

Lab4_1

October 13, 2021

```
[ ]: from sklearn.datasets import make_blobs
from sklearn.cluster import KMeans
from sklearn.metrics.cluster import contingency_matrix
from sklearn.model_selection import cross_val_predict
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

1 Assignment 1

```
[ ]: cluster_a = make_blobs(300, n_features=2, centers=4, cluster_std=0.6,
    ↳center_box = (-10.0, 10.0))
cluster_b = make_blobs(300, n_features=2, centers=4, cluster_std=0.1,
    ↳center_box = (-10.0, 10.0))
cluster_c = make_blobs(300, n_features=2, centers=4, cluster_std=2.5,
    ↳center_box = (-10.0, 10.0))
```

```
[ ]: def get_clustering(cluster, std, random_state=None):
    X, y = cluster
    fig, axs = plt.subplots(1, 11, figsize=(40,40))
    fig.suptitle(f'cluster for standard deviation = {std}, random state =
    ↳{random_state}')
    df = pd.DataFrame(index=[f'actual cluster: {i}' for i in range(4)],
    ↳columns=[[], []])

    sse = []
    for i in range(1, 11):
        kmean = KMeans(n_clusters=i, random_state=random_state)
        y_pred = kmean.fit_predict(X,y)
        centroids = kmean.cluster_centers_
        sse.append(kmean.inertia_)
        axs[i-1].scatter(X[:,0], X[:,1], c= kmean.labels_)
        axs[i-1].scatter(centroids[:,0], centroids[:,1], c=np.unique(kmean.
    ↳labels_), edgecolors='red')
        axs[i-1].set(aspect='equal', title=f'nr. of clusters: {i}')
```

```

cols = [
    [f'predictions for k = {i}'] * i,
    [f'k: {j}' for j in np.arange(i)]]
temp_df = pd.DataFrame(contingency_matrix(y, y_pred), columns=cols,
↳ index=[f'actual cluster: {i}' for i in range(4)])
df = df.merge(temp_df, left_index=True, right_index=True)
print(temp_df)

axs[10].plot(np.arange(1,11), sse)
asp = np.diff(axs[10].get_xlim())[0] / np.diff(axs[10].get_ylim())[0]
axs[10].set(aspect=asp, title=f'SSE for different k')
axs[10].set_xticks(np.arange(1,11))
fig.tight_layout()
fig.set_size_inches(17, 3)

```

```

[ ]: for cluster, std in zip([cluster_a, cluster_b, cluster_c], [0.6, 0.1, 2.5]):
    get_clustering(cluster, std)

```

```

        predictions for k = 1
                k: 0
actual cluster: 0                75
actual cluster: 1                75
actual cluster: 2                75
actual cluster: 3                75
        predictions for k = 2
                k: 0 k: 1
actual cluster: 0                75    0
actual cluster: 1                0    75
actual cluster: 2                0    75
actual cluster: 3                0    75
        predictions for k = 3
                k: 0 k: 1 k: 2
actual cluster: 0                0    75    0
actual cluster: 1                75    0    0
actual cluster: 2                0    0    75
actual cluster: 3                0    0    75
        predictions for k = 4
                k: 0 k: 1 k: 2 k: 3
actual cluster: 0                0    0    75    0
actual cluster: 1                0    0    0    75
actual cluster: 2                75    0    0    0
actual cluster: 3                0    75    0    0
        predictions for k = 5
                k: 0 k: 1 k: 2 k: 3 k: 4
actual cluster: 0                75    0    0    0    0
actual cluster: 1                0    0    0    75    0
actual cluster: 2                0    0    75    0    0

```

actual cluster: 3	0	48	0	0	27			
predictions for k = 6								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5		
actual cluster: 0	0	42	0	0	33	0		
actual cluster: 1	75	0	0	0	0	0		
actual cluster: 2	0	0	0	75	0	0		
actual cluster: 3	0	0	28	0	0	47		
predictions for k = 7								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	
actual cluster: 0	0	35	0	0	0	40	0	
actual cluster: 1	24	0	0	51	0	0	0	
actual cluster: 2	0	0	75	0	0	0	0	
actual cluster: 3	0	0	0	0	31	0	44	
predictions for k = 8								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	39	0	0	36	0	0	0
actual cluster: 1	45	0	0	0	0	0	30	0
actual cluster: 2	0	0	38	0	0	0	0	37
actual cluster: 3	0	0	0	37	0	38	0	0
predictions for k = 9 \								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	44	0	0	0	0	31	0
actual cluster: 1	0	0	32	0	0	43	0	0
actual cluster: 2	0	0	0	42	0	0	0	33
actual cluster: 3	31	0	0	0	21	0	0	0
k: 8								
actual cluster: 0	0							
actual cluster: 1	0							
actual cluster: 2	0							
actual cluster: 3	23							
predictions for k = 10 \								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	25	0	0	0	29	0	21
actual cluster: 1	0	0	36	0	0	0	19	0
actual cluster: 2	38	0	0	0	0	0	0	0
actual cluster: 3	0	0	0	31	44	0	0	0
k: 8 k: 9								
actual cluster: 0	0	0						
actual cluster: 1	20	0						
actual cluster: 2	0	37						
actual cluster: 3	0	0						
predictions for k = 1								
	k: 0							
actual cluster: 0	75							

```

actual cluster: 1          75
actual cluster: 2          75
actual cluster: 3          75
      predictions for k = 2
                k: 0 k: 1
actual cluster: 0          0  75
actual cluster: 1          75  0
actual cluster: 2          0  75
actual cluster: 3          75  0
      predictions for k = 3
                k: 0 k: 1 k: 2
actual cluster: 0          0  0  75
actual cluster: 1          75  0  0
actual cluster: 2          0  75  0
actual cluster: 3          75  0  0
      predictions for k = 4
                k: 0 k: 1 k: 2 k: 3
actual cluster: 0          0  0  75  0
actual cluster: 1          0  0  0  75
actual cluster: 2          0  75  0  0
actual cluster: 3          75  0  0  0
      predictions for k = 5
                k: 0 k: 1 k: 2 k: 3 k: 4
actual cluster: 0          75  0  0  0  0
actual cluster: 1          0  32  0  0  43
actual cluster: 2          0  0  75  0  0
actual cluster: 3          0  0  0  75  0
      predictions for k = 6
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5
actual cluster: 0          0  75  0  0  0  0
actual cluster: 1          53  0  0  0  22  0
actual cluster: 2          0  0  30  0  0  45
actual cluster: 3          0  0  0  75  0  0
      predictions for k = 7
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6
actual cluster: 0          28  0  0  0  47  0  0
actual cluster: 1          0  49  0  0  0  26  0
actual cluster: 2          0  0  45  0  0  0  30
actual cluster: 3          0  0  0  75  0  0  0
      predictions for k = 8
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0          36  0  0  0  0  0  0  39
actual cluster: 1          0  0  0  36  39  0  0  0
actual cluster: 2          0  0  49  0  0  0  26  0
actual cluster: 3          0  42  0  0  0  33  0  0
      predictions for k = 9
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0          43  0  0  0  0  0  32  0

```

\

actual cluster: 1	0	0	0	50	0	0	0	25
actual cluster: 2	0	0	50	0	0	0	0	0
actual cluster: 3	0	27	0	0	20	28	0	0

k: 8

actual cluster: 0	0
actual cluster: 1	0
actual cluster: 2	25
actual cluster: 3	0

predictions for k = 10 \

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	38	0	0	0	0	37	0	0
actual cluster: 1	0	0	0	36	0	0	0	19
actual cluster: 2	0	0	34	0	20	0	0	0
actual cluster: 3	0	50	0	0	0	0	25	0

k: 8 k: 9

actual cluster: 0	0	0
actual cluster: 1	20	0
actual cluster: 2	0	21
actual cluster: 3	0	0

predictions for k = 1

k: 0

actual cluster: 0	75
actual cluster: 1	75
actual cluster: 2	75
actual cluster: 3	75

predictions for k = 2

k: 0 k: 1

actual cluster: 0	75	0
actual cluster: 1	13	62
actual cluster: 2	0	75
actual cluster: 3	52	23

predictions for k = 3

k: 0 k: 1 k: 2

actual cluster: 0	3	0	72
actual cluster: 1	57	15	3
actual cluster: 2	4	71	0
actual cluster: 3	65	0	10

predictions for k = 4

k: 0 k: 1 k: 2 k: 3

actual cluster: 0	0	13	0	62
actual cluster: 1	61	8	6	0
actual cluster: 2	7	0	68	0
actual cluster: 3	29	44	0	2

predictions for k = 5

	k: 0	k: 1	k: 2	k: 3	k: 4
actual cluster: 0	61	0	0	14	0
actual cluster: 1	0	36	1	8	30
actual cluster: 2	0	19	56	0	0
actual cluster: 3	3	9	0	39	24

predictions for k = 6

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5
actual cluster: 0	0	61	0	14	0	0
actual cluster: 1	14	0	0	6	22	33
actual cluster: 2	18	0	43	0	14	0
actual cluster: 3	0	3	0	39	8	25

predictions for k = 7

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6
actual cluster: 0	15	0	0	0	33	27	0
actual cluster: 1	6	13	0	32	0	0	24
actual cluster: 2	0	20	41	0	0	0	14
actual cluster: 3	38	0	0	24	4	0	9

predictions for k = 8

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	17	0	1	0	0	27	30
actual cluster: 1	3	4	32	14	22	0	0	0
actual cluster: 2	34	0	2	0	14	25	0	0
actual cluster: 3	0	15	10	38	8	0	4	0

predictions for k = 9 \

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	19	0	2	30	0	0	24
actual cluster: 1	2	3	20	7	0	18	25	0
actual cluster: 2	34	0	14	0	0	5	0	0
actual cluster: 3	0	8	8	36	3	4	16	0

k: 8

actual cluster: 0	0
actual cluster: 1	0
actual cluster: 2	22
actual cluster: 3	0

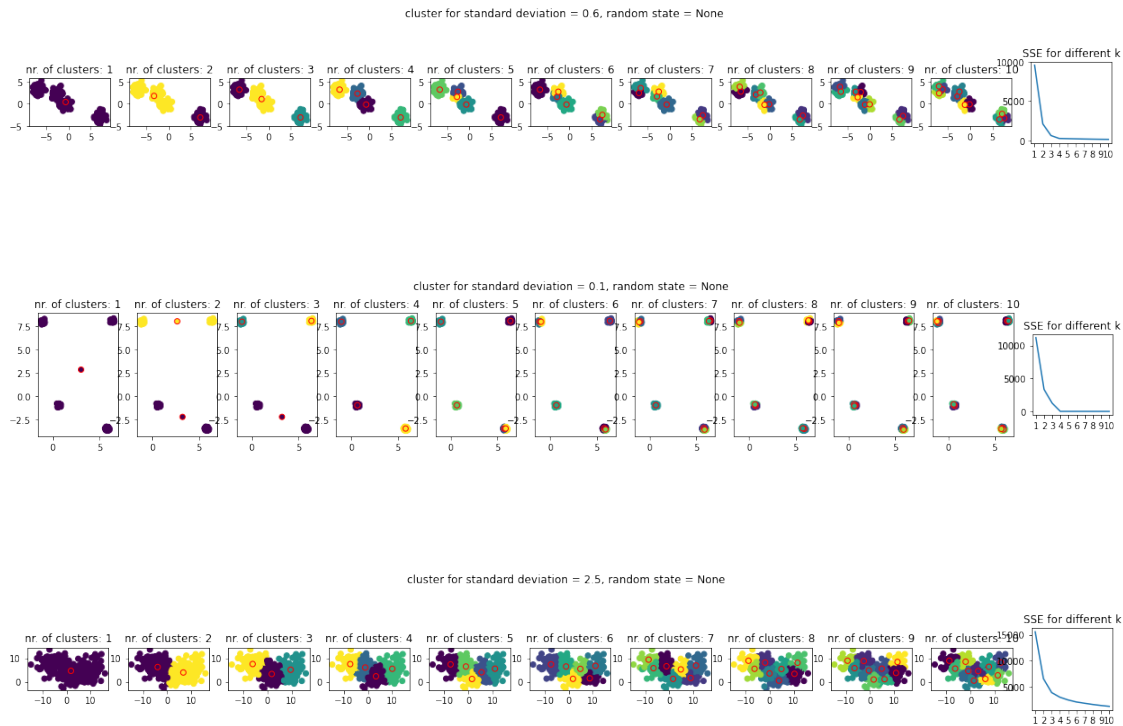
predictions for k = 10 \

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	30	0	0	1	0	0	27	0
actual cluster: 1	0	22	0	2	23	6	0	1
actual cluster: 2	0	3	18	0	0	28	0	21
actual cluster: 3	0	13	0	33	12	1	3	0

k: 8 k: 9

actual cluster: 0	17	0
actual cluster: 1	4	17
actual cluster: 2	0	5

actual cluster: 3 9 4



We can see that the clusters are differently spread out, depending on the standard deviation. It can be stated that the bigger the deviation the bigger the blobs. At standard deviation 2.5 it's actually hard to see the different blobs. Consequently we can see how the SSE behaves differently for the different standard deviations. For deviation = 0.1 the sse curve indicates that there are actually 4 different clusters. An observation that is easily backed up when looking at the plots. Deviation = 0.6 behaves not as clearly. I've run the experiment a couple of times and there are blob distributions that are already linked to a very low SSE for $k=3$. Contrary to this Deviation = 2.5 behaves differently in the sense that 10 clusters perform significantly better than 9 clusters. This indicates that for blobs spread out over a large area we'd need higher number of clusters to minimize the SSE. An insight that is very trivial.

```
[ ]: for cluster, std in zip([cluster_a, cluster_b, cluster_c], [0.6, 0.1, 2.5]):
      for r in [1,5,10]:
          get_clustering(cluster, std, r)
```

predictions for k = 1

k: 0

actual cluster: 0 75

actual cluster: 1 75

actual cluster: 2 75

actual cluster: 3 75

predictions for k = 2

k: 0 k: 1

actual cluster: 0	0	75		
actual cluster: 1	75	0		
actual cluster: 2	75	0		
actual cluster: 3	75	0		

predictions for k = 3

	k: 0	k: 1	k: 2	
actual cluster: 0	0	75	0	
actual cluster: 1	75	0	0	
actual cluster: 2	0	0	75	
actual cluster: 3	0	0	75	

predictions for k = 4

	k: 0	k: 1	k: 2	k: 3	
actual cluster: 0	0	75	0	0	
actual cluster: 1	75	0	0	0	
actual cluster: 2	0	0	0	75	
actual cluster: 3	0	0	75	0	

predictions for k = 5

	k: 0	k: 1	k: 2	k: 3	k: 4	
actual cluster: 0	0	75	0	0	0	
actual cluster: 1	0	0	75	0	0	
actual cluster: 2	75	0	0	0	0	
actual cluster: 3	0	0	0	49	26	

predictions for k = 6

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	
actual cluster: 0	0	31	0	0	44	0	
actual cluster: 1	0	0	75	0	0	0	
actual cluster: 2	0	0	0	75	0	0	
actual cluster: 3	48	0	0	0	0	27	

predictions for k = 7

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	
actual cluster: 0	35	0	0	0	0	40	0	
actual cluster: 1	0	0	28	0	47	0	0	
actual cluster: 2	0	0	0	71	0	0	4	
actual cluster: 3	0	53	0	0	0	0	22	

predictions for k = 8

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7	
actual cluster: 0	0	35	0	0	0	0	40	0	
actual cluster: 1	28	0	0	0	47	0	0	0	
actual cluster: 2	0	0	38	1	0	0	0	36	
actual cluster: 3	0	0	0	32	0	43	0	0	

predictions for k = 9

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7	k: 8	
actual cluster: 0	0	44	0	0	31	0	0	0	0	
actual cluster: 1	0	0	26	0	0	49	0	0	0	
actual cluster: 2	42	0	0	0	0	0	33	0	0	
actual cluster: 3	0	0	0	20	0	0	0	32	0	


```

          k: 8
actual cluster: 0    0
actual cluster: 1    0
actual cluster: 2    0
actual cluster: 3   23
          predictions for k = 10
                                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0              0  45   0   0   0   0   0  30
actual cluster: 1              0   0  32   0   0  19   0   0
actual cluster: 2             42   0   0   0   0   0   0   0
actual cluster: 3              0   0   0  30  23   0  22   0

```

```

          k: 8 k: 9
actual cluster: 0    0    0
actual cluster: 1   24    0
actual cluster: 2    0   33
actual cluster: 3    0    0
          predictions for k = 1
                                k: 0
actual cluster: 0              75
actual cluster: 1              75
actual cluster: 2              75
actual cluster: 3              75
          predictions for k = 2
                                k: 0 k: 1
actual cluster: 0              0   75
actual cluster: 1              75   0
actual cluster: 2              75   0
actual cluster: 3              75   0
          predictions for k = 3
                                k: 0 k: 1 k: 2
actual cluster: 0              0   75   0
actual cluster: 1              75   0   0
actual cluster: 2              0   0   75
actual cluster: 3              0   0   75
          predictions for k = 4
                                k: 0 k: 1 k: 2 k: 3
actual cluster: 0              0   75   0   0
actual cluster: 1              75   0   0   0
actual cluster: 2              0   0   75   0
actual cluster: 3              0   0   0   75
          predictions for k = 5
                                k: 0 k: 1 k: 2 k: 3 k: 4
actual cluster: 0              75   0   0   0   0
actual cluster: 1              0   0   75   0   0
actual cluster: 2              0   75   0   0   0
actual cluster: 3              0   0   0   44  31

```

```

        predictions for k = 6
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5
actual cluster: 0      0  75   0   0   0   0
actual cluster: 1      0   0  30   0   0  45
actual cluster: 2     71   0   0   0   4   0
actual cluster: 3      0   0   0  53  22   0

        predictions for k = 7
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6
actual cluster: 0      0  30   0   0   0   0  45
actual cluster: 1     28   0   0   0   0  47   0
actual cluster: 2      0   0  75   0   0   0   0
actual cluster: 3      0   0   0  37  38   0   0

        predictions for k = 8
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0      0  45   0   0   0  30   0   0
actual cluster: 1      0   0  40   0  35   0   0   0
actual cluster: 2      0   0   0  36   0   0   4  35
actual cluster: 3     53   0   0   0   0   0  22   0

        predictions for k = 9
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0      0  39   0   0   0   0   0   0
actual cluster: 1     40   0   0   0   0   0  35   0
actual cluster: 2      0   0  33   0  42   0   0   0
actual cluster: 3      0   0   0  22   0  23   0  30

        k: 8
actual cluster: 0    36
actual cluster: 1     0
actual cluster: 2     0
actual cluster: 3     0

        predictions for k = 10
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0      0  36   0   0   0  39   0   0
actual cluster: 1     19   0   0   0   0   0  32   0
actual cluster: 2      0   0  38   0   0   0   0   0
actual cluster: 3      0   0   0  23  21   0   0  31

        k: 8 k: 9
actual cluster: 0     0   0
actual cluster: 1     0  24
actual cluster: 2    37   0
actual cluster: 3     0   0

        predictions for k = 1
                k: 0
actual cluster: 0     75
actual cluster: 1     75

```

```

actual cluster: 2          75
actual cluster: 3          75
      predictions for k = 2
                k: 0 k: 1
actual cluster: 0          75    0
actual cluster: 1          0    75
actual cluster: 2          0    75
actual cluster: 3          0    75
      predictions for k = 3
                k: 0 k: 1 k: 2
actual cluster: 0          0    0    75
actual cluster: 1          0    75    0
actual cluster: 2          75    0    0
actual cluster: 3          75    0    0
      predictions for k = 4
                k: 0 k: 1 k: 2 k: 3
actual cluster: 0          0    0    75    0
actual cluster: 1          0    75    0    0
actual cluster: 2          75    0    0    0
actual cluster: 3          0    0    0    75
      predictions for k = 5
                k: 0 k: 1 k: 2 k: 3 k: 4
actual cluster: 0          0    75    0    0    0
actual cluster: 1          0    0    0    75    0
actual cluster: 2          0    0    71    0    4
actual cluster: 3          53    0    0    0    22
      predictions for k = 6
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5
actual cluster: 0          75    0    0    0    0    0
actual cluster: 1          0    0    0    51    24    0
actual cluster: 2          0    0    75    0    0    0
actual cluster: 3          0    32    0    0    0    43
      predictions for k = 7
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6
actual cluster: 0          31    0    0    0    44    0    0
actual cluster: 1          0    0    75    0    0    0    0
actual cluster: 2          0    0    0    38    0    0    37
actual cluster: 3          0    28    0    0    0    47    0
      predictions for k = 8
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0          40    0    0    0    0    0    35    0
actual cluster: 1          0    0    40    0    0    0    0    35
actual cluster: 2          0    0    0    38    0    37    0    0
actual cluster: 3          0    29    0    0    46    0    0    0
      predictions for k = 9
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0          0    25    0    0    21    0    0    0
actual cluster: 1          0    0    49    0    0    0    0    26

```

actual cluster: 2	33	0	0	0	0	0	42	0
actual cluster: 3	0	0	0	32	0	43	0	0

k: 8
 actual cluster: 0 29
 actual cluster: 1 0
 actual cluster: 2 0
 actual cluster: 3 0

predictions for k = 10 \
 k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
 actual cluster: 0 45 0 0 0 0 0 30 0
 actual cluster: 1 0 0 34 0 0 0 0 19
 actual cluster: 2 0 0 0 38 0 37 0 0
 actual cluster: 3 0 22 0 0 22 0 0 0

k: 8 k: 9
 actual cluster: 0 0 0
 actual cluster: 1 0 22
 actual cluster: 2 0 0
 actual cluster: 3 31 0

predictions for k = 1
 k: 0
 actual cluster: 0 75
 actual cluster: 1 75
 actual cluster: 2 75
 actual cluster: 3 75

predictions for k = 2
 k: 0 k: 1
 actual cluster: 0 0 75
 actual cluster: 1 75 0
 actual cluster: 2 0 75
 actual cluster: 3 75 0

predictions for k = 3
 k: 0 k: 1 k: 2
 actual cluster: 0 0 75 0
 actual cluster: 1 75 0 0
 actual cluster: 2 0 0 75
 actual cluster: 3 75 0 0

predictions for k = 4
 k: 0 k: 1 k: 2 k: 3
 actual cluster: 0 0 75 0 0
 actual cluster: 1 75 0 0 0
 actual cluster: 2 0 0 75 0
 actual cluster: 3 0 0 0 75

predictions for k = 5
 k: 0 k: 1 k: 2 k: 3 k: 4

actual cluster: 0	0	0	75	0	0
actual cluster: 1	31	0	0	0	44
actual cluster: 2	0	75	0	0	0
actual cluster: 3	0	0	0	75	0

predictions for k = 6

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5
actual cluster: 0	0	0	75	0	0	0
actual cluster: 1	0	0	0	75	0	0
actual cluster: 2	26	0	0	0	0	49
actual cluster: 3	0	37	0	0	38	0

predictions for k = 7

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6
actual cluster: 0	0	41	0	0	34	0	0
actual cluster: 1	57	0	0	0	0	18	0
actual cluster: 2	0	0	75	0	0	0	0
actual cluster: 3	0	0	0	54	0	0	21

predictions for k = 8

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	75	0	0	0	0	0	0
actual cluster: 1	0	0	0	27	0	0	20	28
actual cluster: 2	0	0	44	0	31	0	0	0
actual cluster: 3	49	0	0	0	0	26	0	0

predictions for k = 9 \

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	33	0	0	0	0	42	0
actual cluster: 1	29	0	0	0	27	0	0	0
actual cluster: 2	0	0	59	0	0	0	0	16
actual cluster: 3	0	0	0	52	0	23	0	0

k: 8

actual cluster: 0	0
actual cluster: 1	19
actual cluster: 2	0
actual cluster: 3	0

predictions for k = 10 \

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	33	0	0	0	0	42	0
actual cluster: 1	29	0	0	0	27	0	0	0
actual cluster: 2	0	0	40	0	0	0	0	14
actual cluster: 3	0	0	0	50	0	25	0	0

k: 8 k: 9

actual cluster: 0	0	0
actual cluster: 1	19	0
actual cluster: 2	0	21
actual cluster: 3	0	0

```

        predictions for k = 1
            k: 0
actual cluster: 0      75
actual cluster: 1      75
actual cluster: 2      75
actual cluster: 3      75
        predictions for k = 2
            k: 0 k: 1
actual cluster: 0      75    0
actual cluster: 1      0    75
actual cluster: 2      75    0
actual cluster: 3      0    75
        predictions for k = 3
            k: 0 k: 1 k: 2
actual cluster: 0      0    0    75
actual cluster: 1      0    75    0
actual cluster: 2      75    0    0
actual cluster: 3      0    75    0
        predictions for k = 4
            k: 0 k: 1 k: 2 k: 3
actual cluster: 0      0    0    75    0
actual cluster: 1      0    0    0    75
actual cluster: 2      75    0    0    0
actual cluster: 3      0    75    0    0
        predictions for k = 5
            k: 0 k: 1 k: 2 k: 3 k: 4
actual cluster: 0      0    0    75    0    0
actual cluster: 1      0    0    0    38    37
actual cluster: 2      75    0    0    0    0
actual cluster: 3      0    75    0    0    0
        predictions for k = 6
            k: 0 k: 1 k: 2 k: 3 k: 4 k: 5
actual cluster: 0      0    0    75    0    0    0
actual cluster: 1      33    0    0    0    42    0
actual cluster: 2      0    75    0    0    0    0
actual cluster: 3      0    0    0    29    0    46
        predictions for k = 7
            k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6
actual cluster: 0      0    75    0    0    0    0    0
actual cluster: 1      0    0    0    34    0    41    0
actual cluster: 2      0    0    25    0    50    0    0
actual cluster: 3      29    0    0    0    0    0    46
        predictions for k = 8
            k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0      0    29    0    0    46    0    0    0
actual cluster: 1      0    0    0    53    0    22    0    0
actual cluster: 2      0    0    42    0    0    0    33    0
actual cluster: 3      34    0    0    0    0    0    0    41

```

	predictions for k = 9								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7	
actual cluster: 0	0	34	0	0	41	0	0	0	
actual cluster: 1	22	0	0	0	0	0	53	0	
actual cluster: 2	0	0	18	0	0	29	0	0	
actual cluster: 3	0	0	0	50	0	0	0	25	

	k: 8
actual cluster: 0	0
actual cluster: 1	0
actual cluster: 2	28
actual cluster: 3	0

	predictions for k = 10								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7	
actual cluster: 0	0	0	37	0	0	0	38	0	
actual cluster: 1	0	0	0	29	0	26	0	0	
actual cluster: 2	0	49	0	0	26	0	0	0	
actual cluster: 3	40	0	0	0	0	0	0	22	

	k: 8	k: 9
actual cluster: 0	0	0
actual cluster: 1	20	0
actual cluster: 2	0	0
actual cluster: 3	0	13

	predictions for k = 1	
	k: 0	
actual cluster: 0	75	
actual cluster: 1	75	
actual cluster: 2	75	
actual cluster: 3	75	

	predictions for k = 2	
	k: 0	k: 1
actual cluster: 0	0	75
actual cluster: 1	75	0
actual cluster: 2	0	75
actual cluster: 3	75	0

	predictions for k = 3		
	k: 0	k: 1	k: 2
actual cluster: 0	0	75	0
actual cluster: 1	75	0	0
actual cluster: 2	0	0	75
actual cluster: 3	75	0	0

	predictions for k = 4			
	k: 0	k: 1	k: 2	k: 3
actual cluster: 0	0	75	0	0
actual cluster: 1	0	0	0	75

actual cluster: 2	0	0	75	0				
actual cluster: 3	75	0	0	0				
predictions for k = 5								
	k: 0	k: 1	k: 2	k: 3	k: 4			
actual cluster: 0	75	0	0	0	0			
actual cluster: 1	0	46	0	0	29			
actual cluster: 2	0	0	75	0	0			
actual cluster: 3	0	0	0	75	0			
predictions for k = 6								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5		
actual cluster: 0	0	0	75	0	0	0		
actual cluster: 1	0	48	0	0	27	0		
actual cluster: 2	30	0	0	0	0	45		
actual cluster: 3	0	0	0	75	0	0		
predictions for k = 7								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	
actual cluster: 0	0	0	41	0	0	0	34	
actual cluster: 1	0	0	0	36	39	0	0	
actual cluster: 2	0	33	0	0	0	42	0	
actual cluster: 3	75	0	0	0	0	0	0	
predictions for k = 8								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	26	0	0	0	0	49	0	0
actual cluster: 1	0	0	0	37	0	0	18	20
actual cluster: 2	0	0	75	0	0	0	0	0
actual cluster: 3	0	22	0	0	53	0	0	0
predictions for k = 9								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	33	0	0	42	0	0	0
actual cluster: 1	0	0	0	44	0	0	31	0
actual cluster: 2	0	0	44	0	0	0	0	31
actual cluster: 3	30	0	0	0	0	25	0	0
k: 8								
actual cluster: 0	0							
actual cluster: 1	0							
actual cluster: 2	0							
actual cluster: 3	20							
predictions for k = 10								
	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	0	28	0	47	0	0	0
actual cluster: 1	0	29	0	0	0	26	0	20
actual cluster: 2	34	0	0	0	0	0	0	0
actual cluster: 3	0	0	0	39	0	0	13	0
k: 8 k: 9								


