# * Hierarchical Clustering [HC]
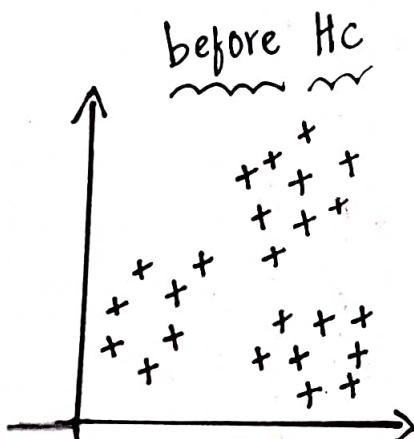
↓
Combine data points which are very close by distance

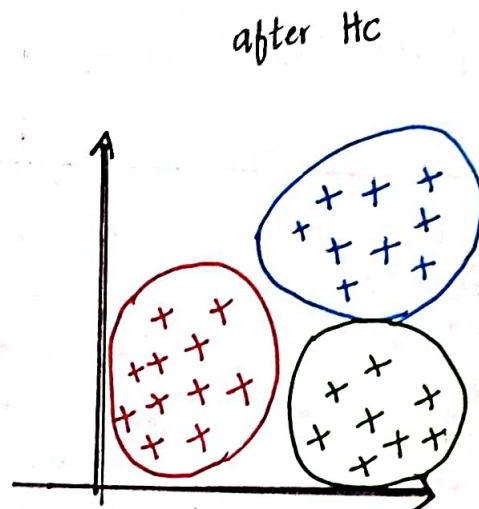Ex:-   Dogs

```
        Dogs                    ] ⇒ Hierachy
       /  |  \  \
   Kids friendly Security Professional
    /\    /\      /\         /\
```

* Hierarchical clustering algorithms build a hierarchy of clusters where " Each node is a Cluster " Consists of clusters of its daughter nodes.
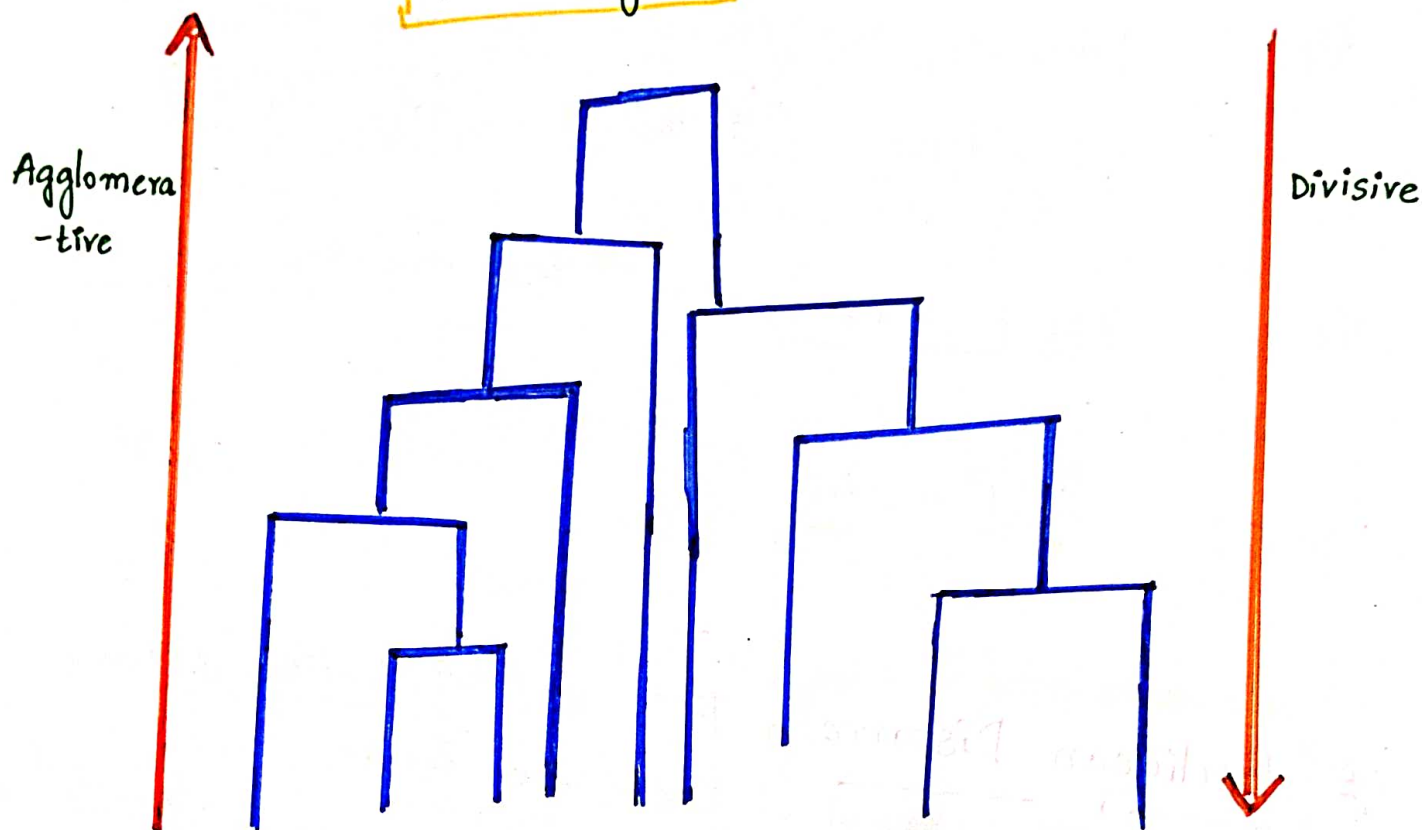
What "HC" does ?

before HC                              after Hc



HC ