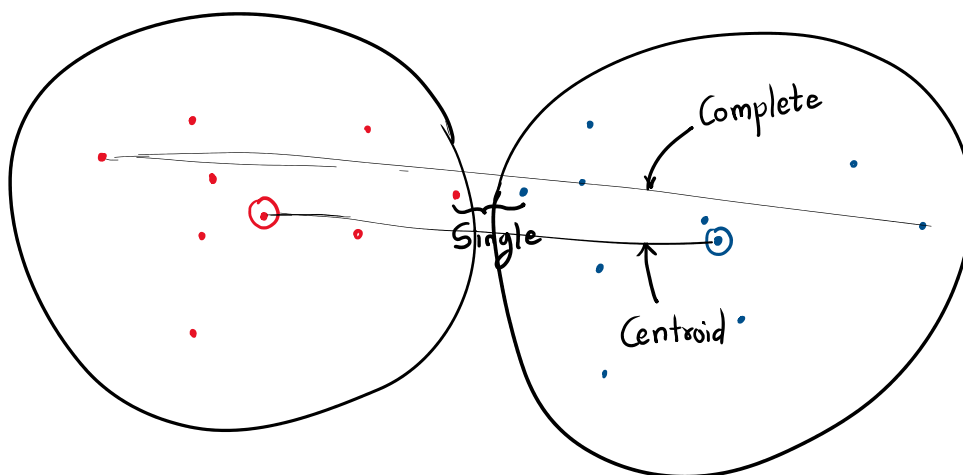
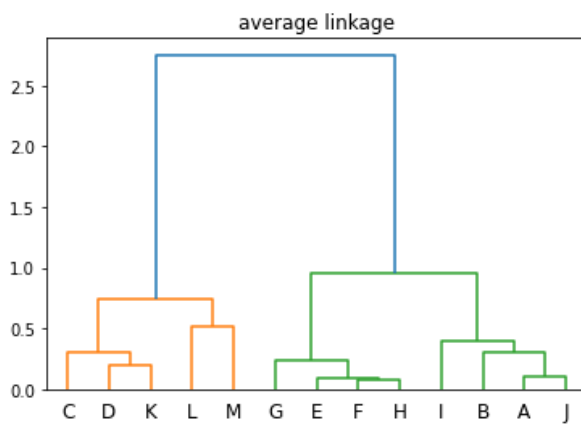
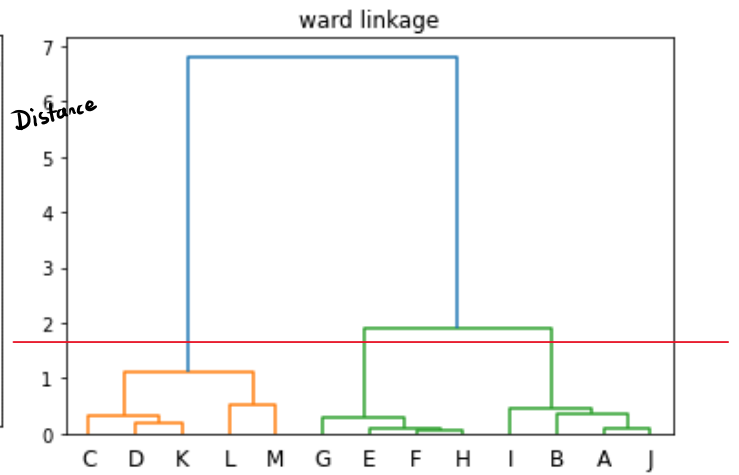
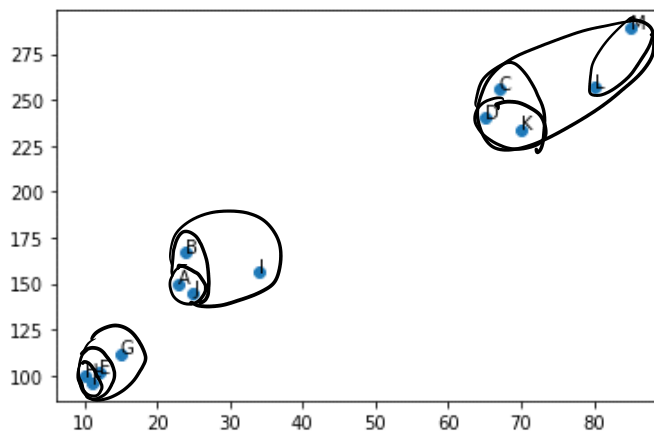


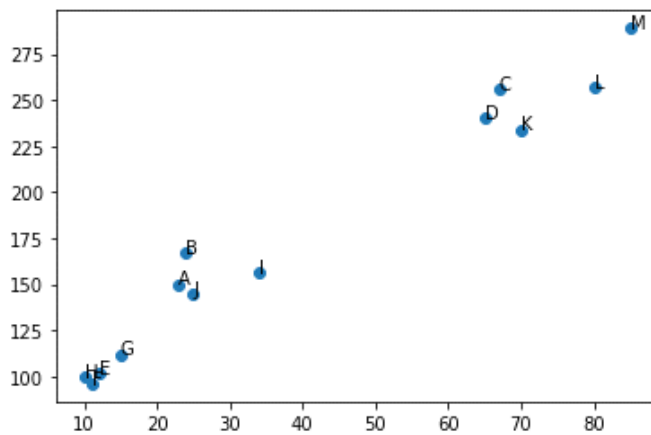
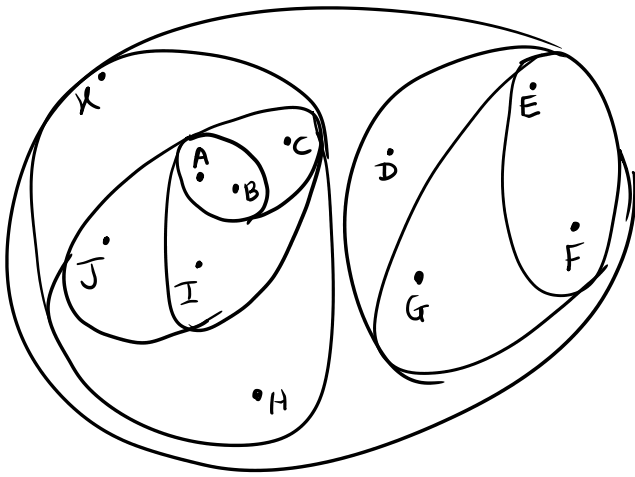
# Clustering