```

actual cluster: 0    0    0
actual cluster: 1    0    0
actual cluster: 2    0   41
actual cluster: 3   23    0
      predictions for k = 1
                k: 0
actual cluster: 0           75
actual cluster: 1           75
actual cluster: 2           75
actual cluster: 3           75
      predictions for k = 2
                k: 0 k: 1
actual cluster: 0           0   75
actual cluster: 1          62   13
actual cluster: 2          75    0
actual cluster: 3          23   52
      predictions for k = 3
                k: 0 k: 1 k: 2
actual cluster: 0           0   74    1
actual cluster: 1          15    3   57
actual cluster: 2          71    0    4
actual cluster: 3           0   12   63
      predictions for k = 4
                k: 0 k: 1 k: 2 k: 3
actual cluster: 0          62    0    0   13
actual cluster: 1           0   61    6    8
actual cluster: 2           0   12   63    0
actual cluster: 3           2   27    0   46
      predictions for k = 5
                k: 0 k: 1 k: 2 k: 3 k: 4
actual cluster: 0          61    0    0    0   14
actual cluster: 1           0   36   30    1    8
actual cluster: 2           0   19    0   56    0
actual cluster: 3           3    9   24    0   39
      predictions for k = 6
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5
actual cluster: 0           0   62    0    0   13    0
actual cluster: 1          33    0   23    0    6   13
actual cluster: 2           0    0   15   36    0   24
actual cluster: 3          22    2    8    0   43    0
      predictions for k = 7
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6
actual cluster: 0          13    0    0    0   30    0   32
actual cluster: 1           6   22   33    0    0   14    0
actual cluster: 2           0   15    0   36    0   24    0
actual cluster: 3          42    8   22    0    0    0    3
      predictions for k = 8
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7

```

actual cluster: 0	21	0	0	5	0	0	27	22
actual cluster: 1	0	14	32	6	0	22	0	1
actual cluster: 2	0	18	0	0	43	14	0	0
actual cluster: 3	0	0	21	40	0	8	5	1

predictions for k = 9 \

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	28	0	0	0	22	20	5	0
actual cluster: 1	0	20	25	21	1	0	5	0
actual cluster: 2	0	4	0	15	0	0	0	32
actual cluster: 3	4	5	19	8	1	0	38	0

k: 8

actual cluster: 0	0
actual cluster: 1	3
actual cluster: 2	24
actual cluster: 3	0

predictions for k = 10 \

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	18	0	1	29	0	0	27	0
actual cluster: 1	4	4	4	0	26	0	0	18
actual cluster: 2	0	24	0	0	0	18	0	7
actual cluster: 3	9	0	33	3	16	0	0	10

k: 8 k: 9

actual cluster: 0	0	0
actual cluster: 1	18	1
actual cluster: 2	5	21
actual cluster: 3	4	0

predictions for k = 1

k: 0

actual cluster: 0	75
actual cluster: 1	75
actual cluster: 2	75
actual cluster: 3	75

predictions for k = 2

k: 0 k: 1

actual cluster: 0	75	0
actual cluster: 1	13	62
actual cluster: 2	0	75
actual cluster: 3	52	23

predictions for k = 3

k: 0 k: 1 k: 2

actual cluster: 0	72	0	3
actual cluster: 1	3	11	61
actual cluster: 2	0	71	4
actual cluster: 3	10	0	65

```

        predictions for k = 4
                k: 0 k: 1 k: 2 k: 3
actual cluster: 0      0  62   0  13
actual cluster: 1     61   0   6   8
actual cluster: 2      9   0  66   0
actual cluster: 3     27   2   0  46

        predictions for k = 5
                k: 0 k: 1 k: 2 k: 3 k: 4
actual cluster: 0     62   0  13   0   0
actual cluster: 1      0  36   8   1  30
actual cluster: 2      0  19   0  56   0
actual cluster: 3      3   9  41   0  22

        predictions for k = 6
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5
actual cluster: 0     13   0   0  62   0   0
actual cluster: 1      6  14  33   0   0  22
actual cluster: 2      0  24   0   0  36  15
actual cluster: 3     43   0  22   2   0   8

        predictions for k = 7
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6
actual cluster: 0     31   0  33   0   0  11   0
actual cluster: 1      0  23   1   0  33   5  13
actual cluster: 2      0  15   0  36   0   0  24
actual cluster: 3      0   8   9   0  22  36   0

        predictions for k = 8
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0     29   0   0  28   0   0  18   0
actual cluster: 1      0   3  22   0  29   0   4  17
actual cluster: 2      0  26  15   0   2  32   0   0
actual cluster: 3      0   0   9   3  10   0  15  38

        predictions for k = 9
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0     27   0   0   0  30   0  17   1
actual cluster: 1      0   2  25  20   0   0   4   6
actual cluster: 2      0  34   0  14   0  22   0   0
actual cluster: 3      4   0  16   8   0   0  10  33

        k: 8
actual cluster: 0      0
actual cluster: 1     18
actual cluster: 2      5
actual cluster: 3      4

        predictions for k = 10
                k: 0 k: 1 k: 2 k: 3 k: 4 k: 5 k: 6 k: 7
actual cluster: 0      0  28   0   1   0   0  29   0
actual cluster: 1     22   0   2   4   5   0   0  22
actual cluster: 2      0   0  25   0  24  18   0   3

```

```
actual cluster: 3      8      3      0      34      2      0      0      13
```

actual cluster: 2	0	20	0	0	40	15	0
actual cluster: 3	8	0	38	21	0	8	0

predictions for k = 8

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	17	29	0	28	0	0	1
actual cluster: 1	14	4	0	0	0	21	26	10
actual cluster: 2	18	0	0	43	0	14	0	0
actual cluster: 3	0	12	0	0	3	8	14	38

predictions for k = 9

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	20	0	1	30	0	0	0
actual cluster: 1	25	3	2	7	0	0	20	18
actual cluster: 2	0	0	34	0	0	22	14	5
actual cluster: 3	16	8	0	36	3	0	8	4

k: 8

actual cluster: 0	24
actual cluster: 1	0
actual cluster: 2	0
actual cluster: 3	0

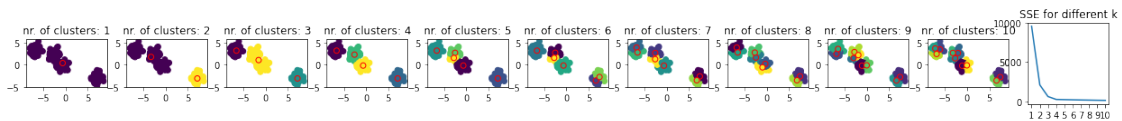
predictions for k = 10

	k: 0	k: 1	k: 2	k: 3	k: 4	k: 5	k: 6	k: 7
actual cluster: 0	0	28	0	0	5	0	0	20
actual cluster: 1	17	0	22	1	5	25	0	0
actual cluster: 2	5	0	4	21	0	0	18	0
actual cluster: 3	4	4	11	0	39	16	0	0

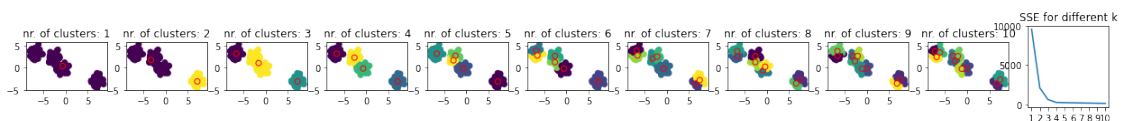
k: 8 k: 9

actual cluster: 0	22	0
actual cluster: 1	1	4
actual cluster: 2	0	27
actual cluster: 3	1	0

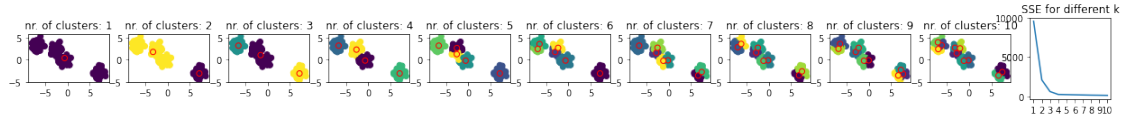
cluster for standard deviation = 0.6, random state = 1



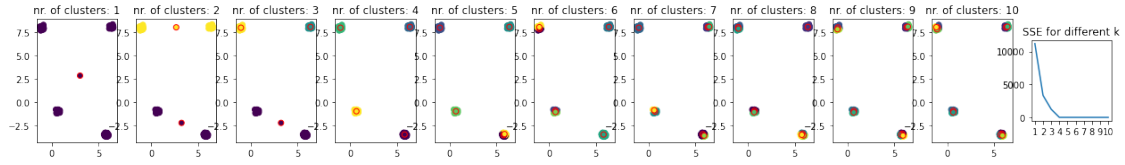
cluster for standard deviation = 0.6, random state = 5



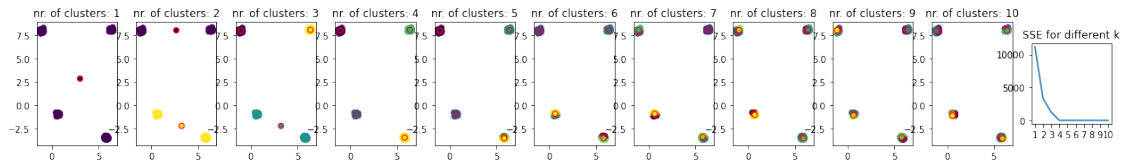
cluster for standard deviation = 0.6, random state = 10



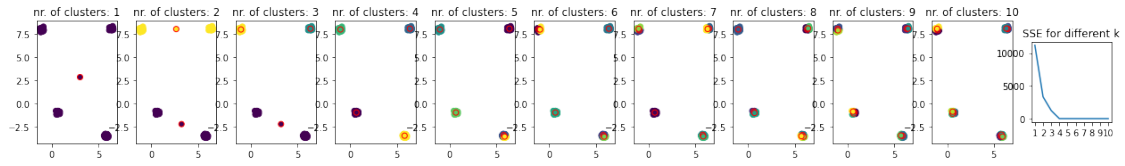
cluster for standard deviation = 0.1, random state = 1



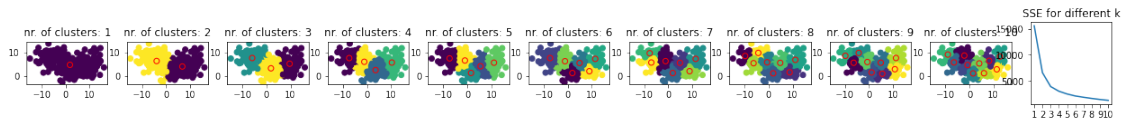
cluster for standard deviation = 0.1, random state = 5



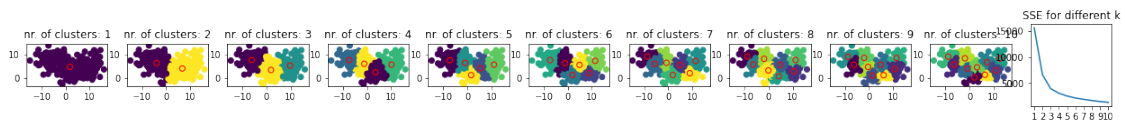
cluster for standard deviation = 0.1, random state = 10

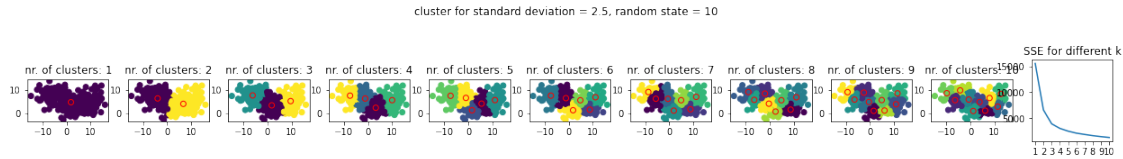


cluster for standard deviation = 2.5, random state = 1



cluster for standard deviation = 2.5, random state = 5





The implementation of k-means relies on randomness in order to find the first set of centroids. Accordingly each different random state produces a different result. However, these differences are only apparent after the actual number of clusters is reached. Meaning that the algorithm produces similar results for $k \leq 4$. with $k > 4$ the solutions start to differgate. It's not supprising that this is the case, given that the data is distributed into 4 blobs. However these changes do not influence the SSE, which seems to be constant for the different random states. One possible way to chane that would be to have the algorithm start at k fixed locations in the grid provided by the min-max values on all axes. This would ensure that the results stay the same.

2 Assignment 2

