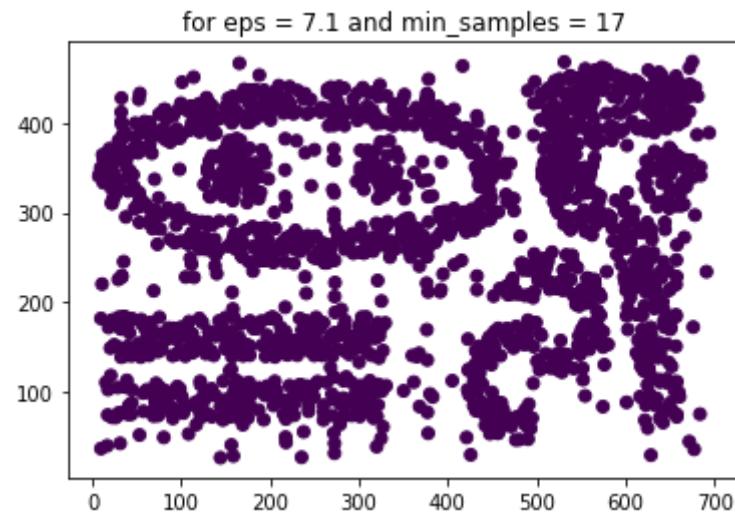
for eps = 7.1 and min_samples = 3

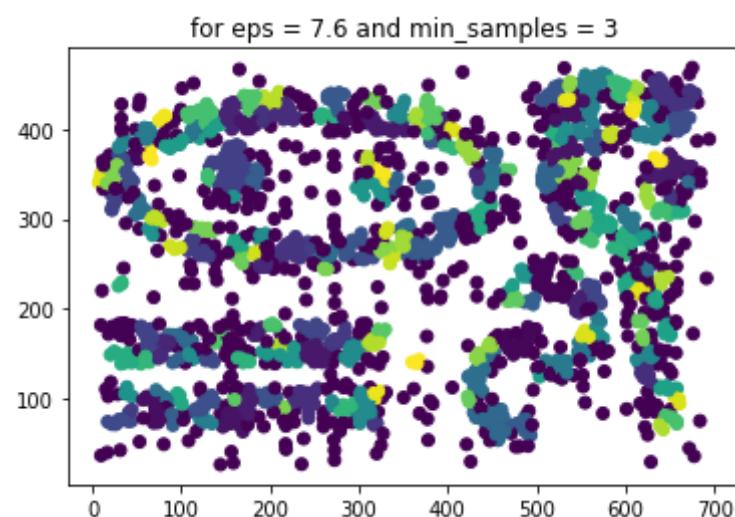
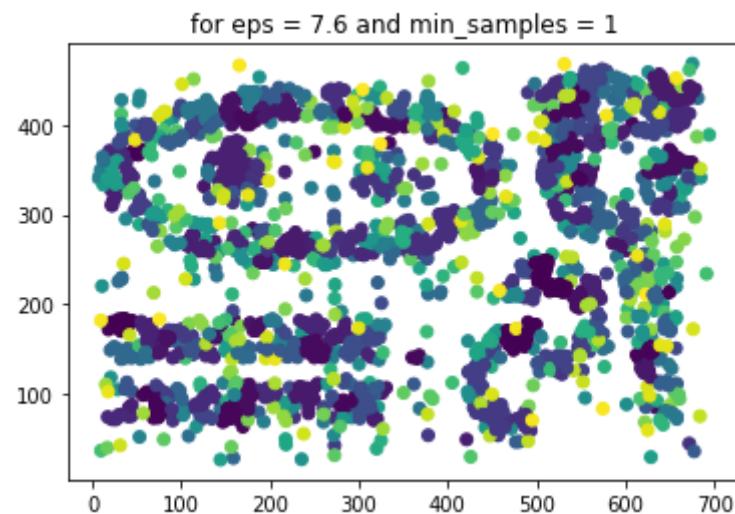


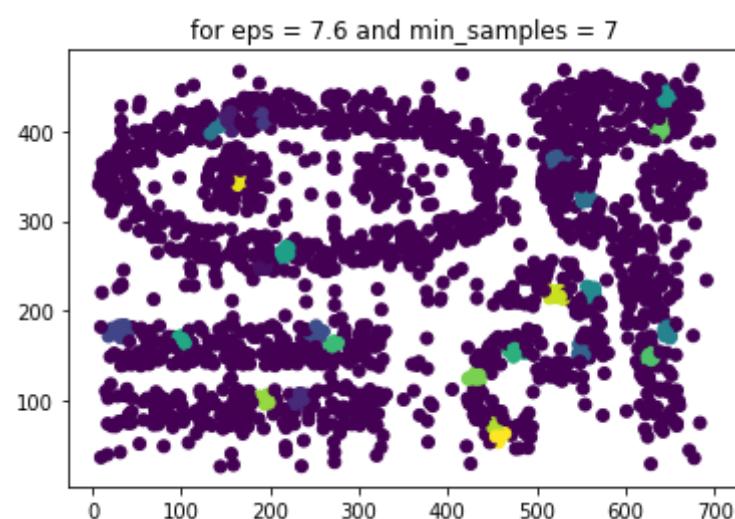
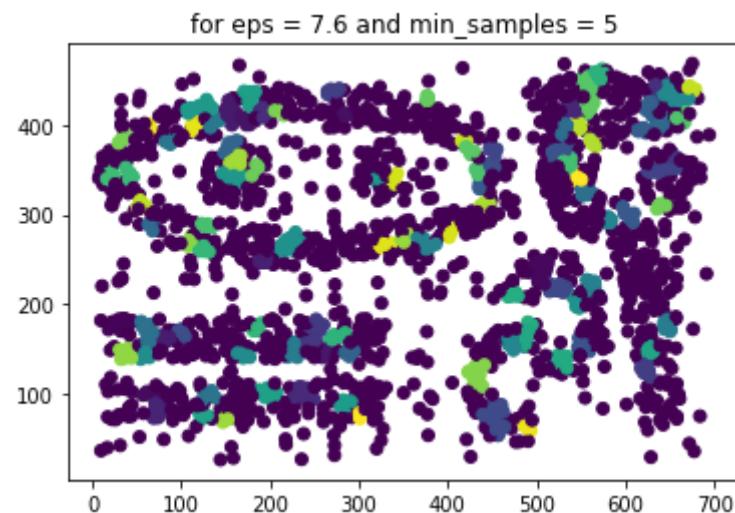


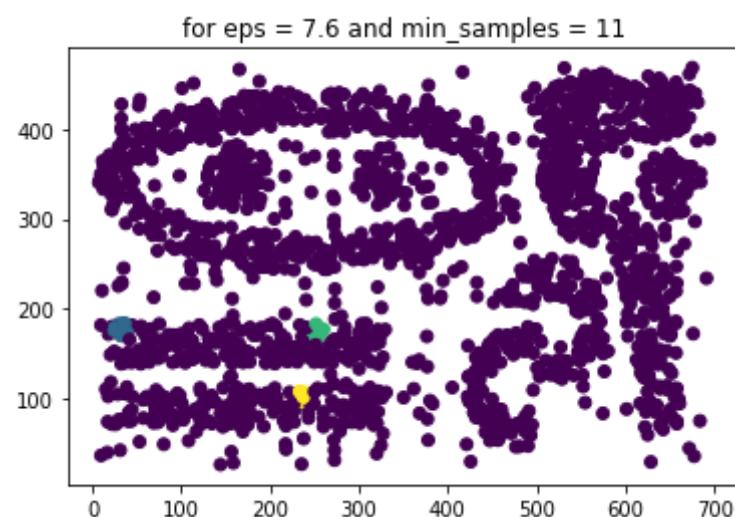
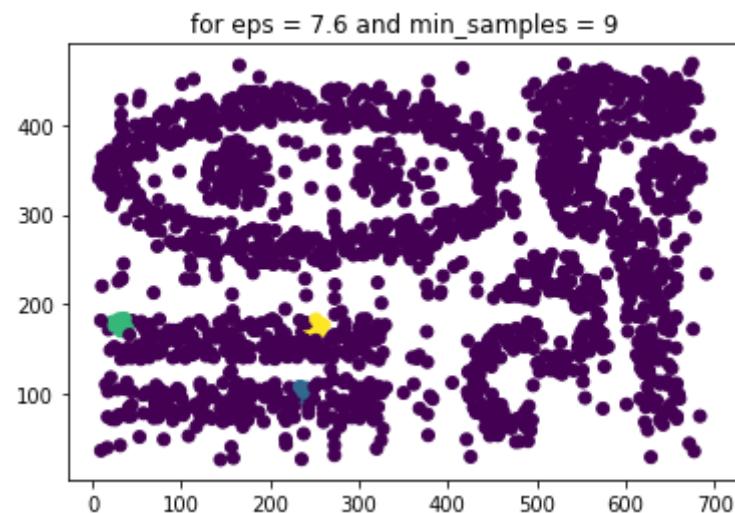


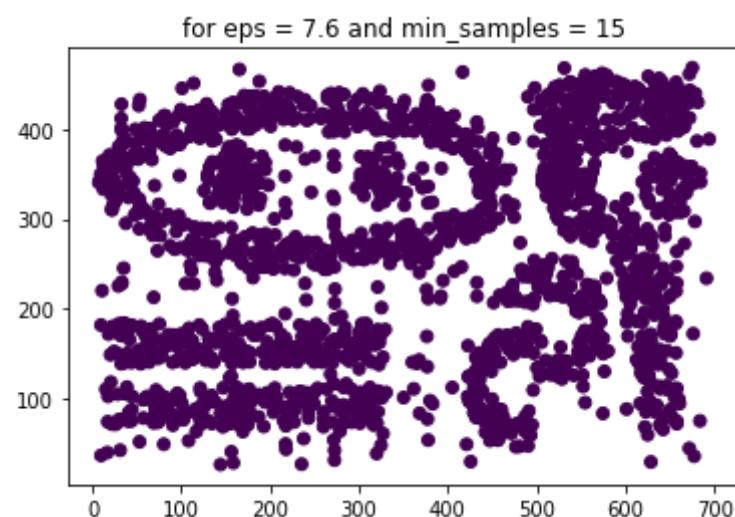
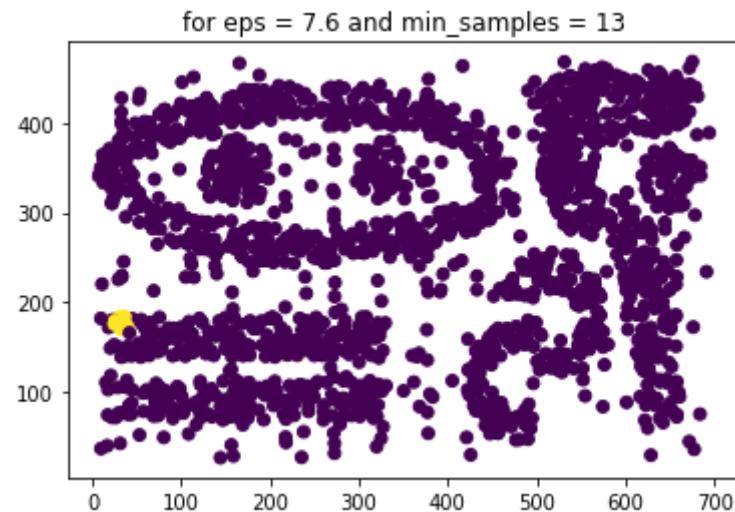


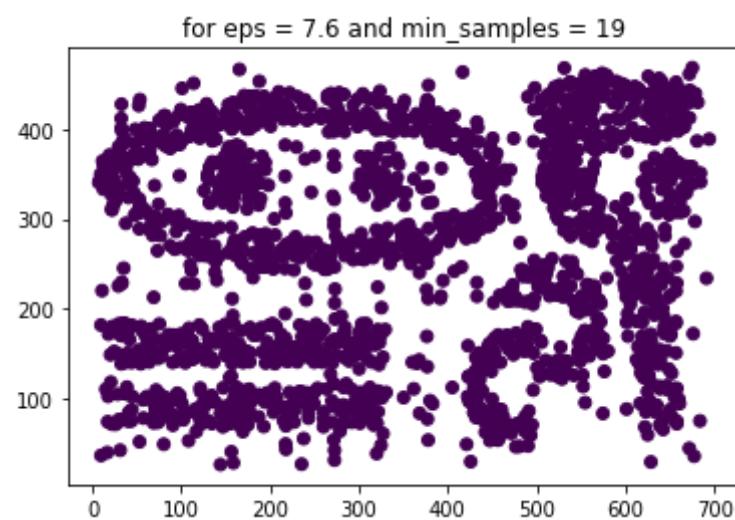
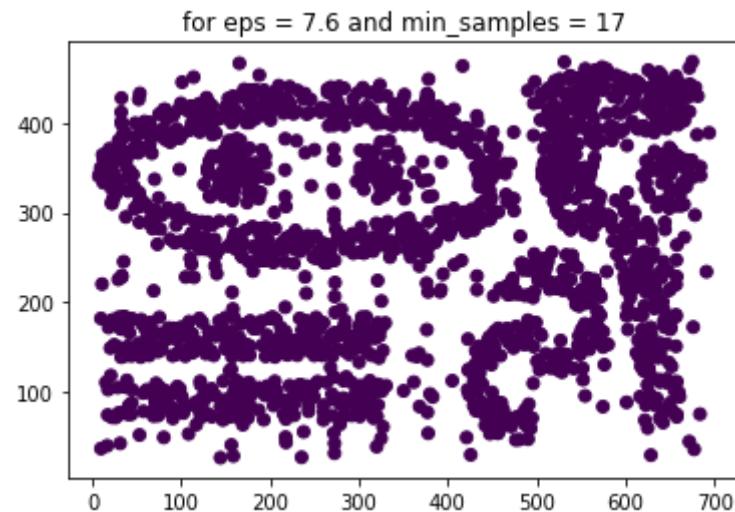


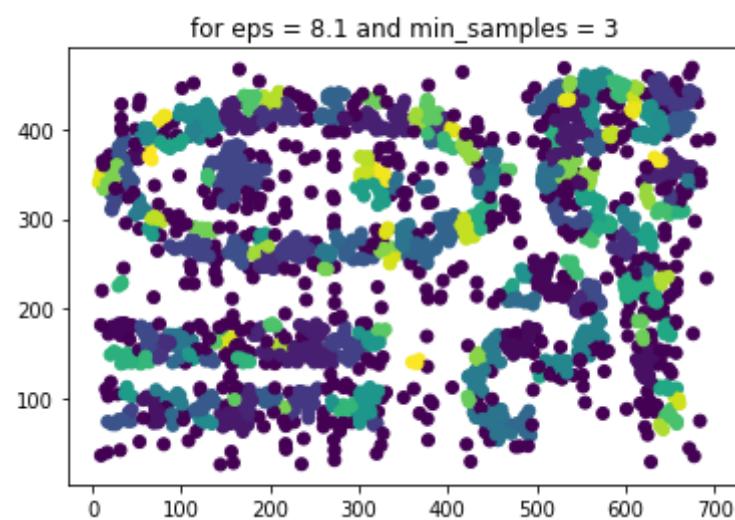
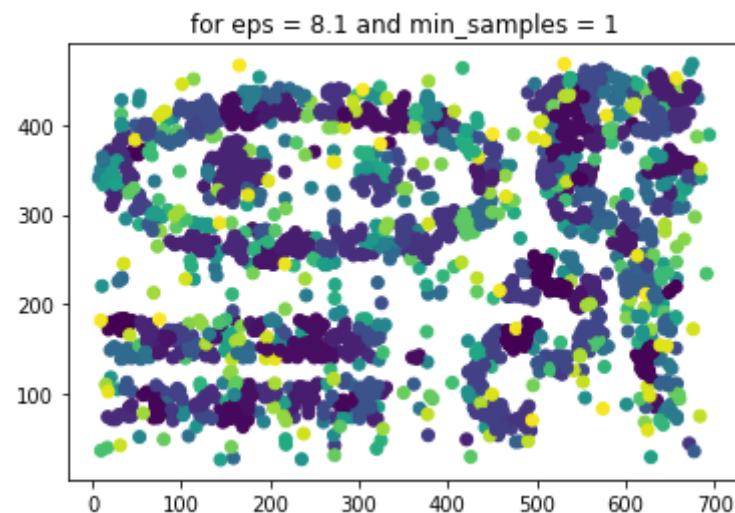


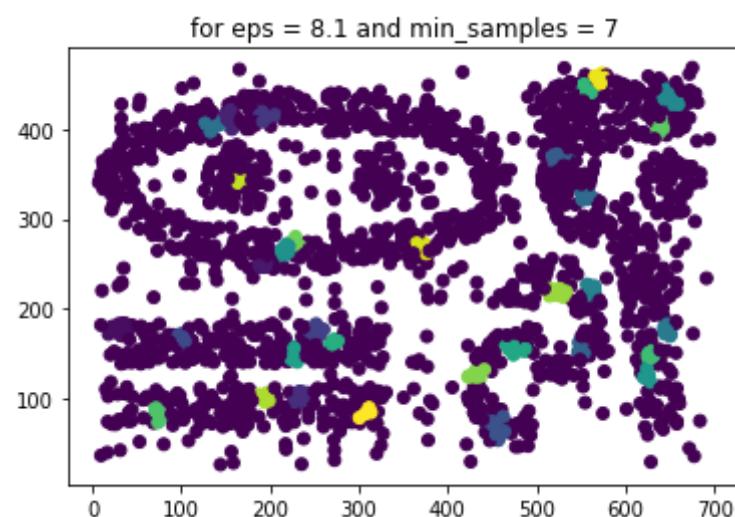
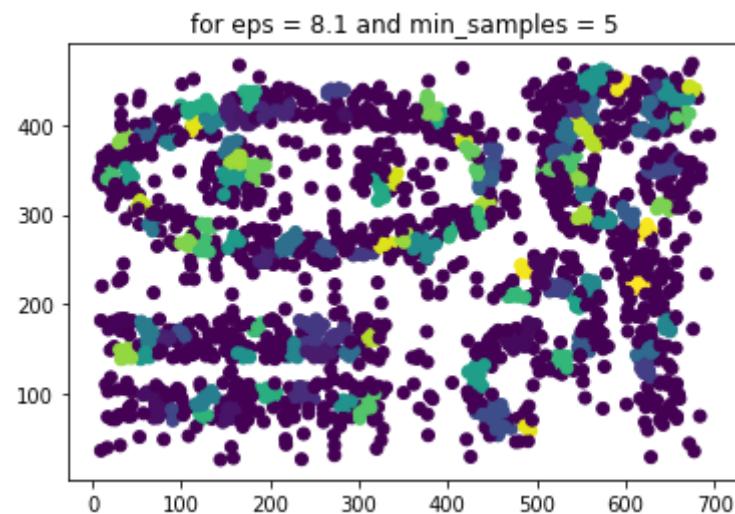


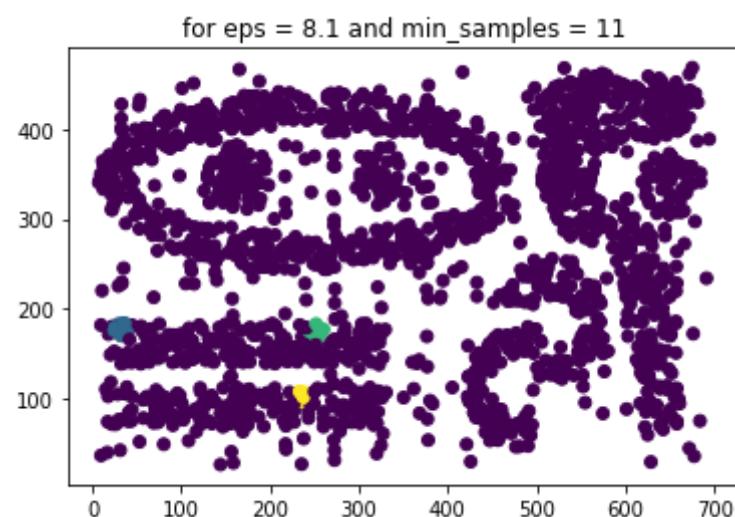
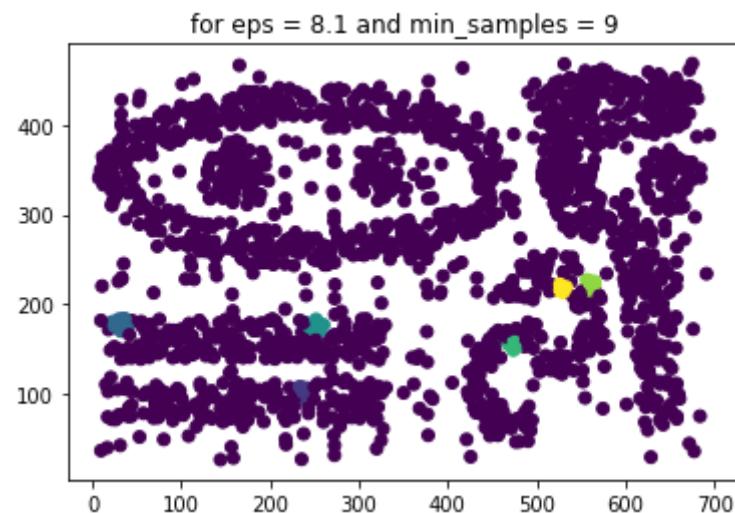


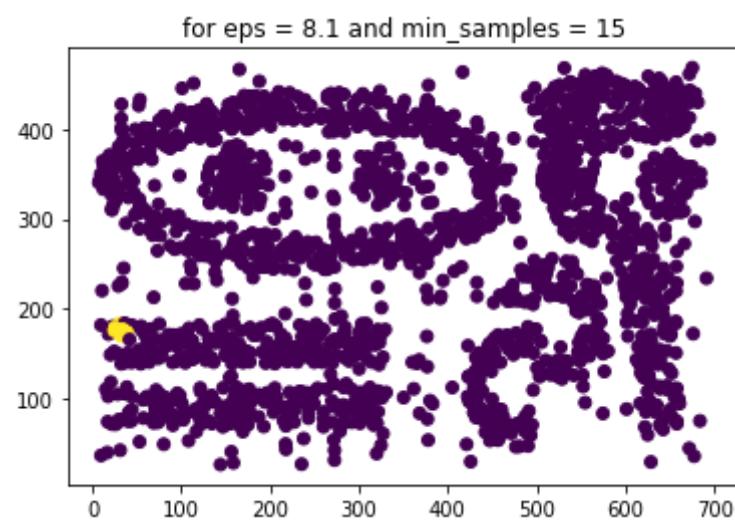
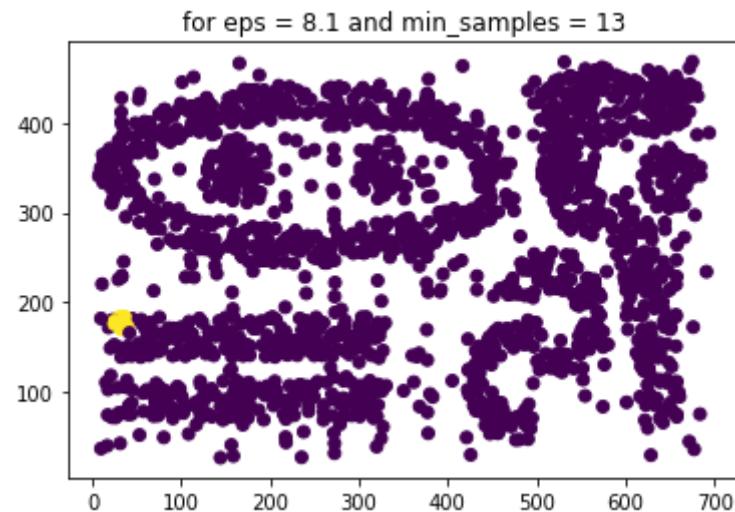


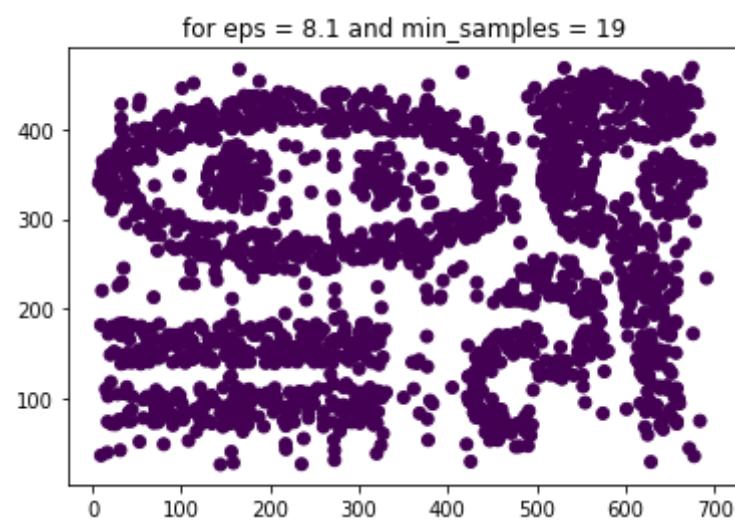
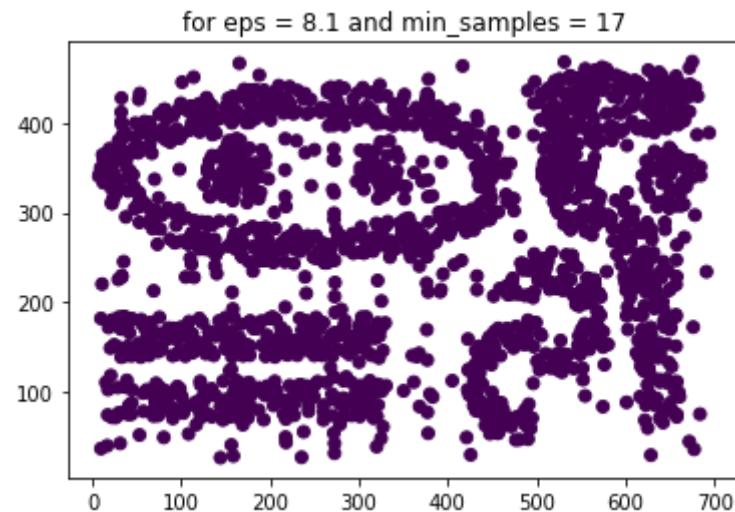


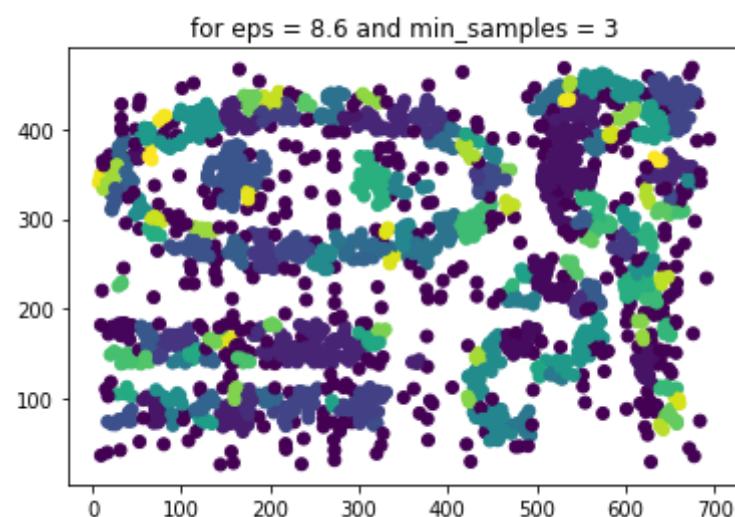
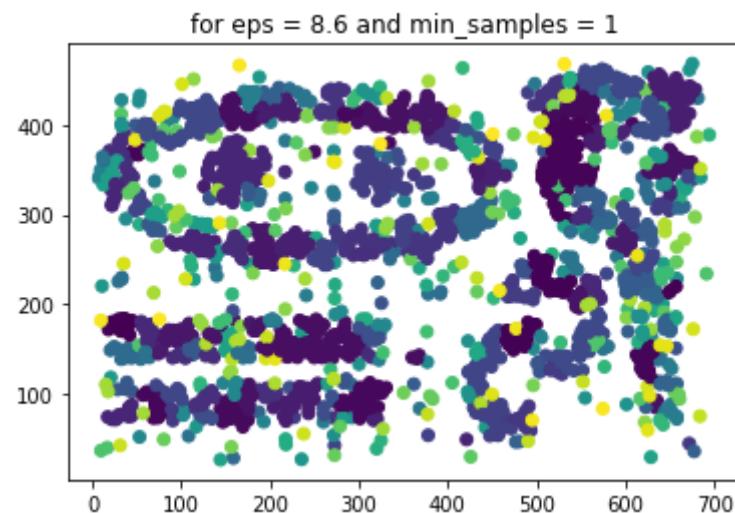


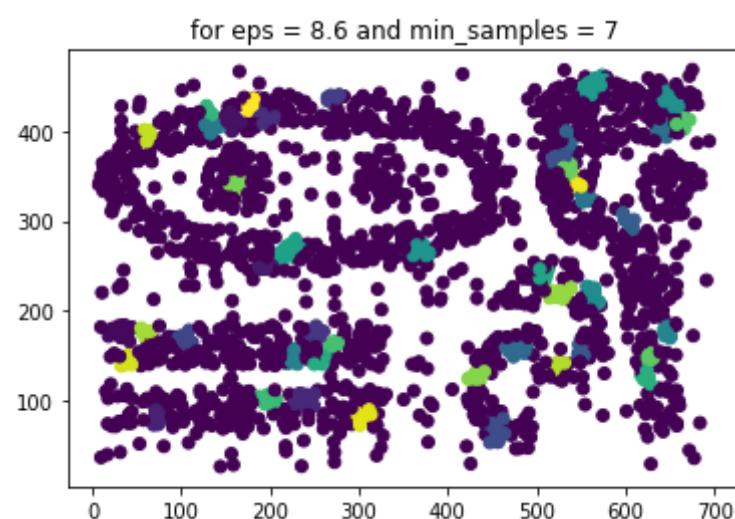
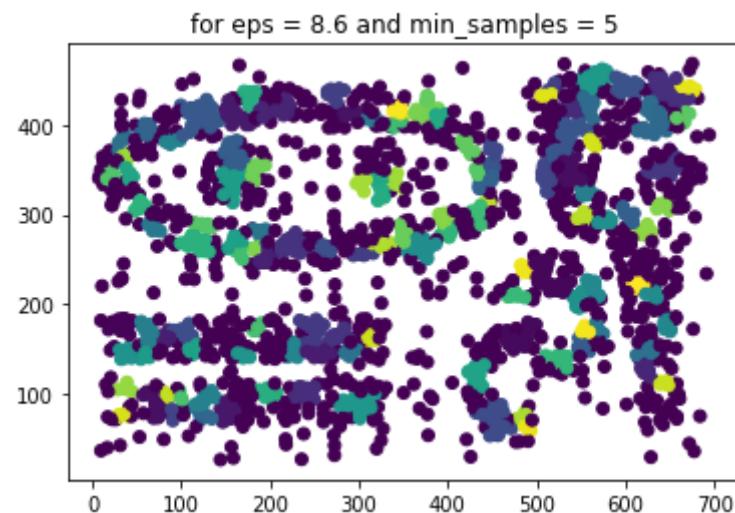


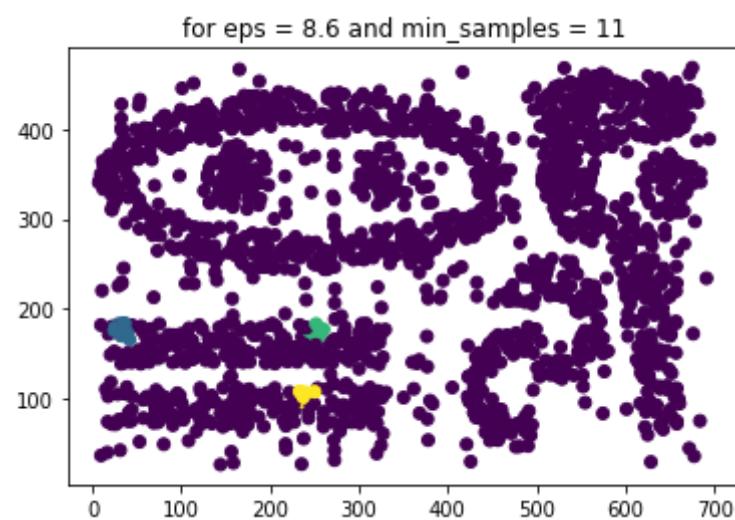
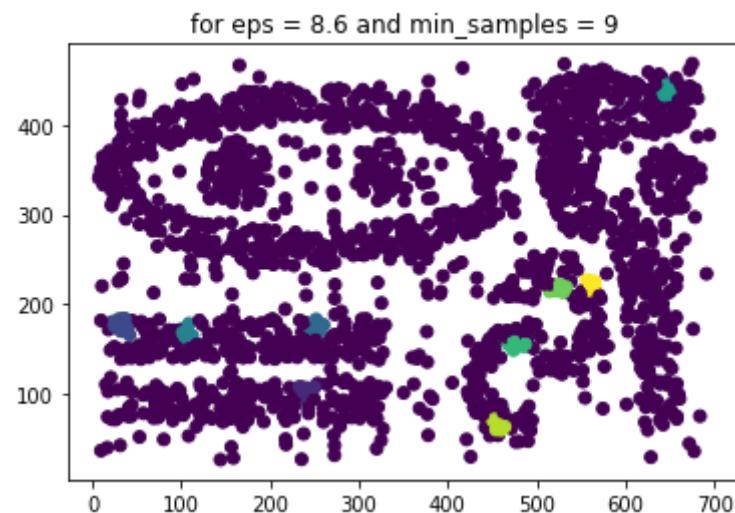


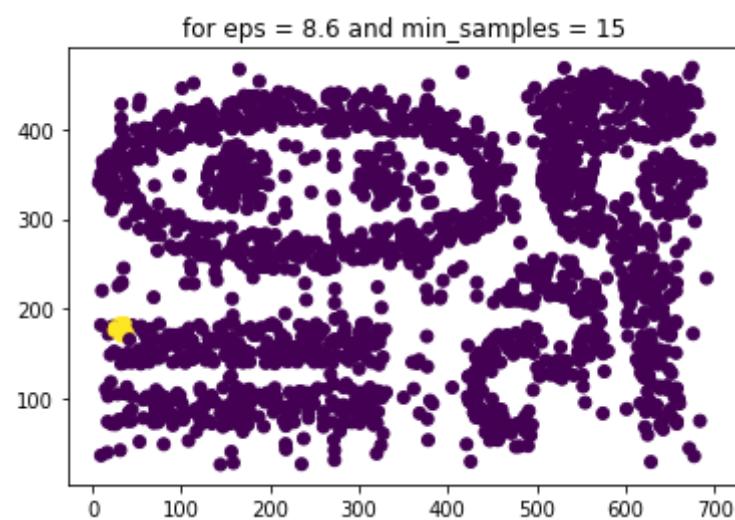
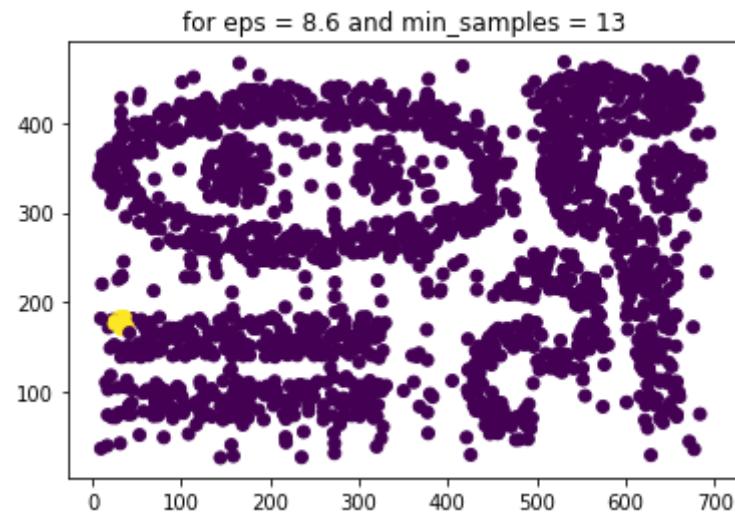


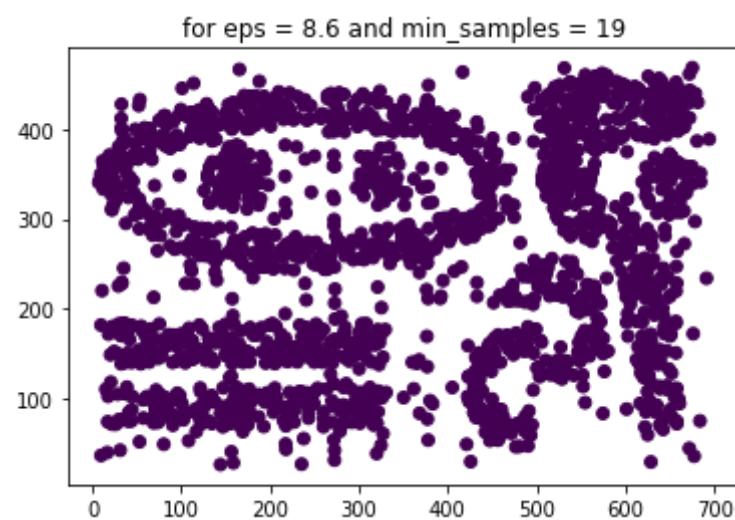
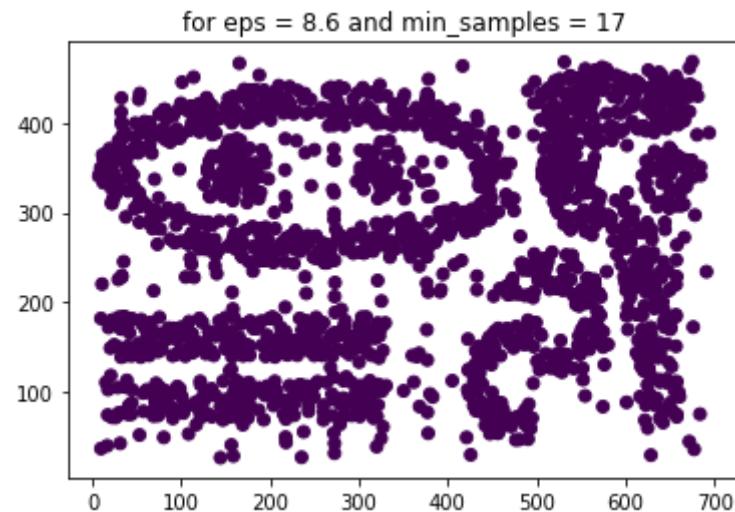


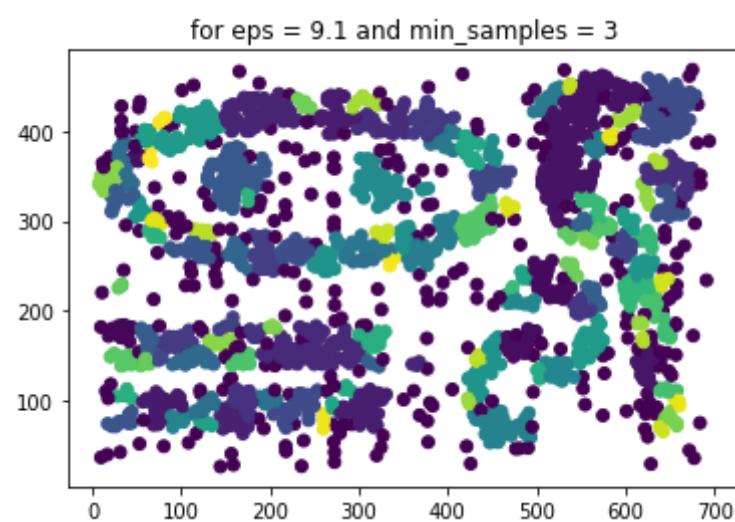
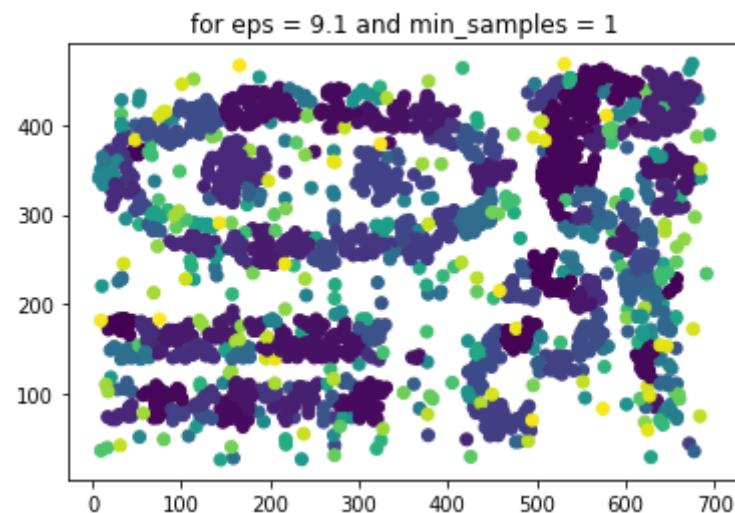


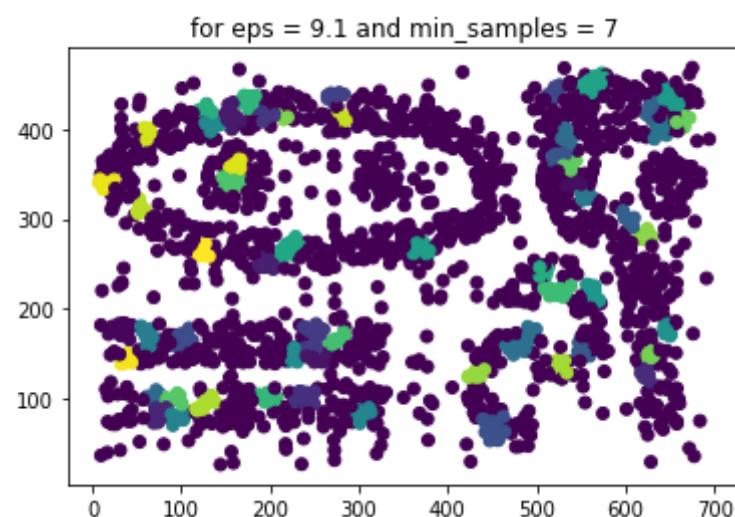
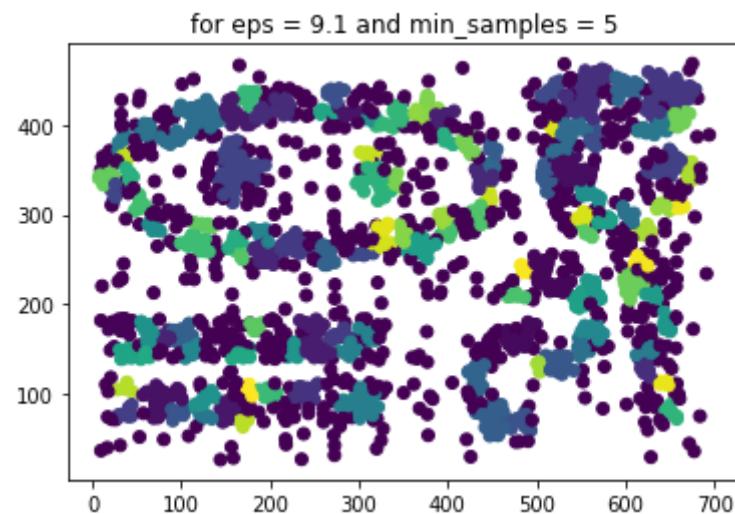


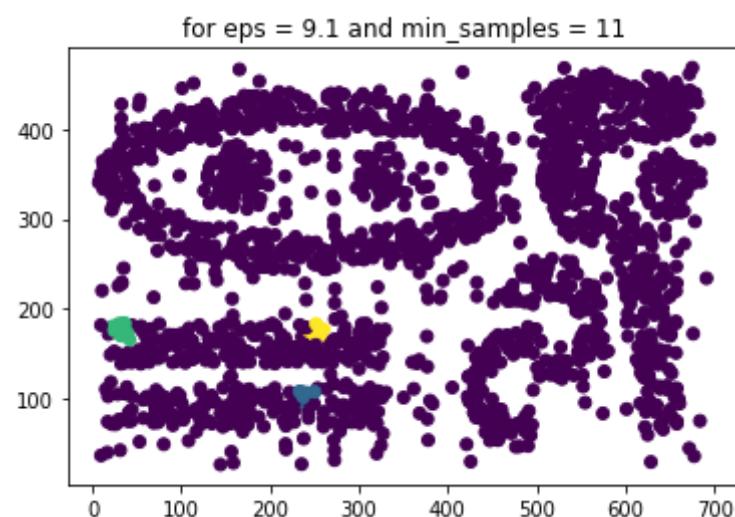
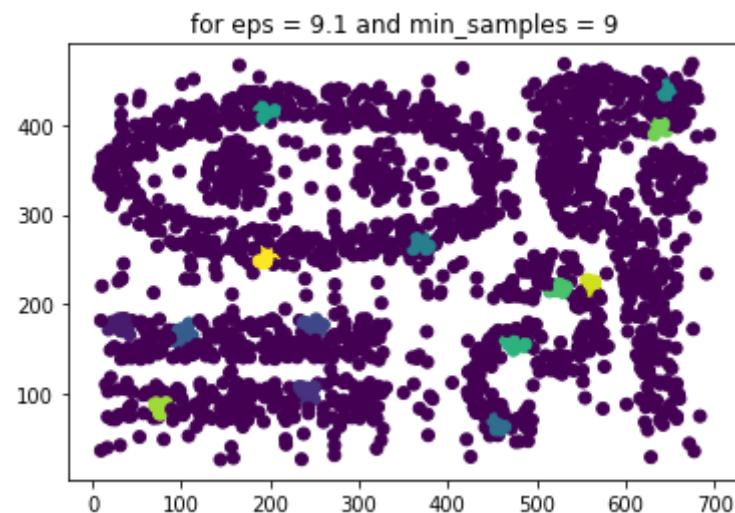


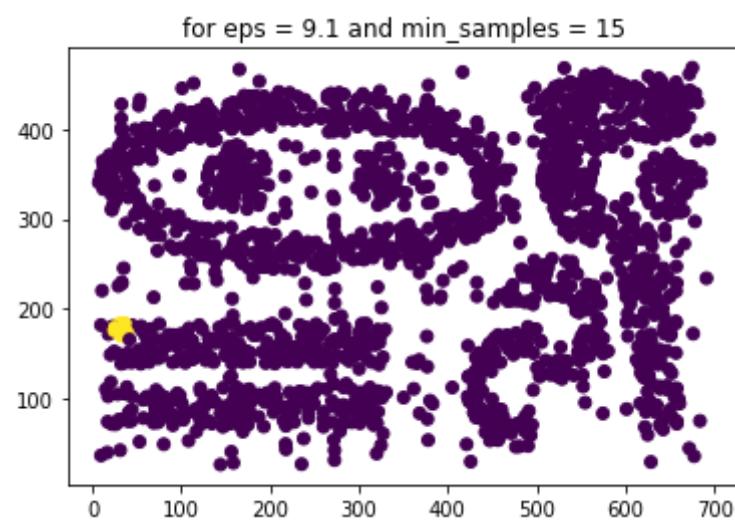
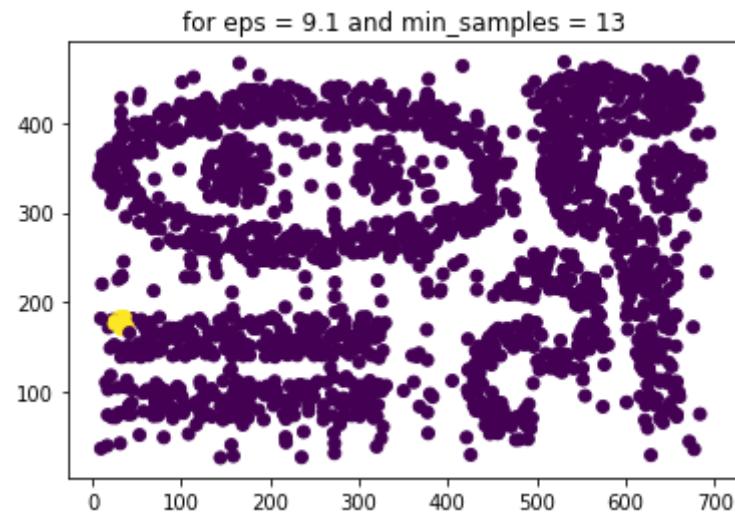


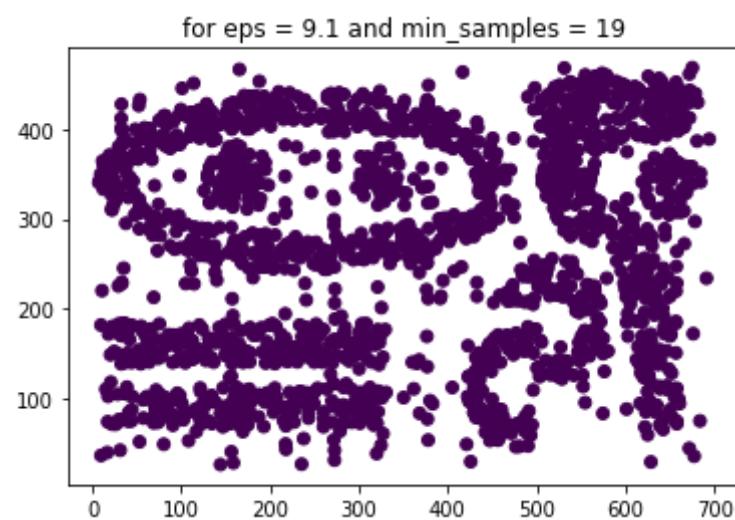
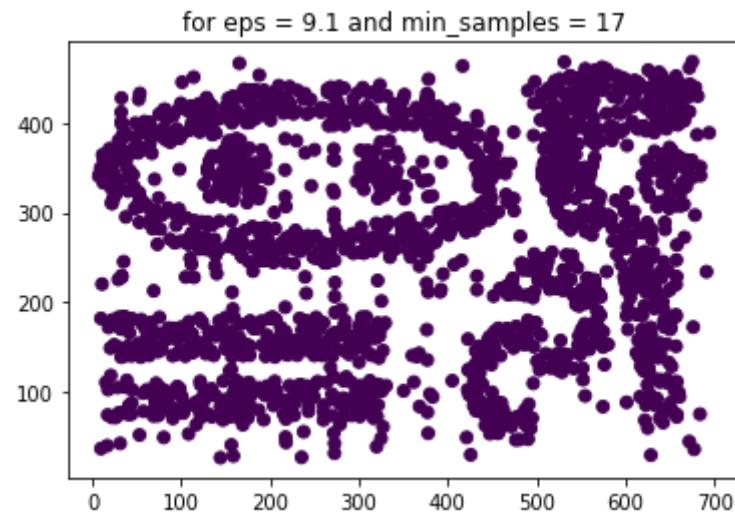


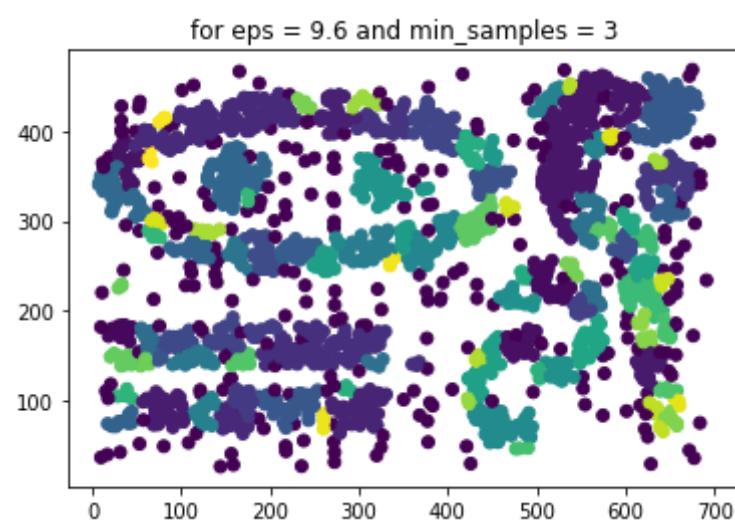
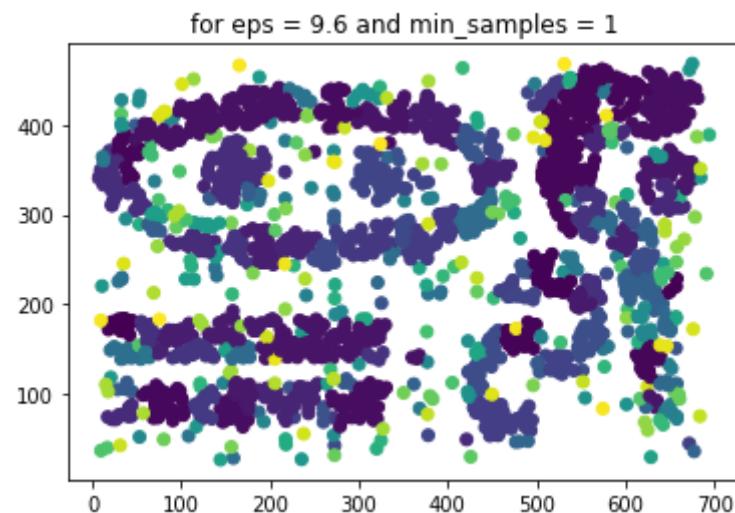


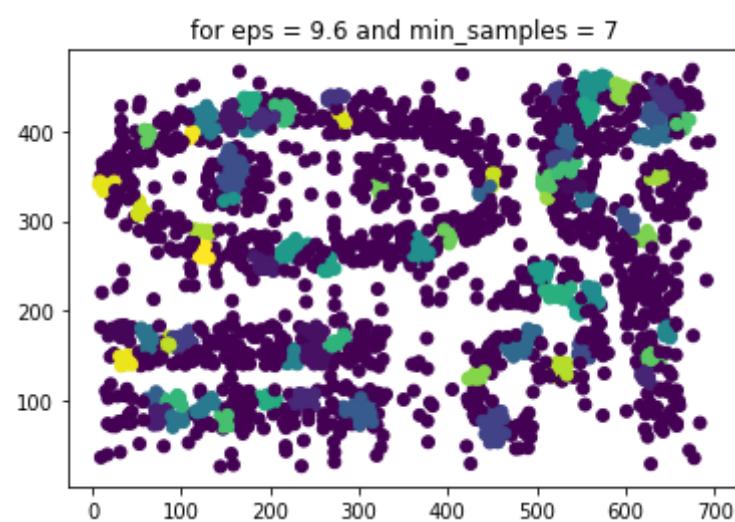
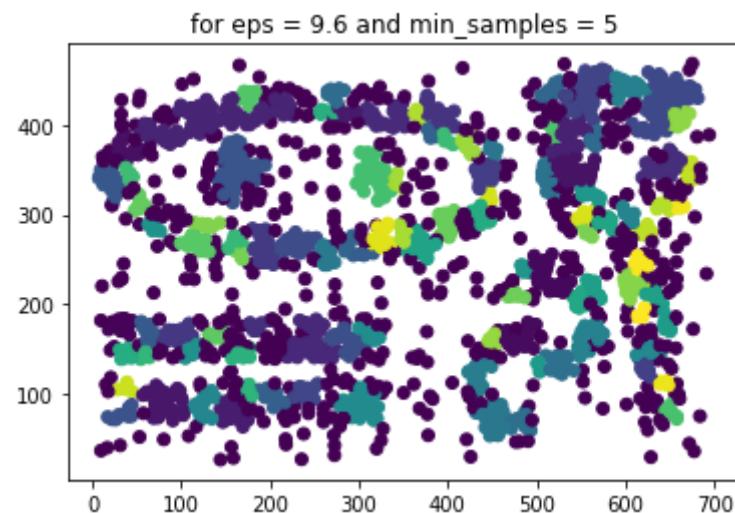


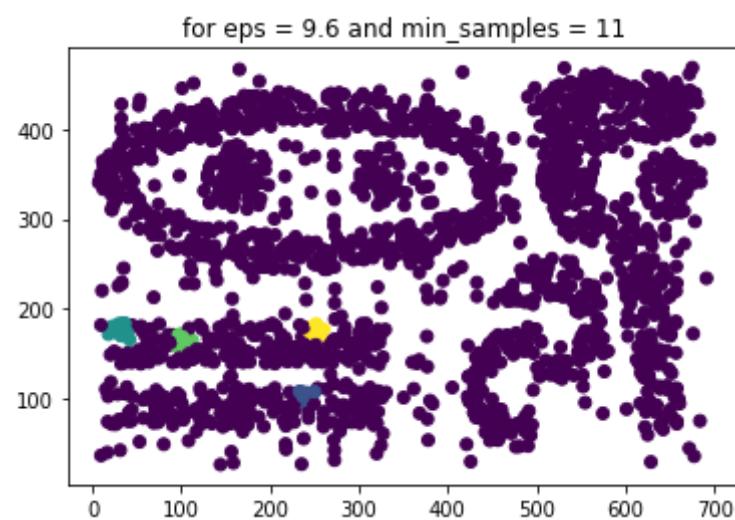
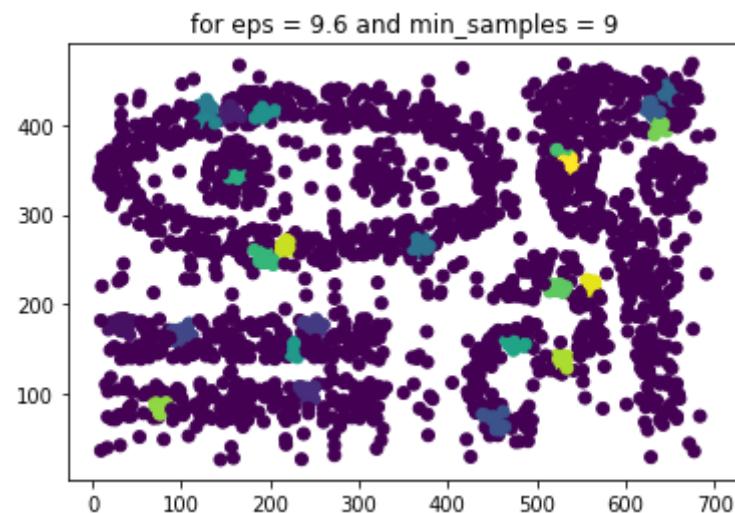


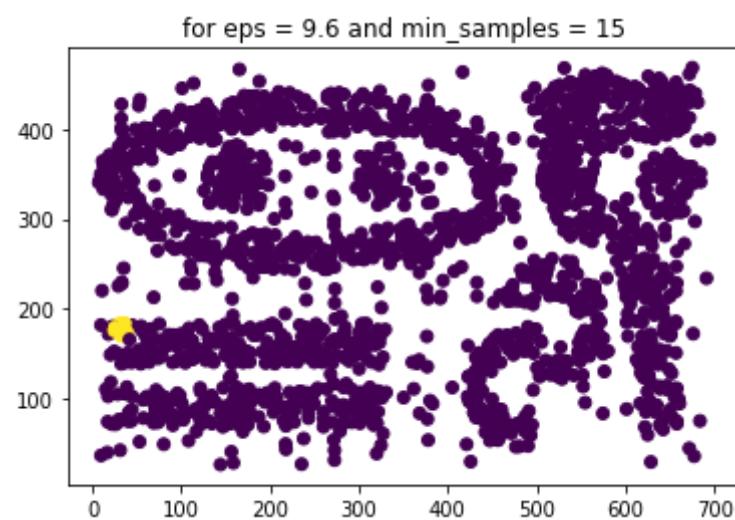
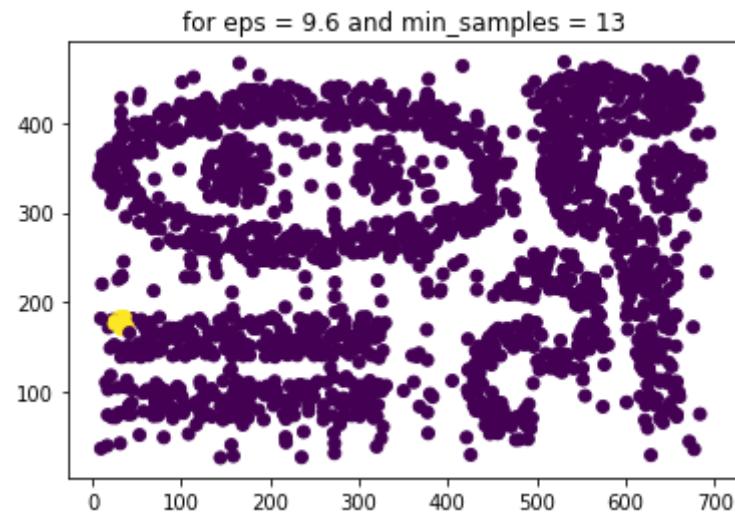


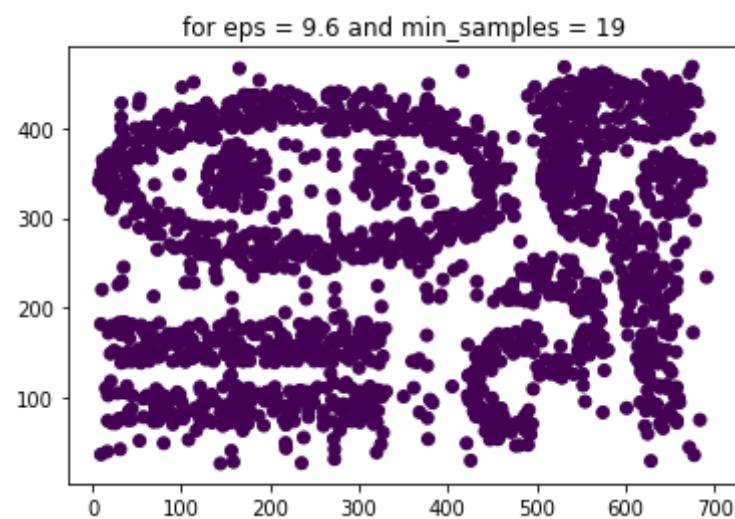
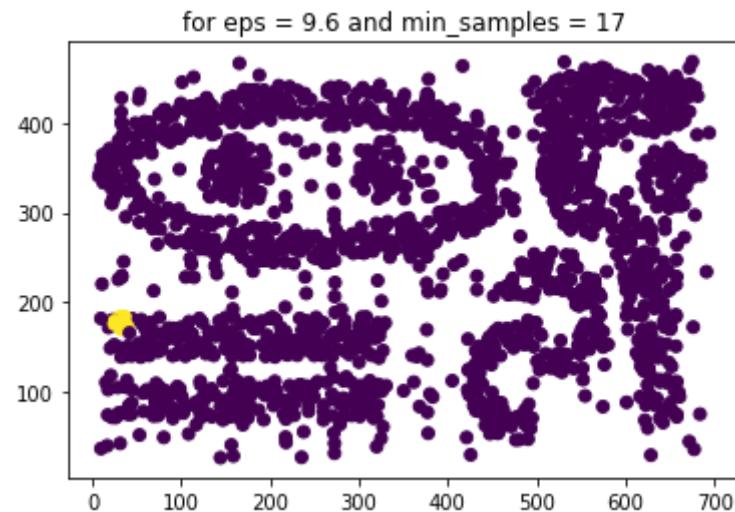


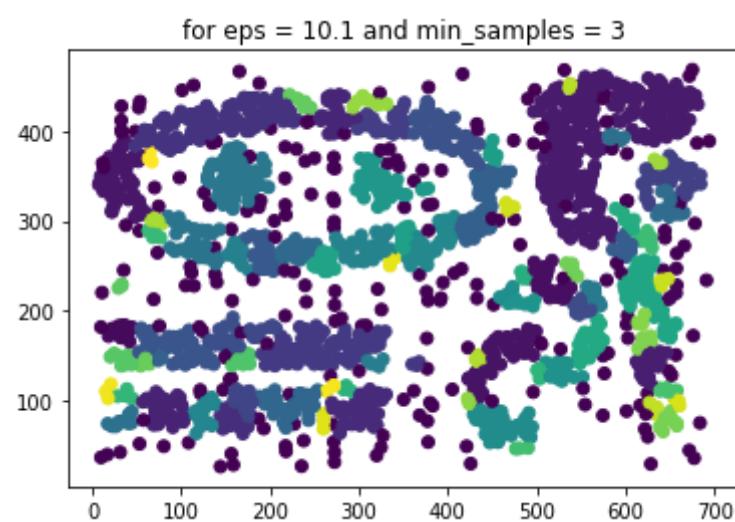
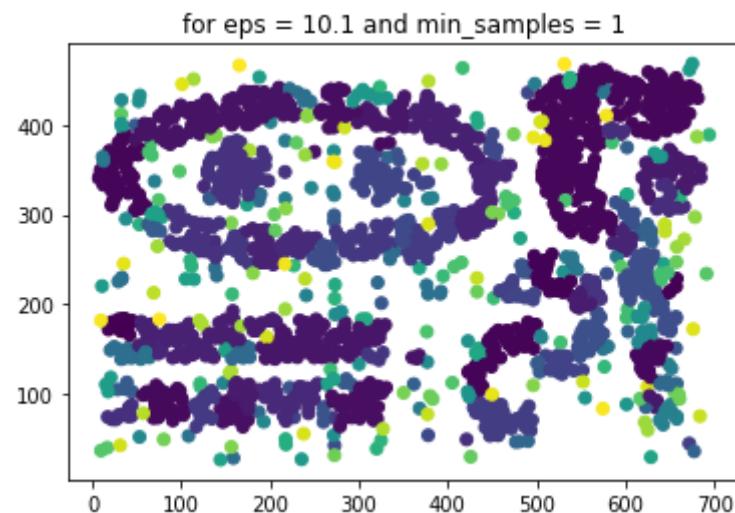


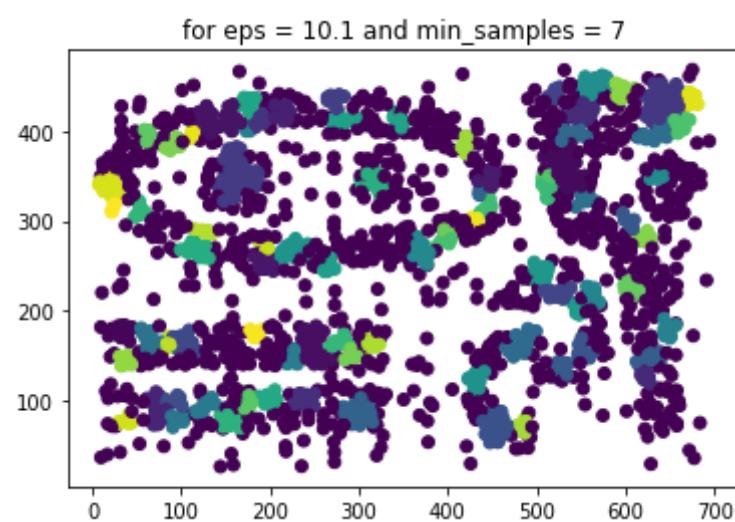
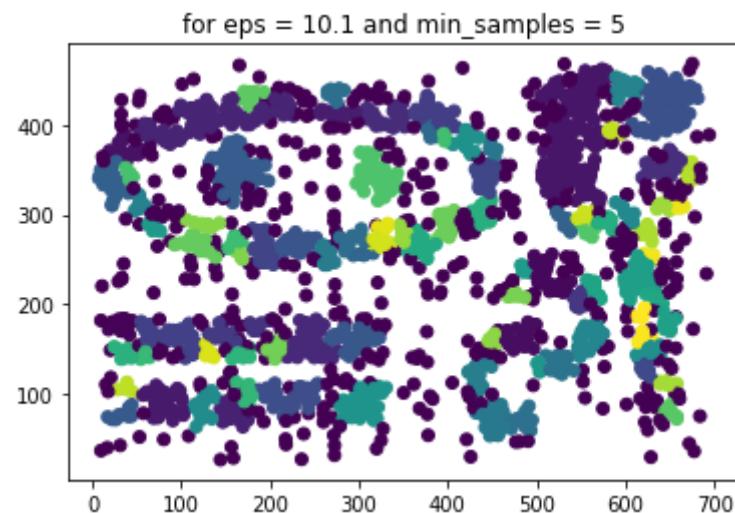


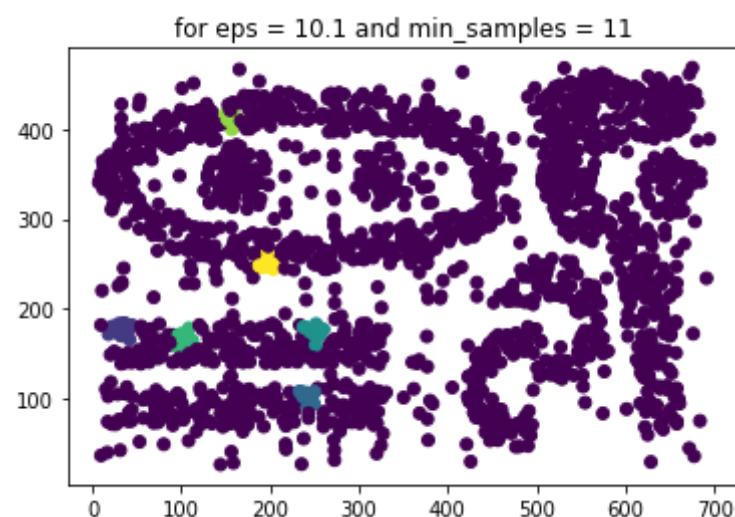
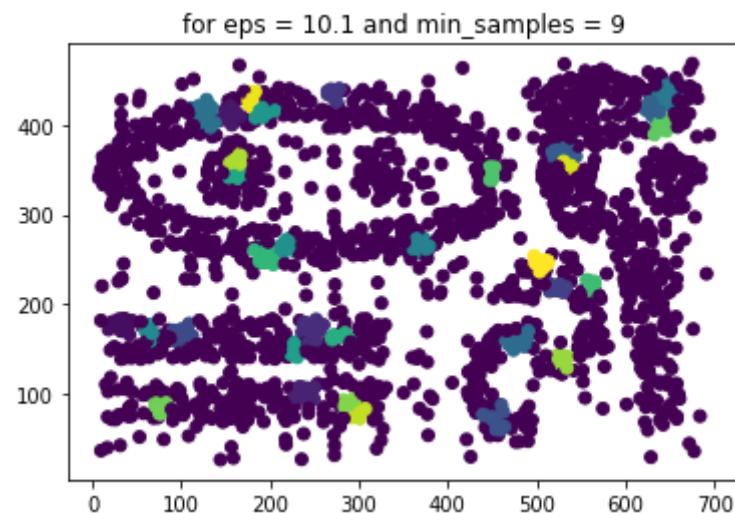


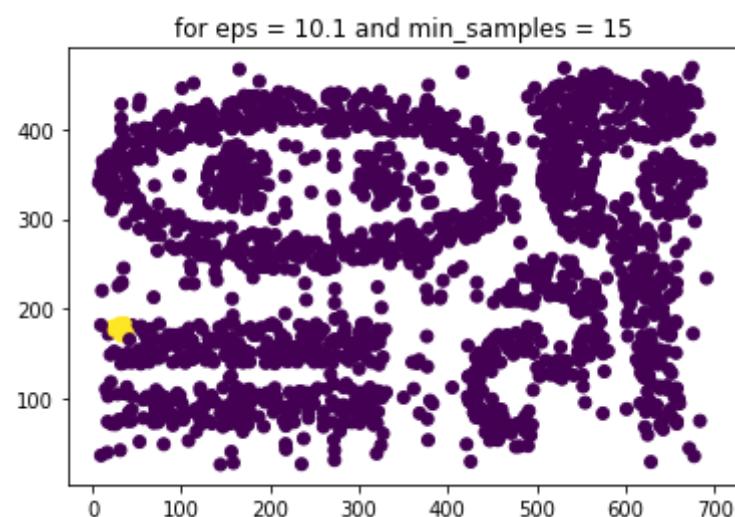
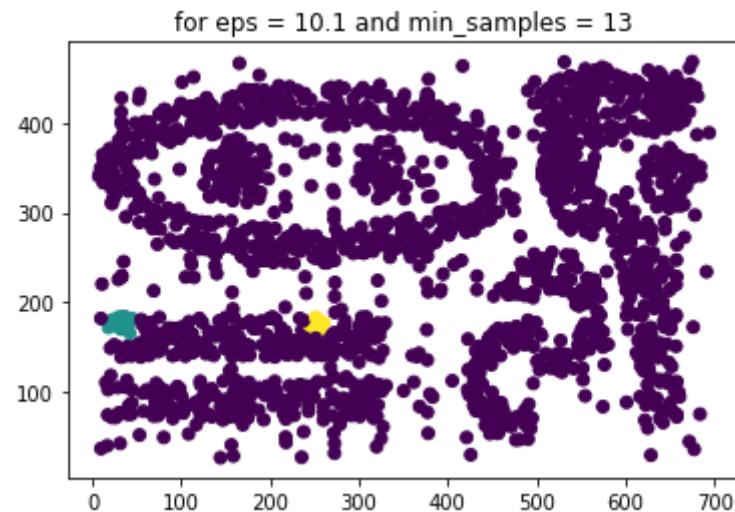


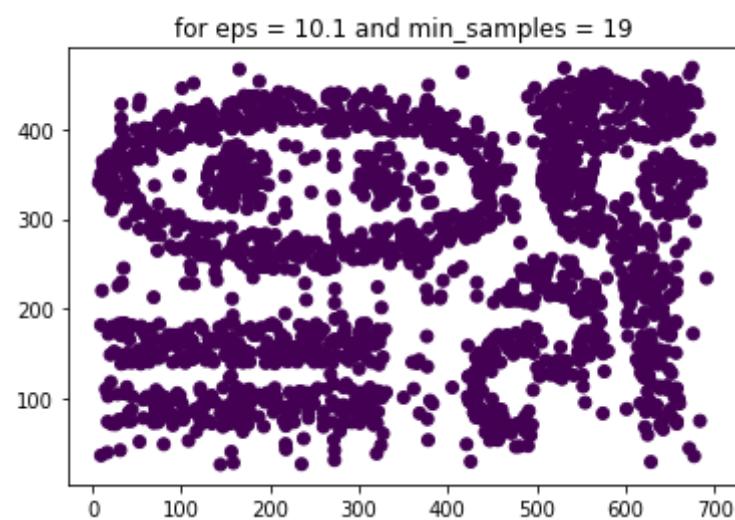
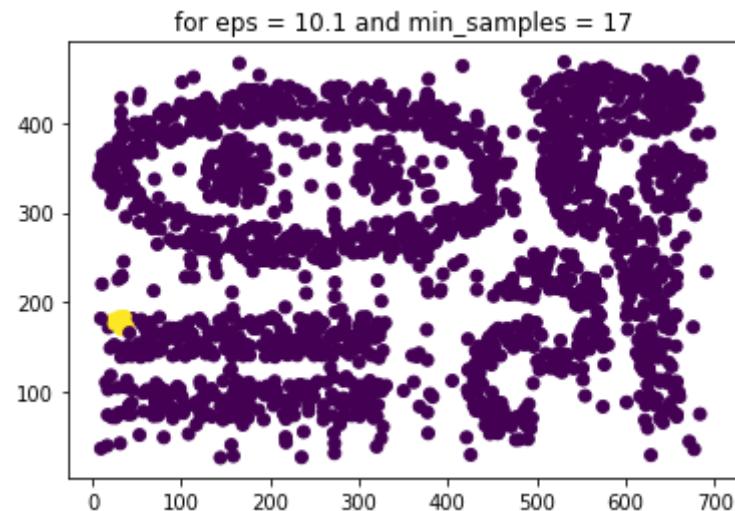


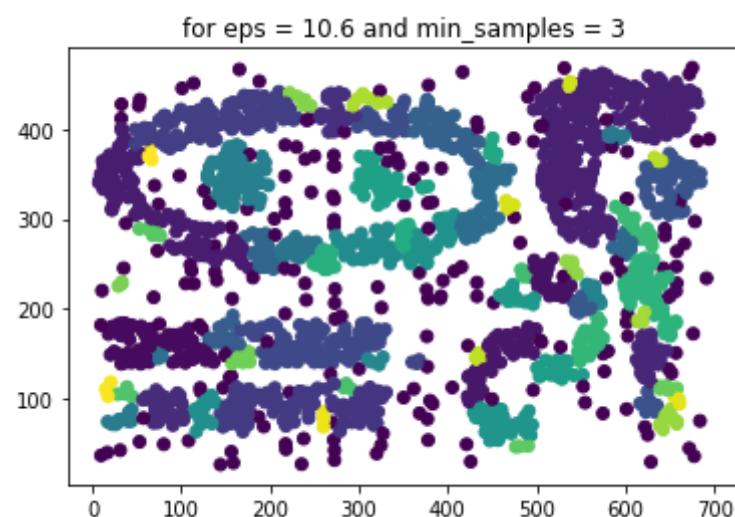
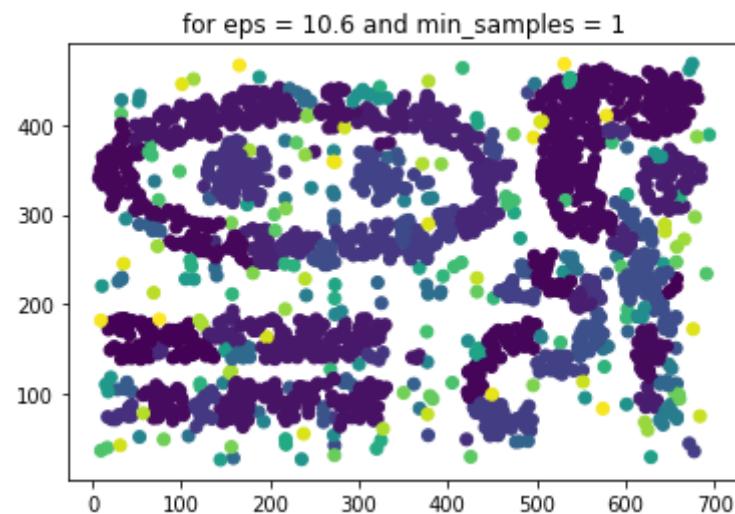


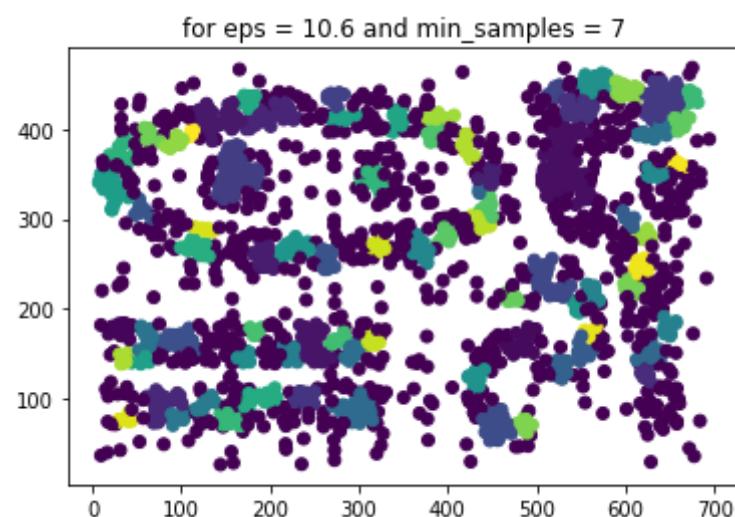
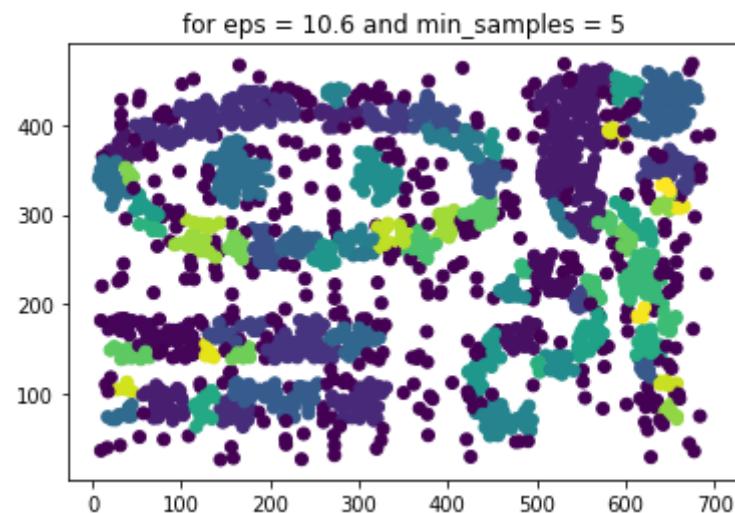


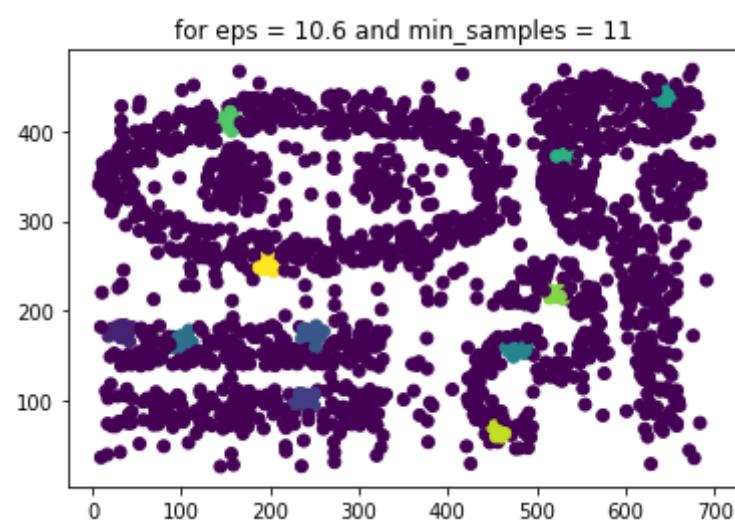
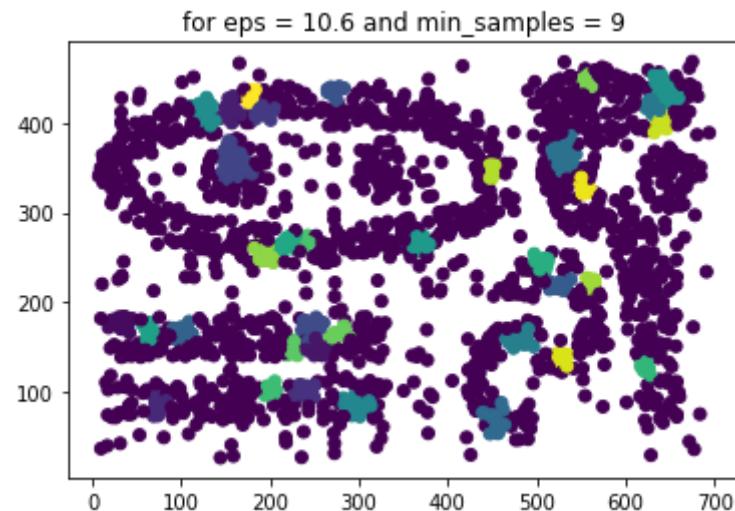


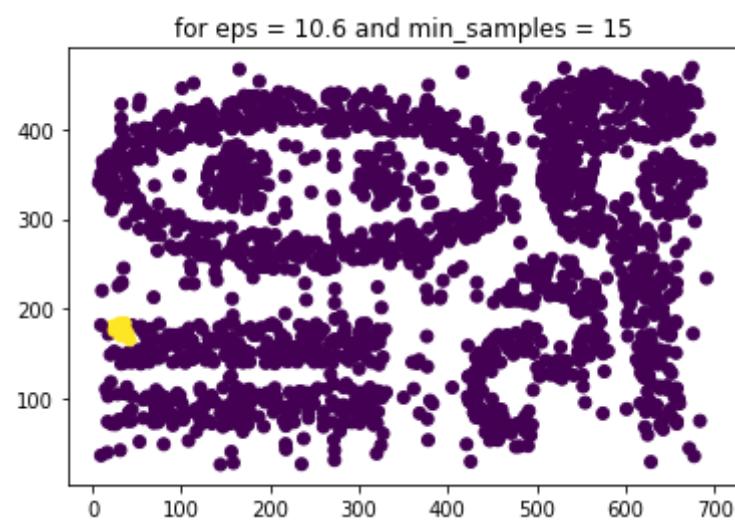
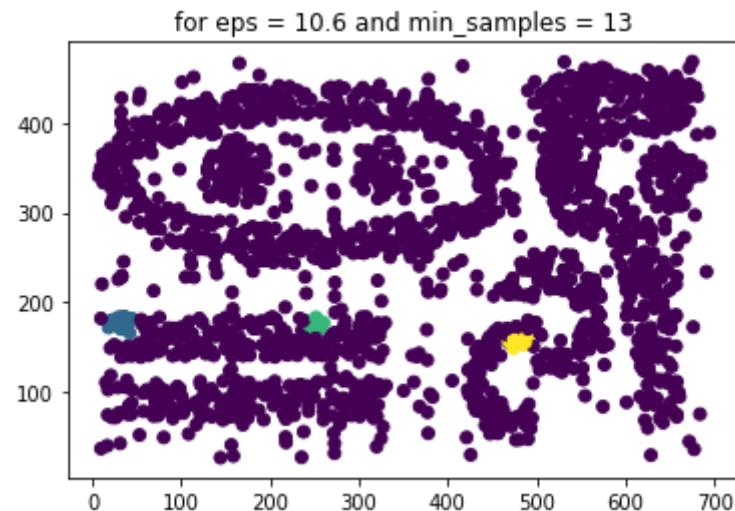


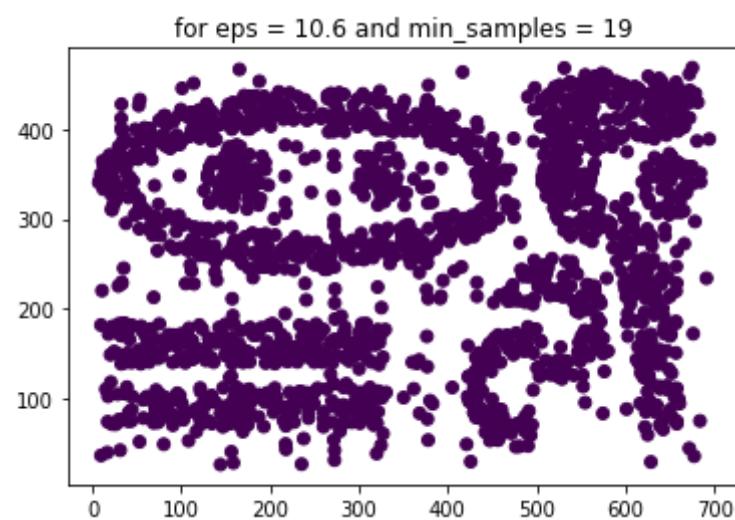
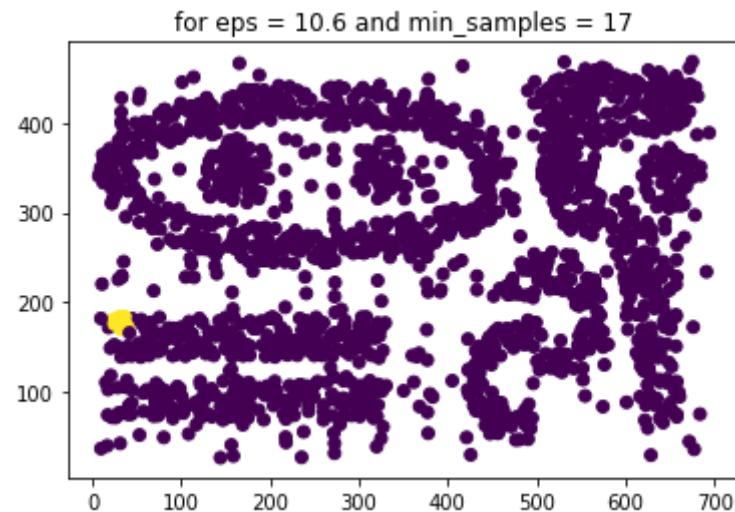


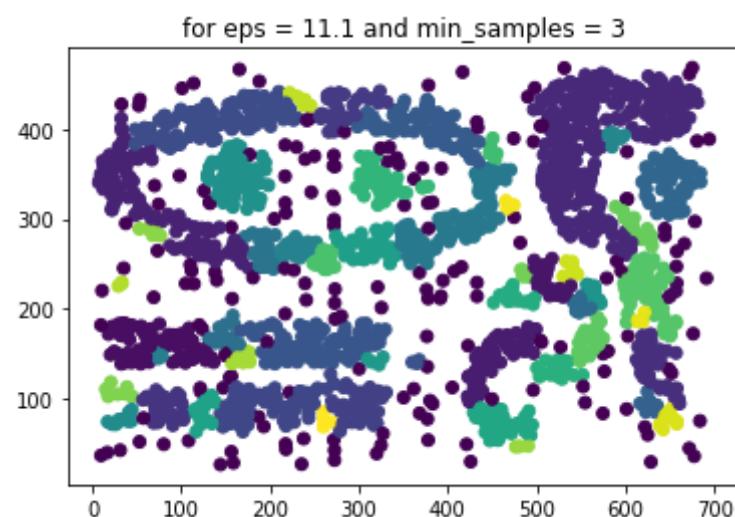
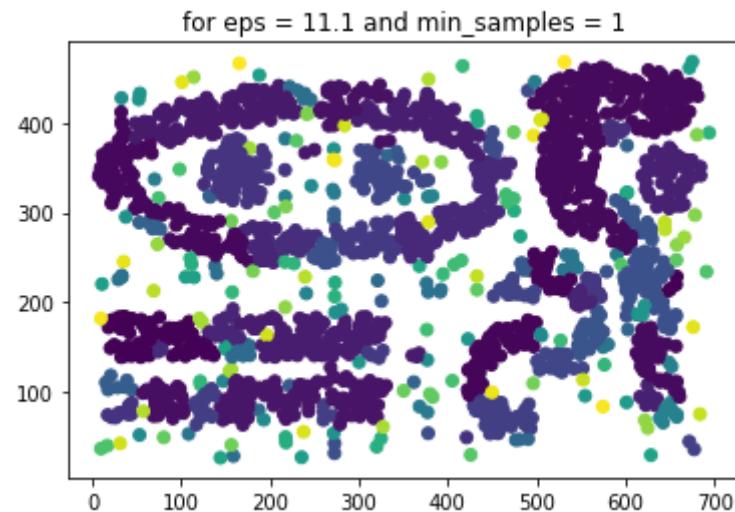


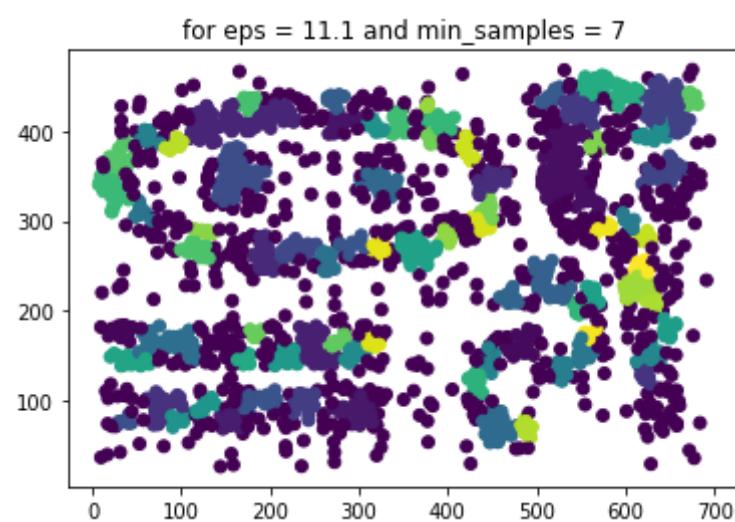
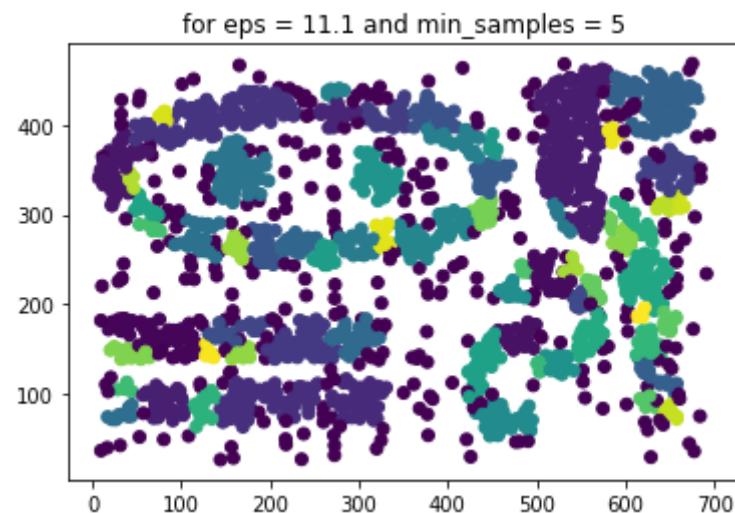


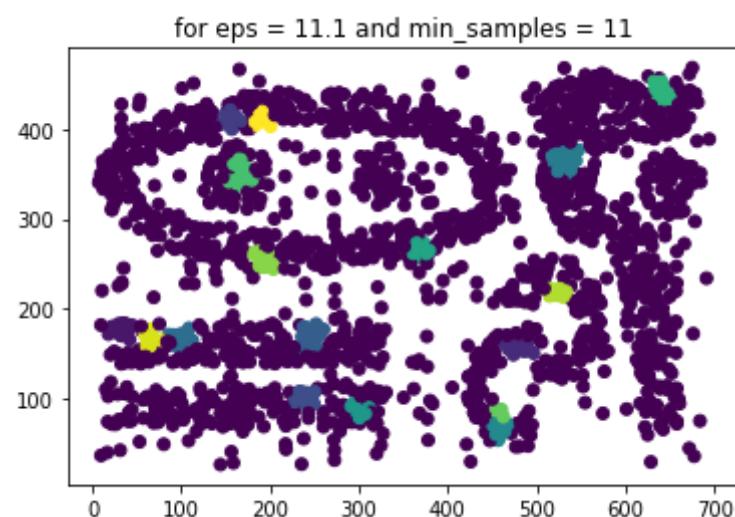
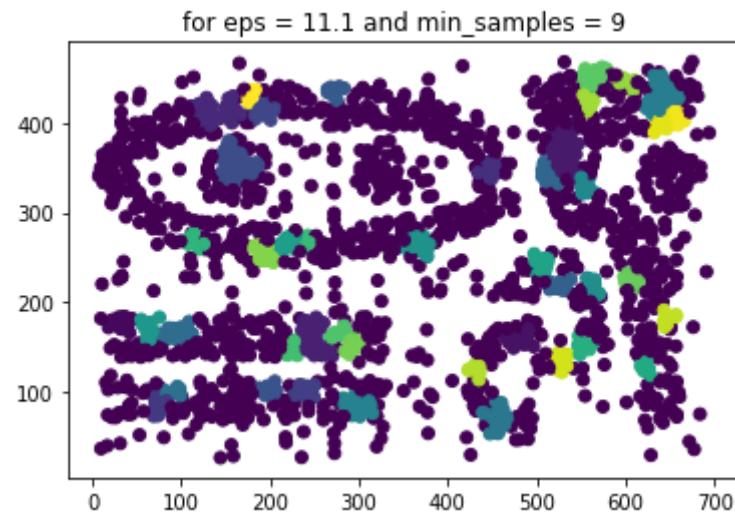


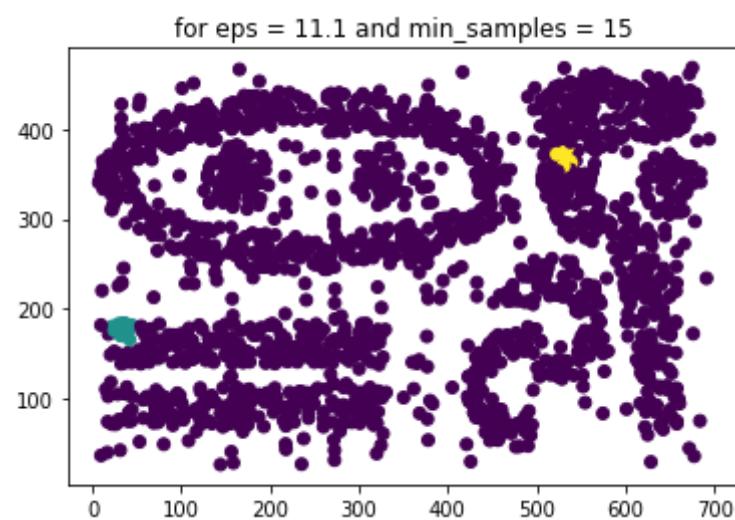
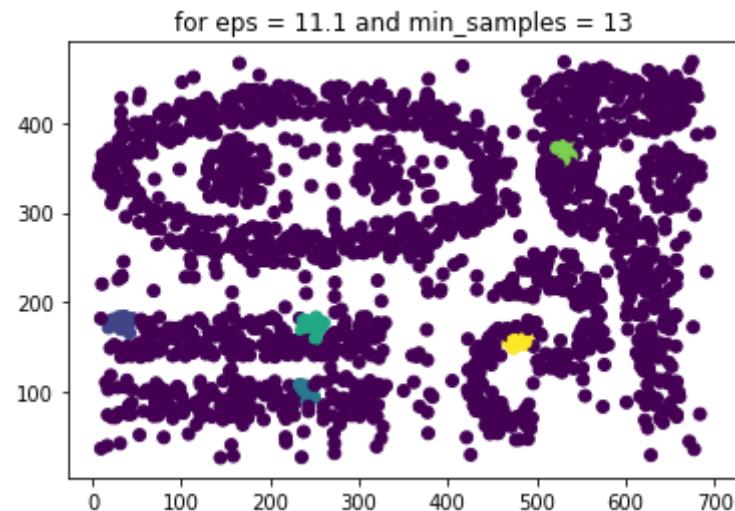


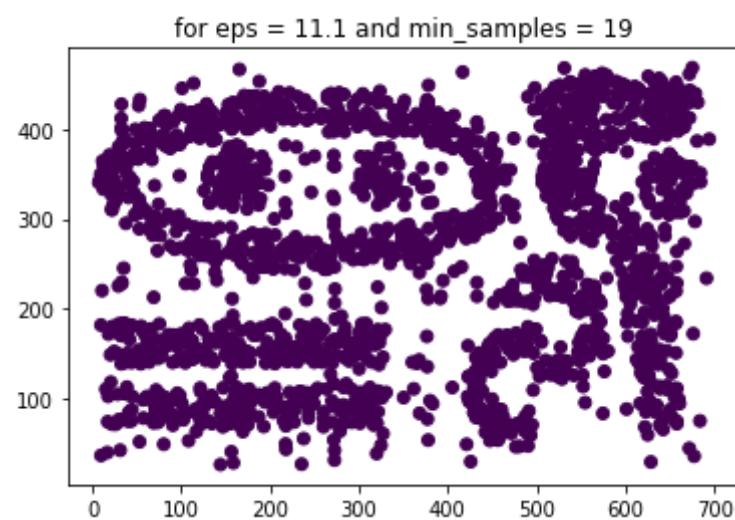
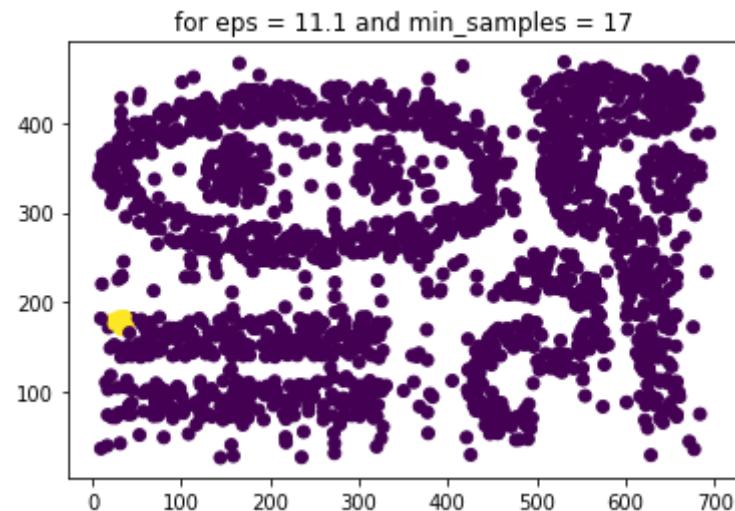


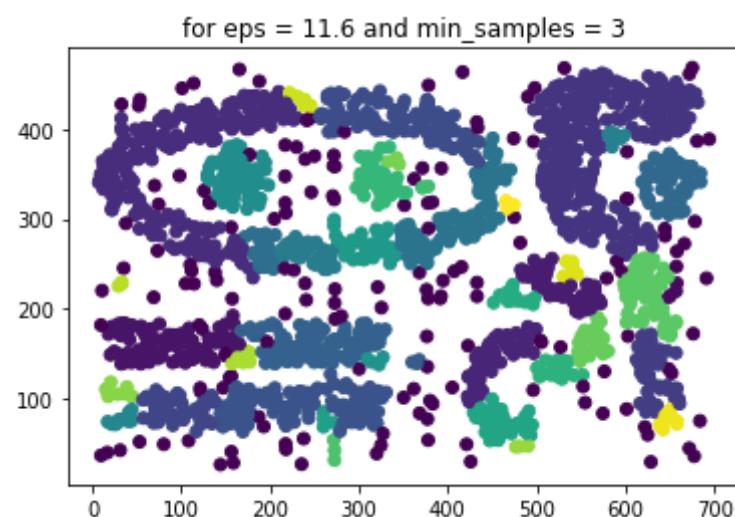
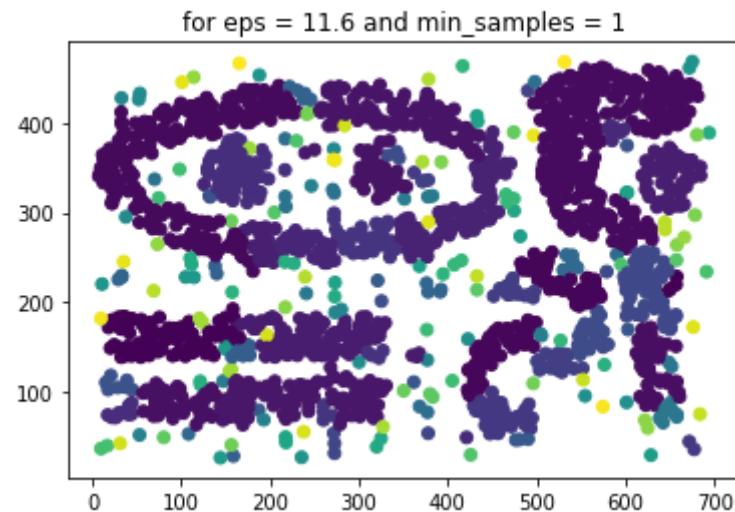


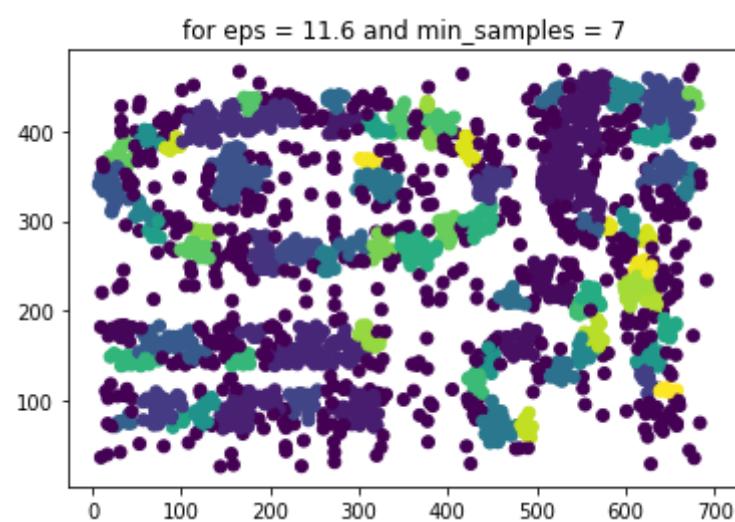
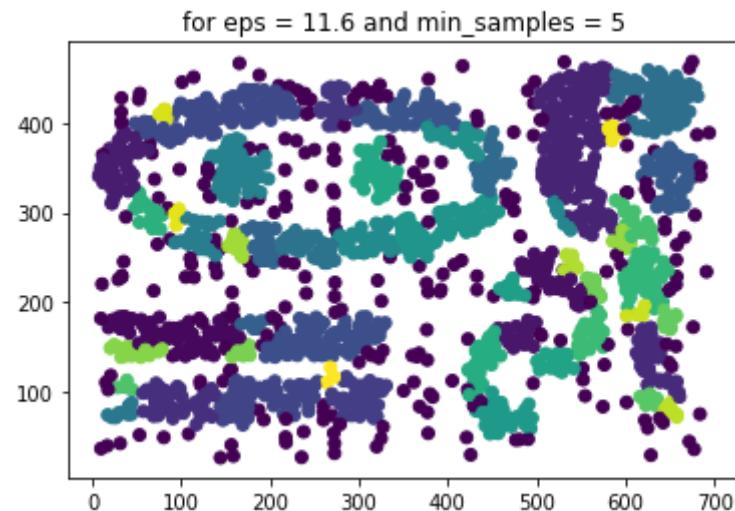


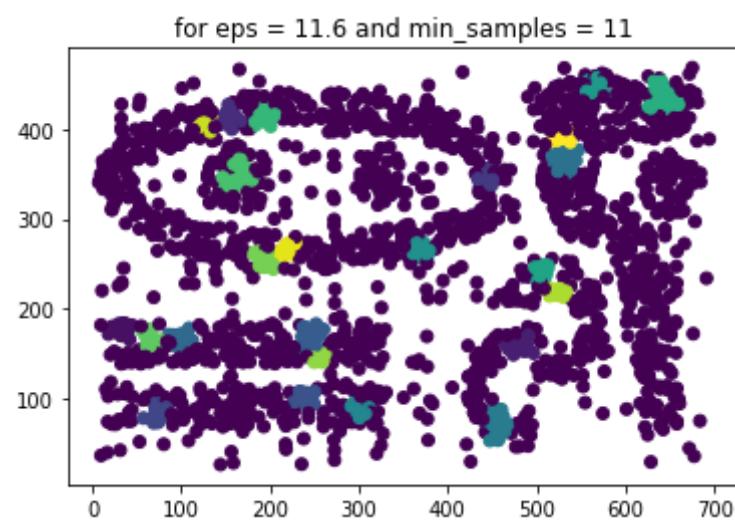
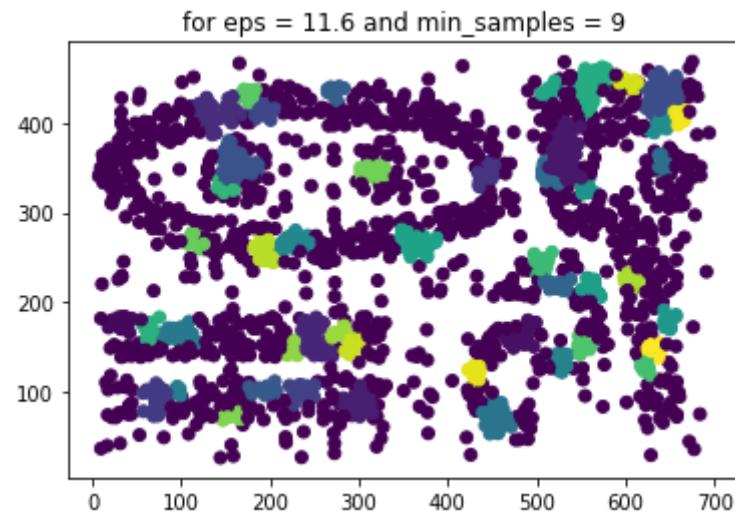


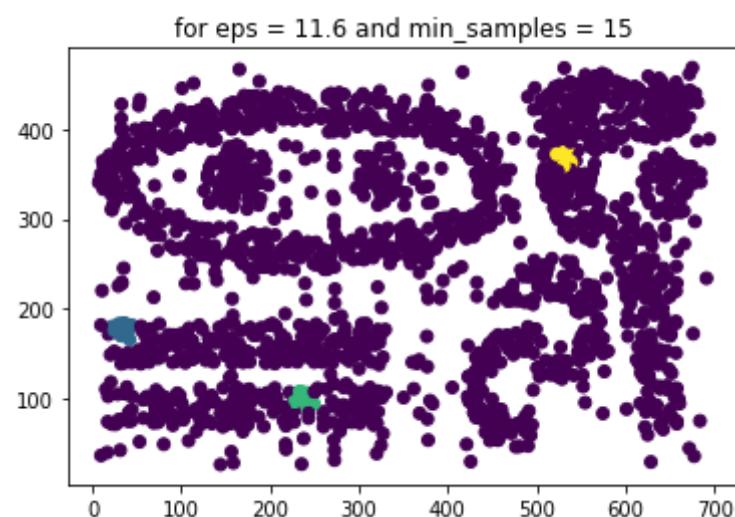
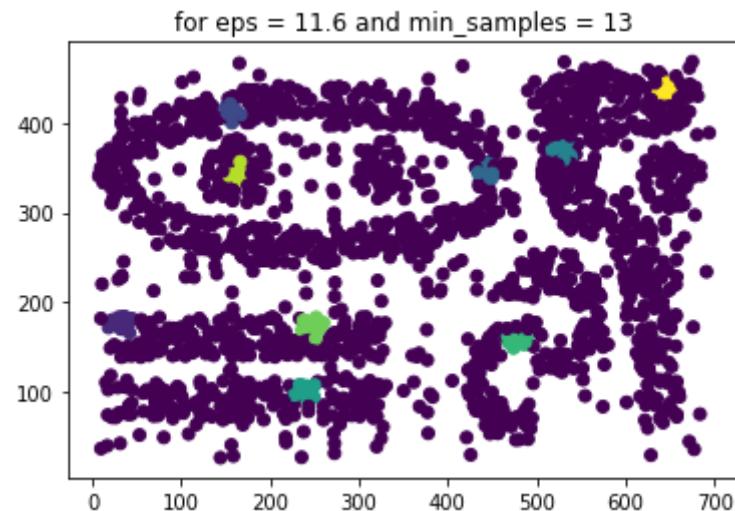


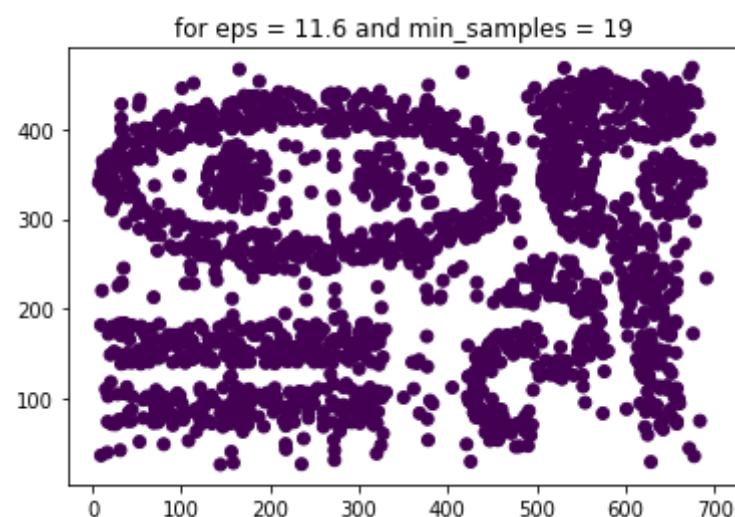
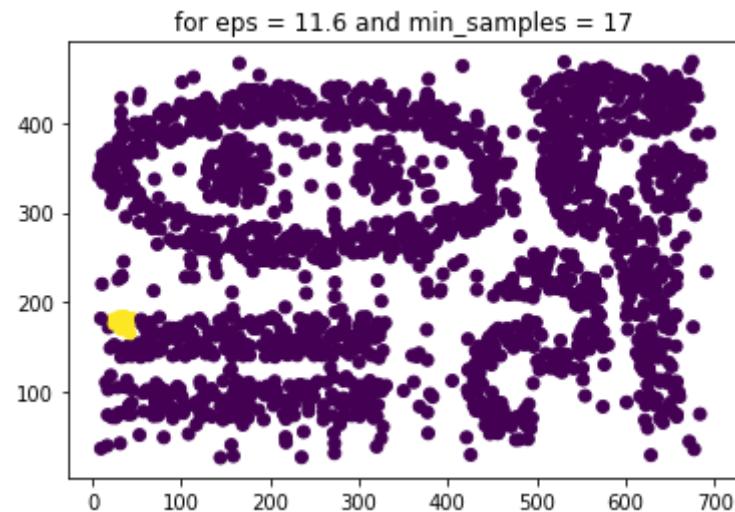


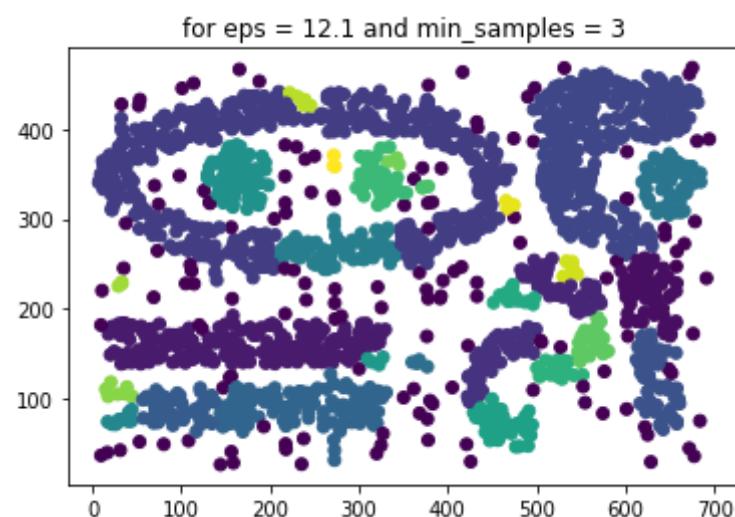
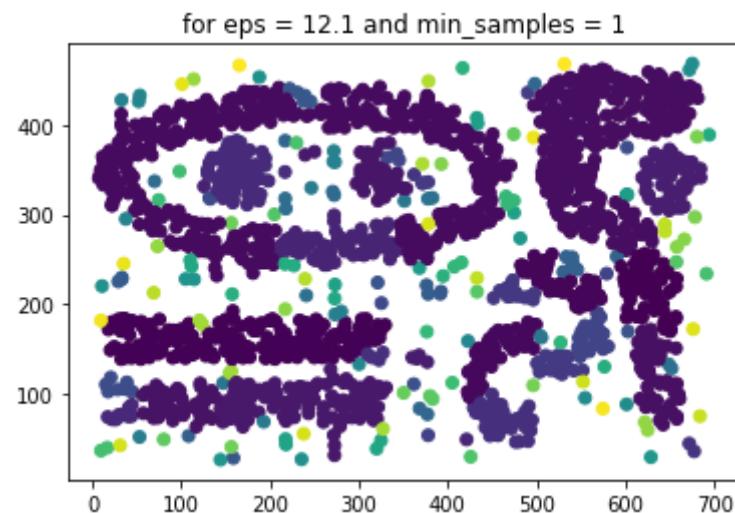


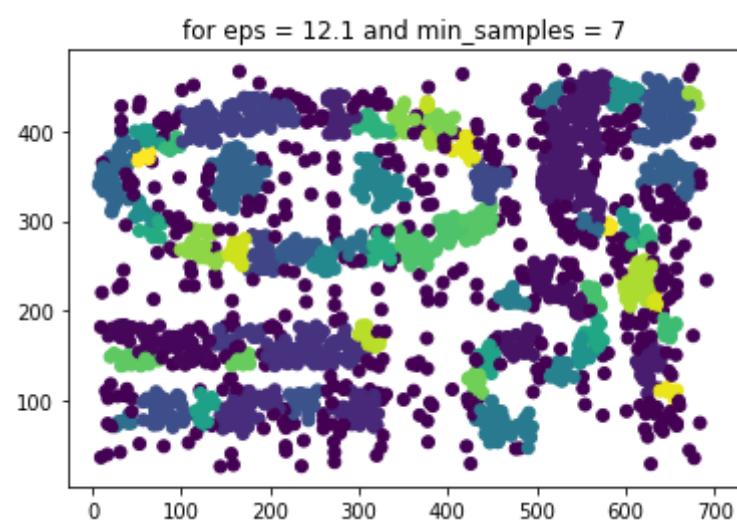
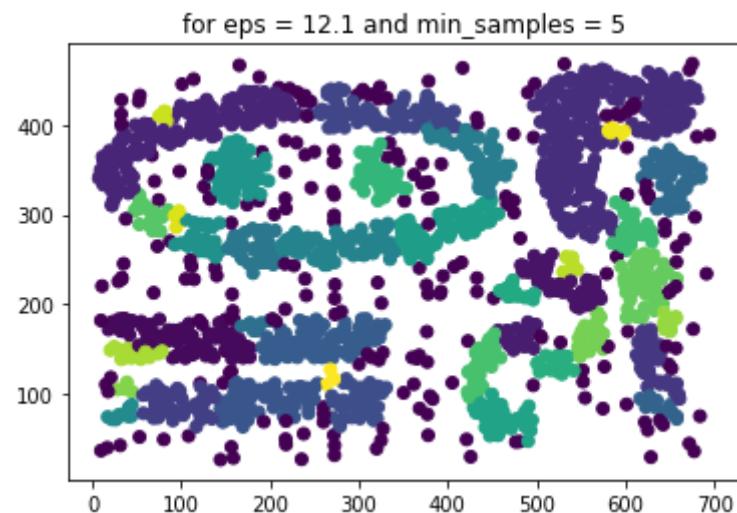