06 April 2024 14:06



ward : based on ANOVA

Average

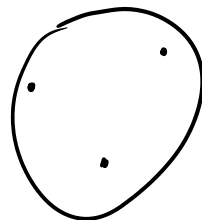


$$\text{sil coef} = \frac{b-a}{\max(a,b)}$$

b: avg nearest cluster distance

a: avg intra-cluster distance

$$\text{Sil score} = \text{avg} \left( \text{Sil coef of all point} \right)$$



Datasets / milk.csv

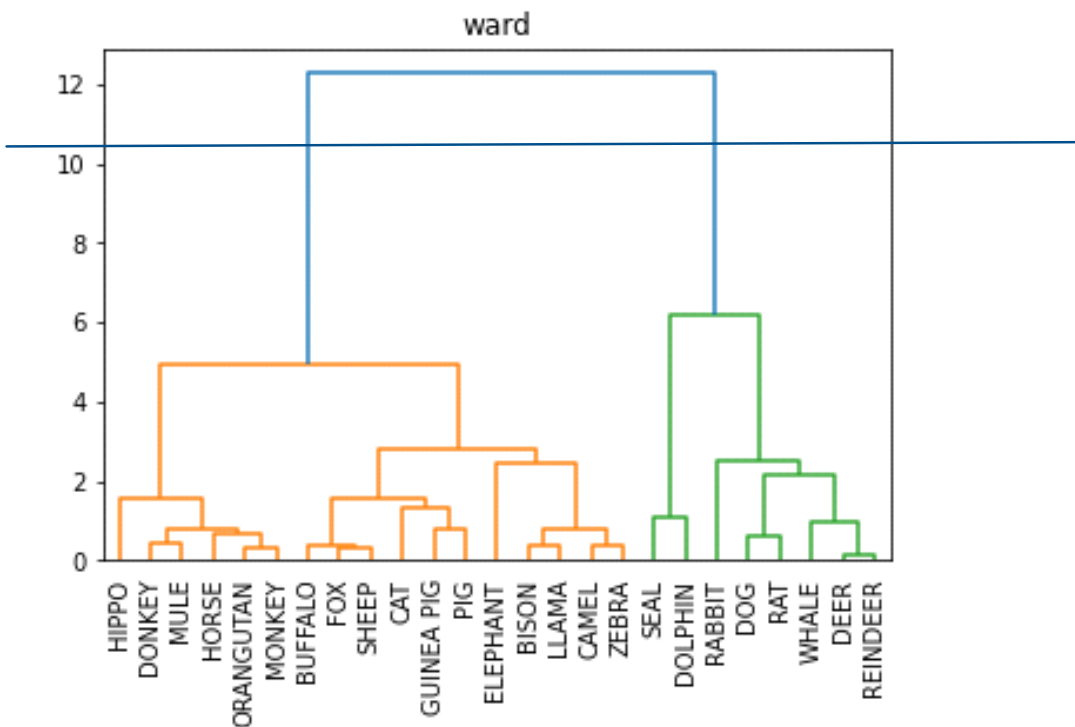
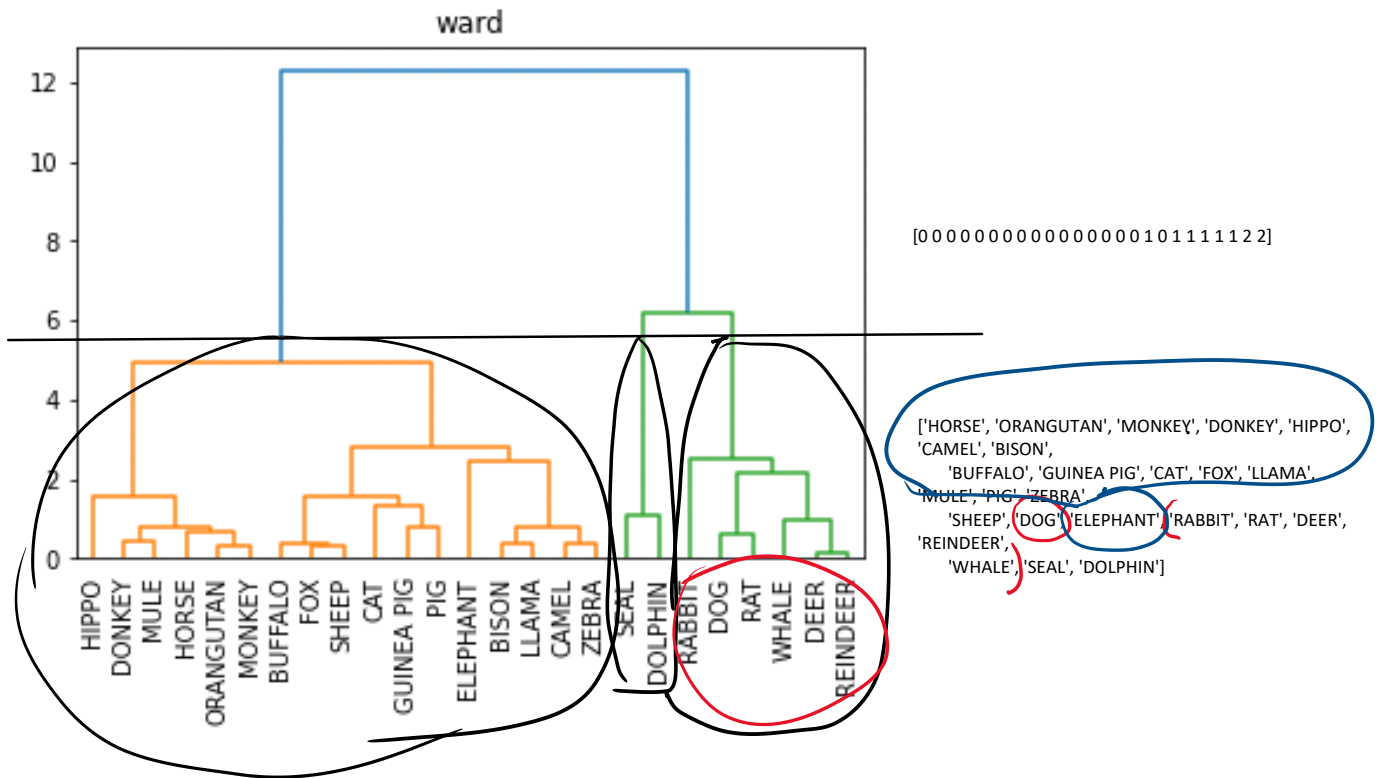
index\_col = 0