```
[ ]: data = pd.read_csv('vertebrate.csv')
data
```

```
[ ]:
```

	Name	Warm-blooded	Gives Birth	Aquatic Creature	\
0	human	1	1	0	
1	python	0	0	0	
2	salmon	0	0	1	
3	whale	1	1	1	
4	frog	0	0	1	
5	komodo	0	0	0	
6	bat	1	1	0	
7	pigeon	1	0	0	
8	cat	1	1	0	
9	leopard shark	0	1	1	
10	turtle	0	0	1	
11	penguin	1	0	1	
12	porcupine	1	1	0	
13	eel	0	0	1	
14	salamander	0	0	1	

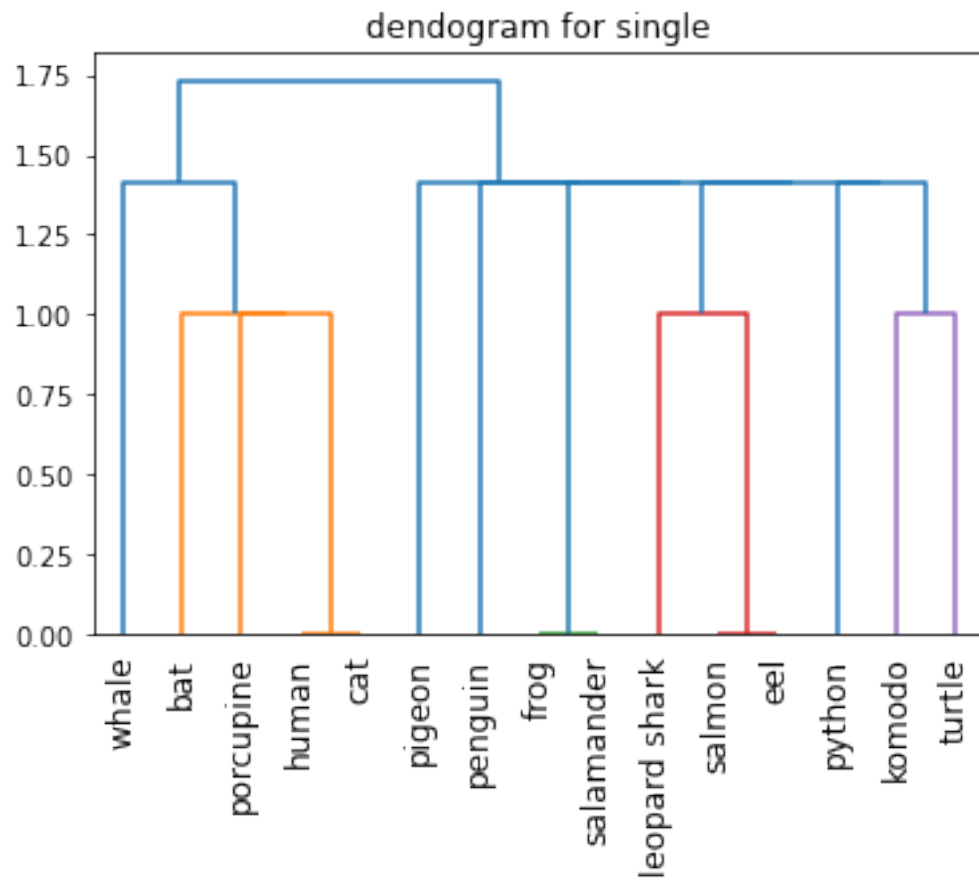
	Aerial Creature	Has Legs	Hibernates	Class
0	0	1	0	mammals
1	0	0	1	reptiles
2	0	0	0	fishes
3	0	0	0	mammals
4	0	1	1	amphibians

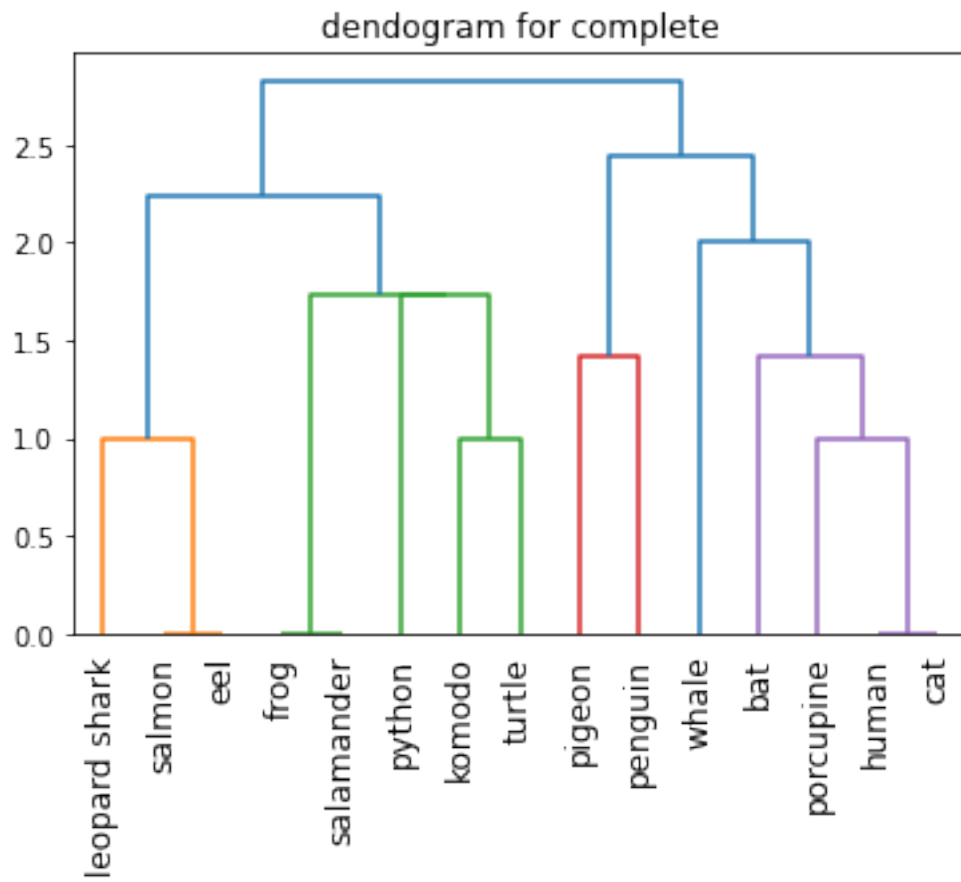
5	0	1	0	reptiles
6	1	1	1	mammals
7	1	1	0	birds
8	0	1	0	mammals
9	0	0	0	fishes
10	0	1	0	reptiles
11	0	1	0	birds
12	0	1	1	mammals
13	0	0	0	fishes
14	0	1	1	amphibians

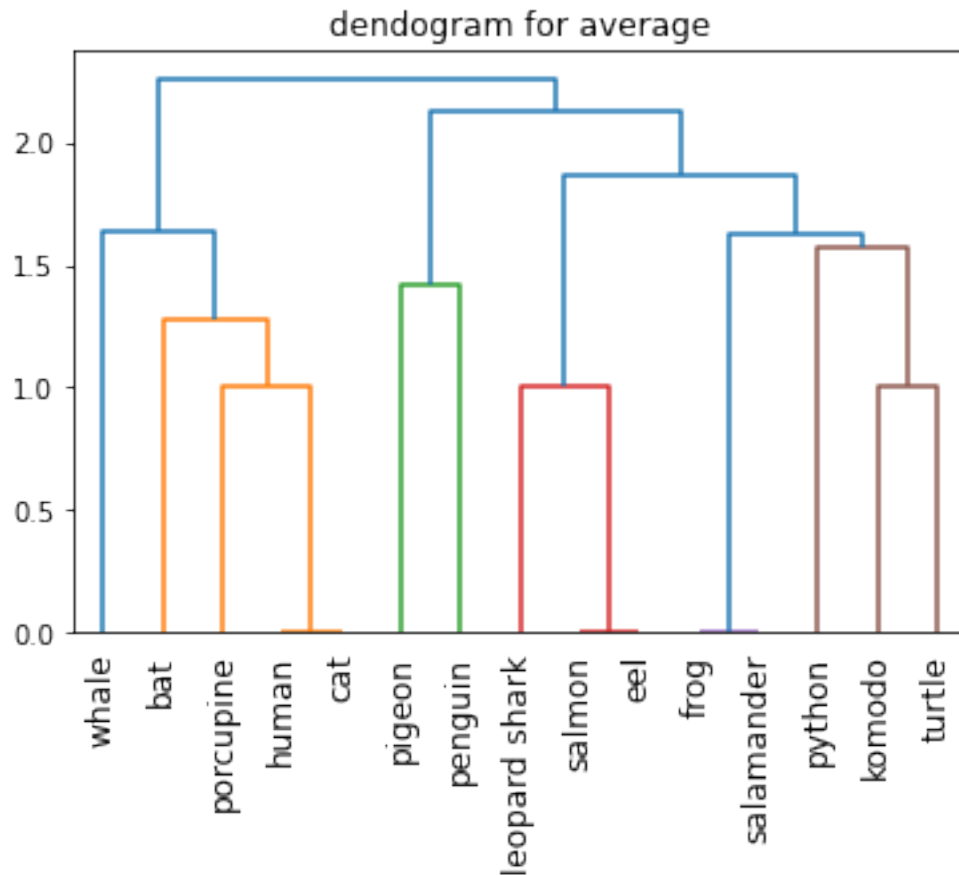
```
[ ]: data = data.convert_dtypes()
data['Class'] = data['Class'].astype(object)
data = pd.get_dummies(data)
```

```
[ ]: from scipy.cluster.hierarchy import dendrogram, single, complete, average
from scipy.spatial.distance import pdist
```

```
[ ]: distances = pdist(data.iloc[:,1:-1].astype(int).values)
for i in [single, complete, average]:
    Z = i(distances)
    dendrogram(Z, labels= data.iloc[:,0].values)
    plt.xticks(rotation=90)
    plt.title(f'dendrogram for {i.__name__}')
    plt.show()
```





I'd argue that the representation provided by 'average' is the best, simply because the distinction between mammals and non-mammals is the most natural to me. This difference is also given in the 'single' dendrogram. However, due to the fact that the rest is clustered into one major subclass is somewhat off in this solution. Comparing 'complete' and 'average' I'd still argue for 'average', simply because birds are closer to reptiles (given that both produce eggs) than to mammals. Generally speaking it's obvious that the best (most natural) solution depends highly on the metric used to evaluate this.

3 Assignment 3

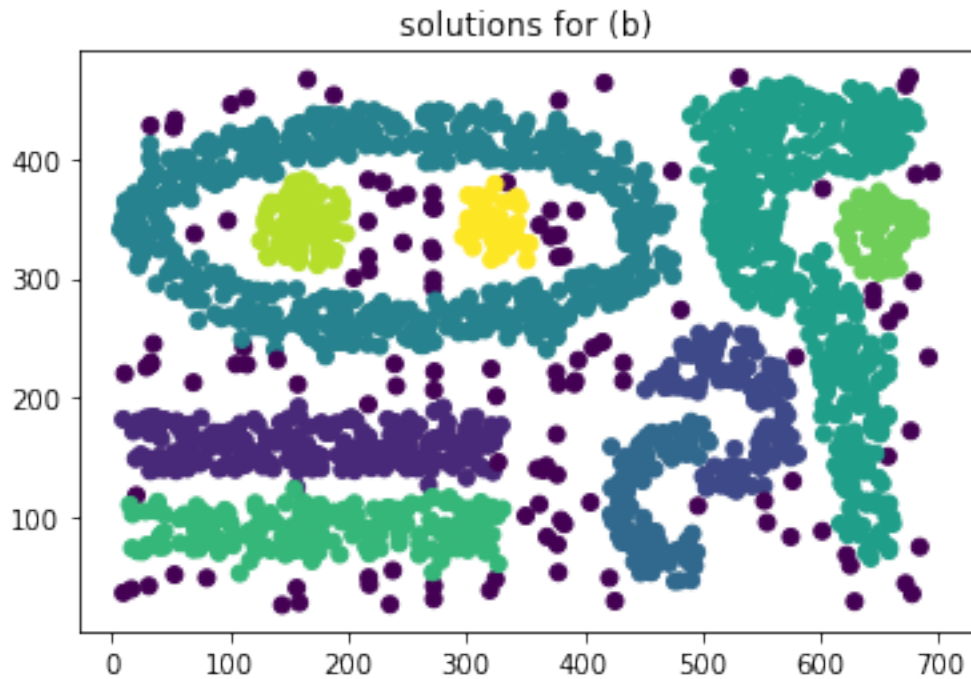
```
[ ]: from sklearn.cluster import DBSCAN
```