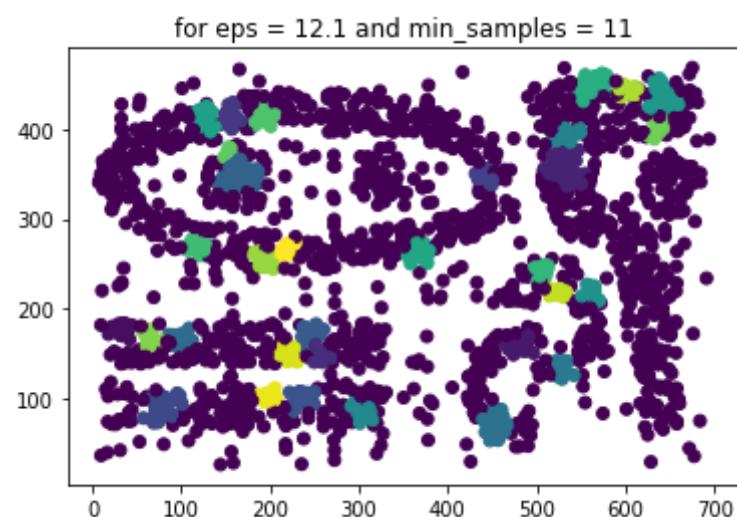
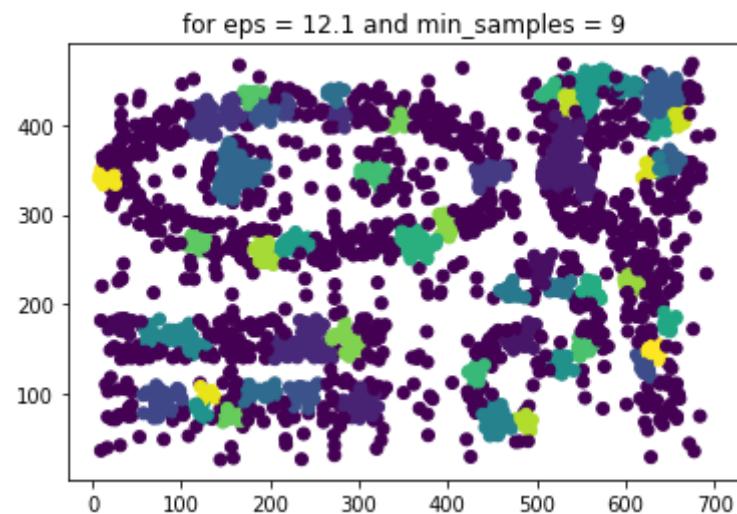


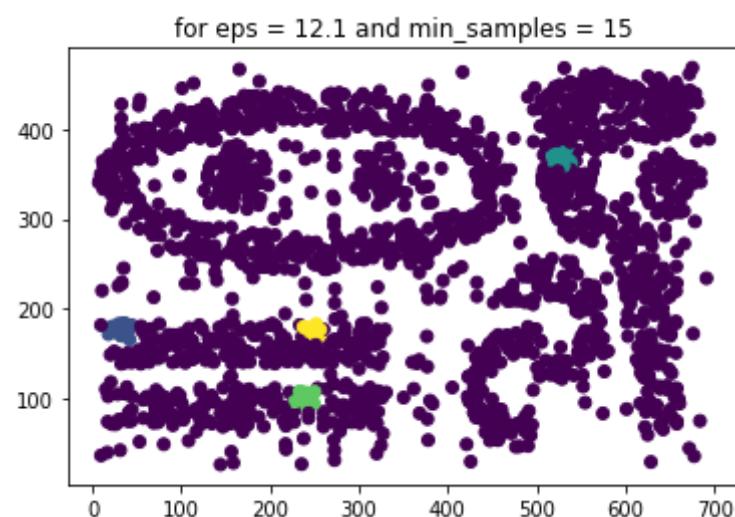
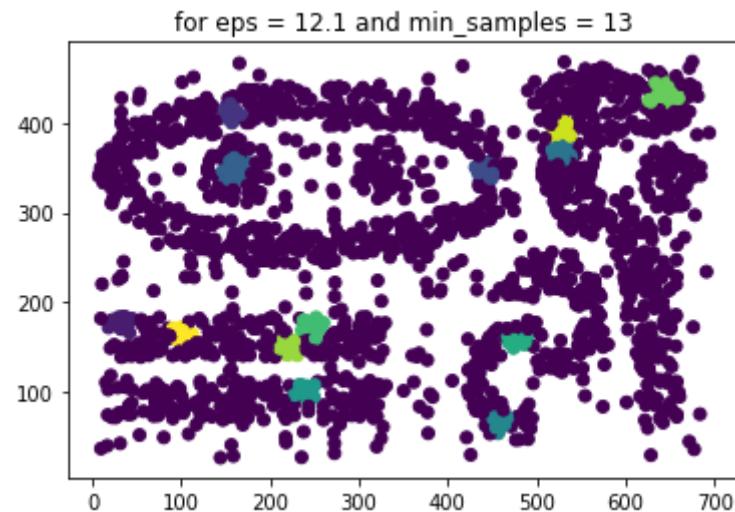


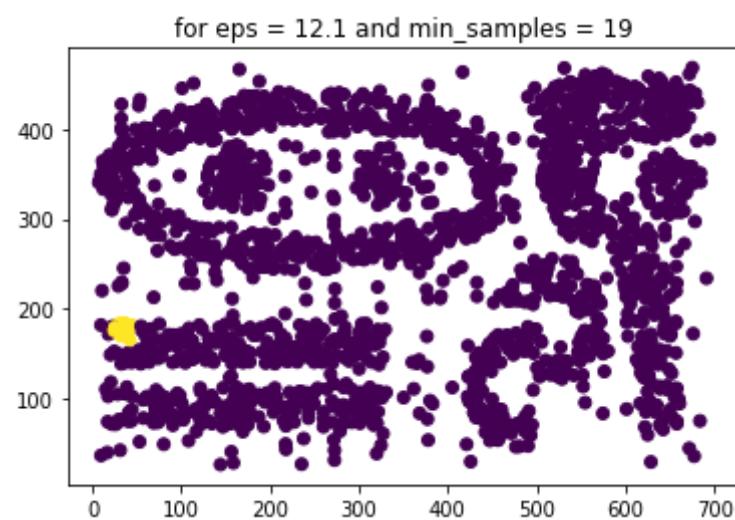
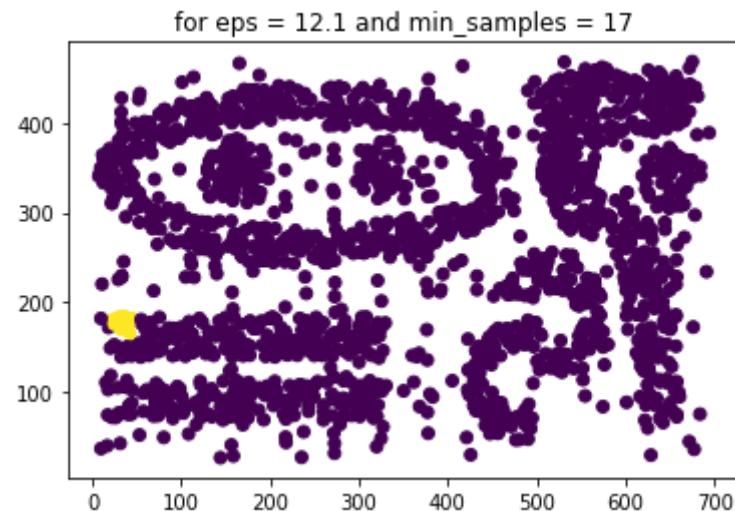


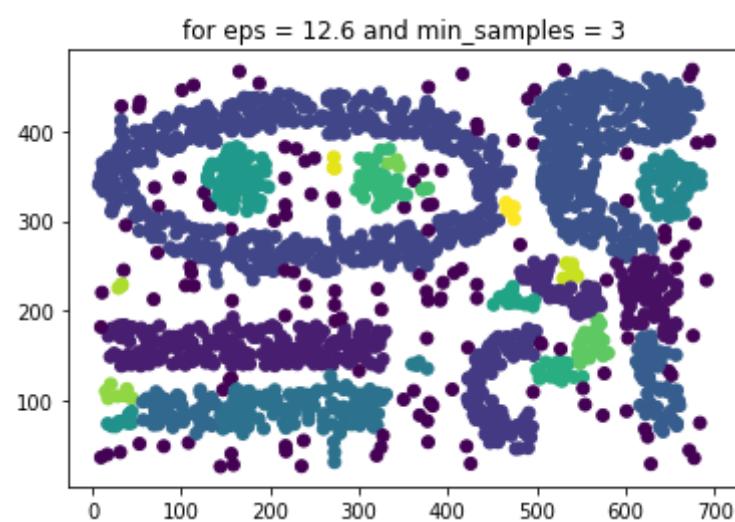
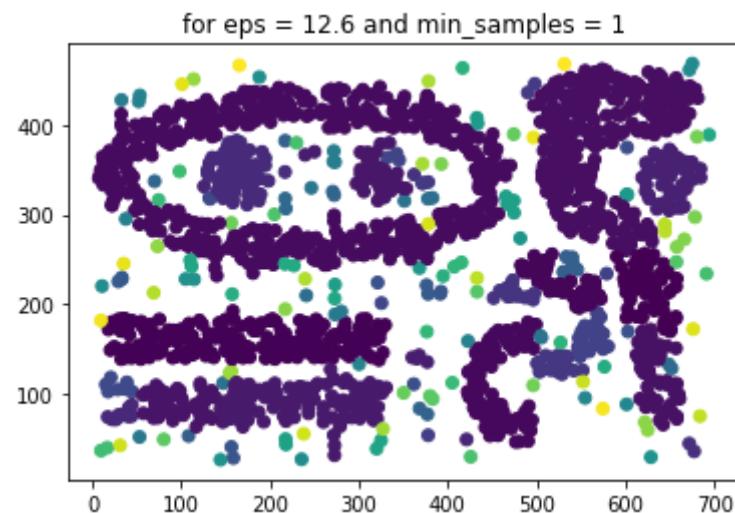


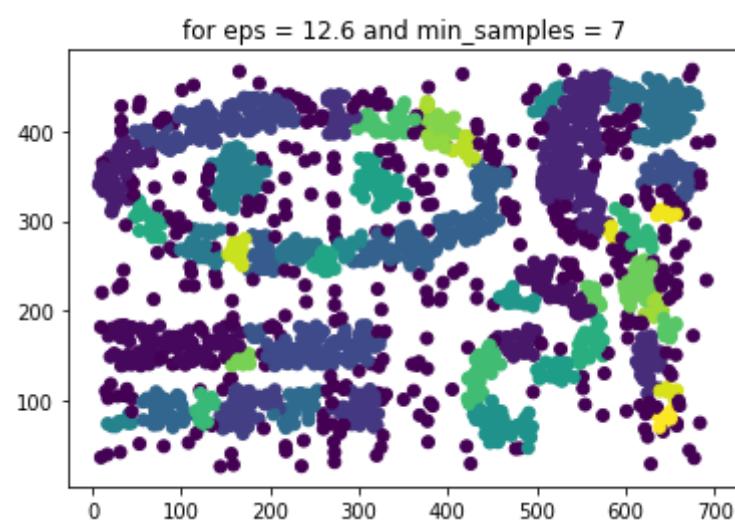
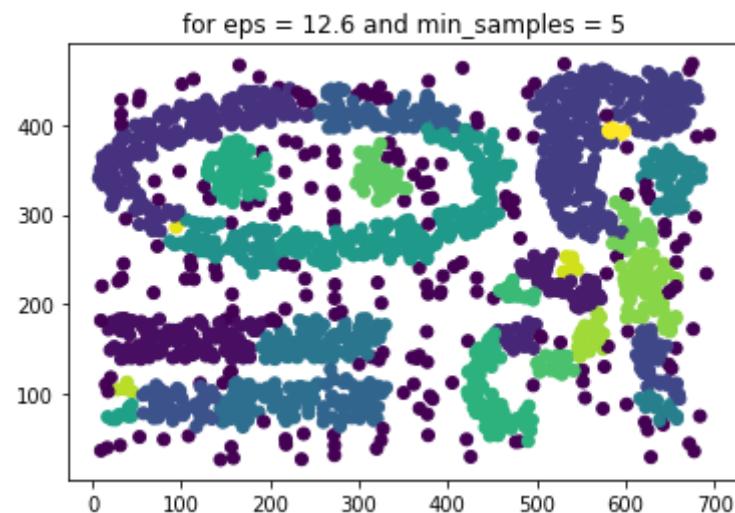


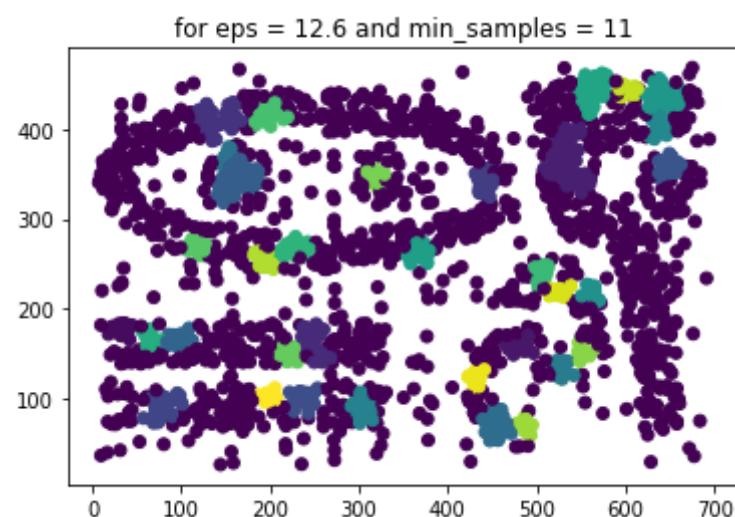
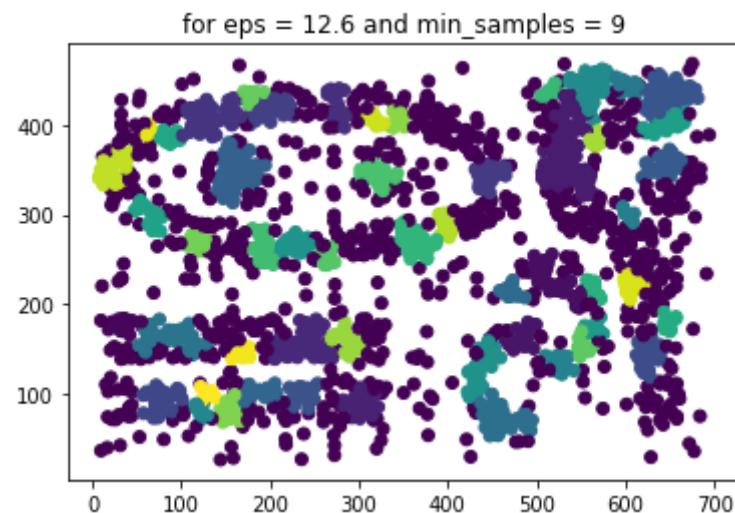


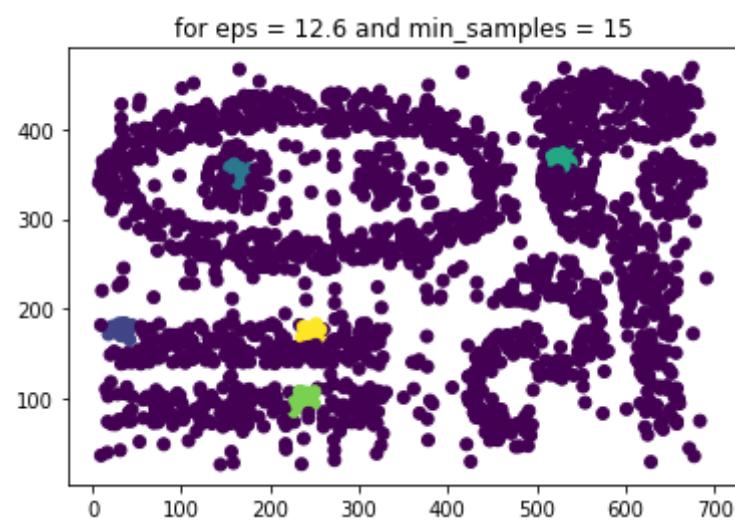
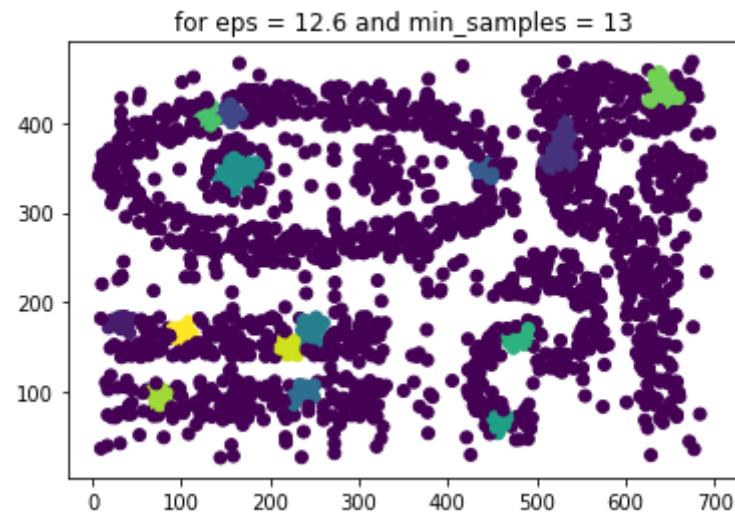


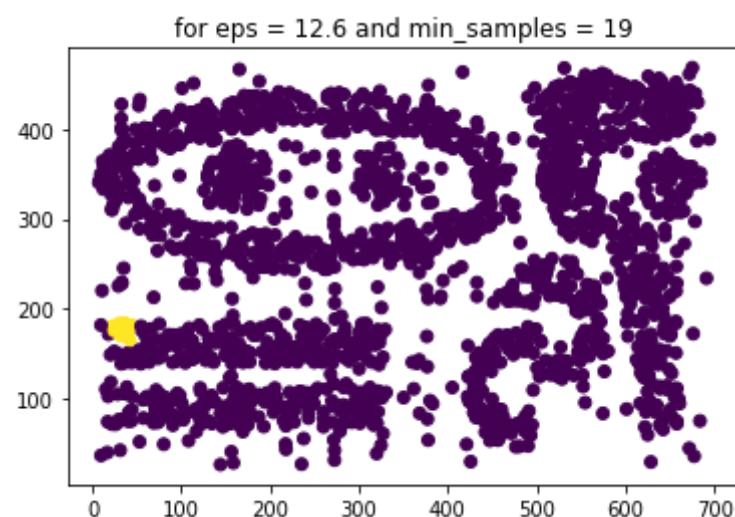
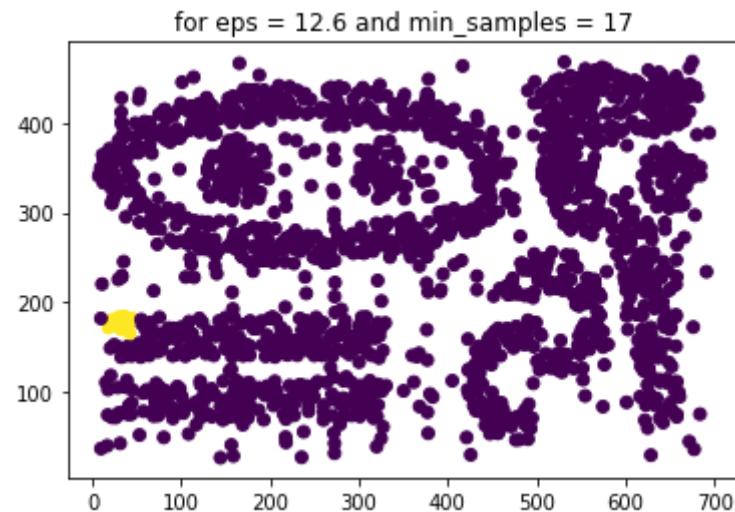


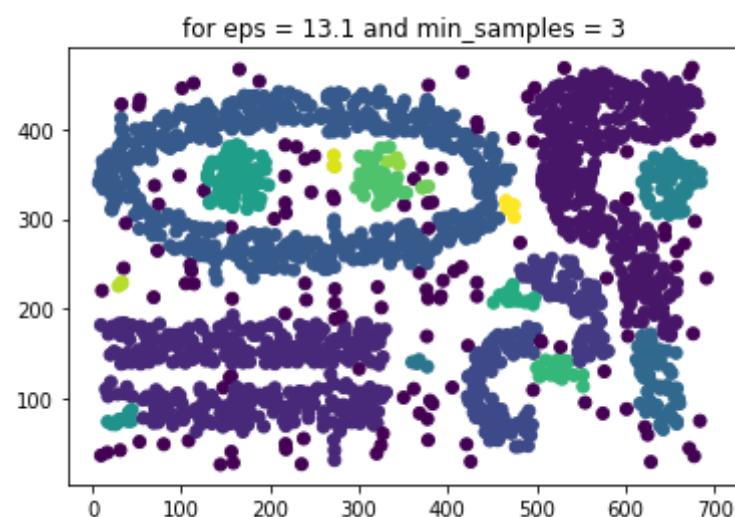
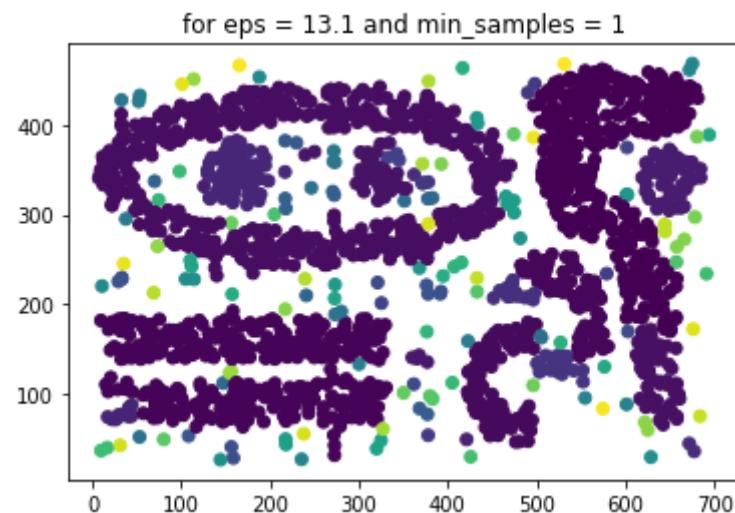


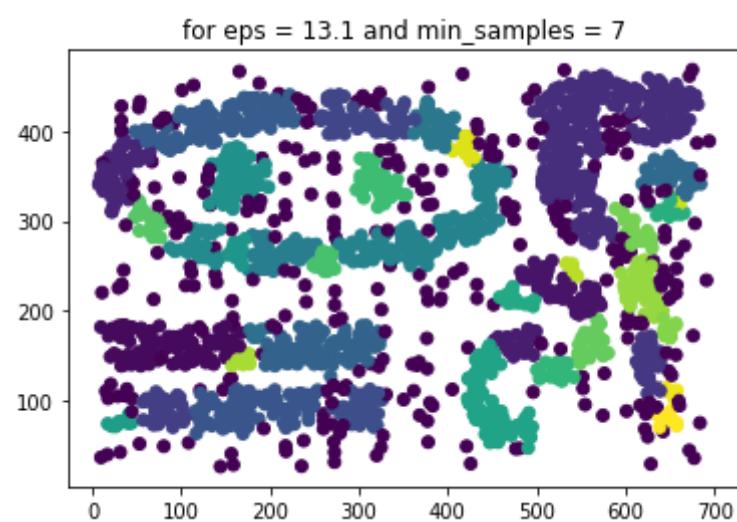
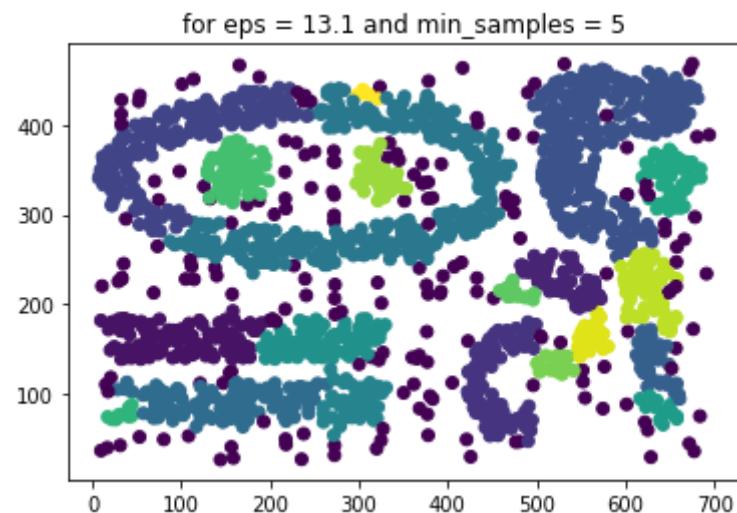


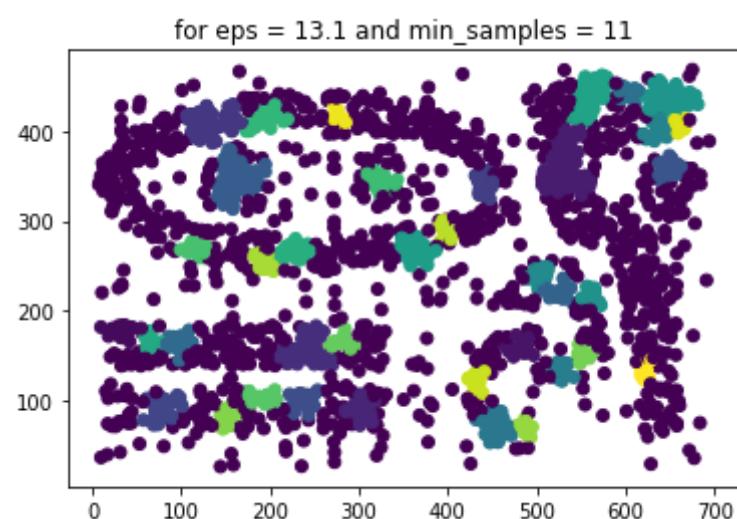
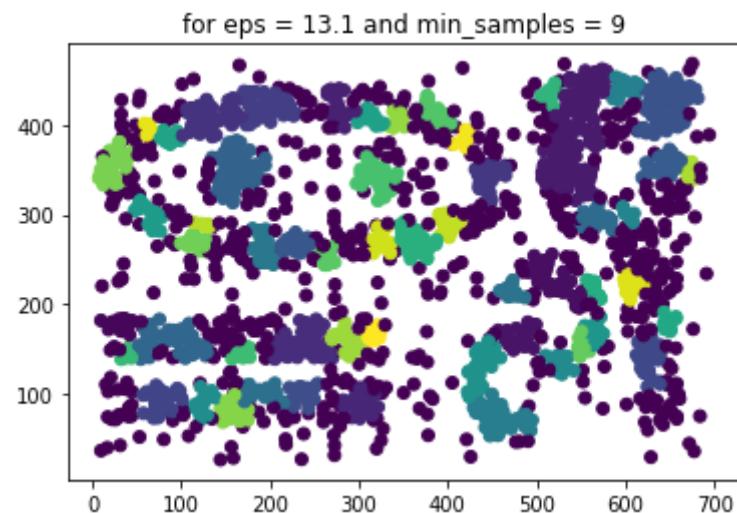


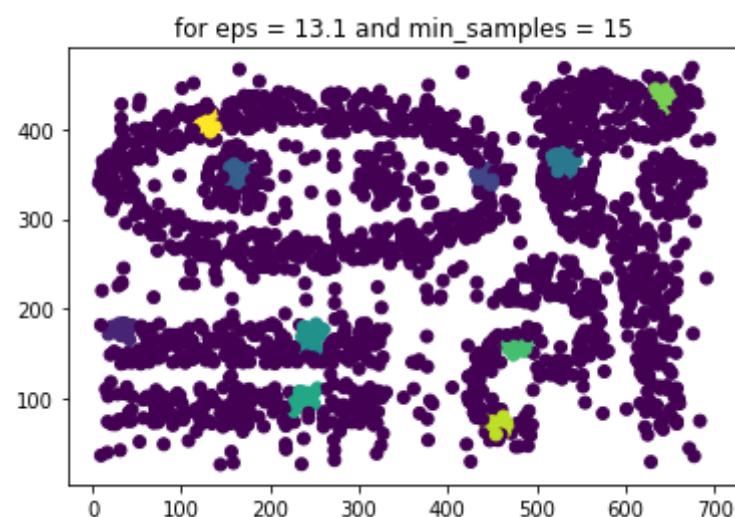
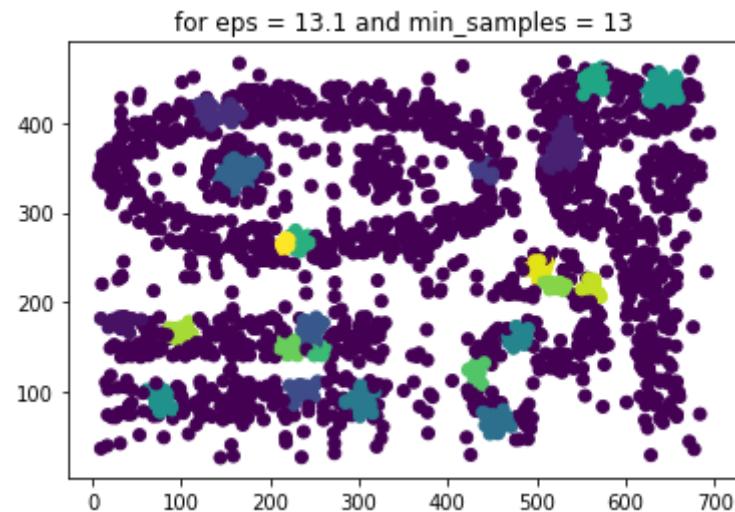


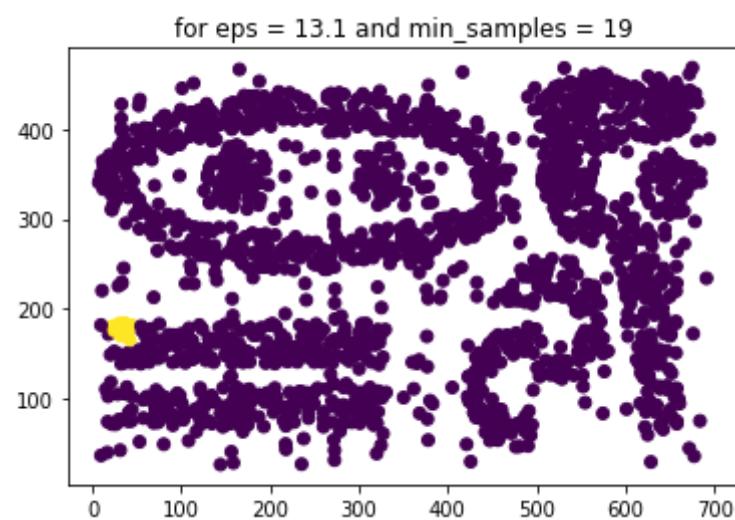
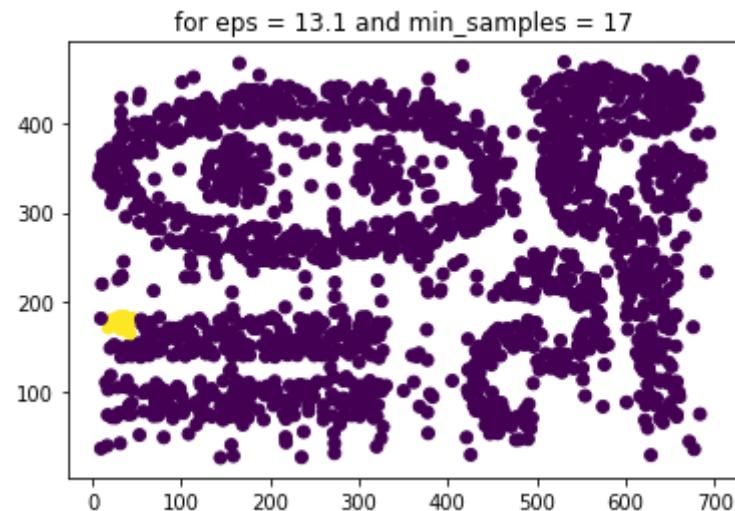


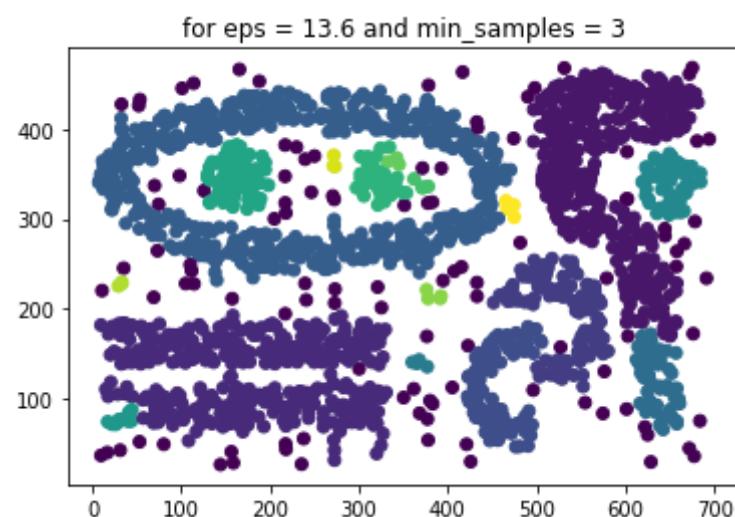
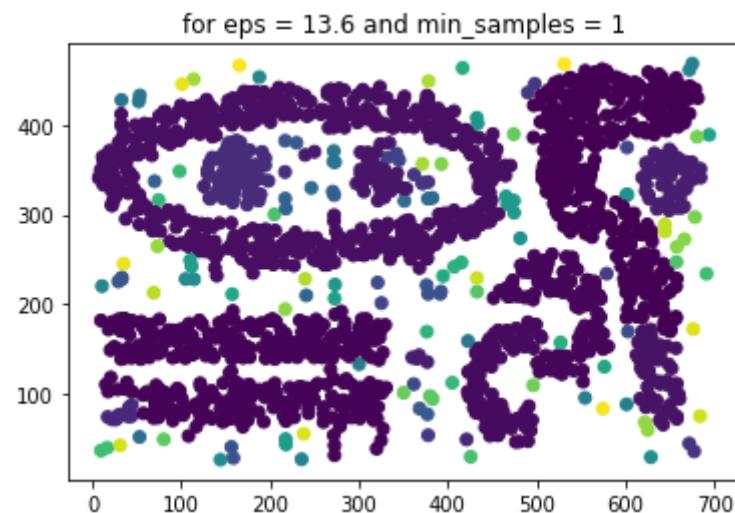


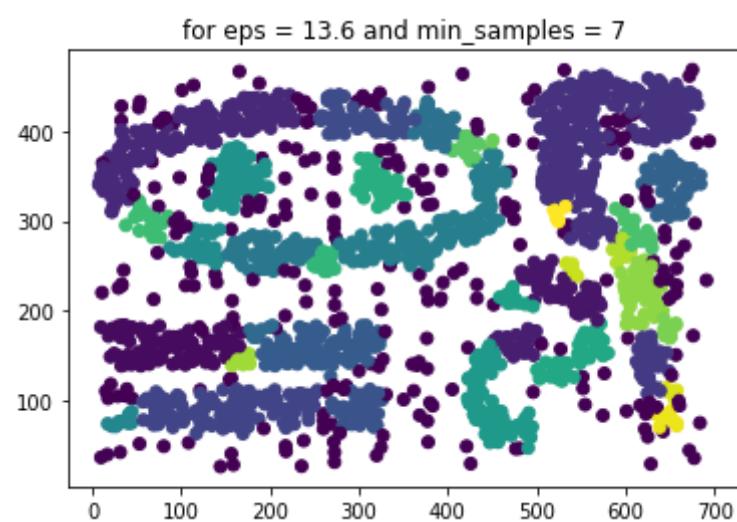
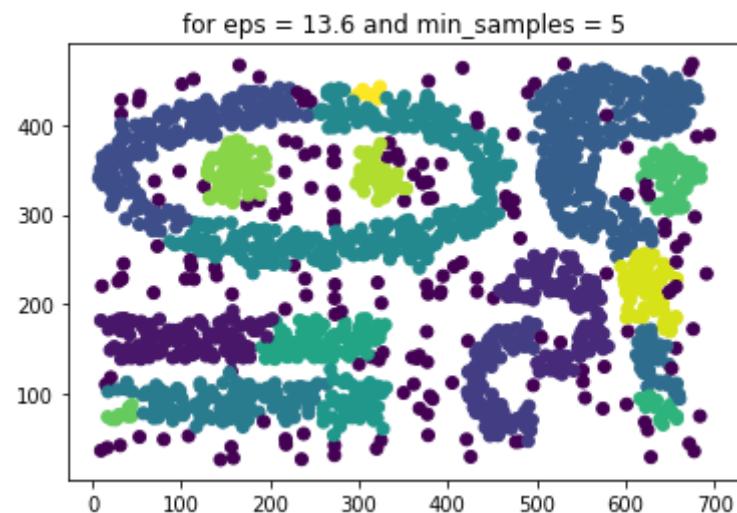