1) Dendrogram

└─ clusters ─ 2 3 4 5 6

1> Dendrogram

2> Try clusters = 2, 3, 4, 5, 6



Datasets / nutrient.csv

index\_col = 0

1) Dendrogram

2) Try clusters = 2, 3, 4, 5, 6

Datasets / Protein.csv

index\_col = 0

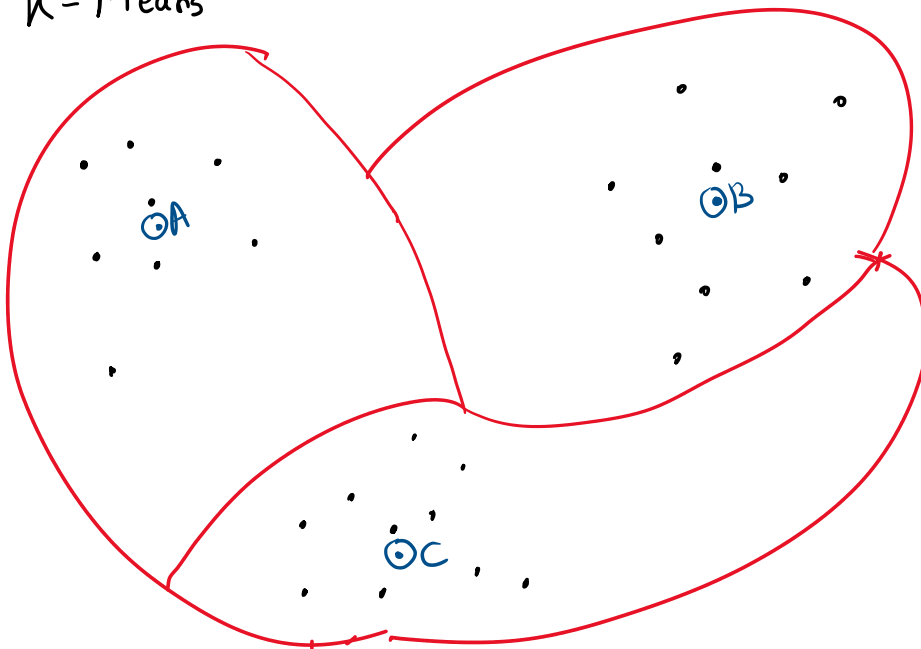
1) Dendrogram

2) Try clusters = 2, 3, 4, 5, 6

$x_1, y_1$        $x_2, y_2$

$$\begin{matrix} \odot \\ \cdot \end{matrix} \begin{matrix} x_3, y_3 \\ \cdot \end{matrix} \begin{matrix} \odot \\ \cdot \end{matrix} \begin{matrix} x_4, y_4 \\ \cdot \end{matrix} \left( \begin{matrix} x_1 + x_2 + x_3 \\ + x_4 \\ 4 \end{matrix}, \begin{matrix} y_1 + y_2 + y_3 \\ + y_4 \\ 4 \end{matrix} \right)$$