```
[ ]: data = pd.read_csv('chameleon.csv')
```

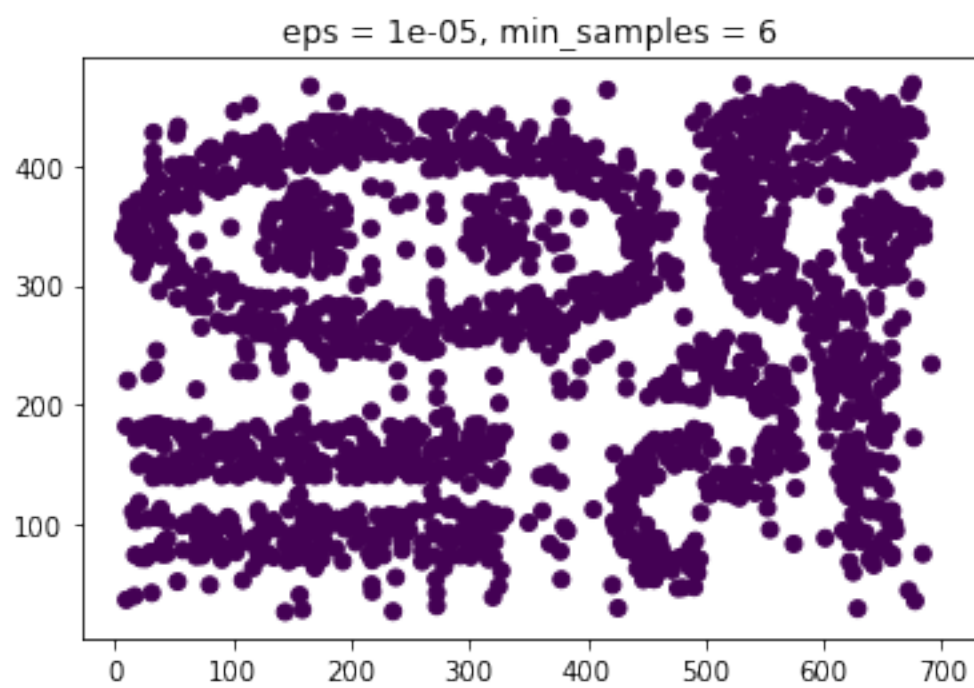
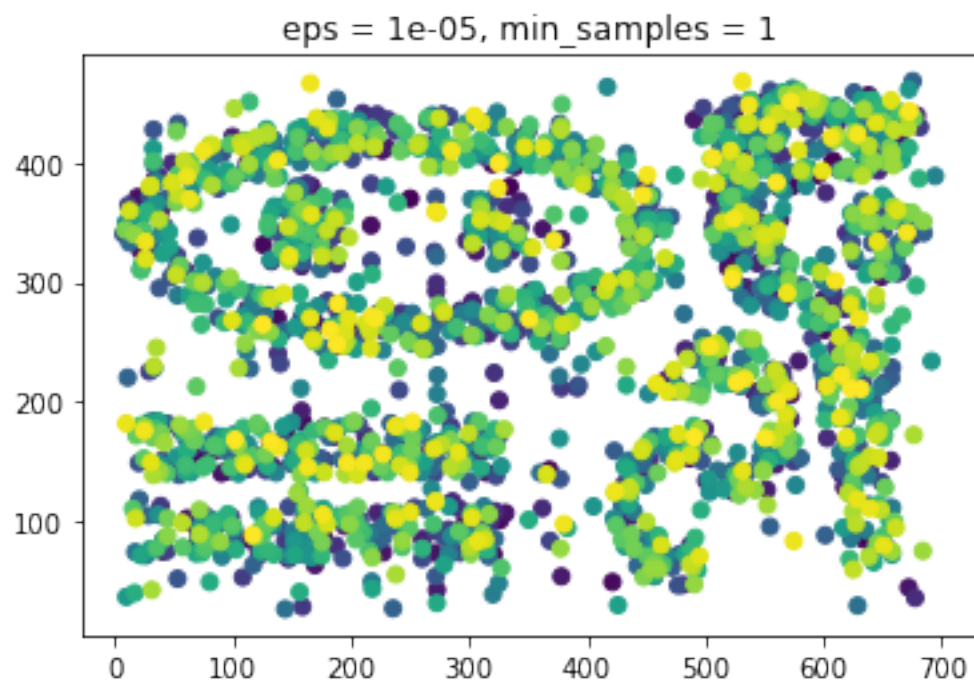
```
[ ]: dbscan = DBSCAN(eps=15.5, min_samples=5)
      class_pred = dbscan.fit_predict(data)
      data['labels'] = class_pred
```