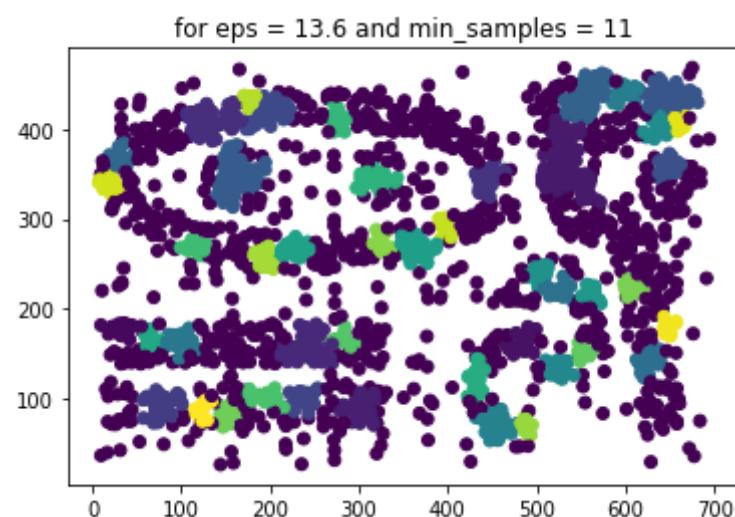
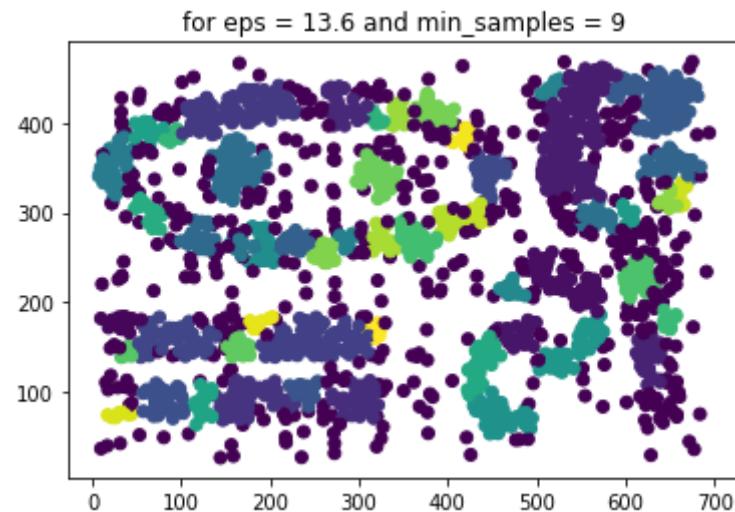


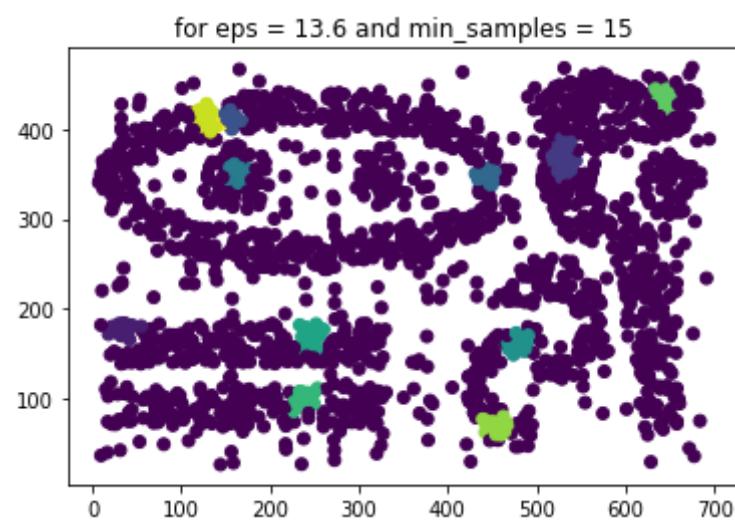
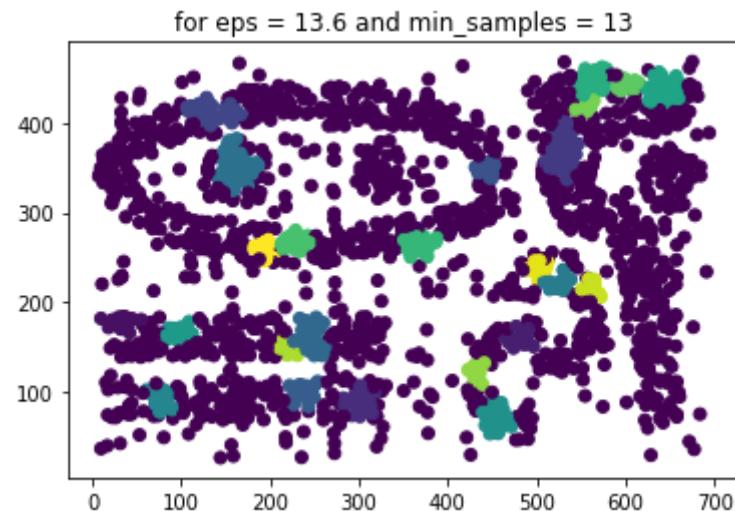


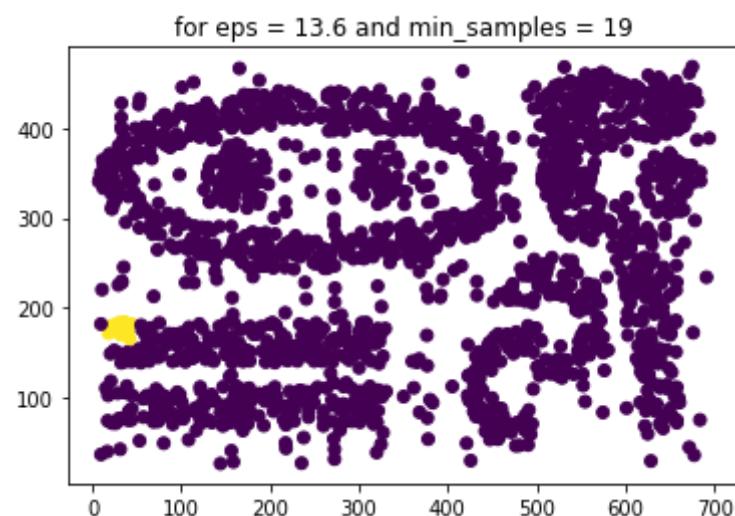
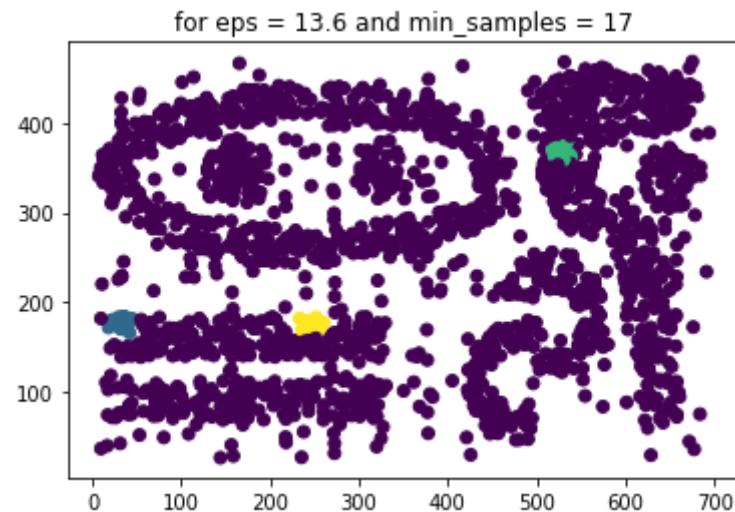


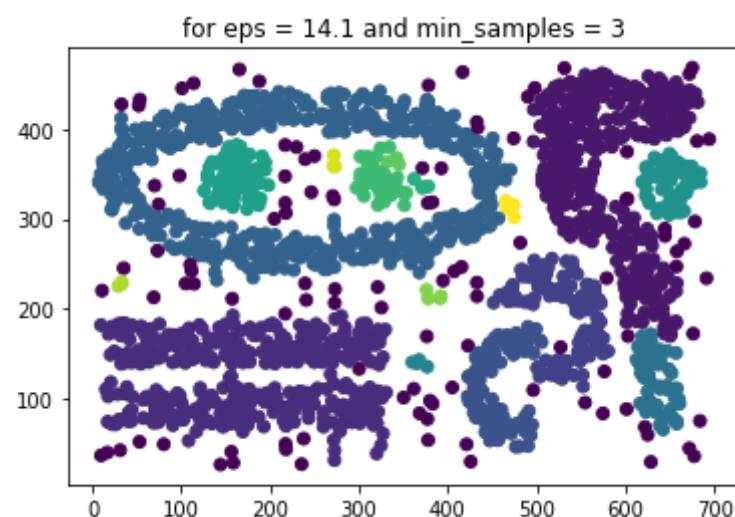
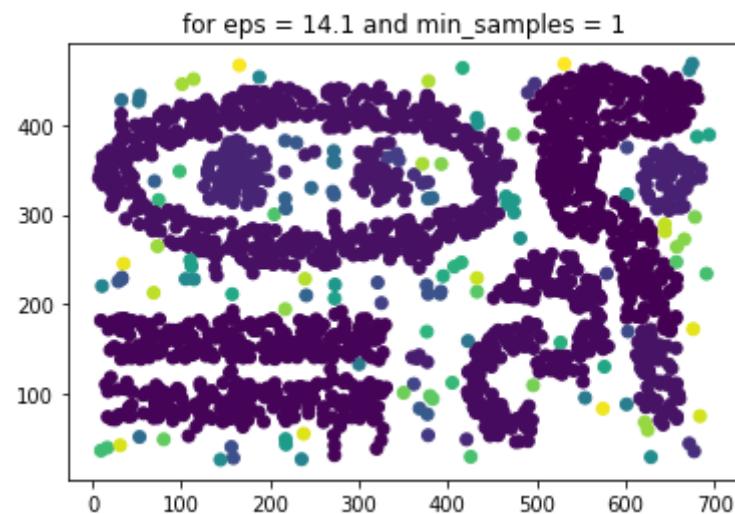


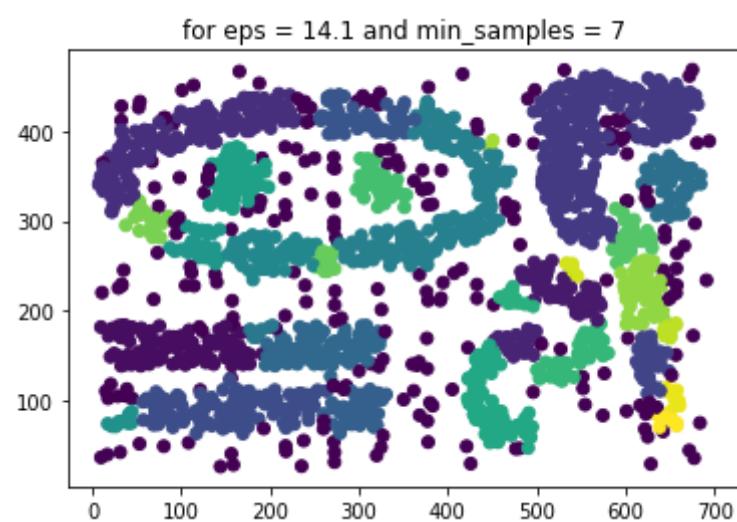
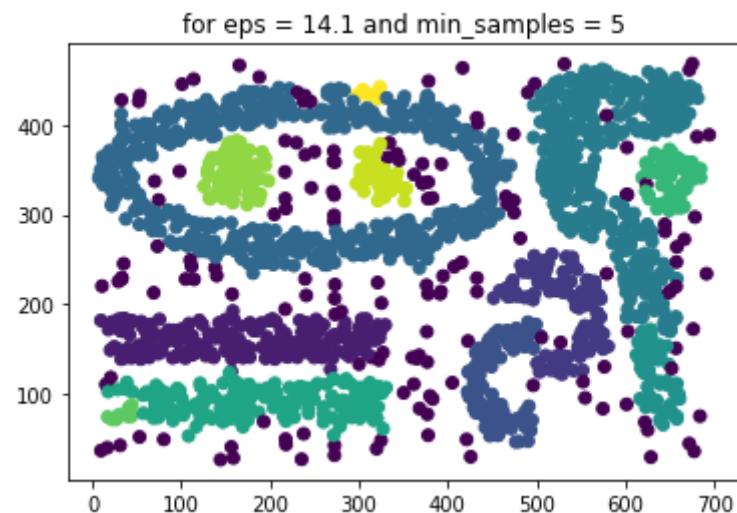


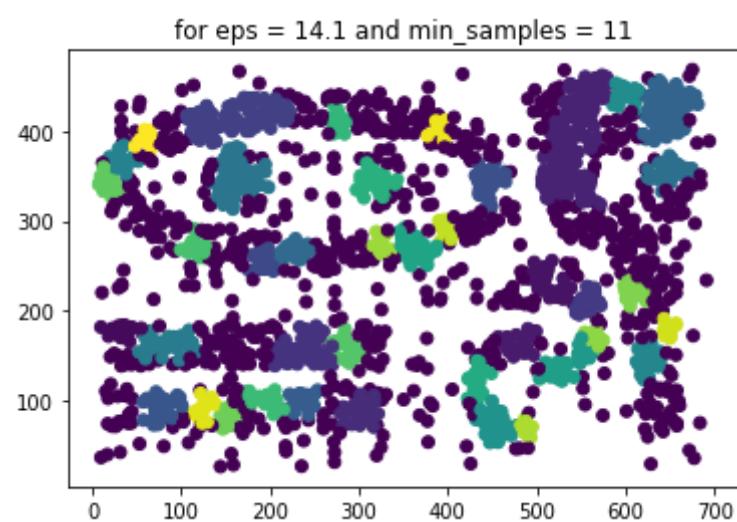
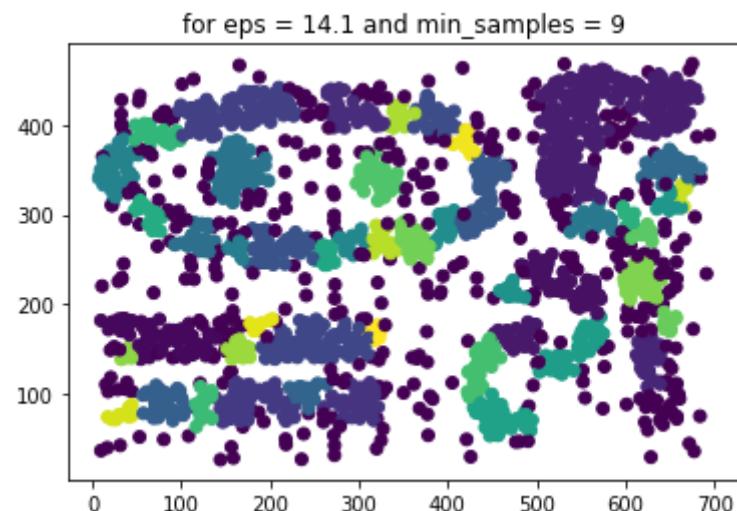


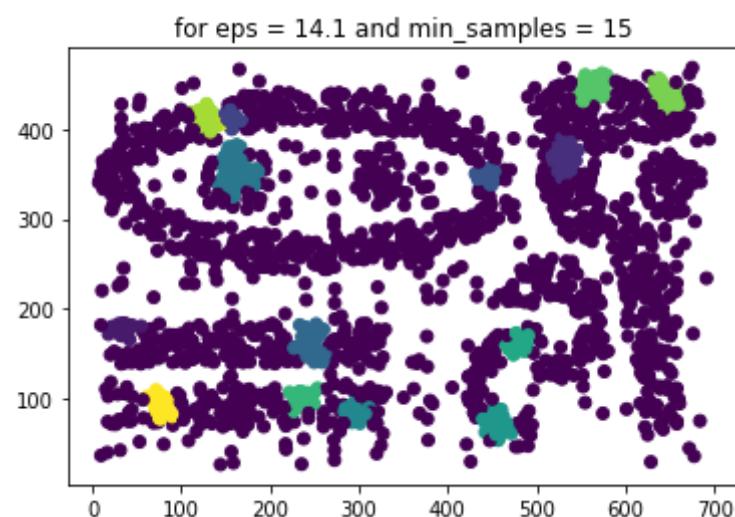
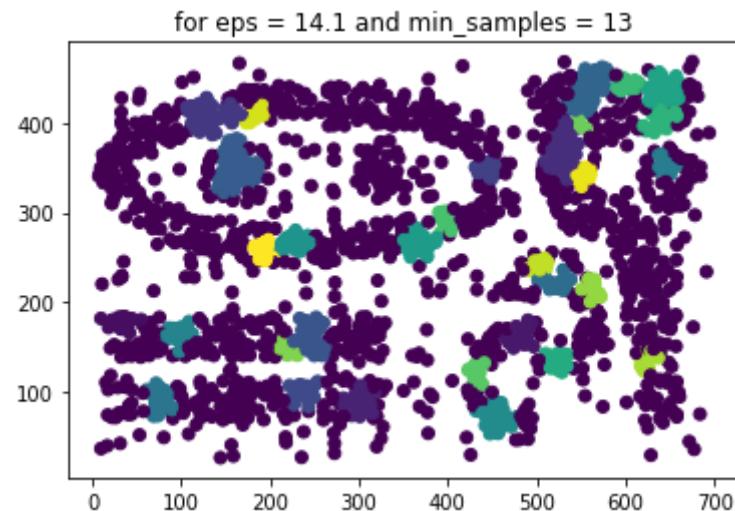


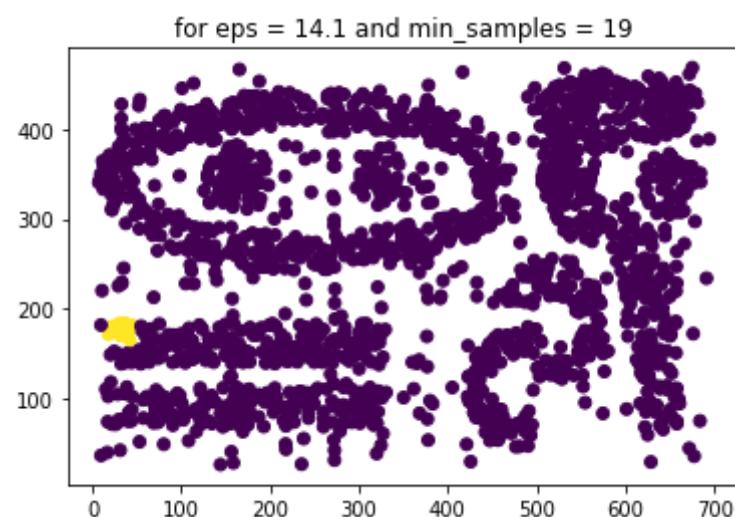
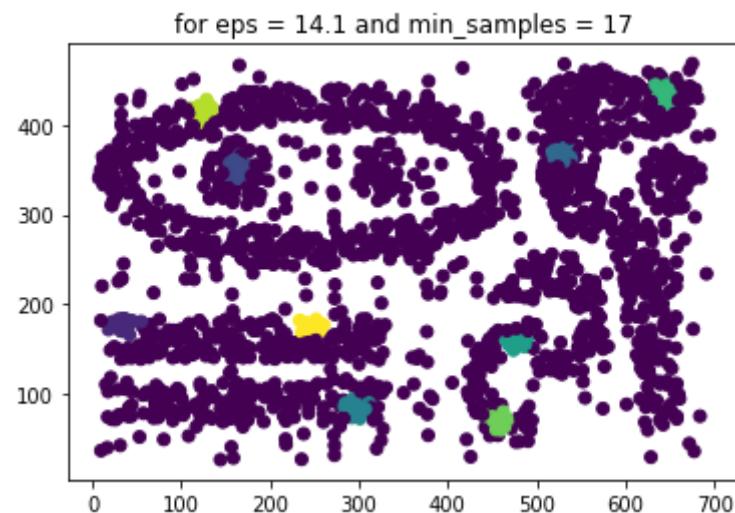


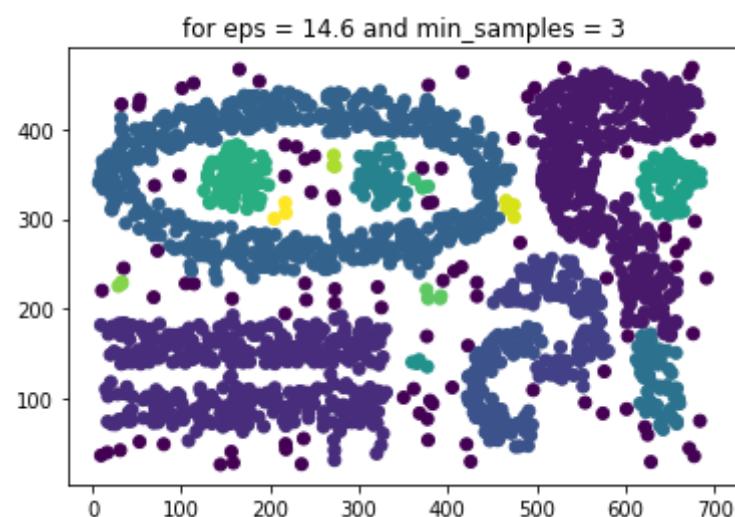
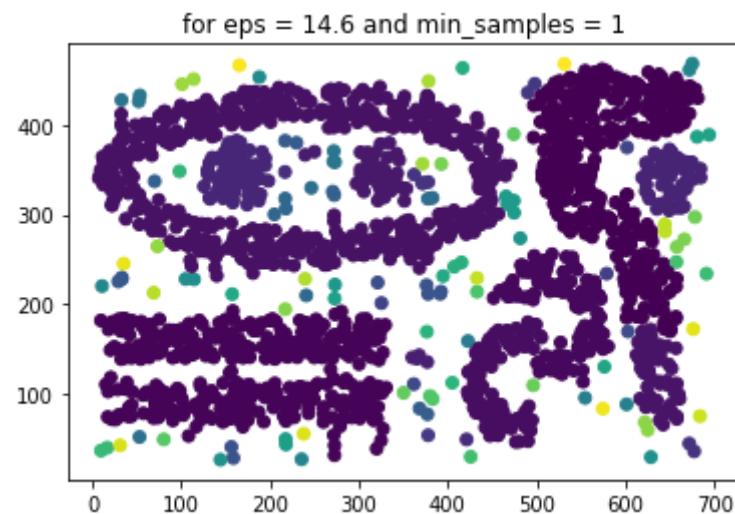


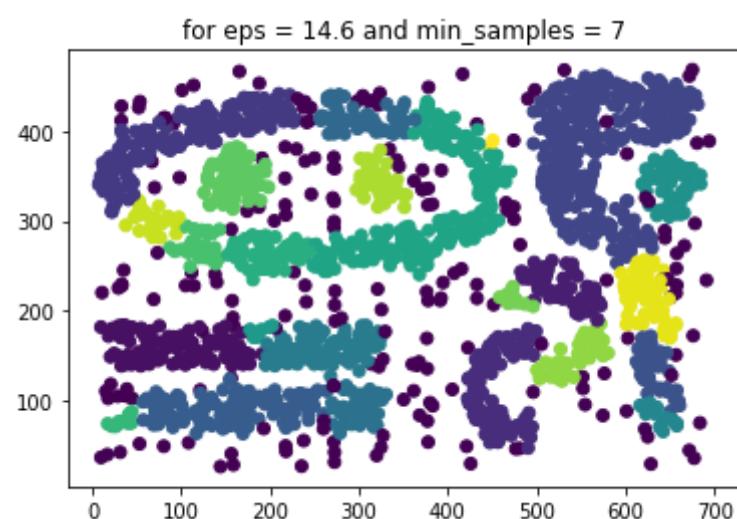
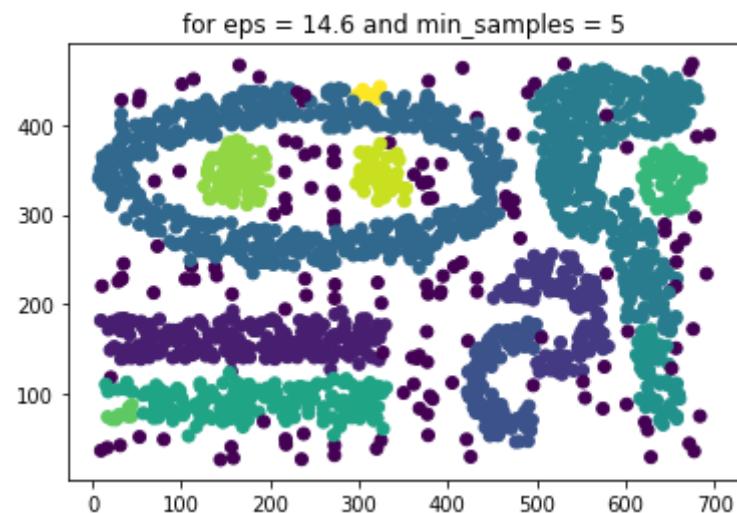


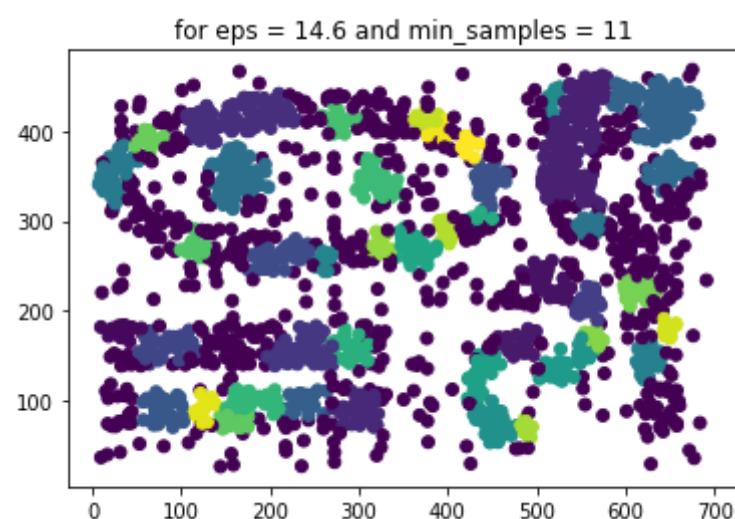
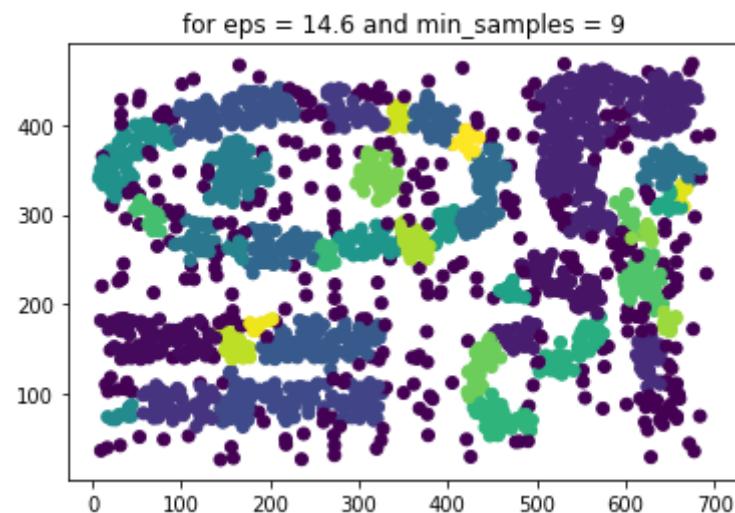


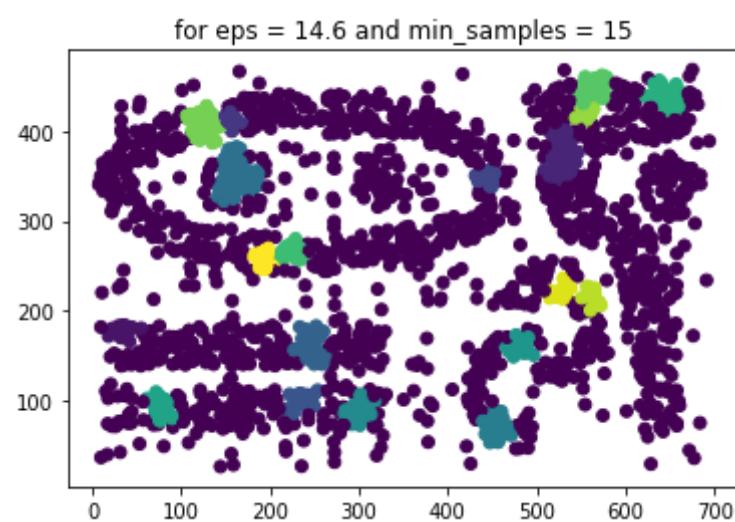
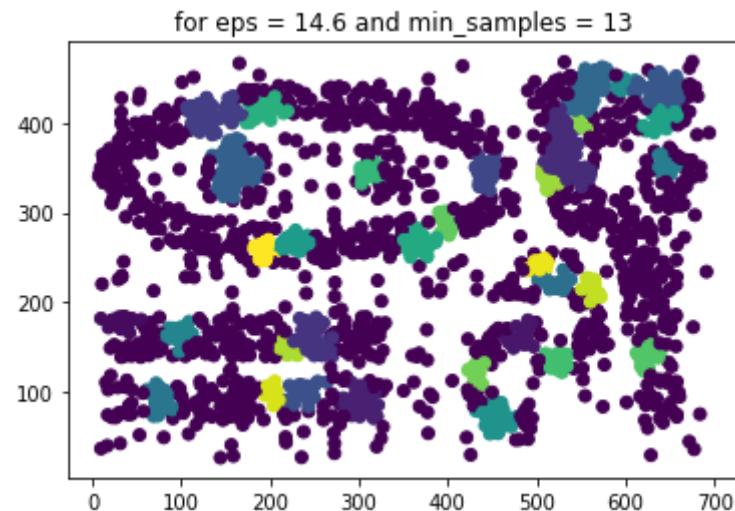


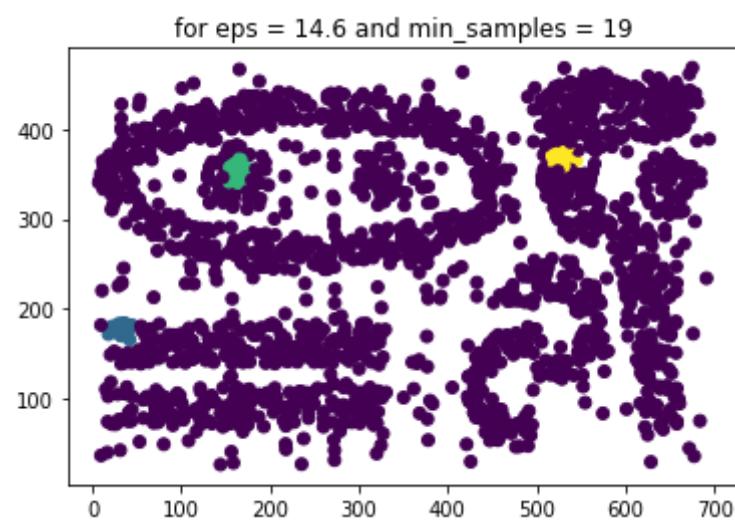
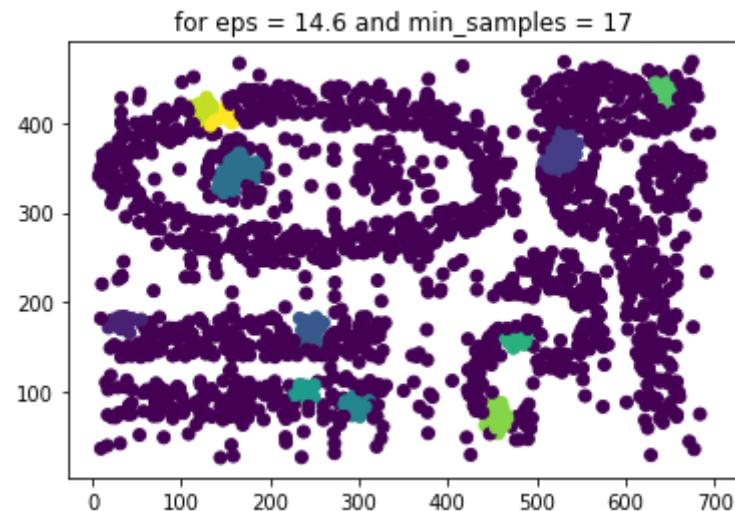


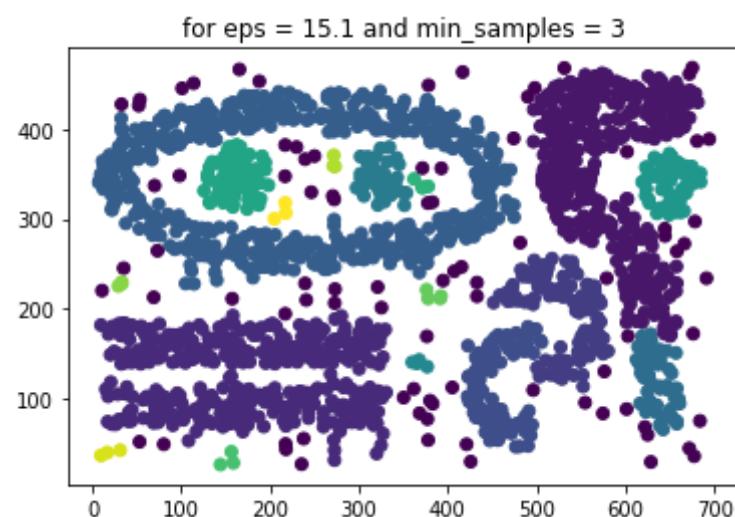
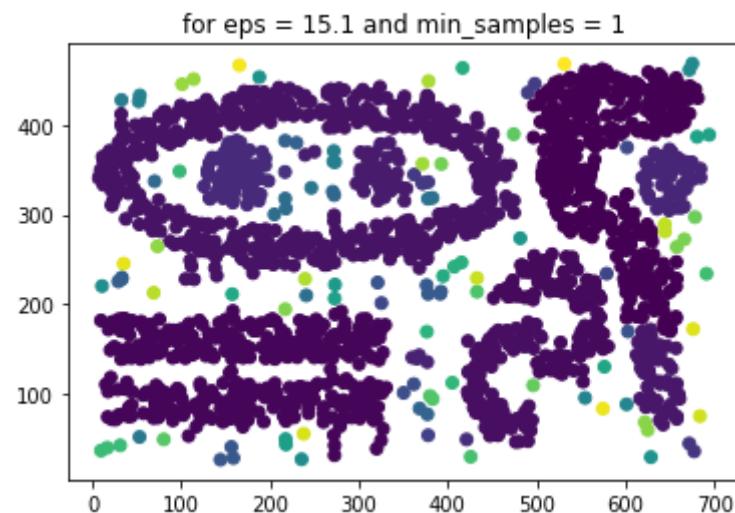


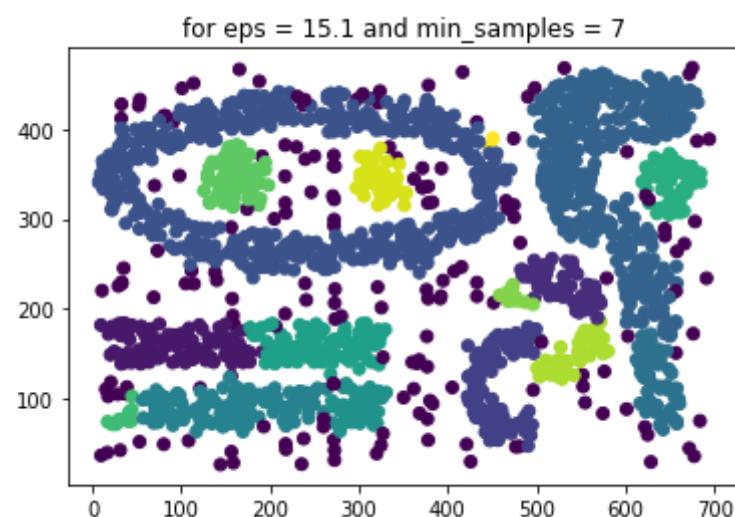
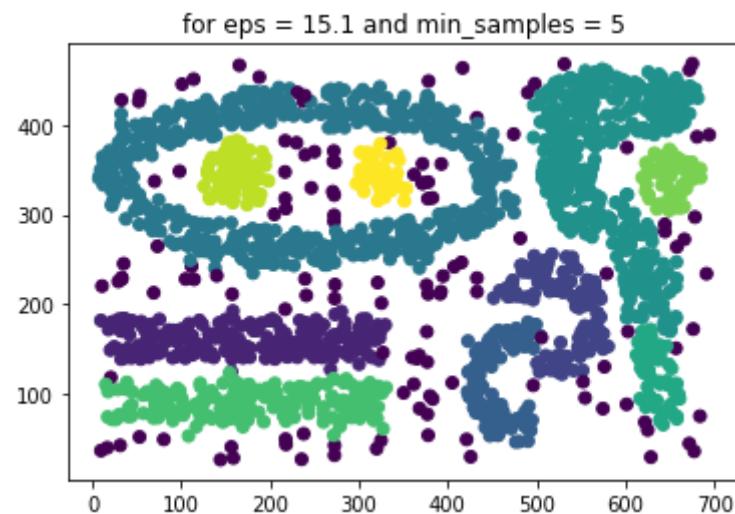


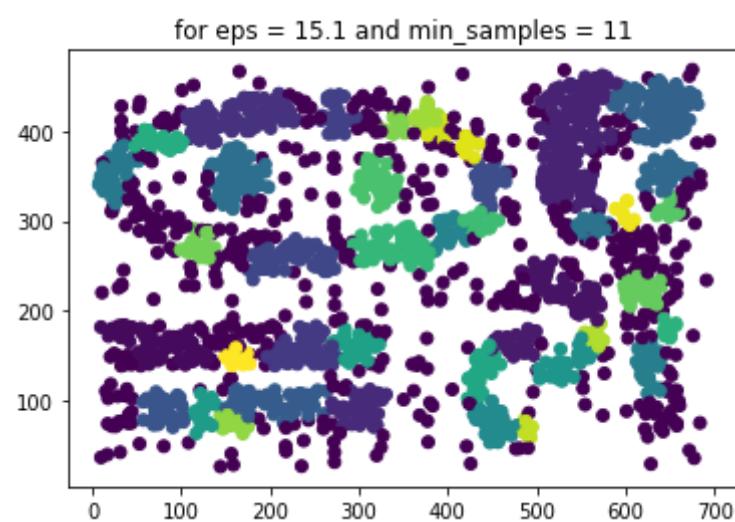
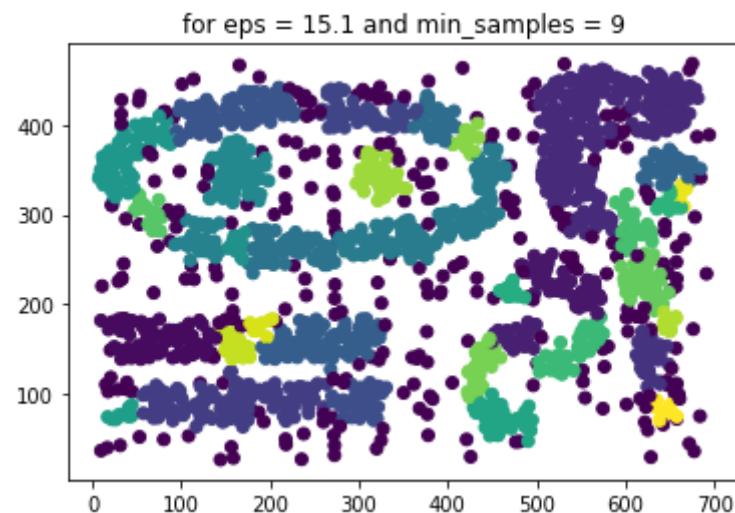