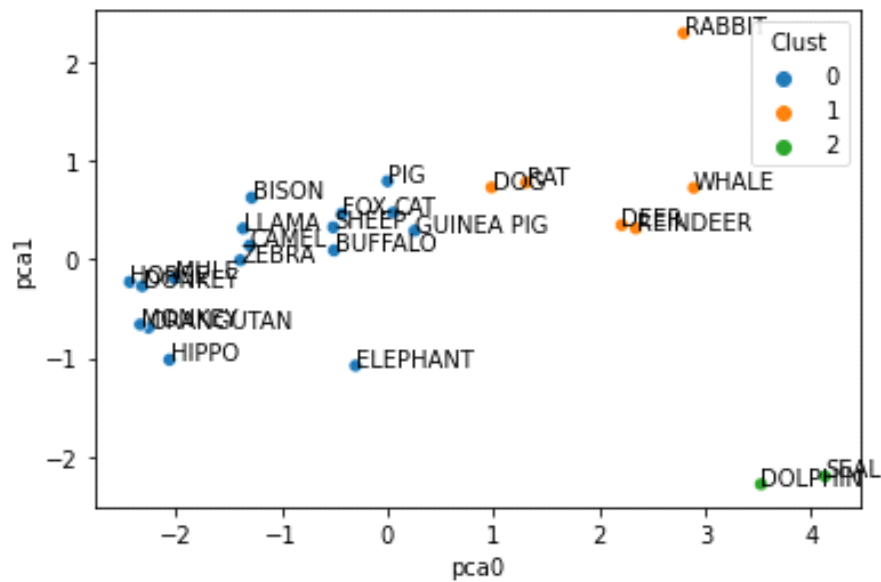
K-Means





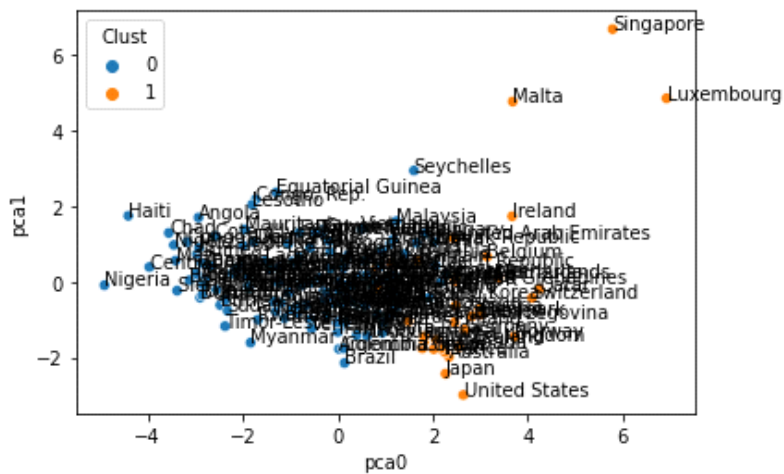
	X1	X2
Name		
A	-0.620863	-0.411108
B	-0.584506	-0.154456
C	0.978838	1.189194
D	0.906124	0.947639
E	-1.020788	-1.135773
F	-1.057145	-1.226356
G	-0.911718	-0.984801
H	-1.093502	-1.165967
I	-0.220938	-0.320525
J	-0.548149	-0.486594
K	1.087909	0.857056
L	1.451477	1.204291
M	1.633261	1.687401

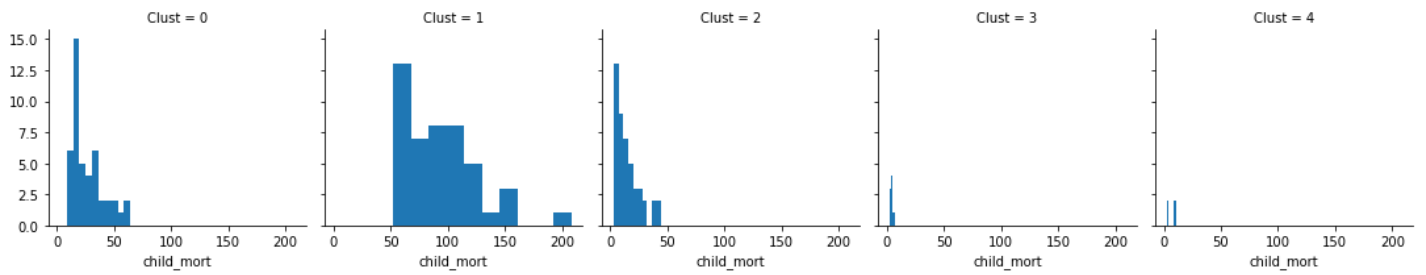
```
In [15]: print(clust.labels_)
[0 0 1 1 0 0 0 0 0 0 1 1 1]
```



Ward's Distance

<https://www.statisticshowto.com/wards-method/#:~:text=Calculate%20the%20distance%20between%20each,of%20squares%20from%20Step%204.>





```
In [97]: clust_data.groupby('Clust')['child_mort'].mean()
Out[97]:
Clust
0    26.917778
1    94.180435
2    15.160870
3     4.300000
4     6.937500
```