```
[ ]: plt.scatter(data['x'], data['y'], c=data['labels'])
plt.title('solutions for (b)')
```

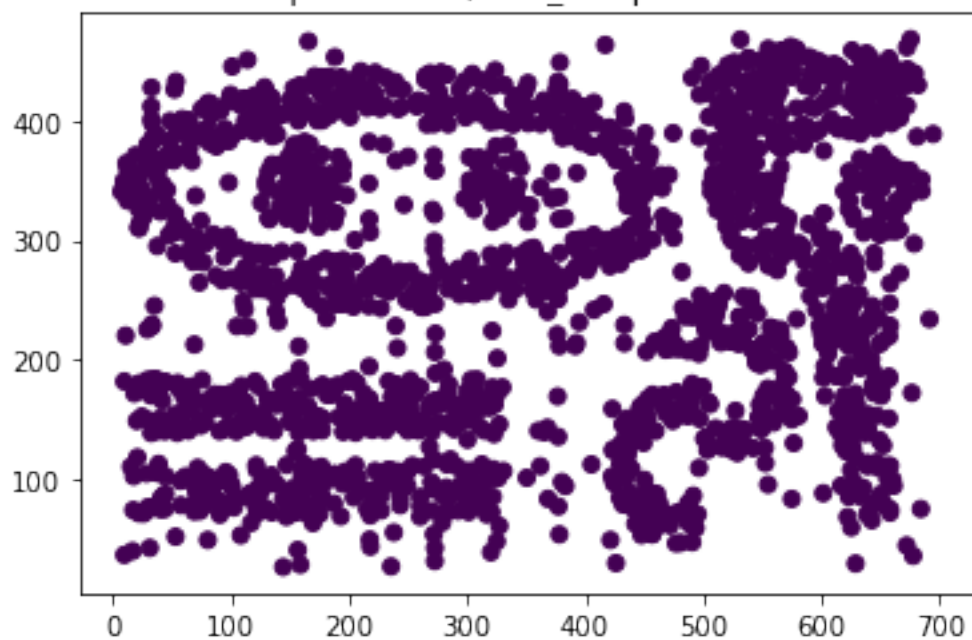
```
[ ]: Text(0.5, 1.0, 'solutions for (b)')
```



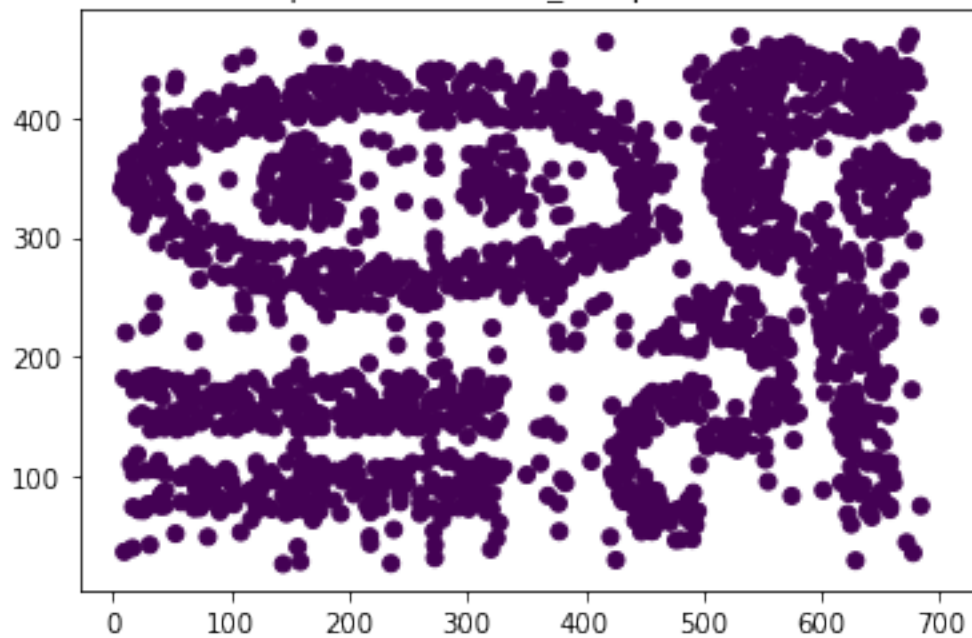
```
[ ]: for eps in np.hstack((1e-5,np.arange(5,21,5))):
    for min_sample in range(1,22,5):
        dbscan = DBSCAN(eps=eps, min_samples=min_sample)
        class_pred = dbscan.fit_predict(data)
        data['labels'] = class_pred
        plt.scatter(data['x'], data['y'], c=data['labels'])
        plt.title(f'eps = {eps}, min_samples = {min_sample}')
        plt.show()
```



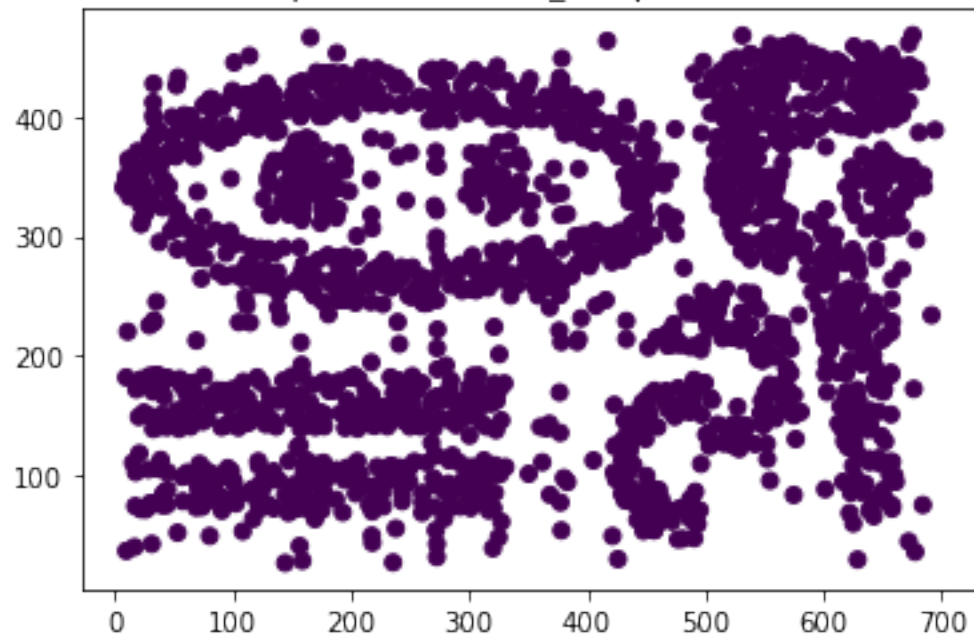
eps = 1e-05, min_samples = 11



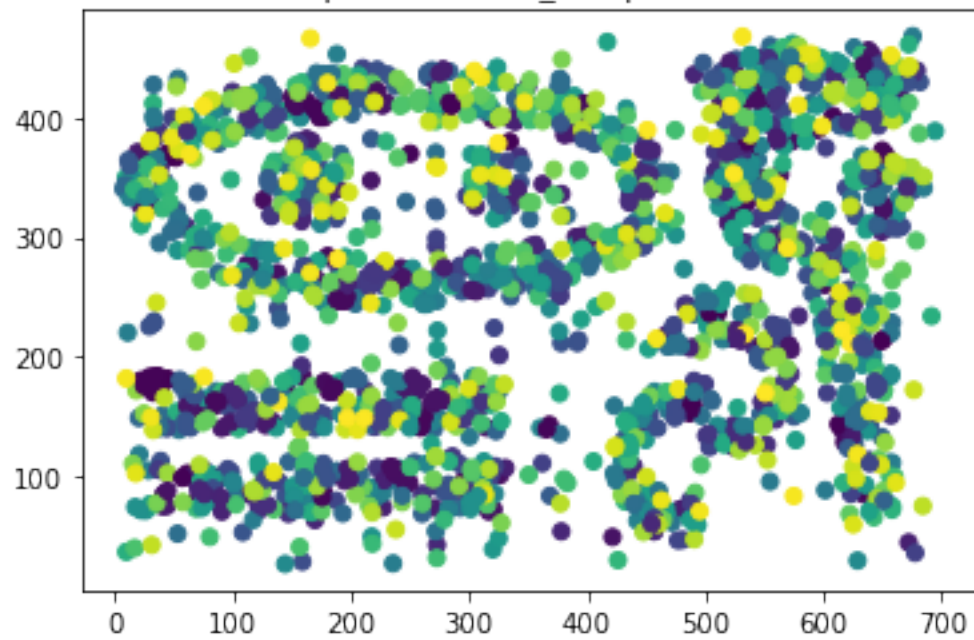
eps = 1e-05, min_samples = 16



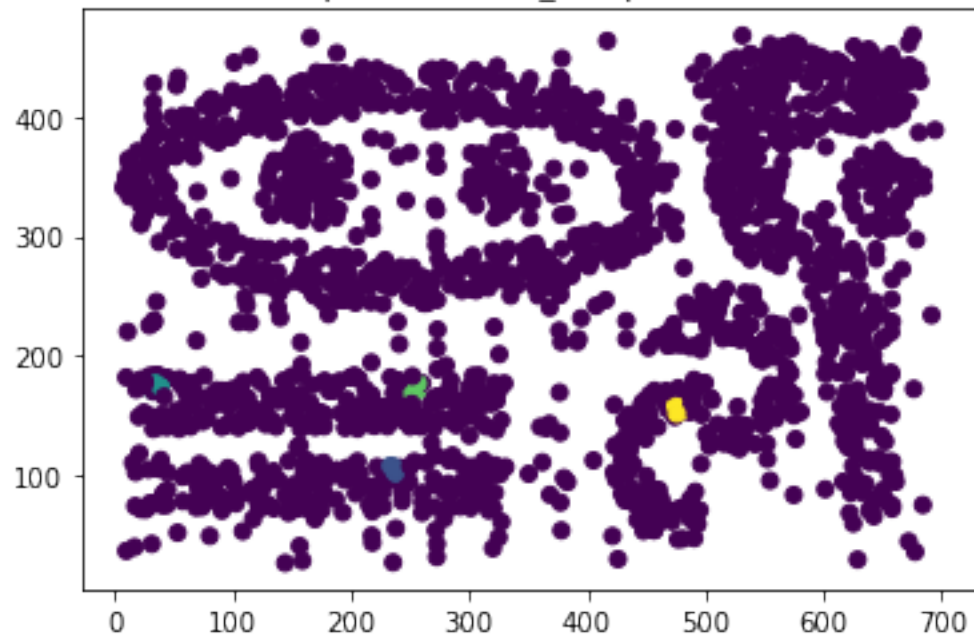
eps = 1e-05, min_samples = 21



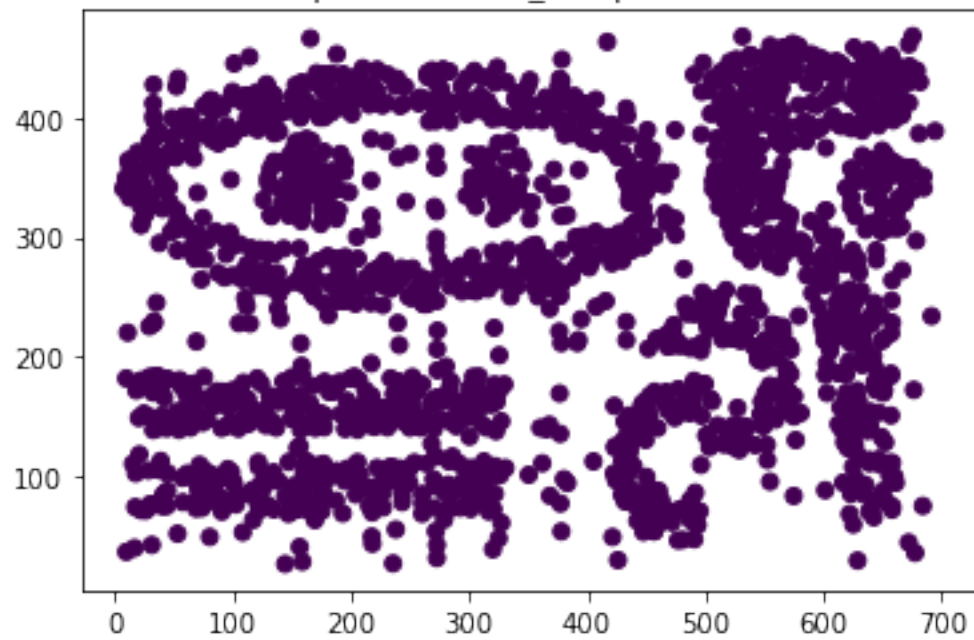
eps = 5.0, min_samples = 1