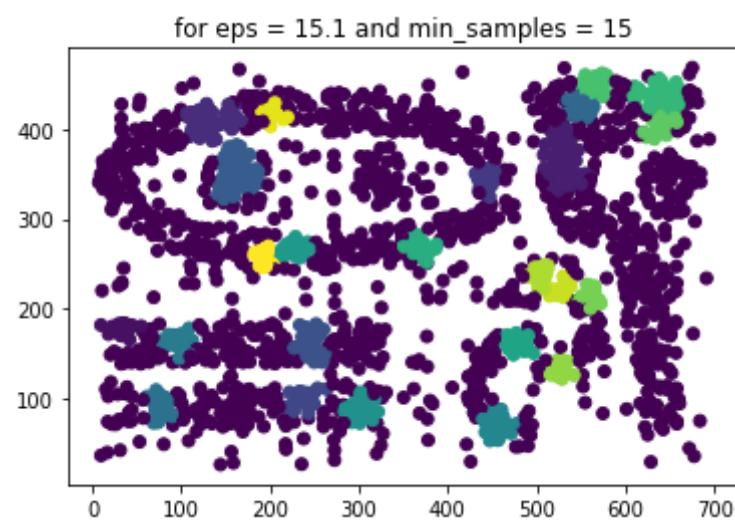
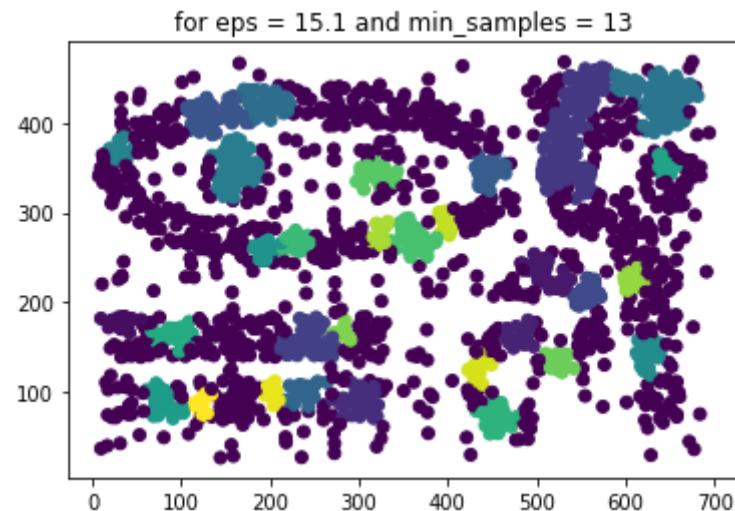


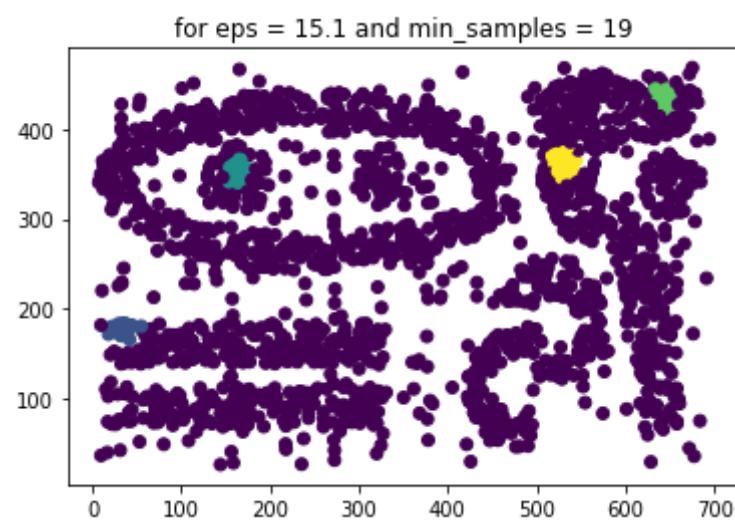
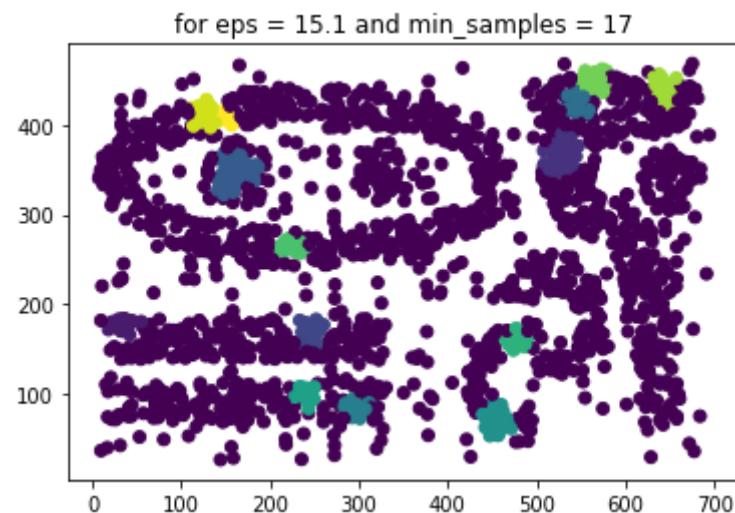


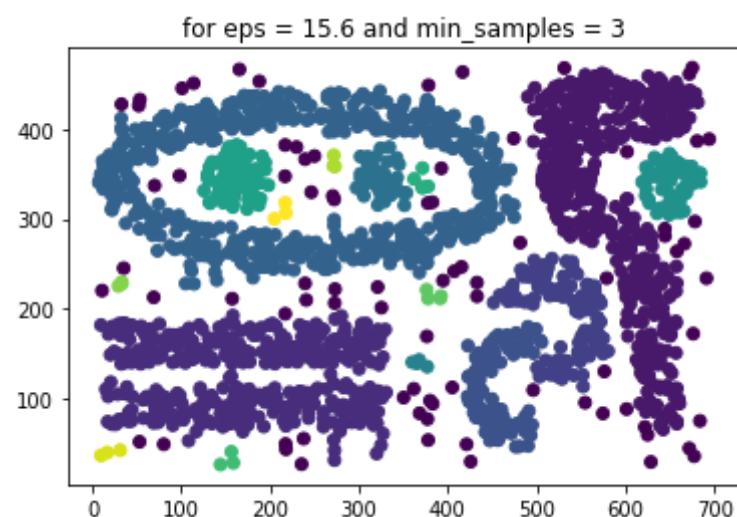
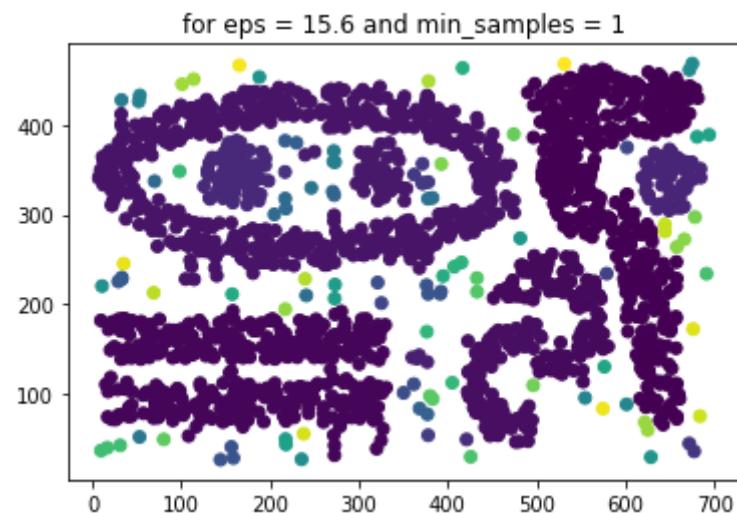


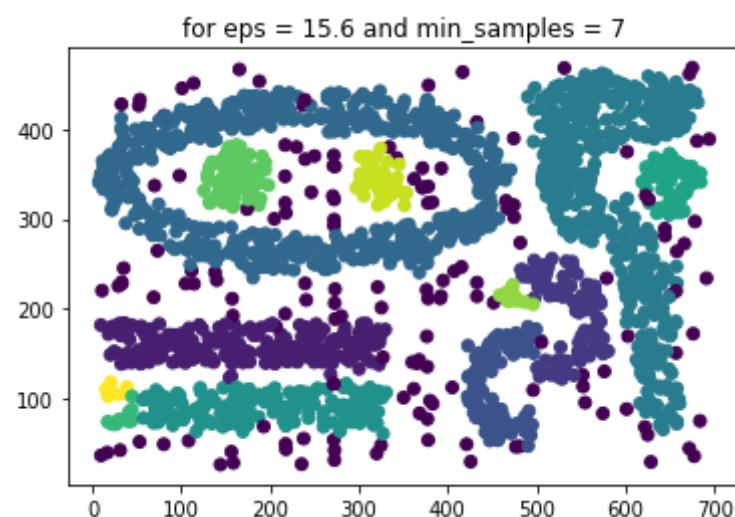
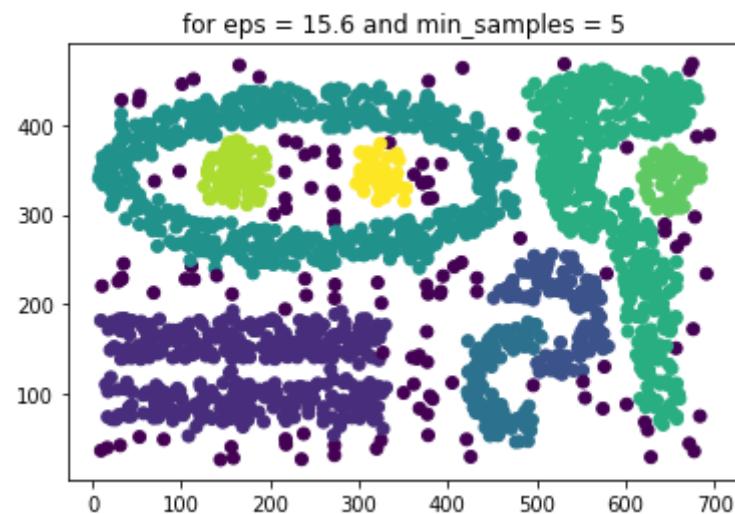


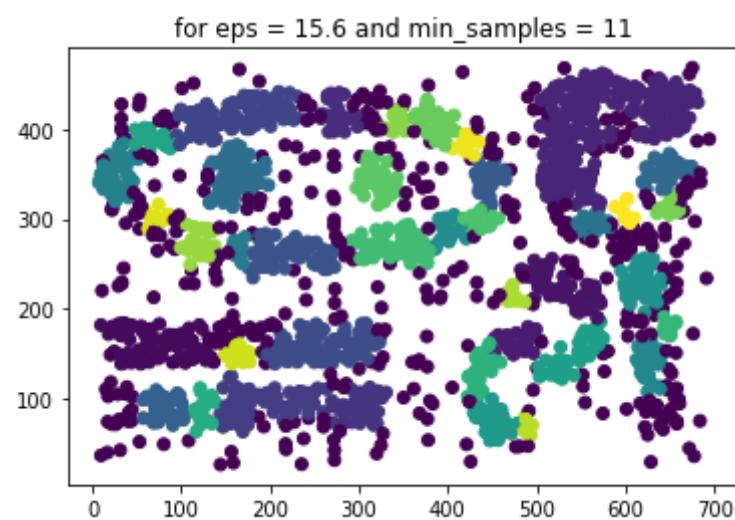
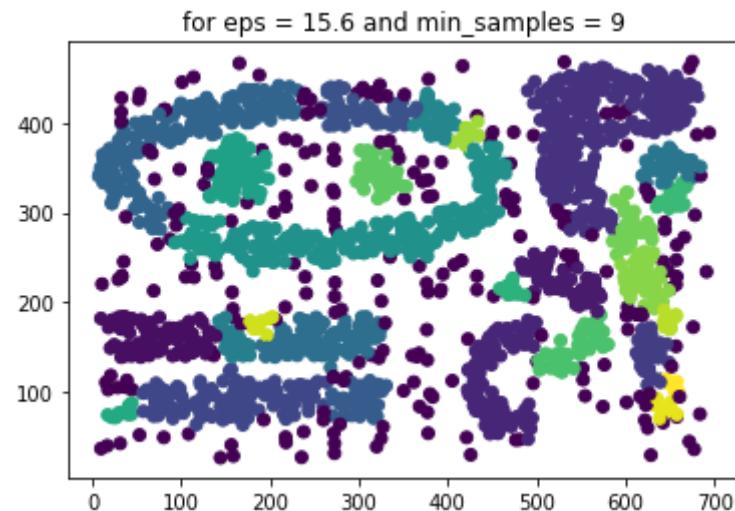


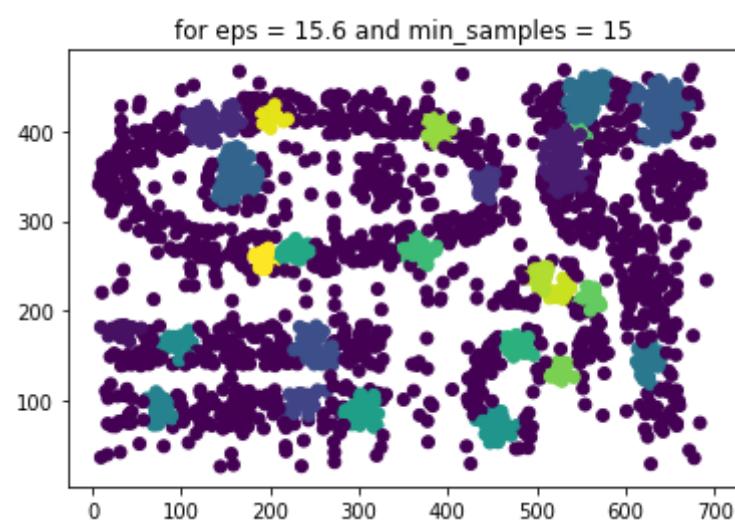
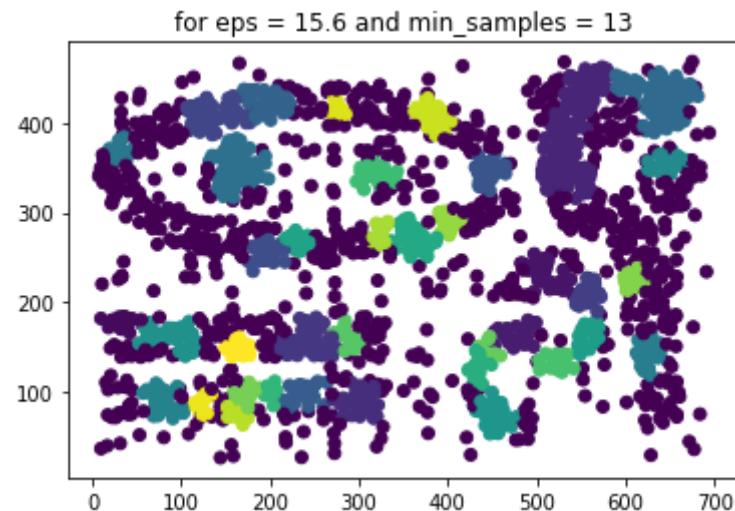


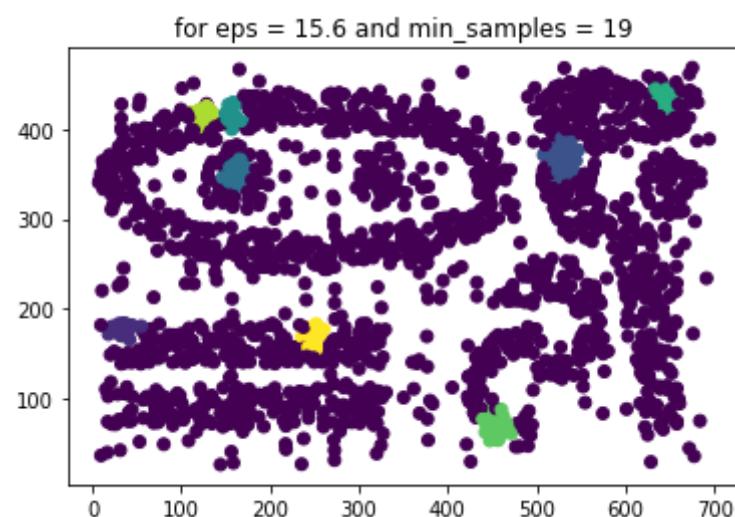
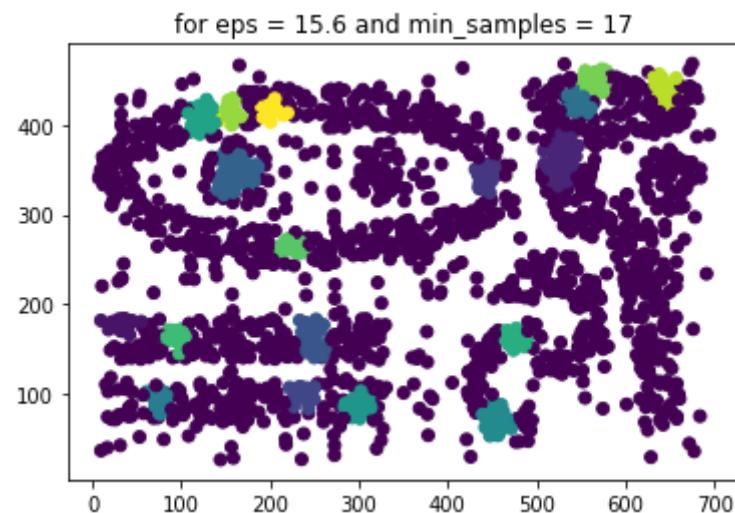


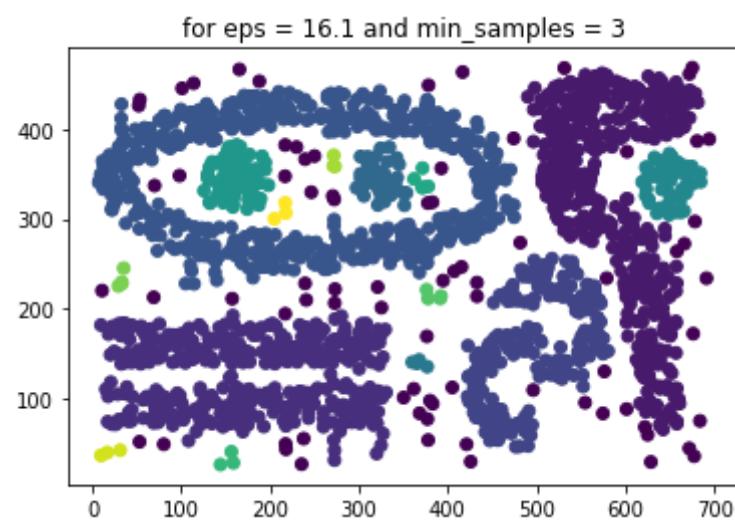
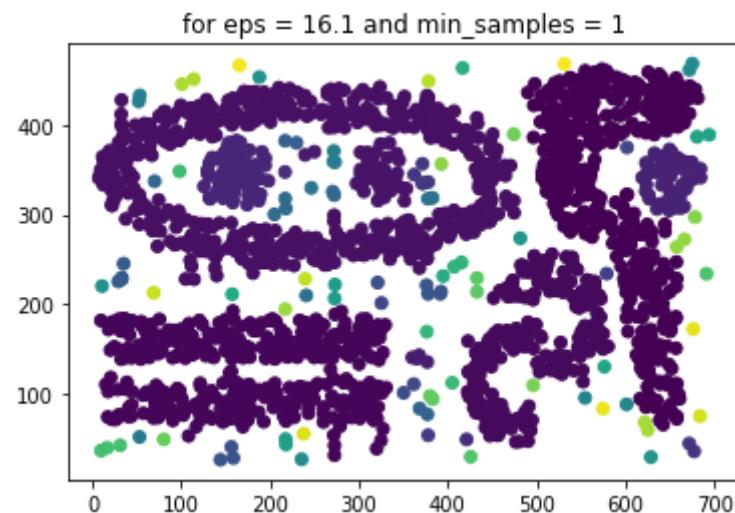


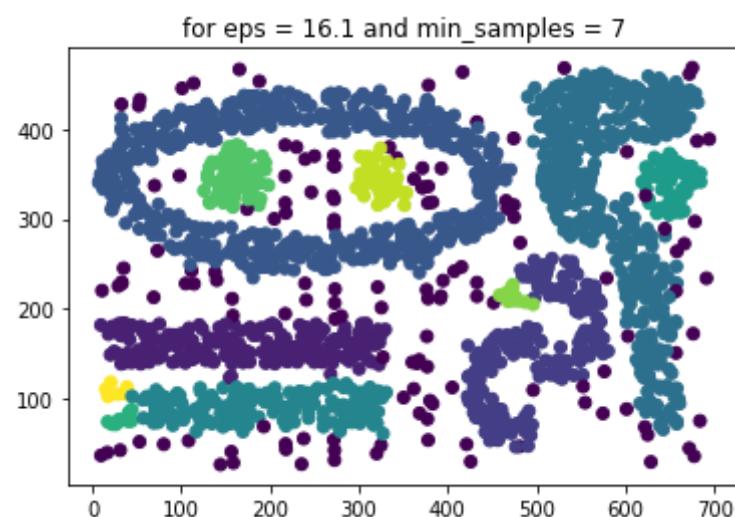
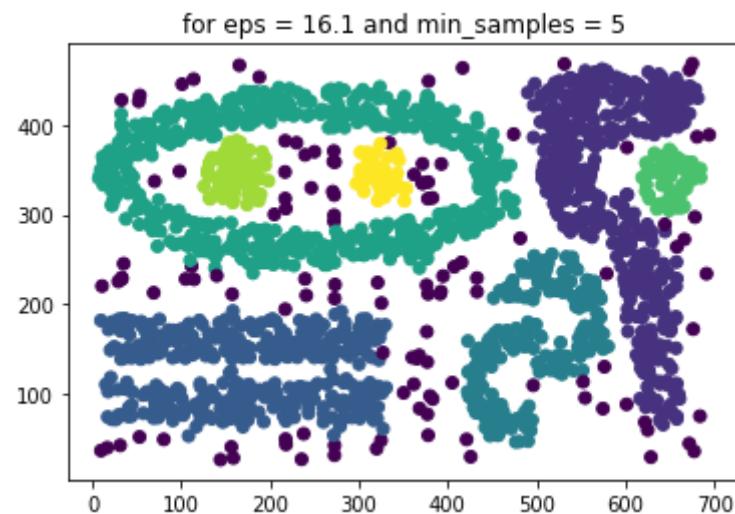


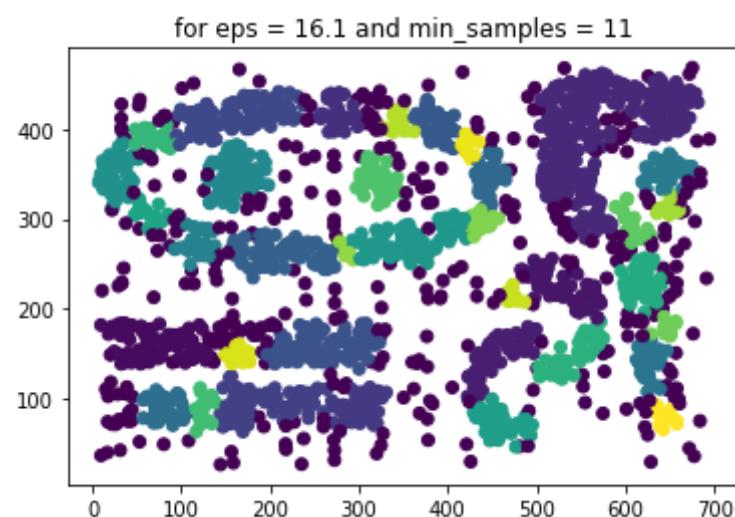
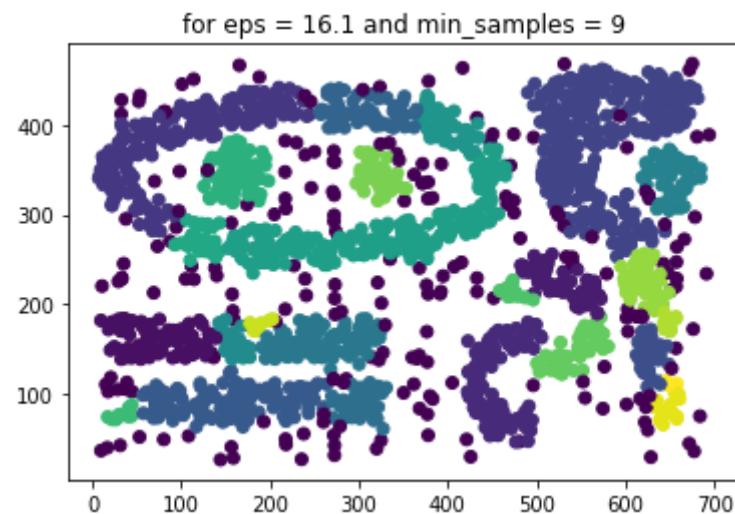


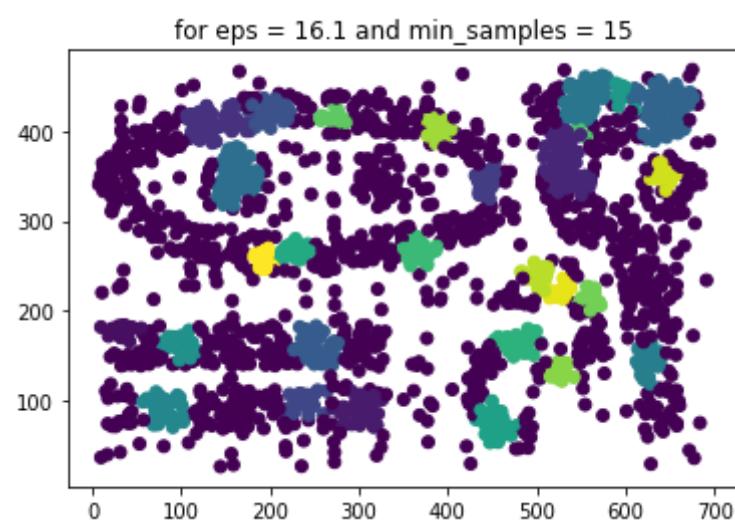
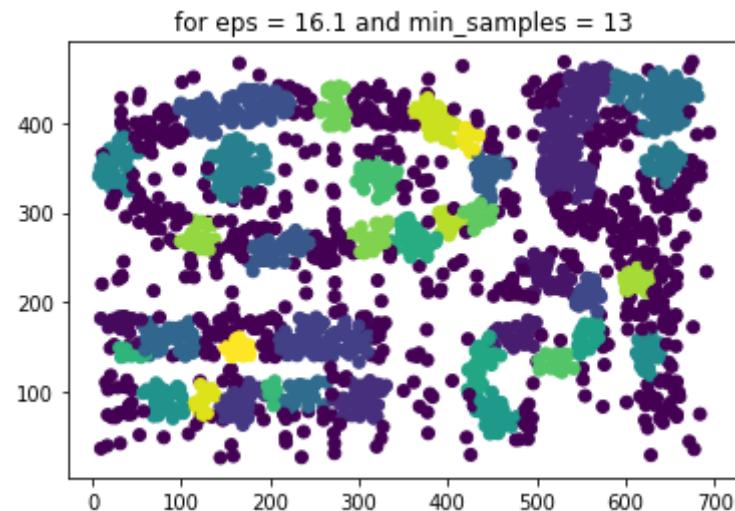


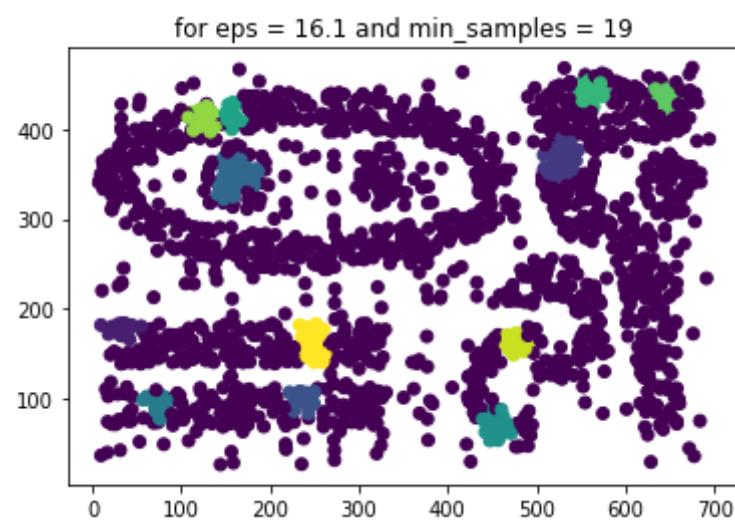
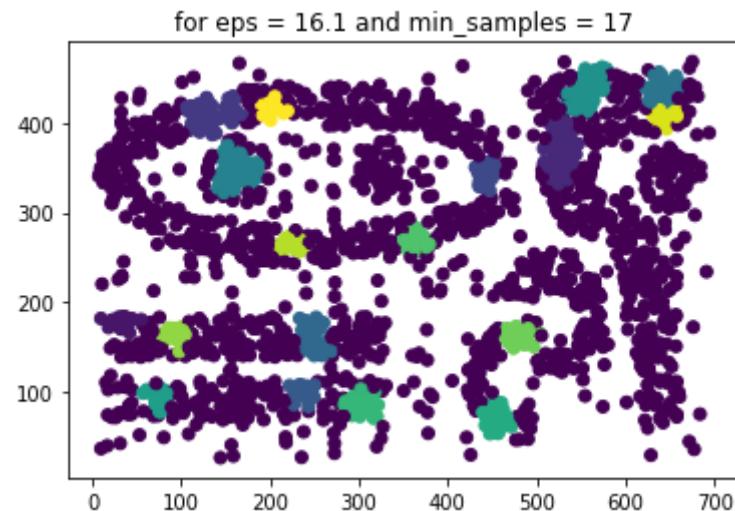


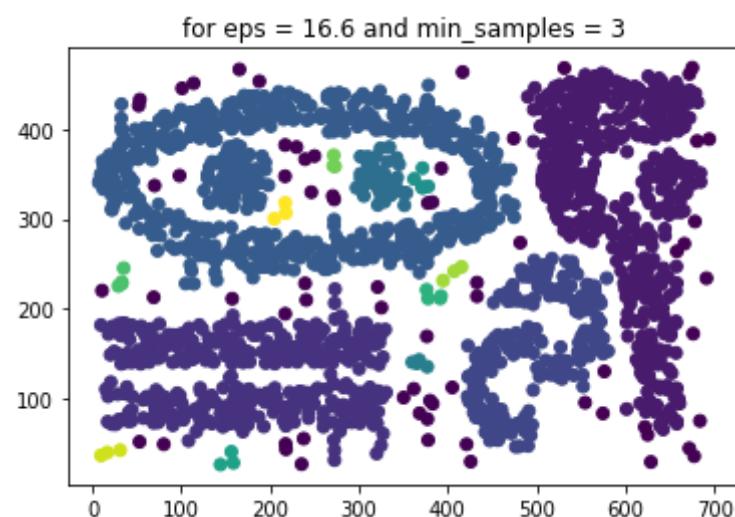
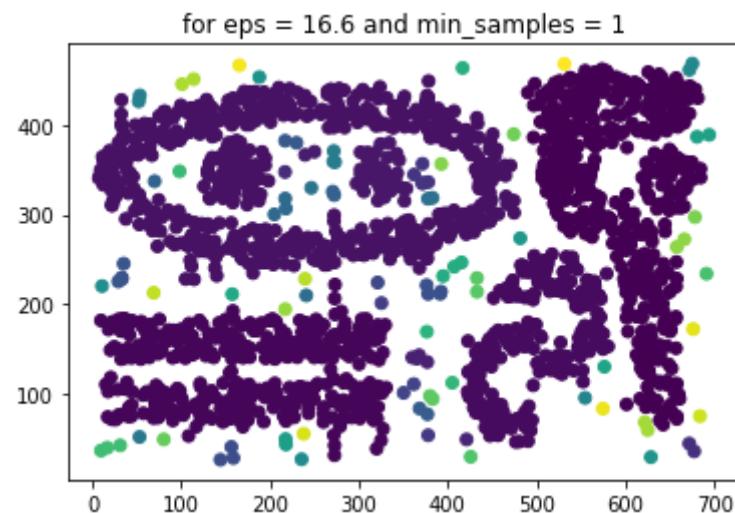


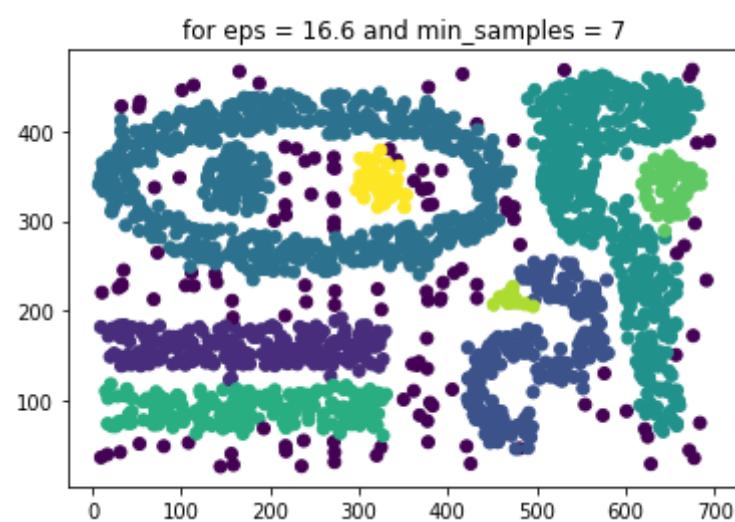
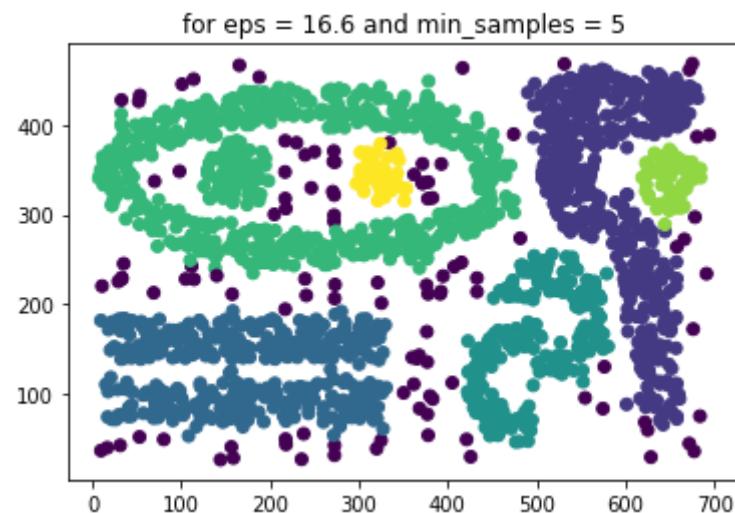


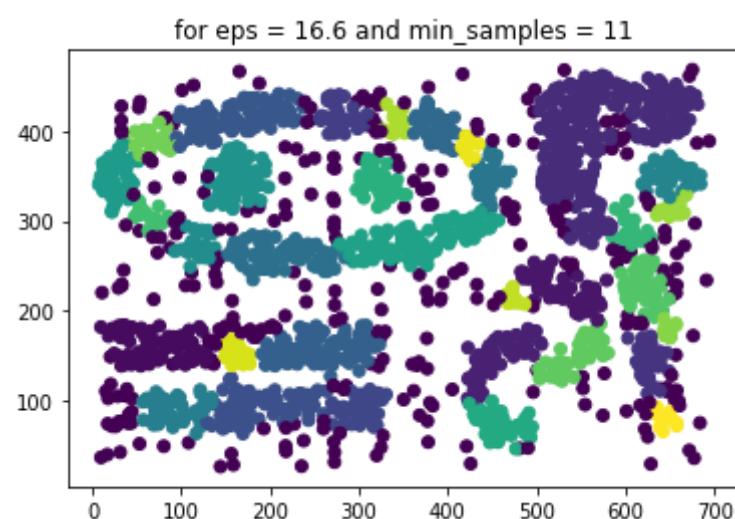
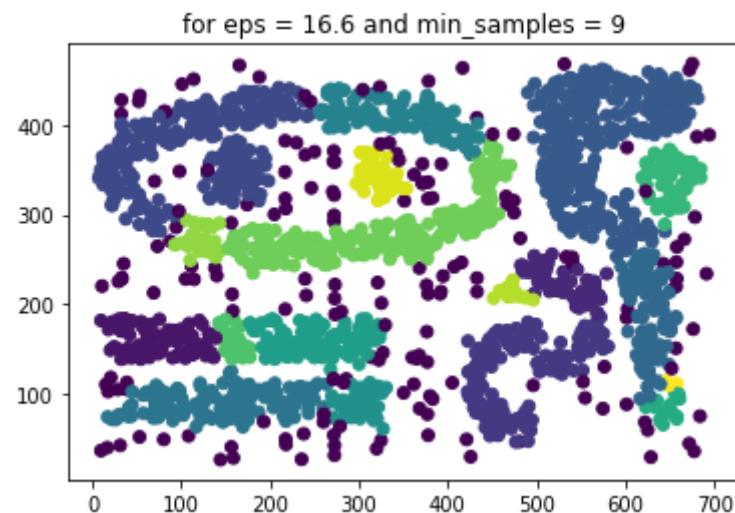


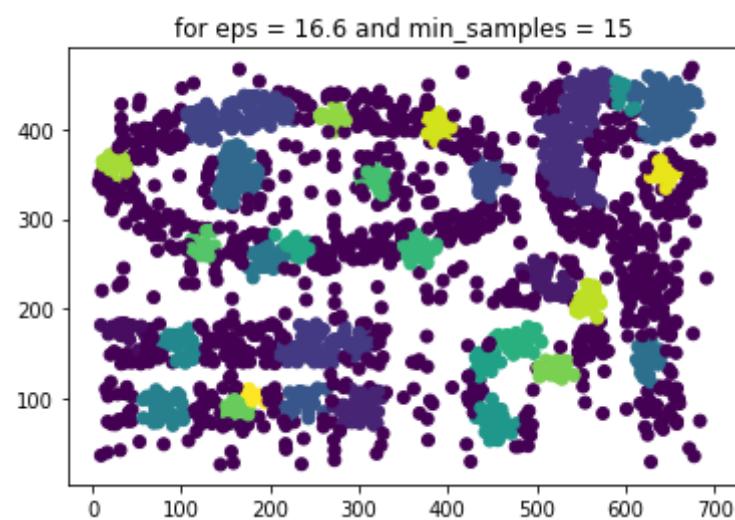
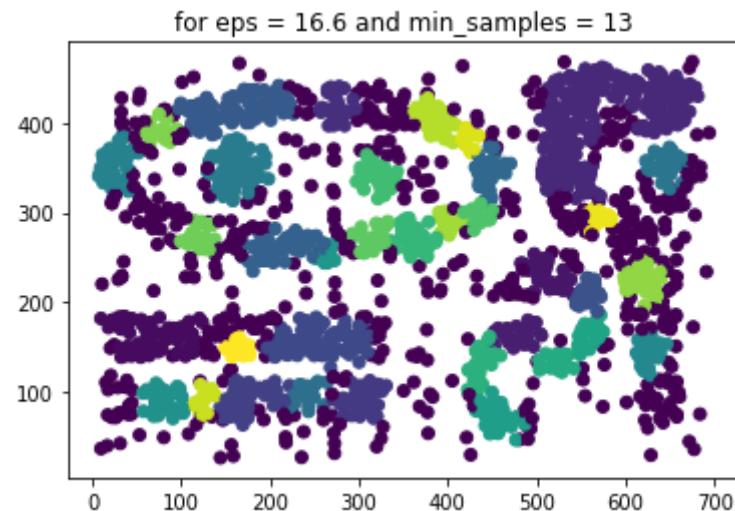


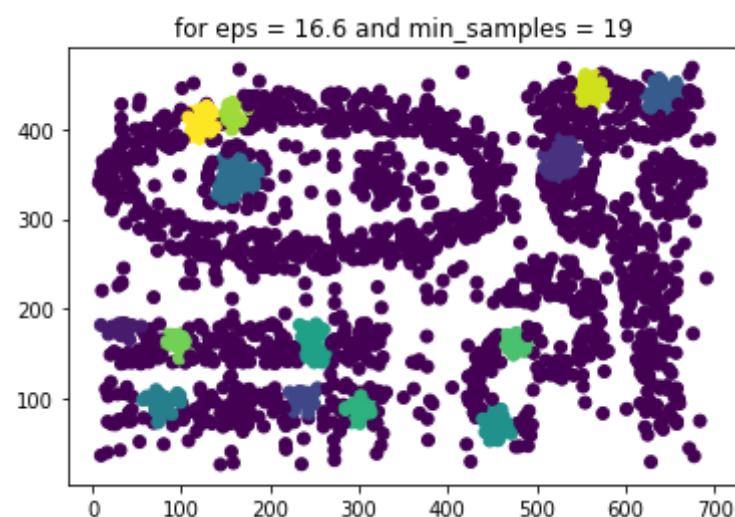
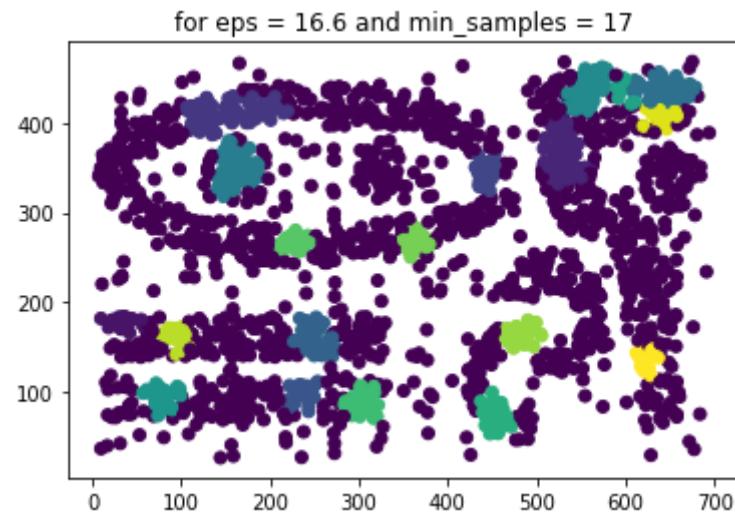


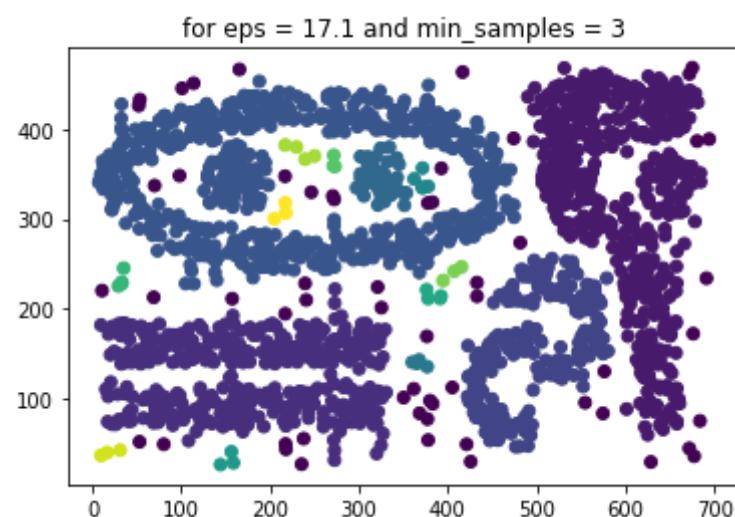
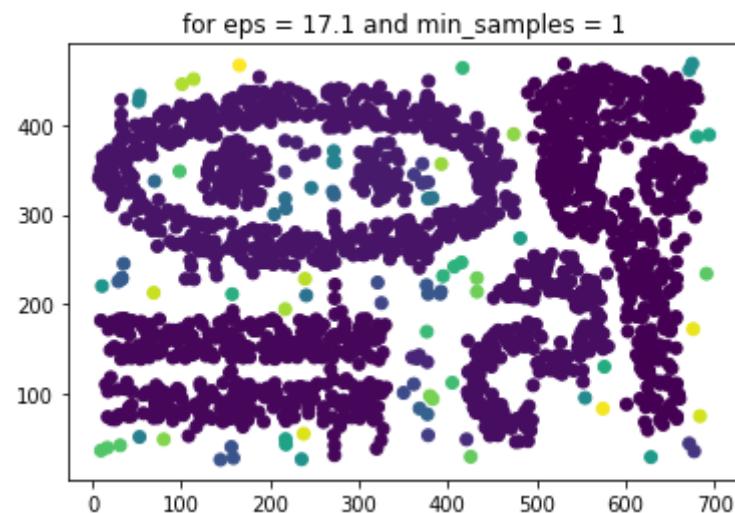


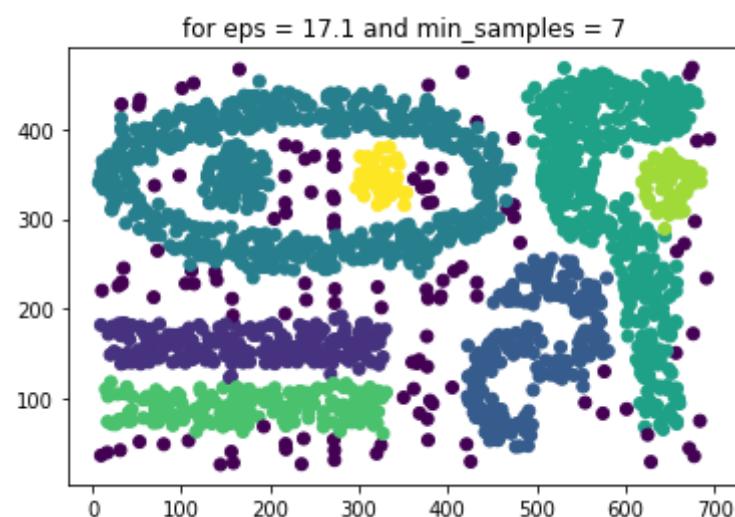
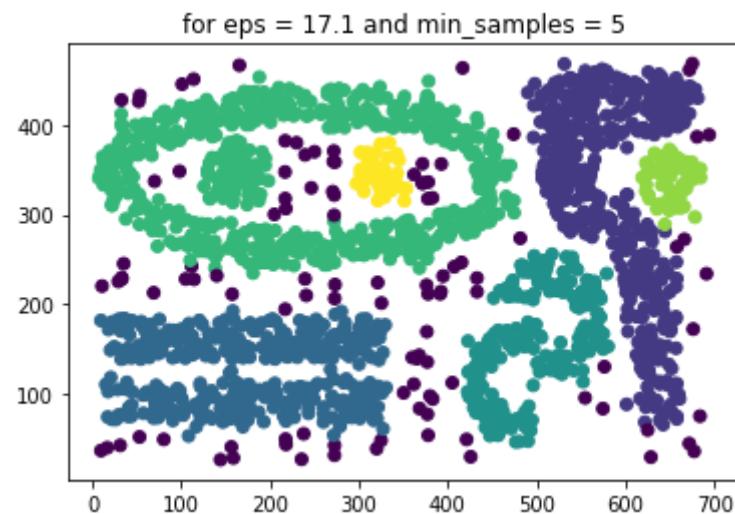


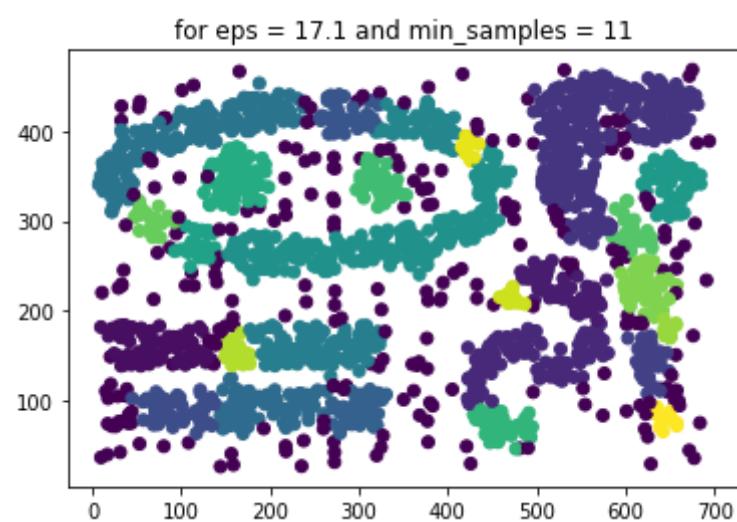
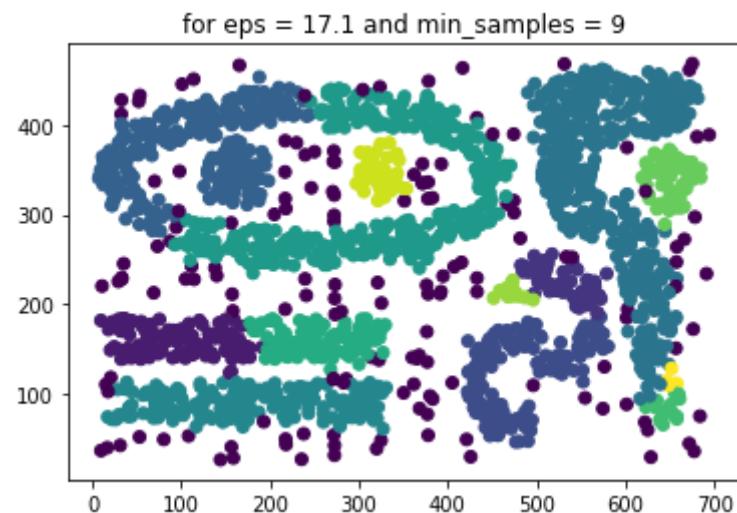


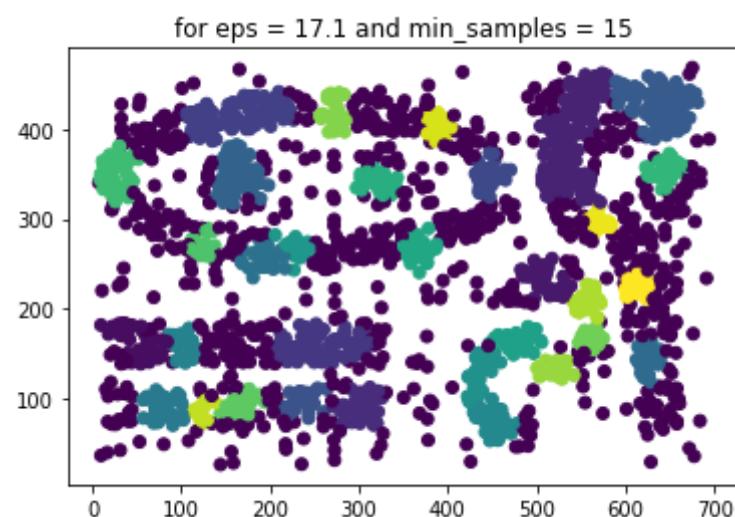
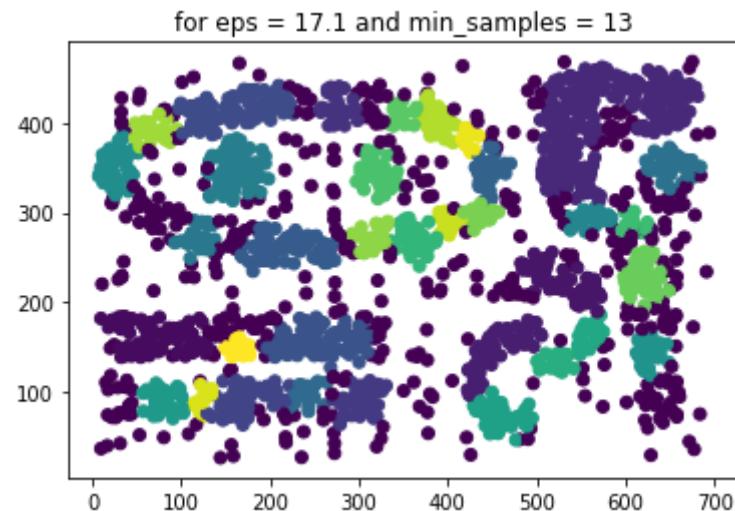


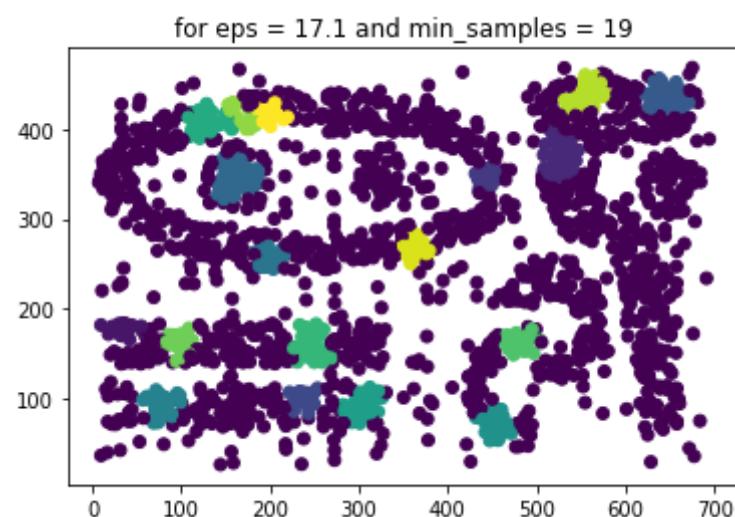
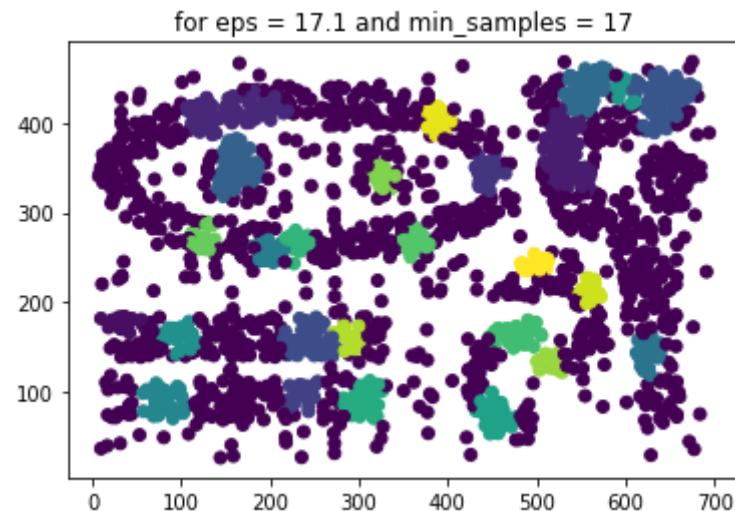


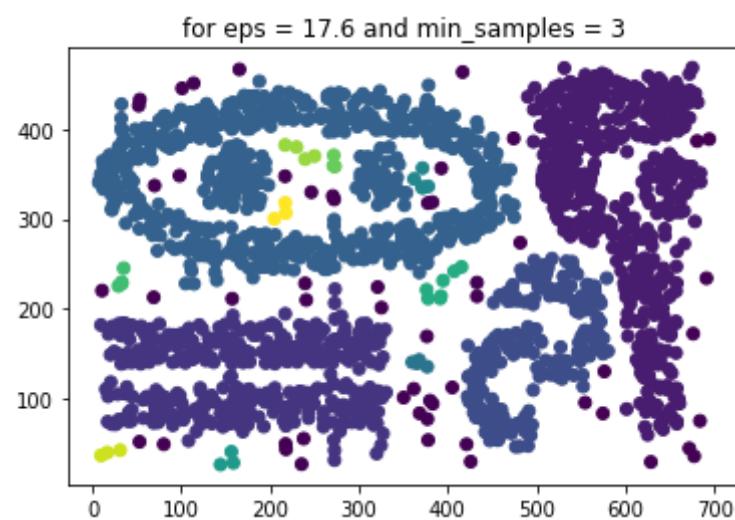
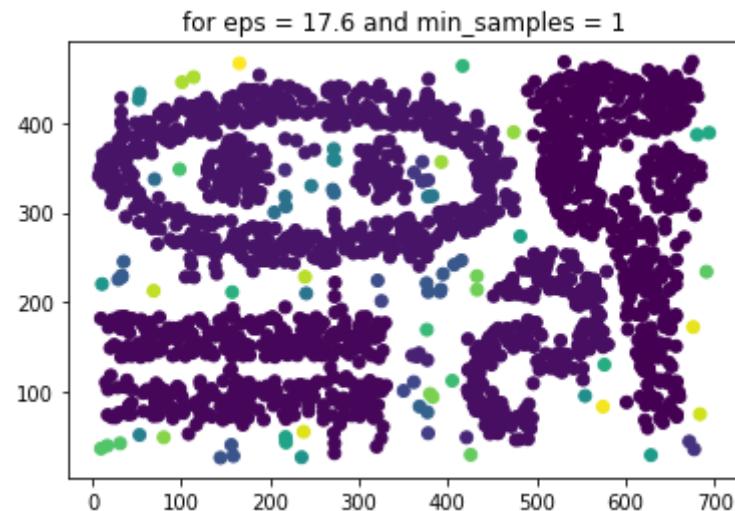


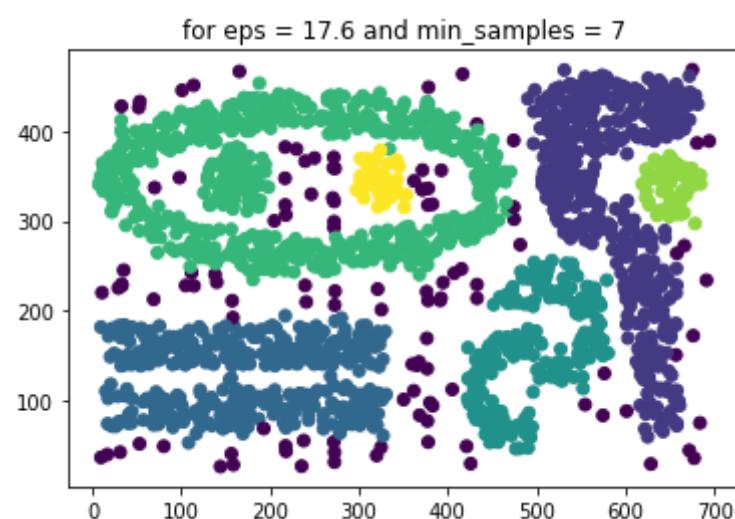
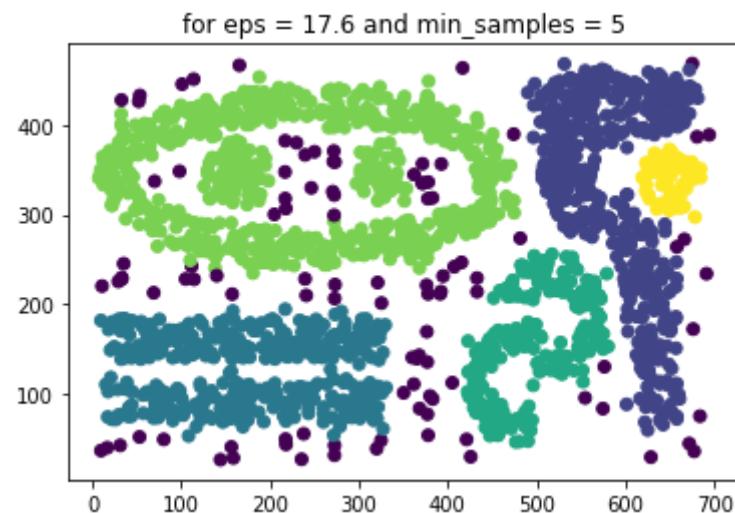


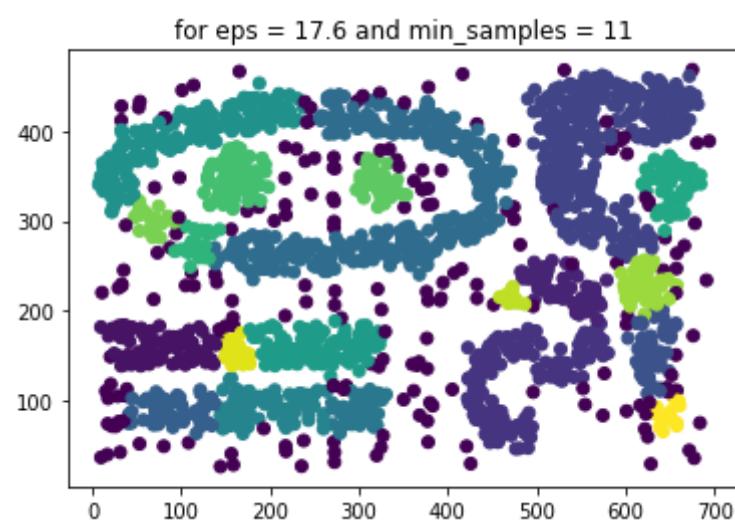
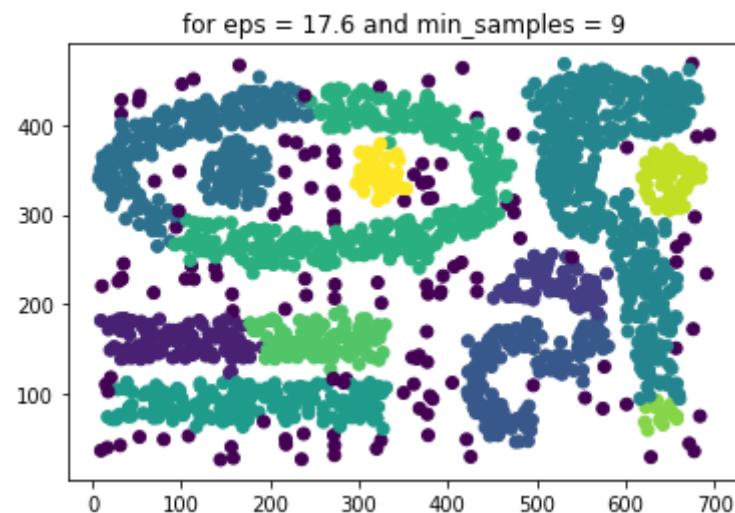


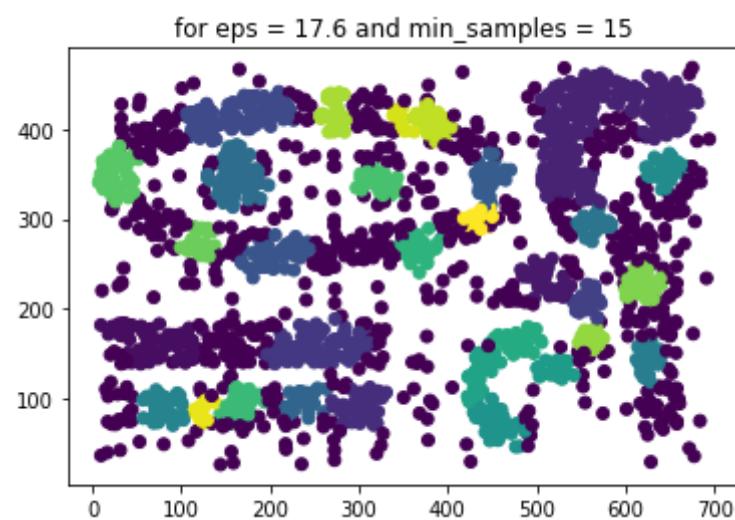
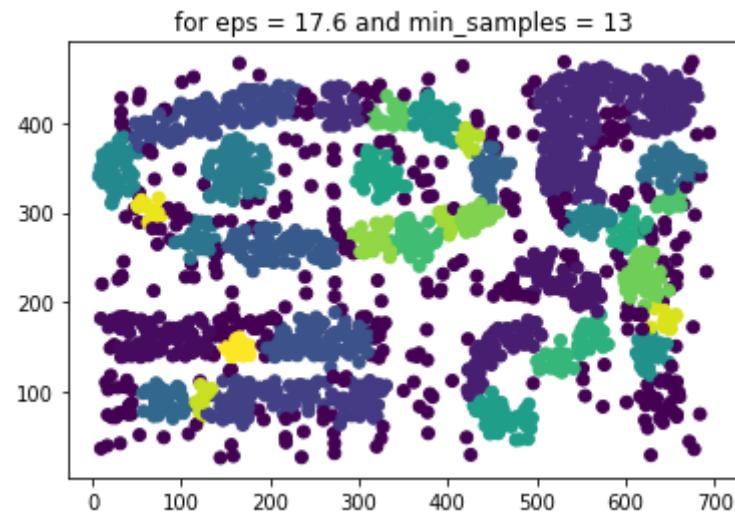


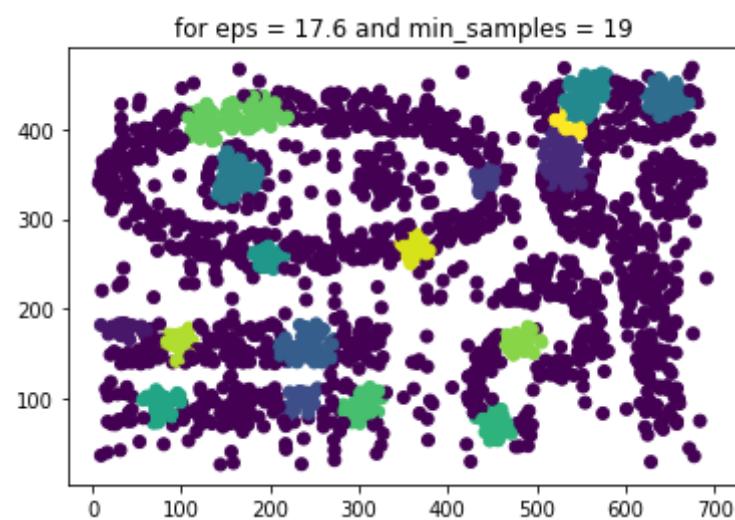
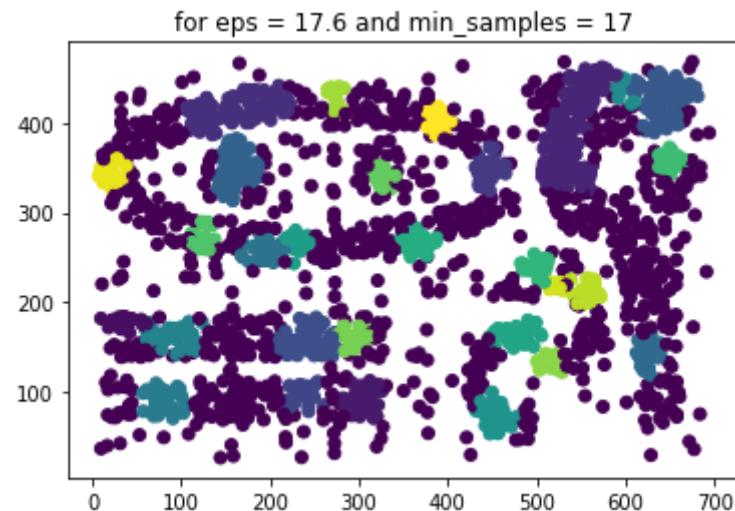


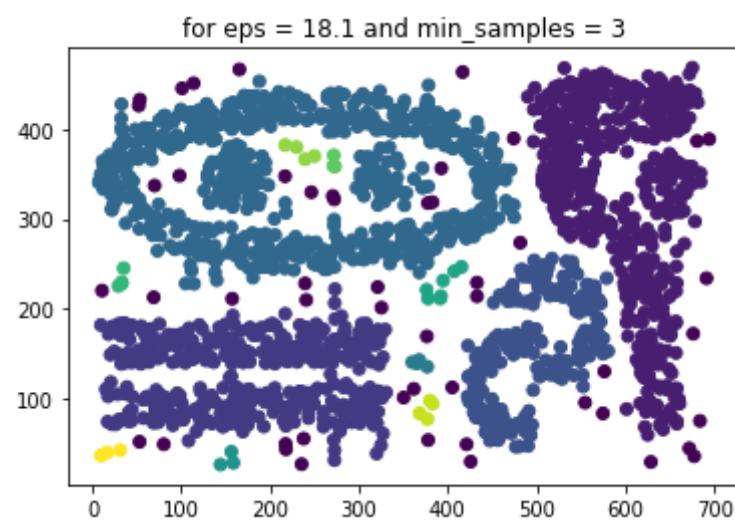
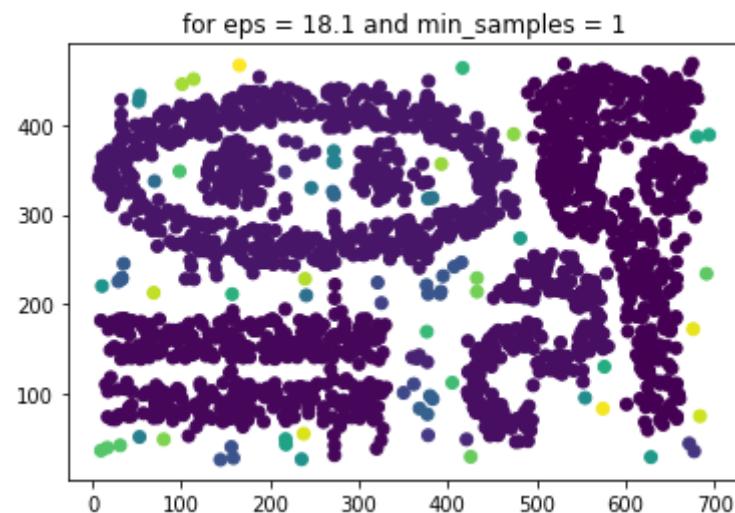


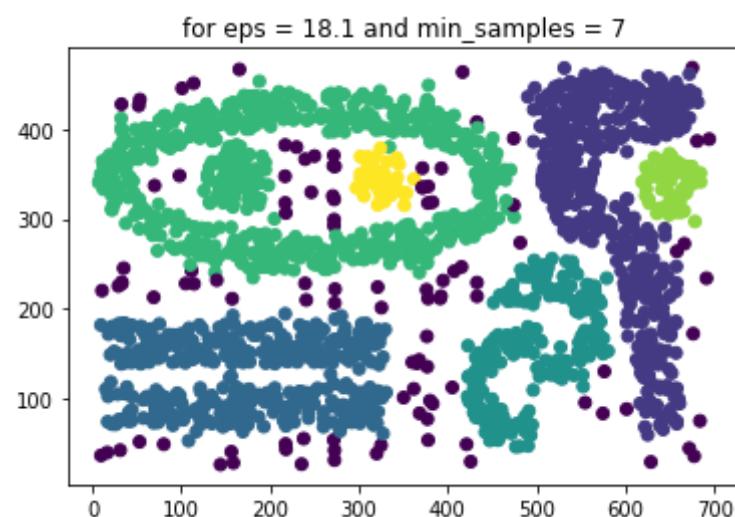
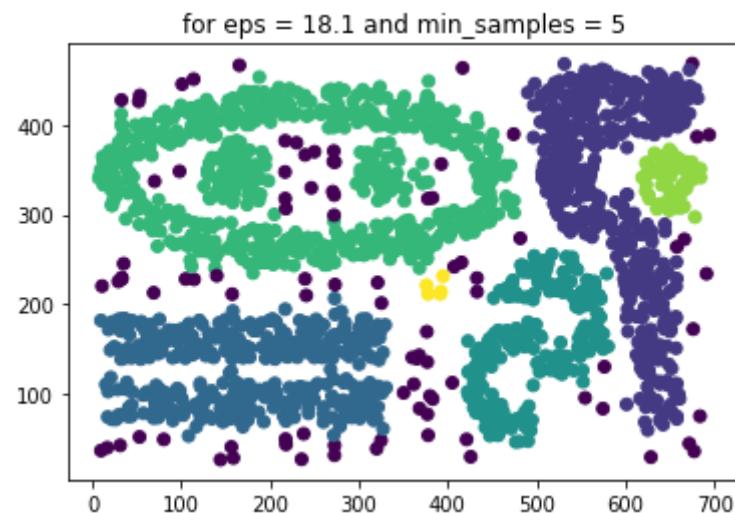


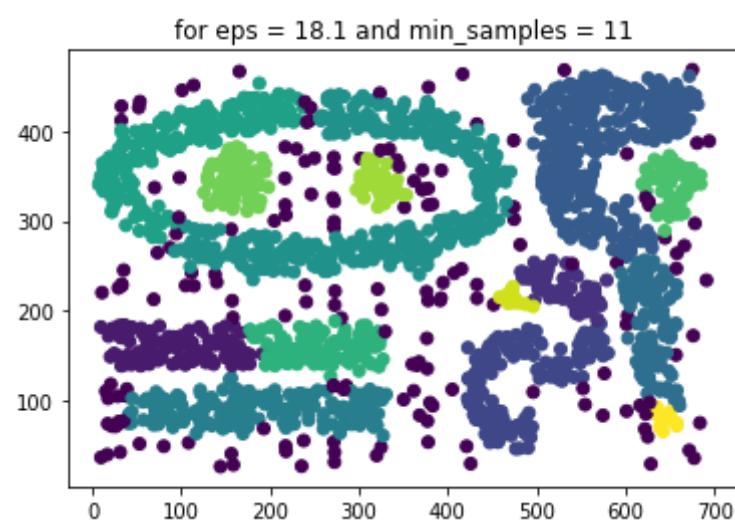
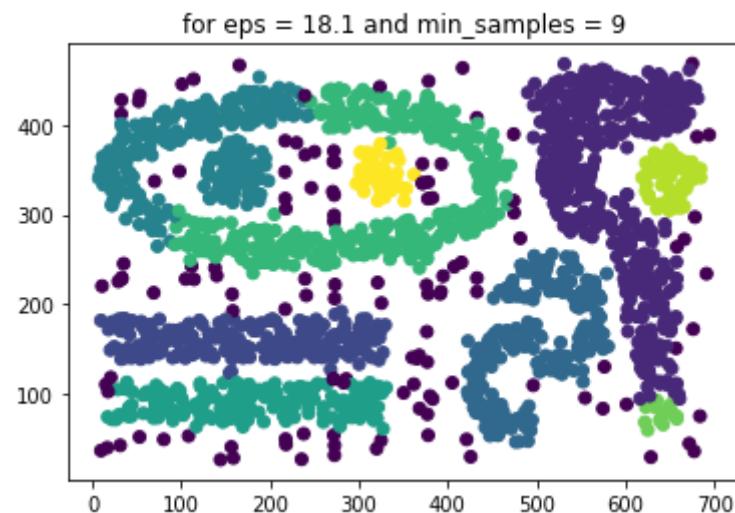


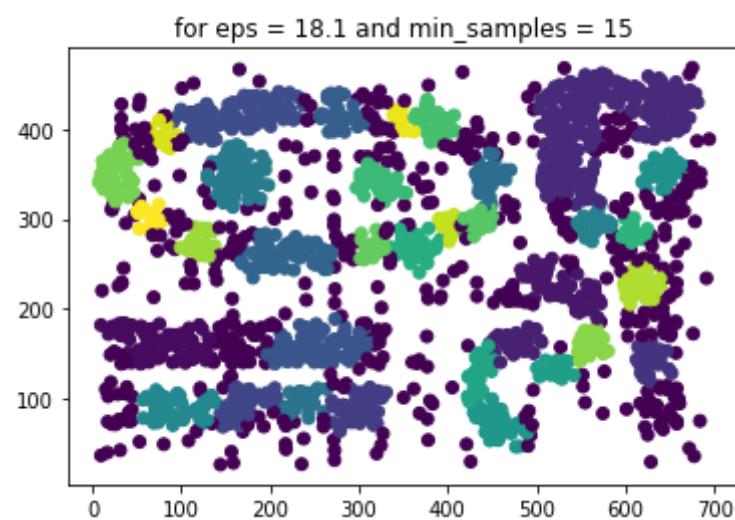
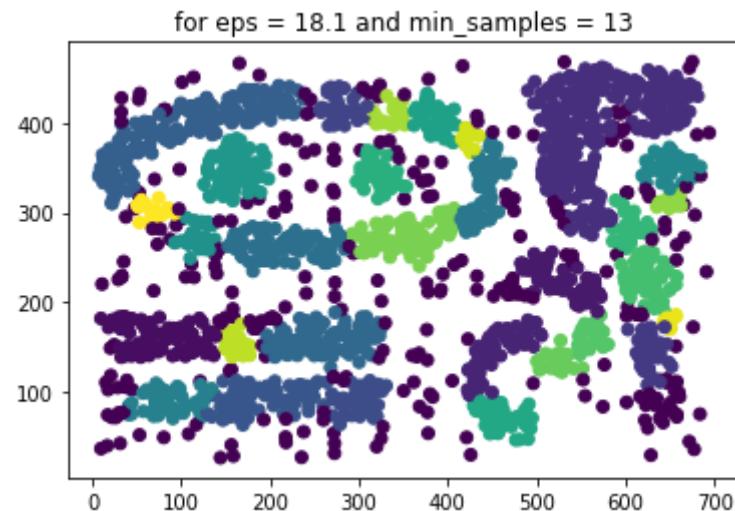


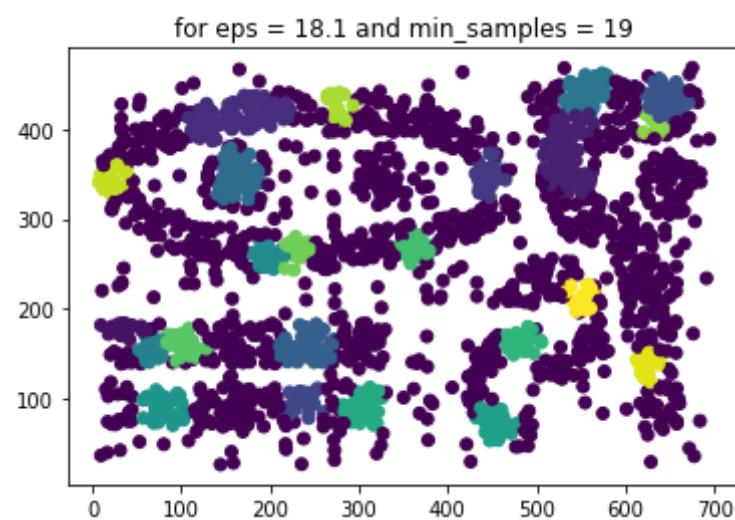
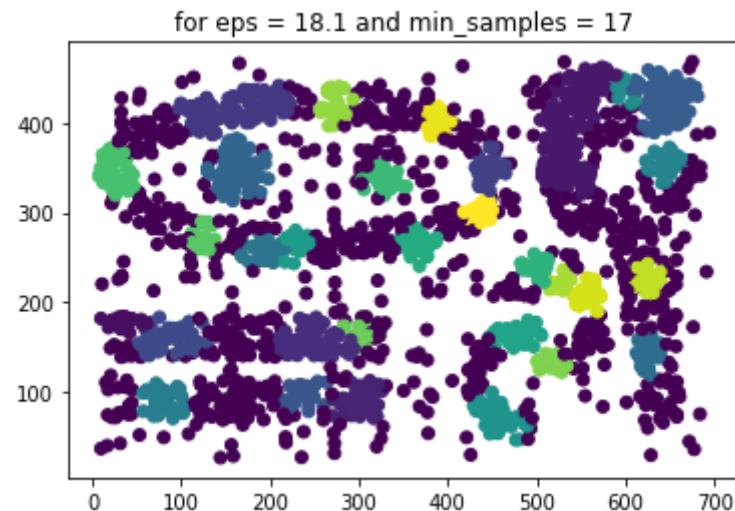