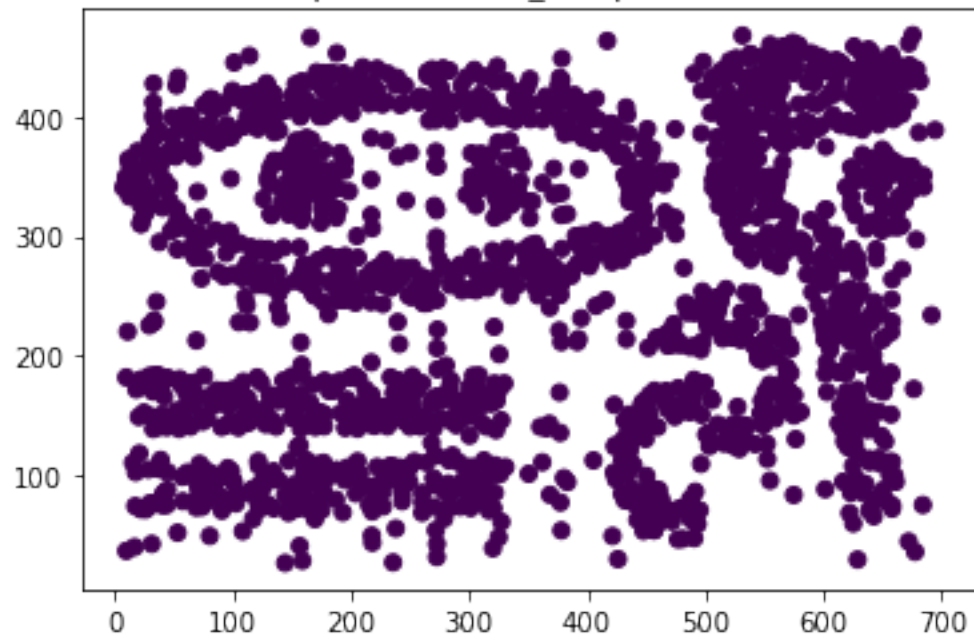
eps = 5.0, min_samples = 6



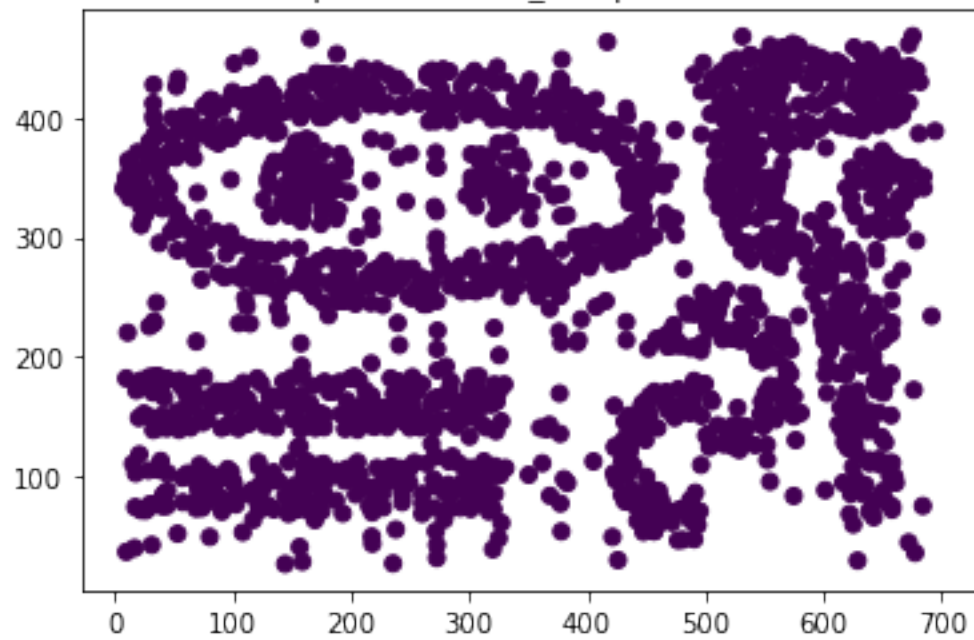
eps = 5.0, min_samples = 11



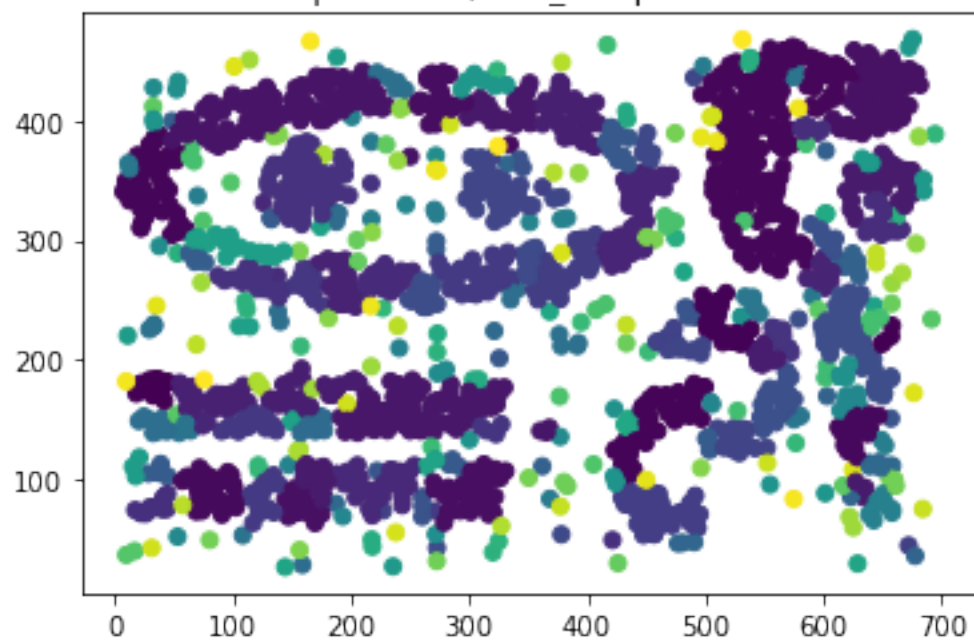
eps = 5.0, min_samples = 16



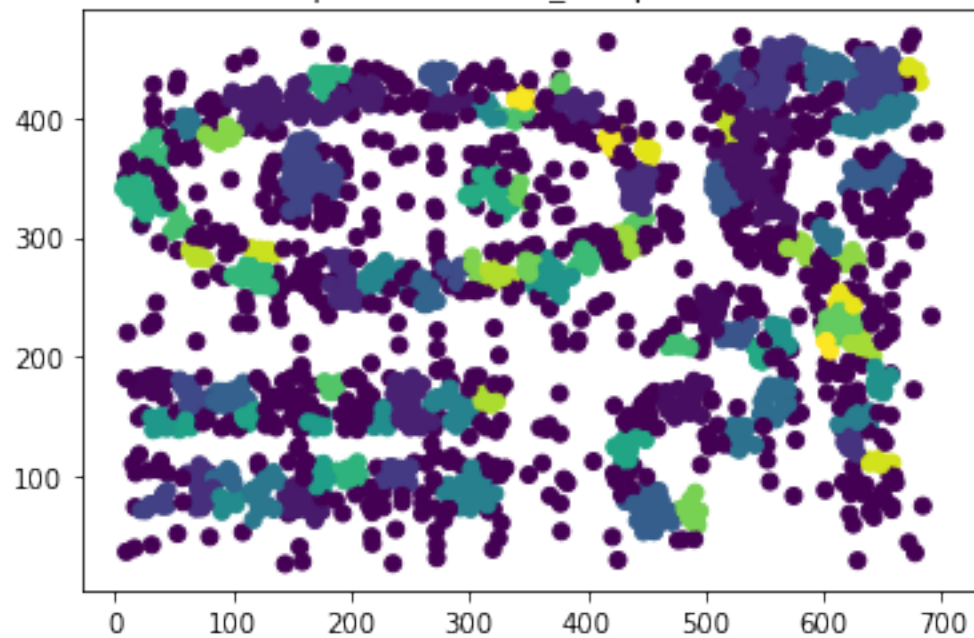
eps = 5.0, min_samples = 21



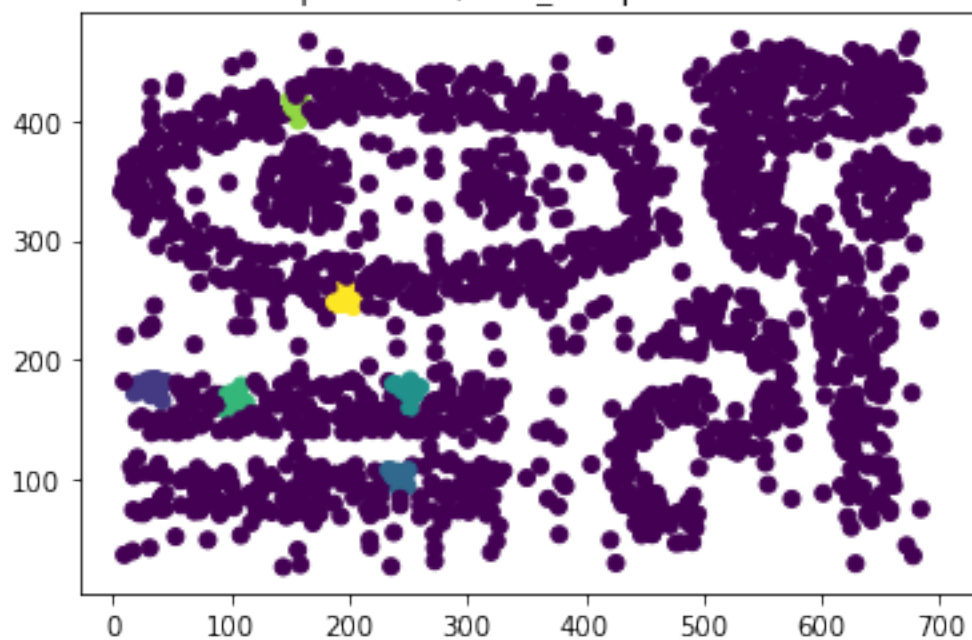
eps = 10.0, min_samples = 1



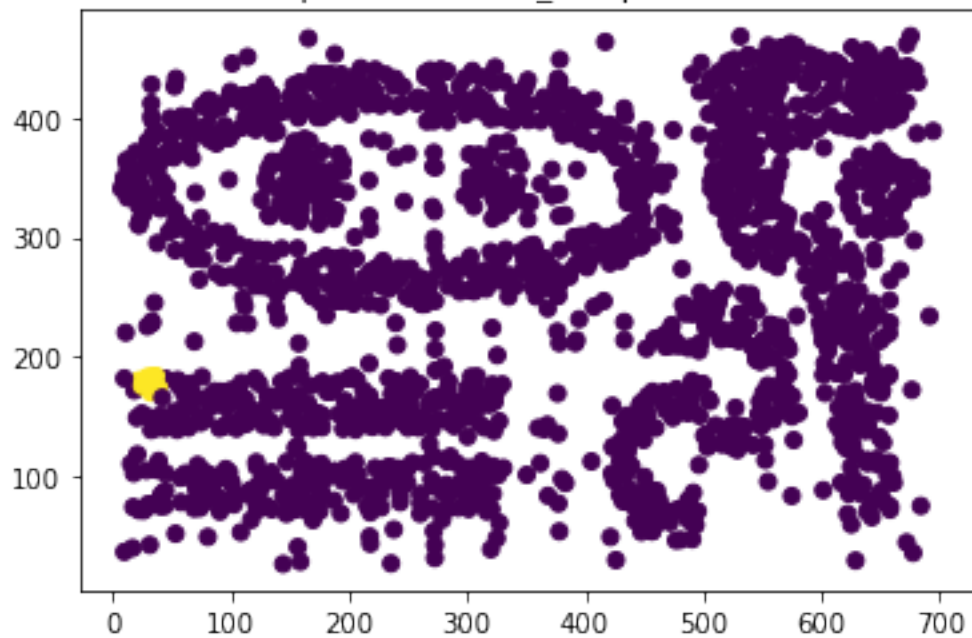
eps = 10.0, min_samples = 6

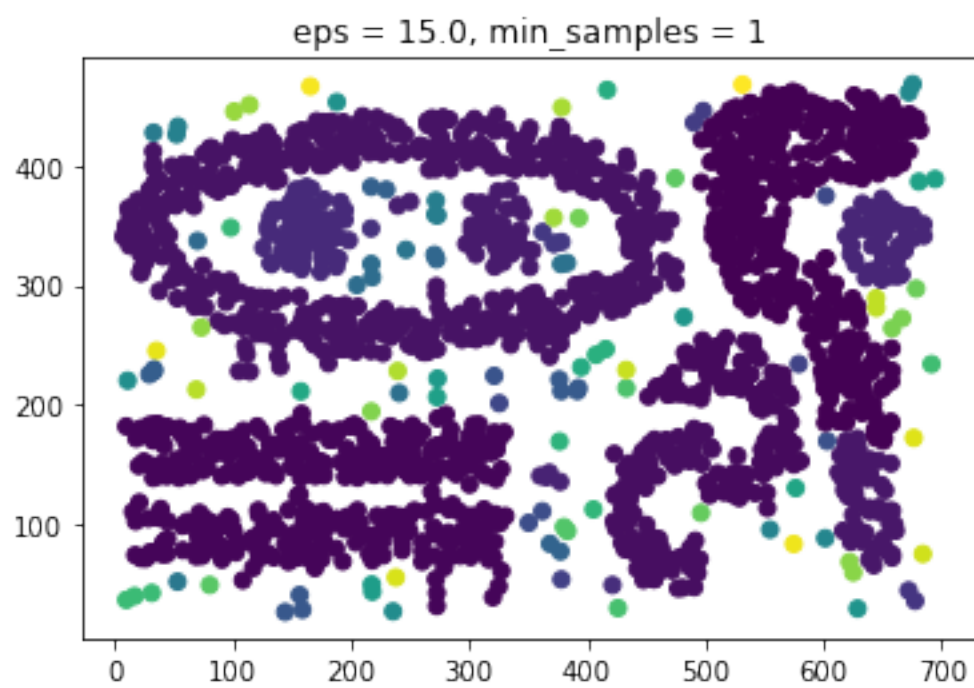
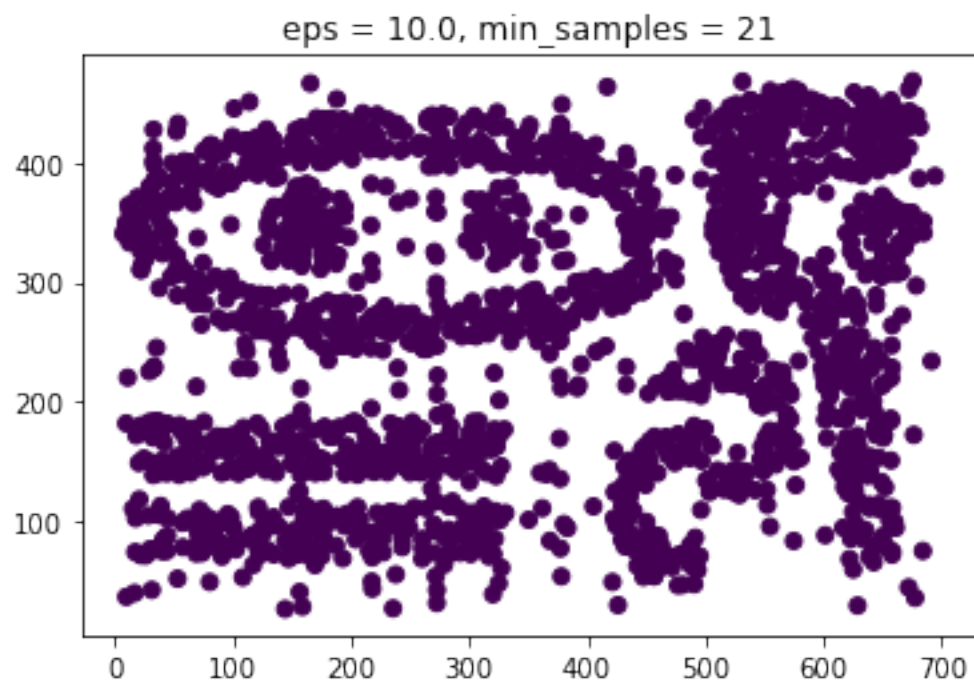


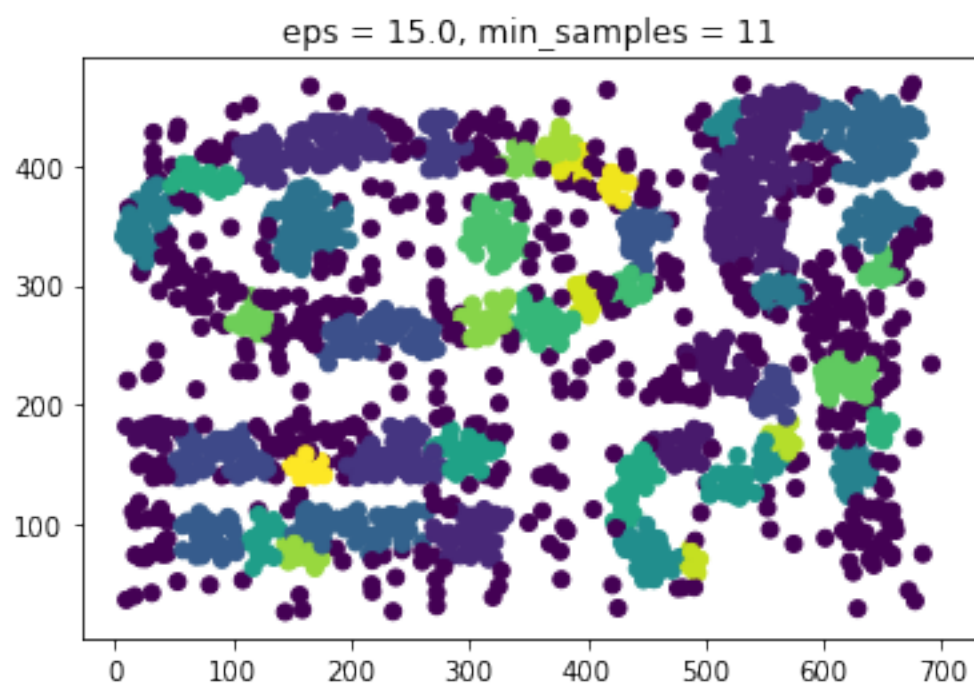
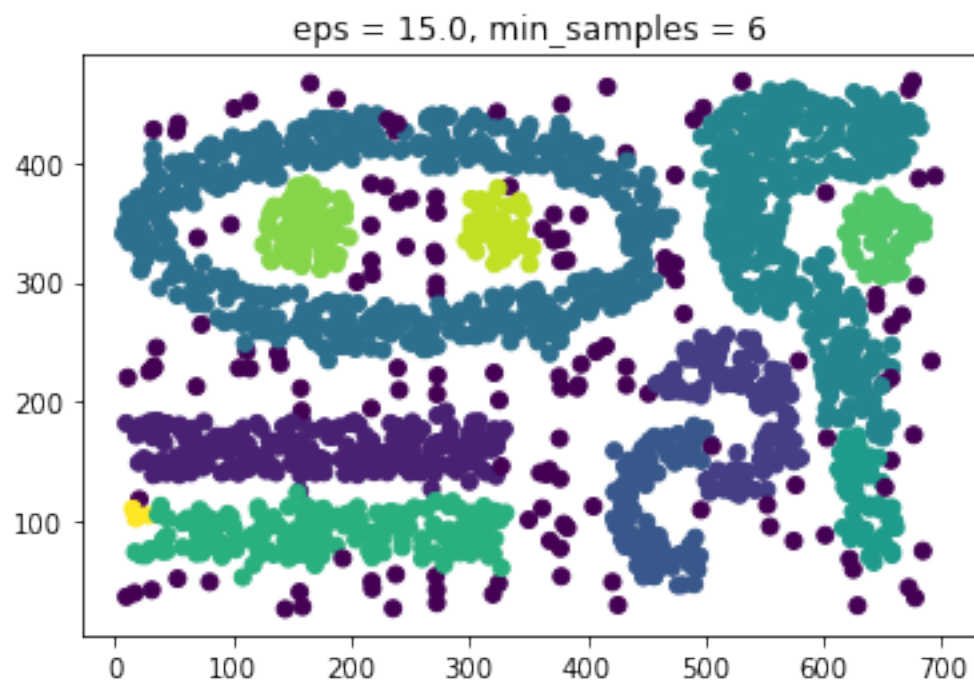
eps = 10.0, min_samples = 11

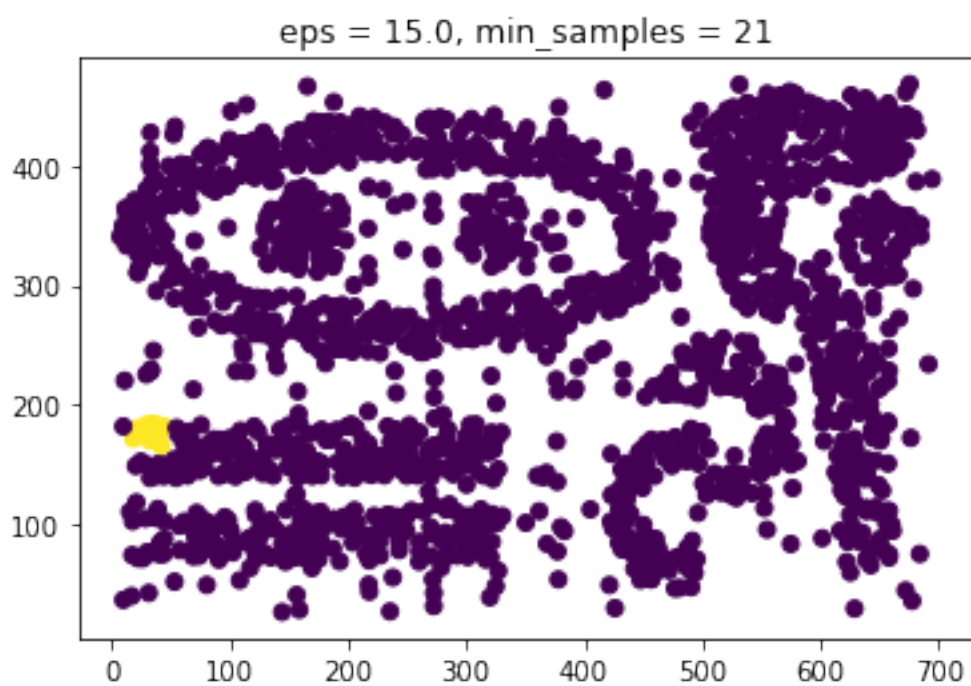
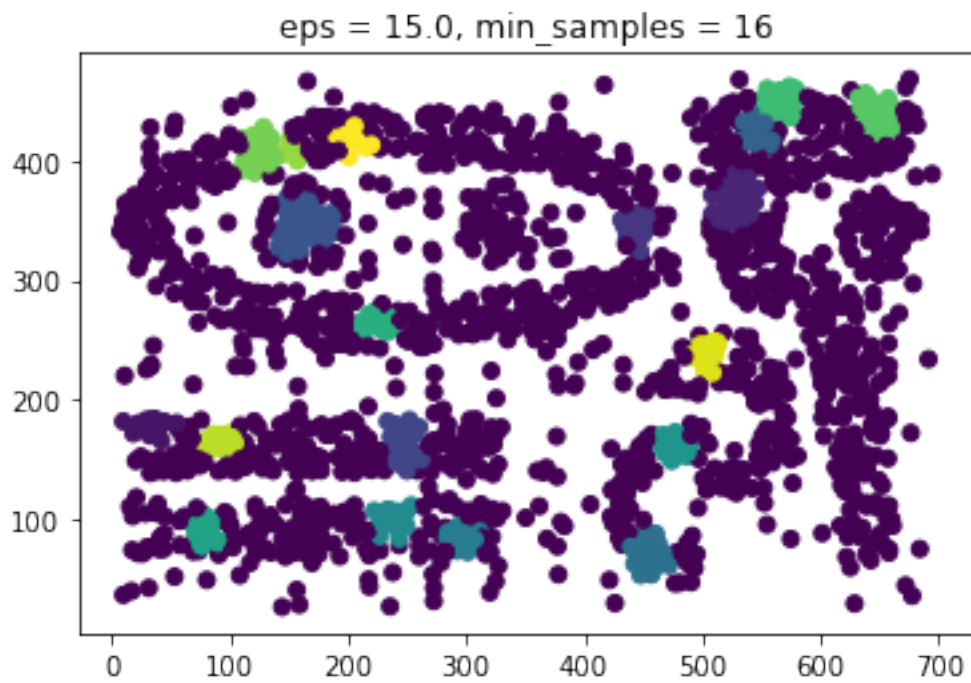


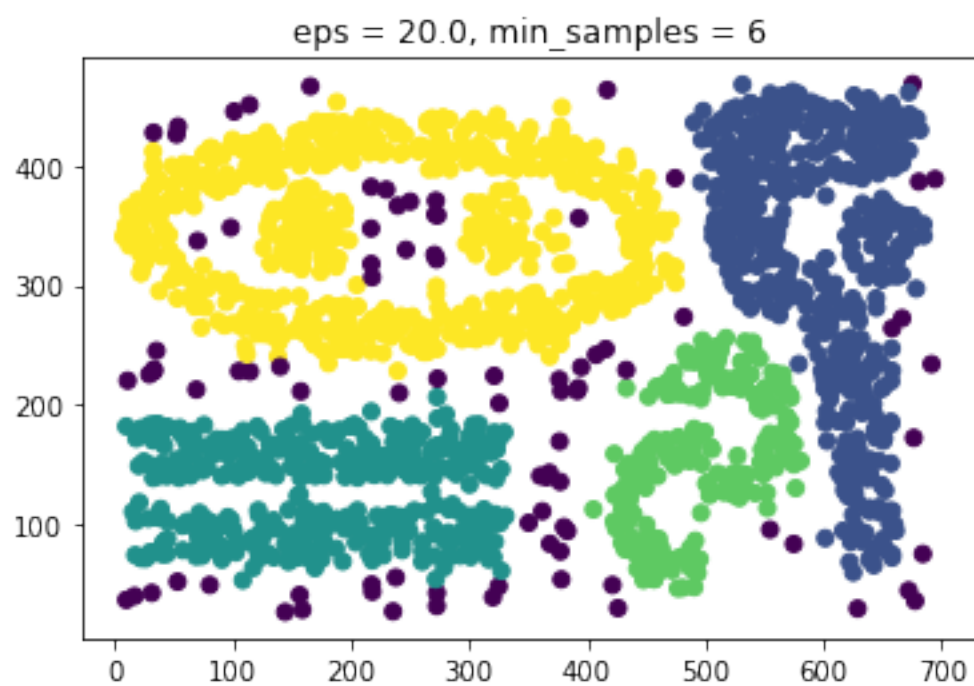
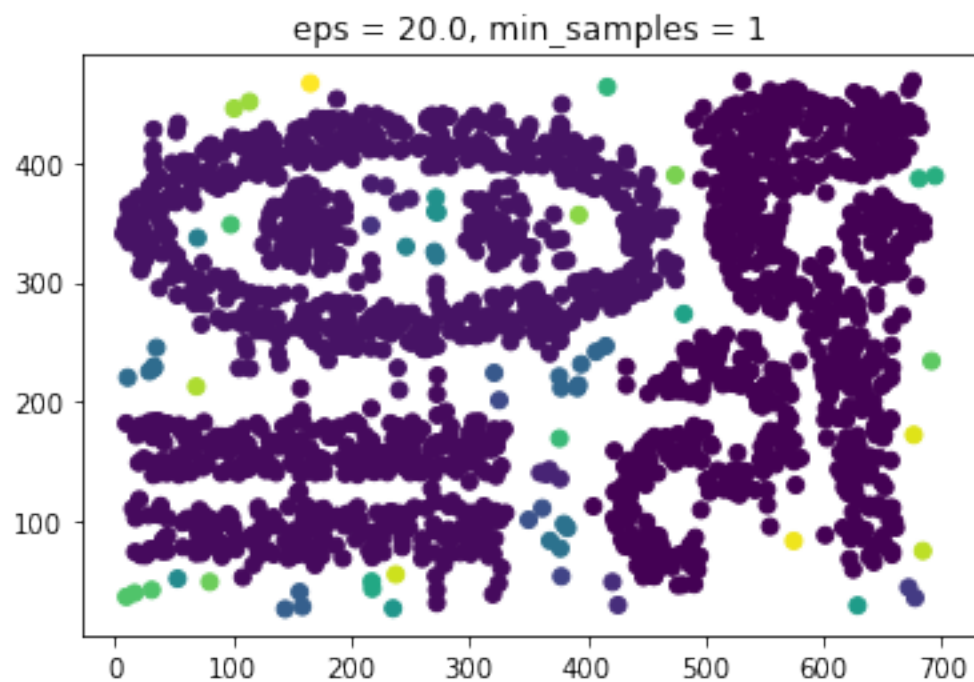
eps = 10.0, min_samples = 16



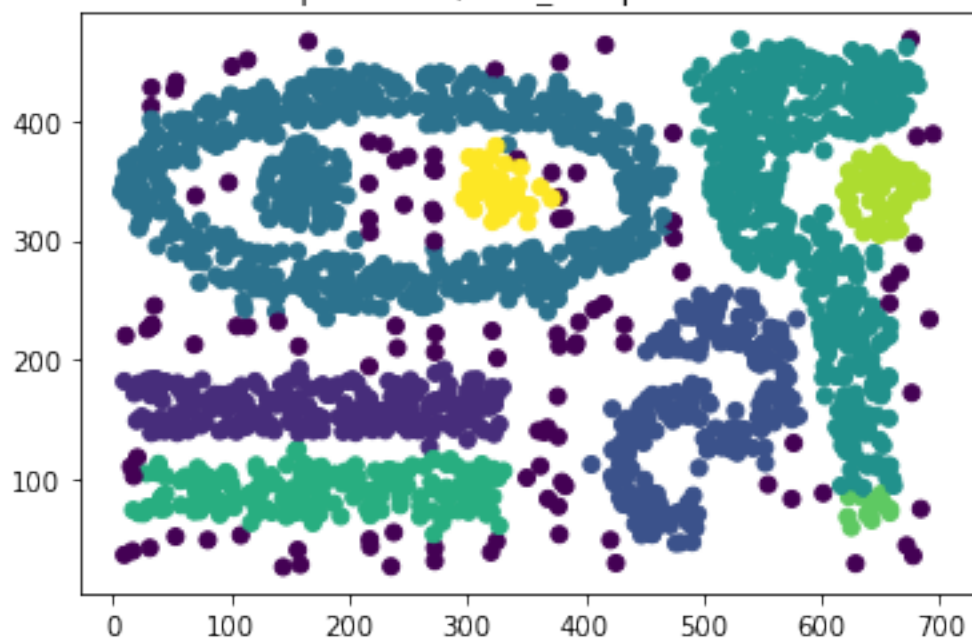




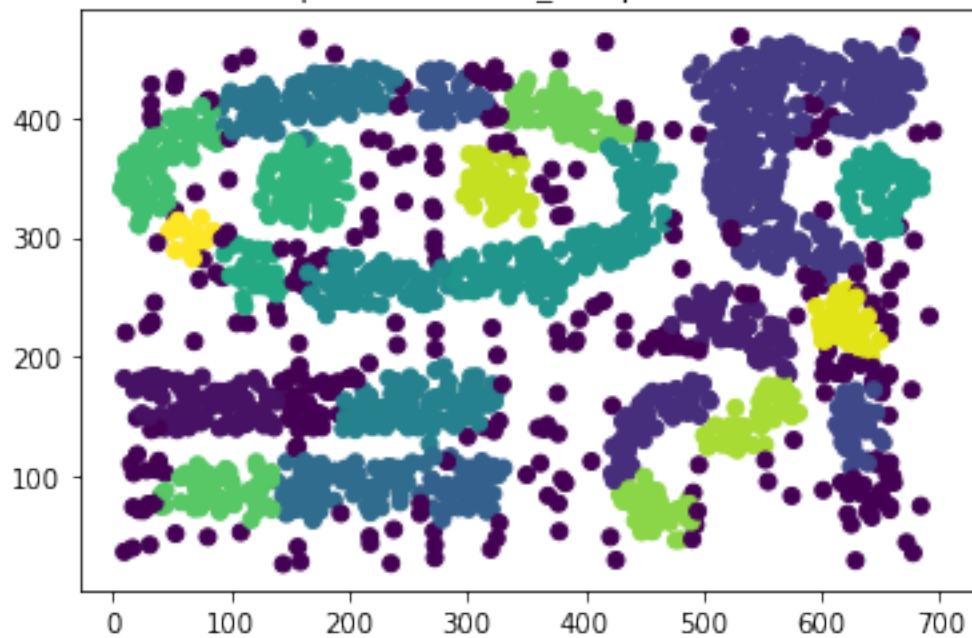