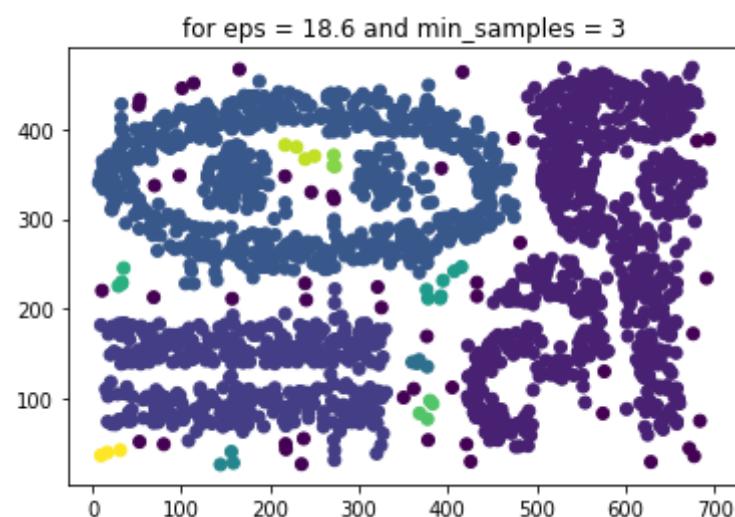
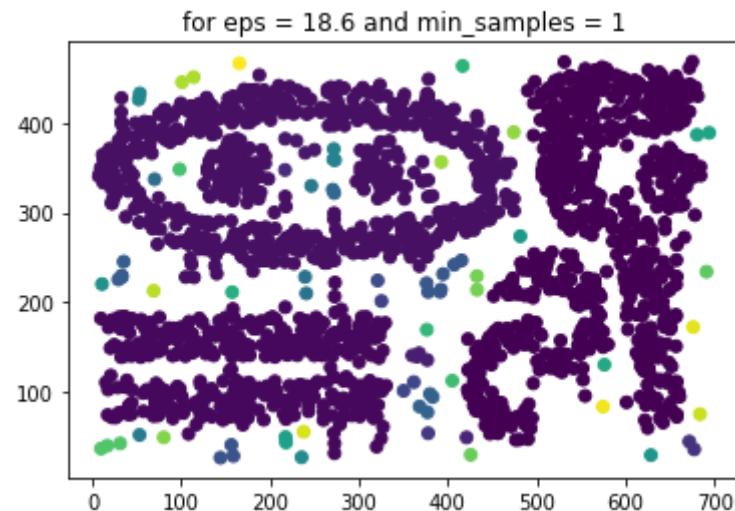


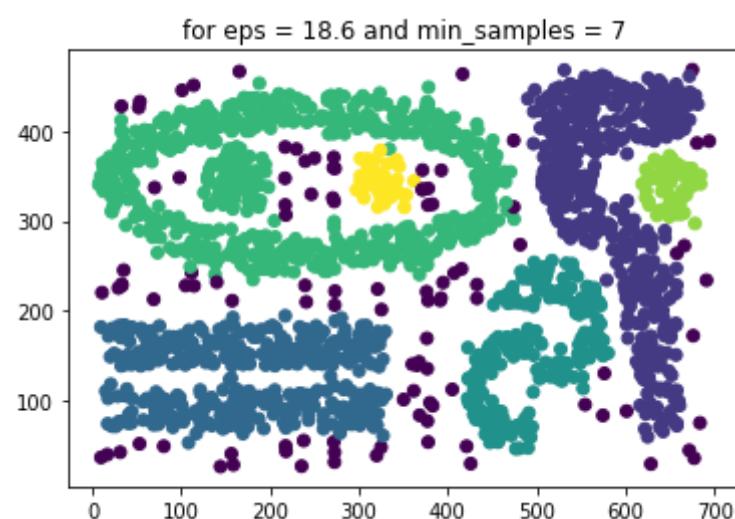
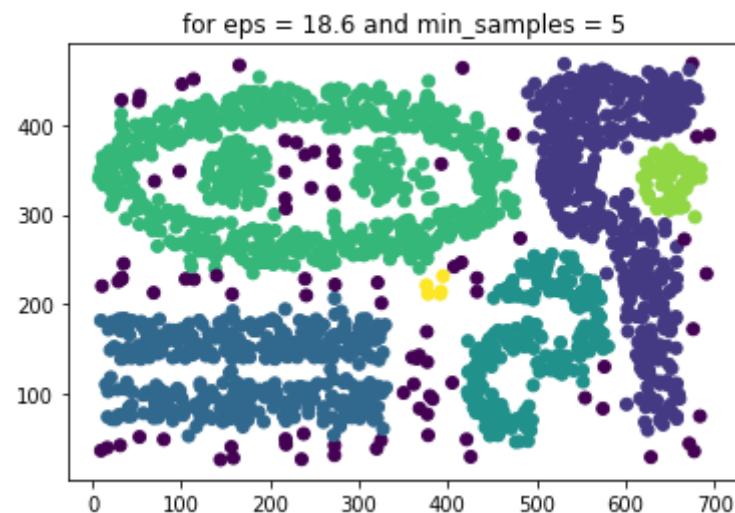


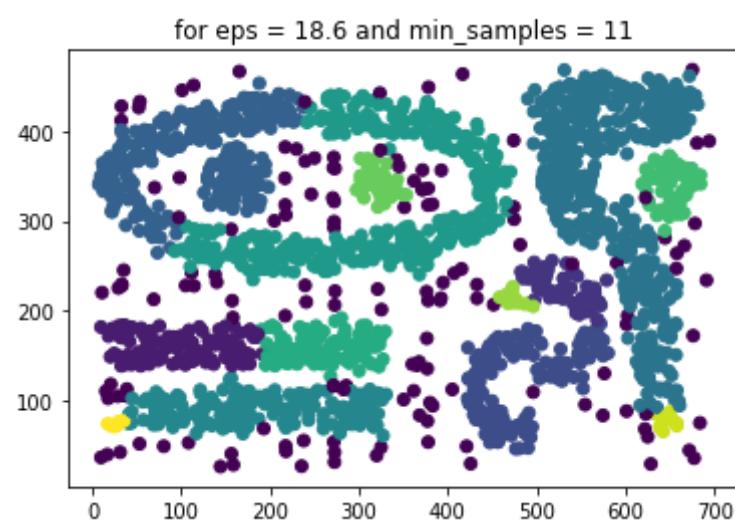
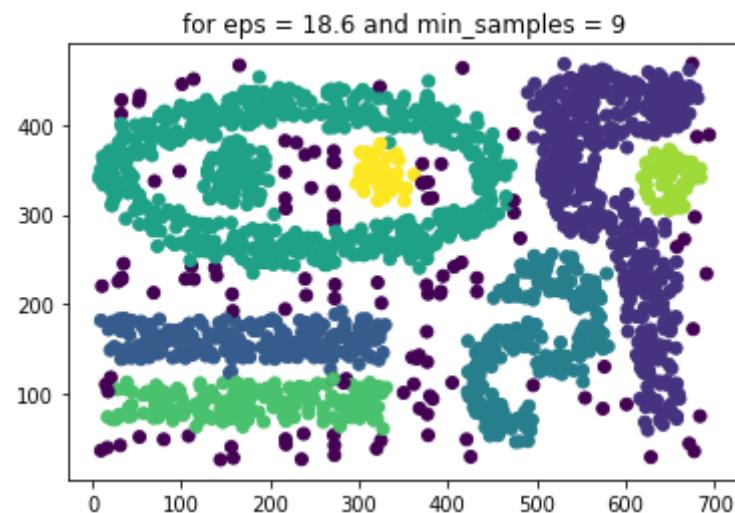


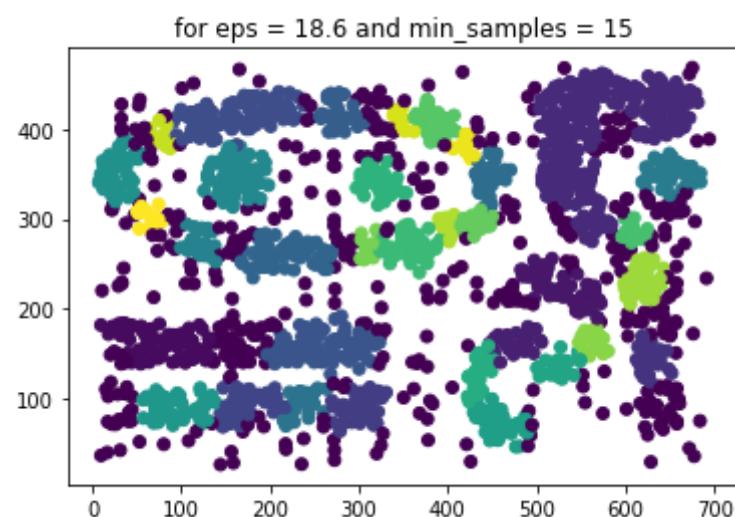
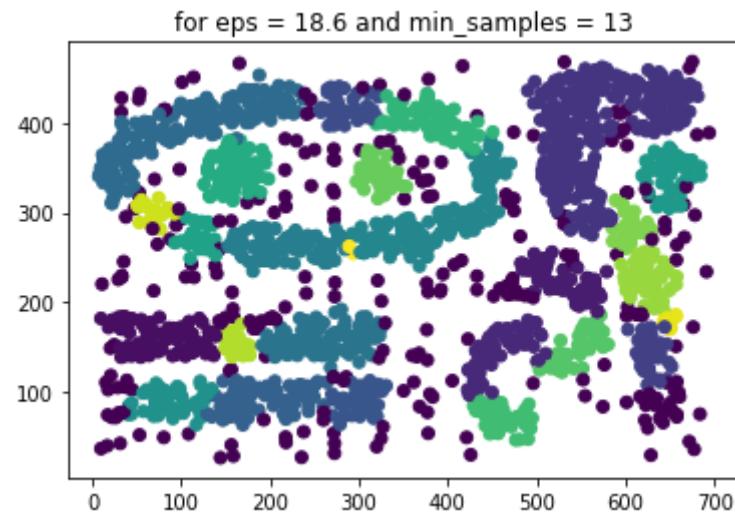


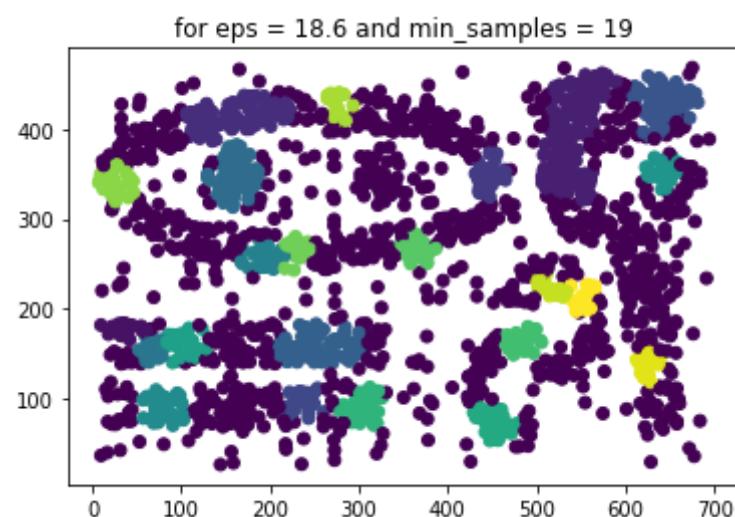
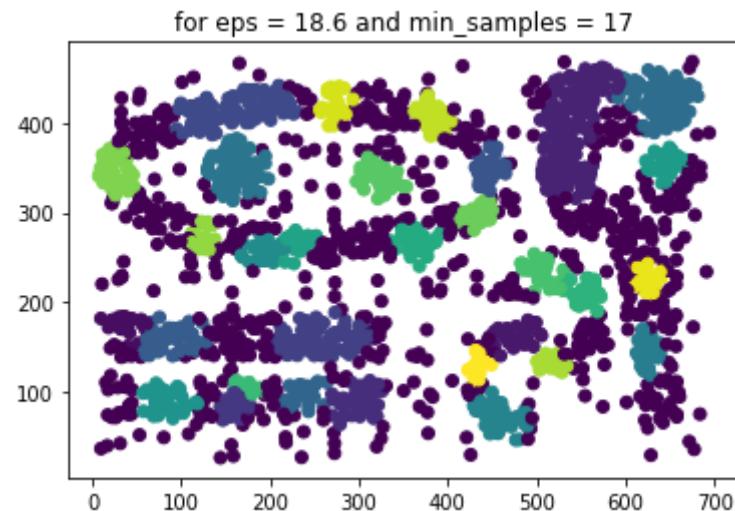


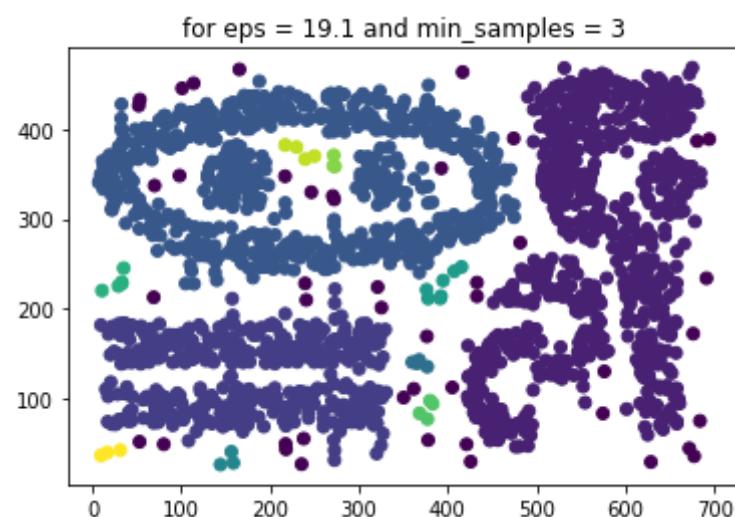
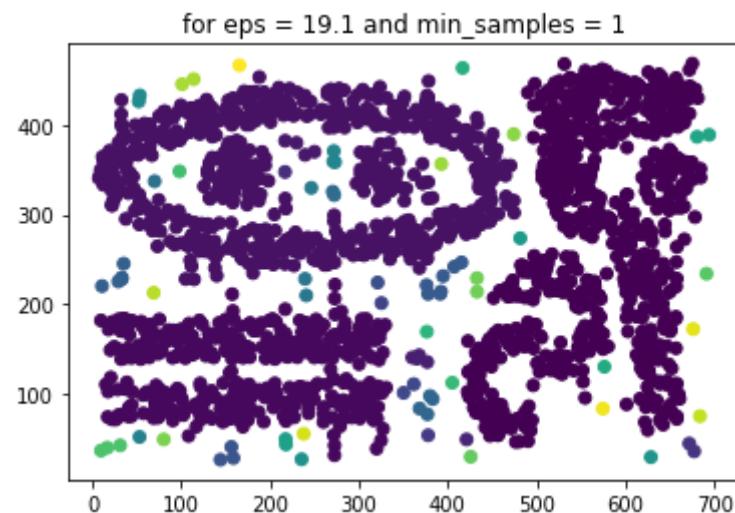


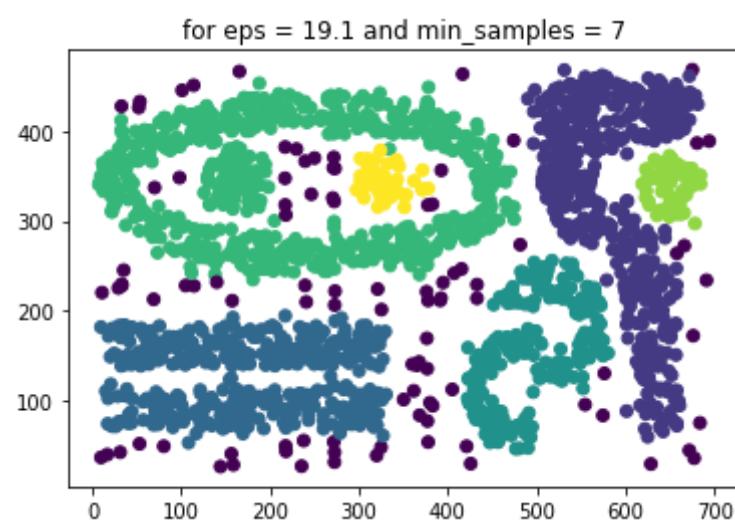
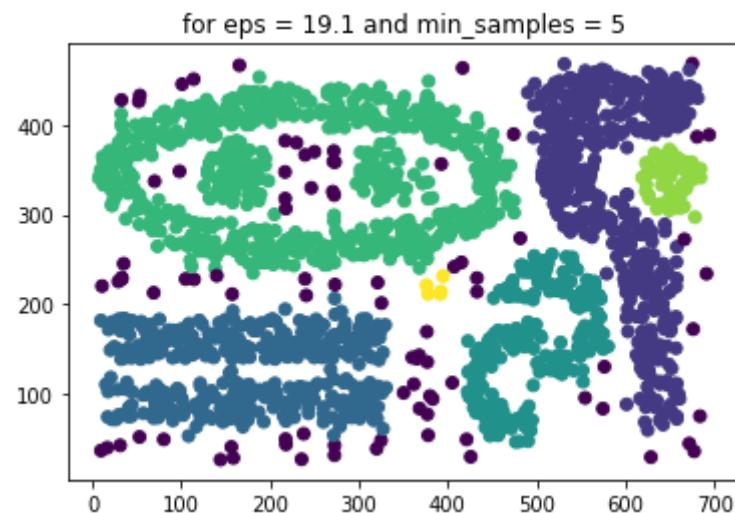


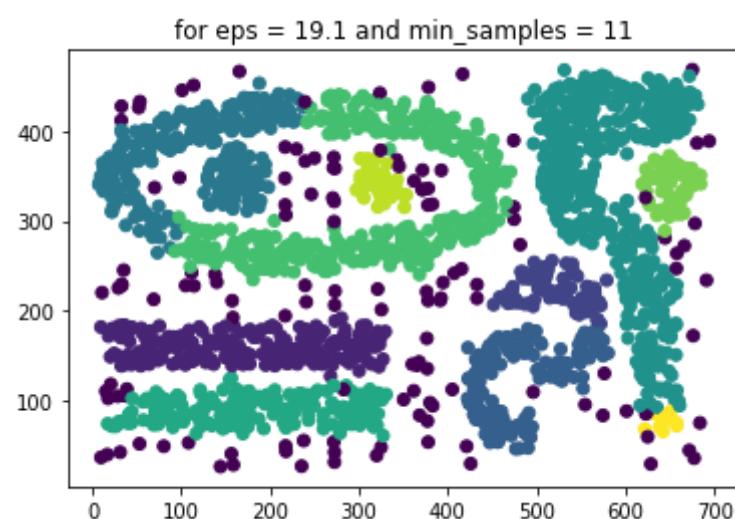
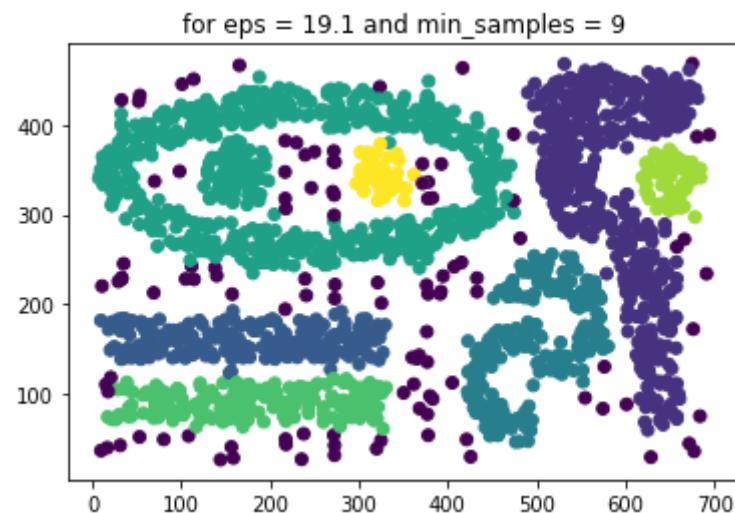


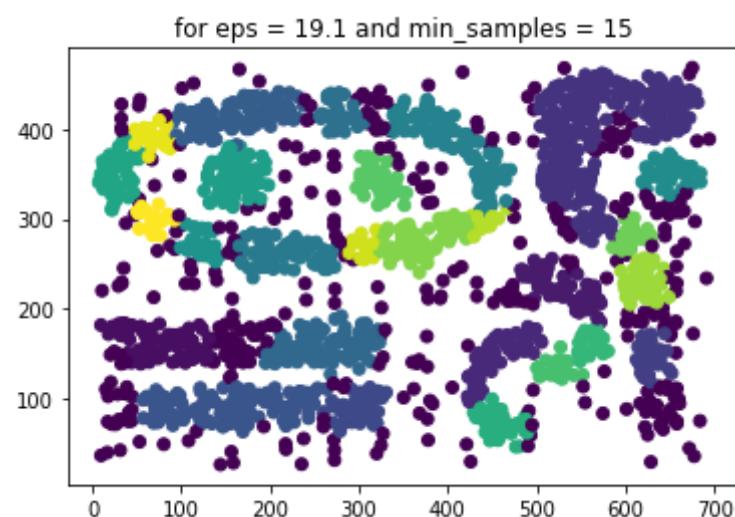
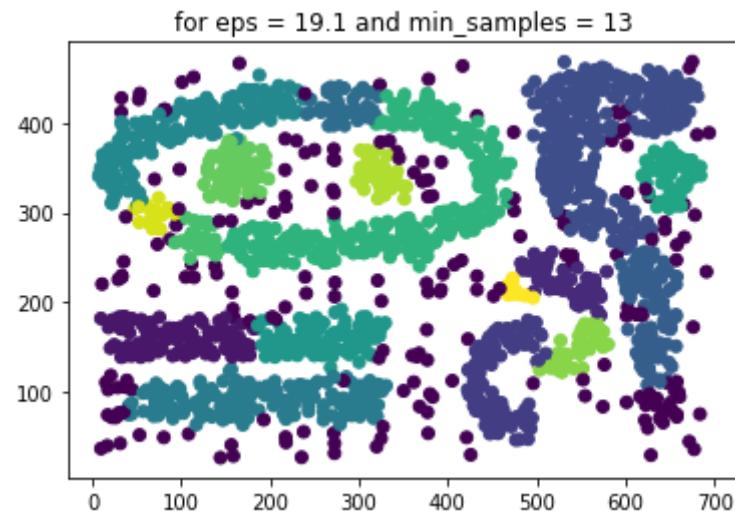


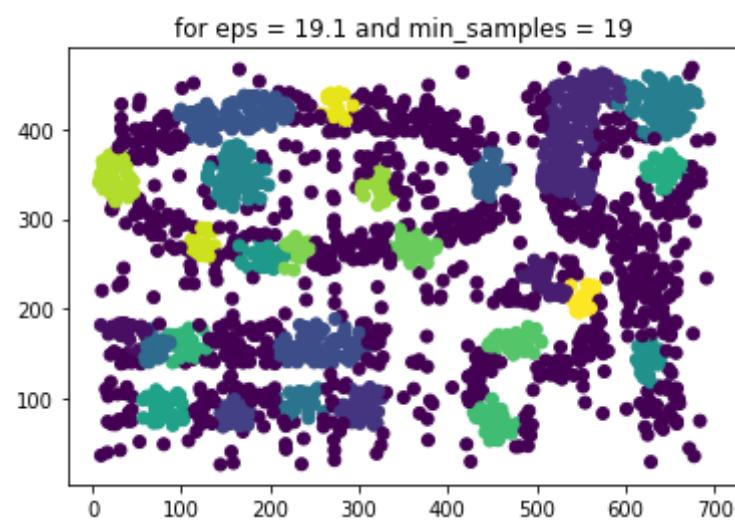
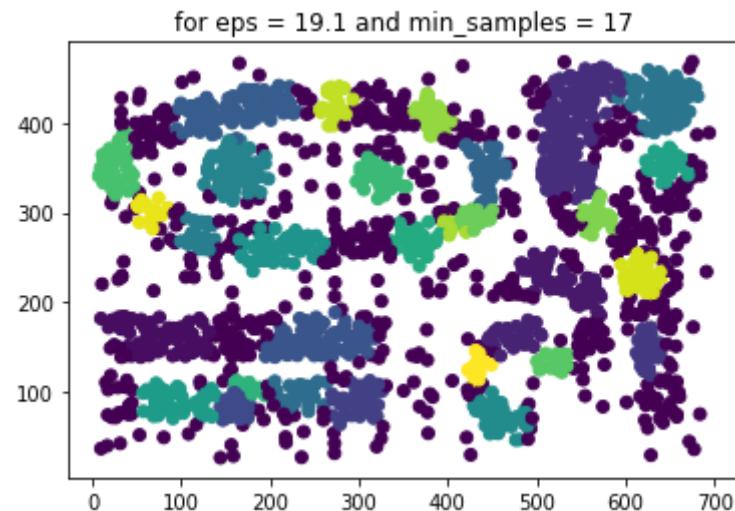


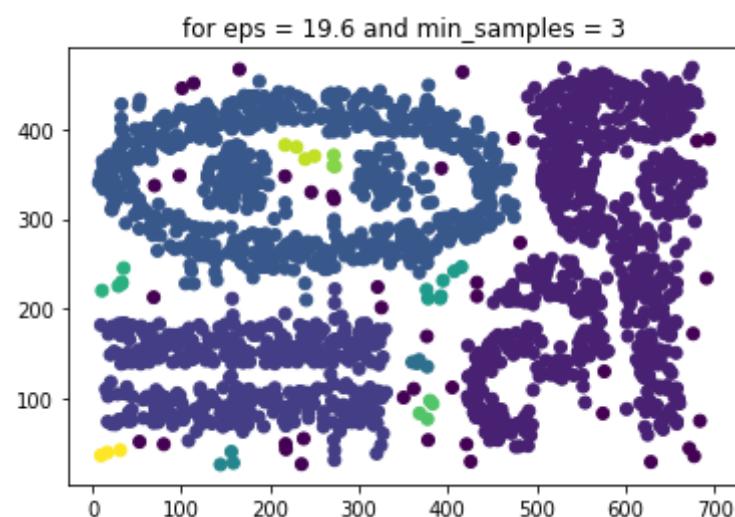
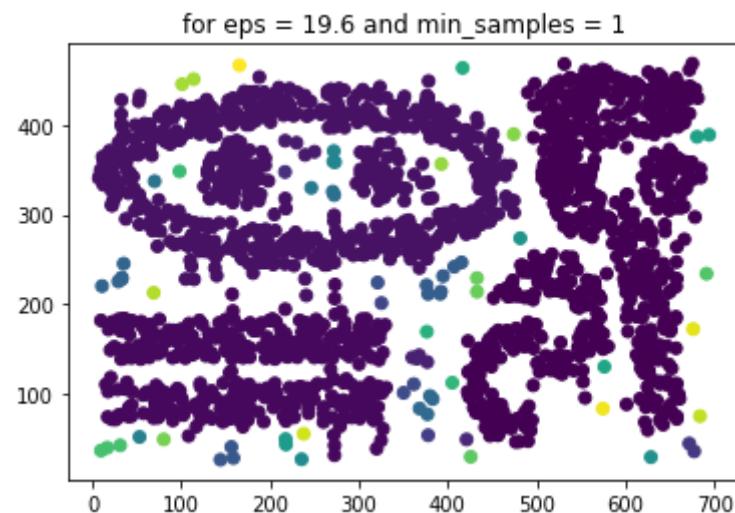


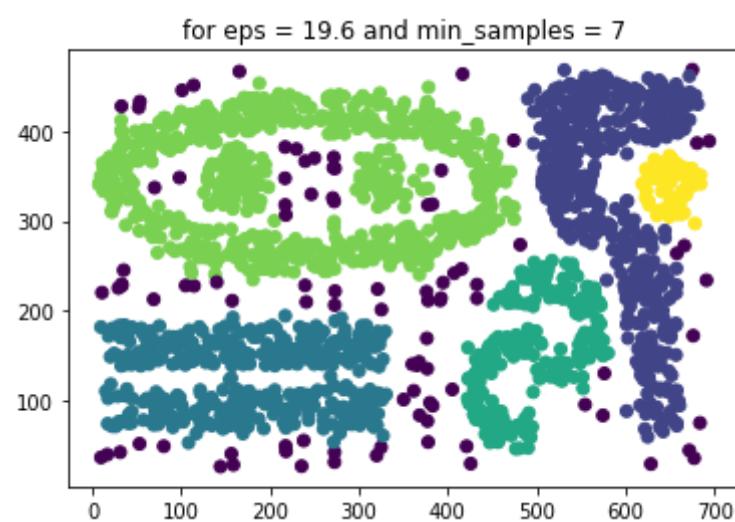
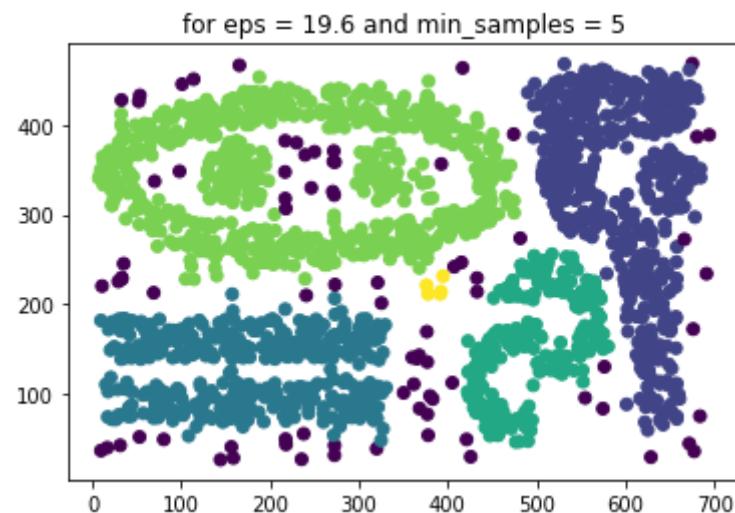


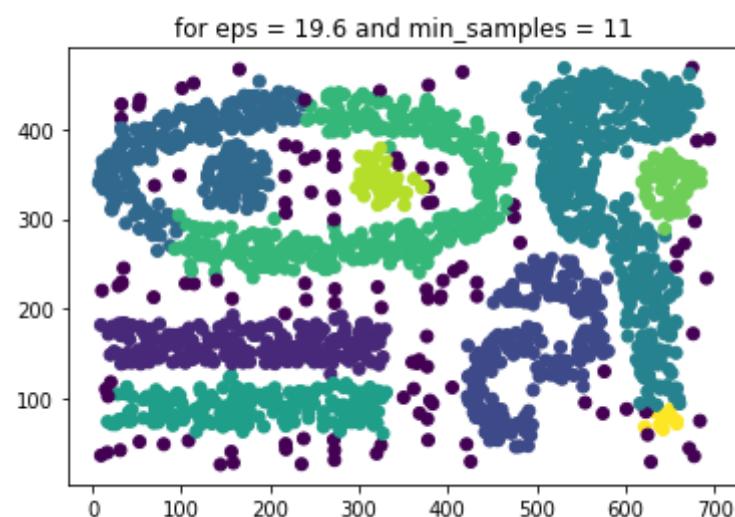
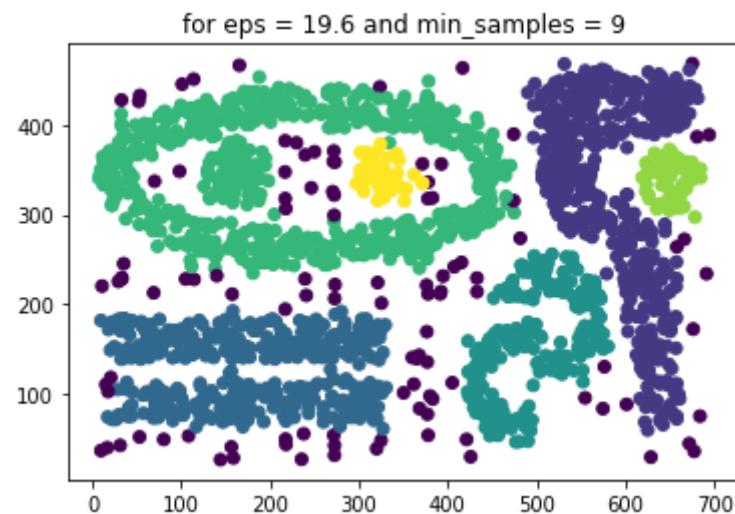


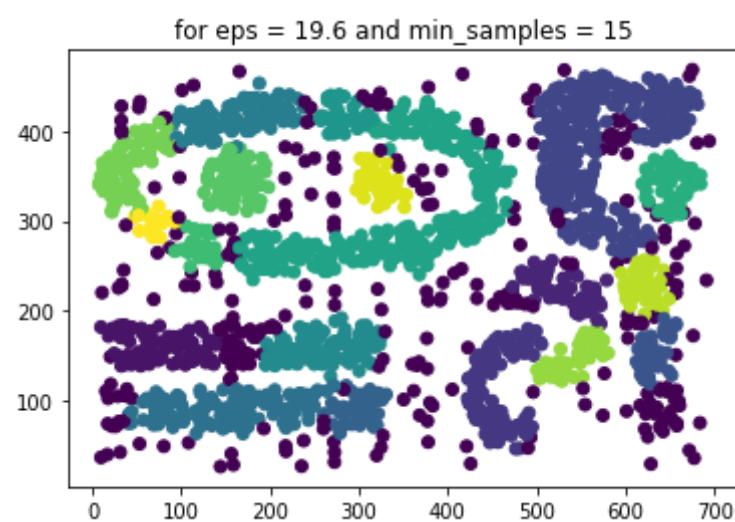
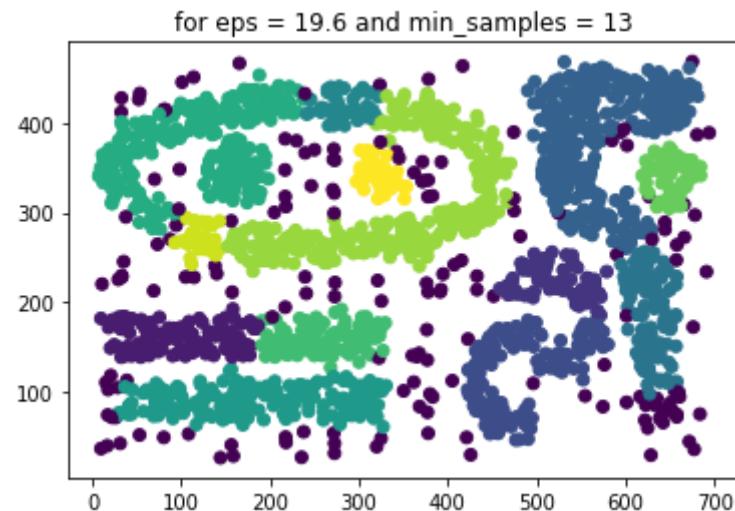


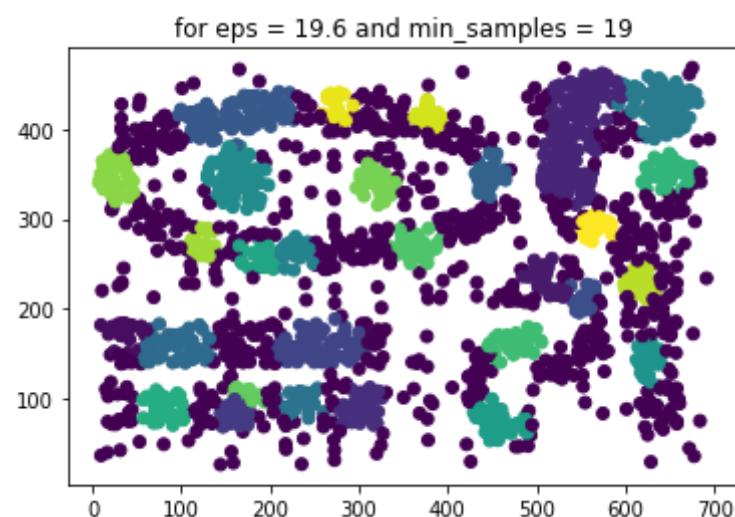
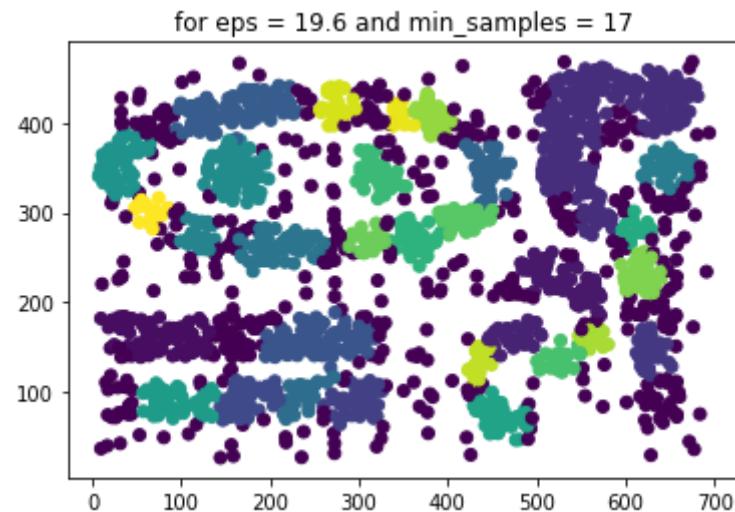












eps : float, optional The maximum distance between two samples for one to be considered as in the neighborhood of the other. This is not a maximum bound on the distances of points within a cluster. This is the most important DBSCAN parameter to choose appropriately for your data set and distance function.

min_samples : int, optional The number of samples (or total weight) in a neighborhood for a point to be considered as a core point. This includes the point itself.

Small eps results in small clusters - we can't have any good generalization then. Where min samples is much bigger than eps, dbscan is not effective in differentiating between clusters

As we could observe, if min samples is significantly bigger than eps, algorithm can't find appropriate clusters - data are too sparse for this setting of metaparameters.

As we can see, the best results (most natural clusters) were achieved for eps's from the range of 15.1-16.1 and for min-samples from the range of 5-7.