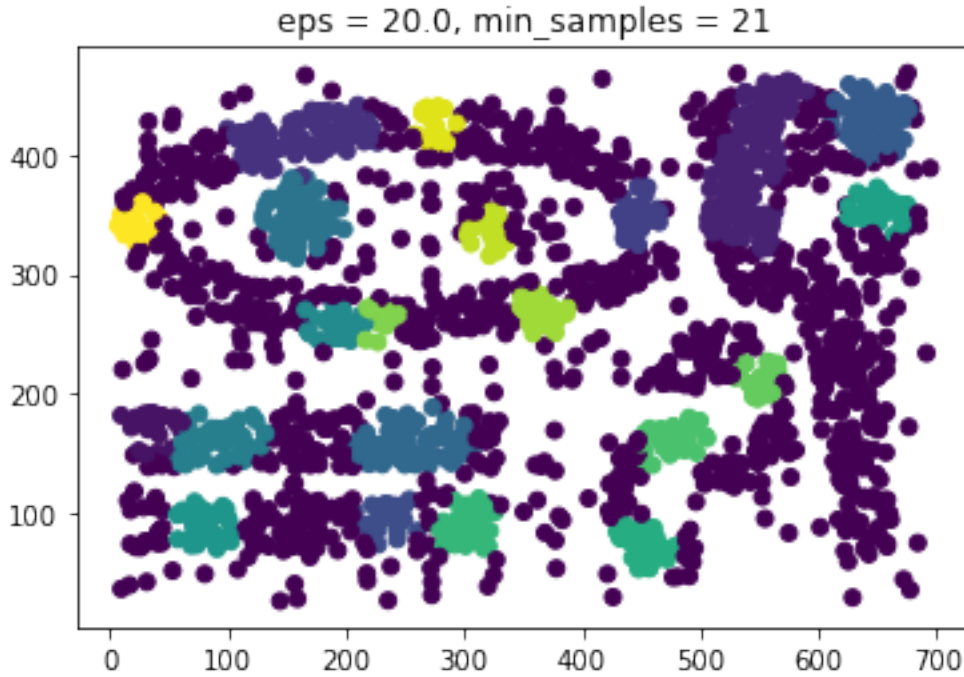


eps = 20.0, min_samples = 11



eps = 20.0, min_samples = 16





The result provided in (b) seems to be the best. However, some solutions in (c) can be considered good as well. especially $\text{eps}=15$ and $\text{min_samples}=6$ (unsurprisingly as this is quite close to the parameters used in (b)), $\text{eps}=20$ and $\text{min_samples}=6$ and $\text{eps}=20$ and $\text{min_samples}=11$. Generally speaking we get an increased number of clusters with increasing min_samples , which peaks at some number, after which no points in the neighborhood of another point can be found. Eps on the other hand seems to have a severe influence on the size of the obtained clusters. Which is backed up by the documentation which describes eps as [The maximum distance between two samples for one to be considered as in the neighborhood of the other](#)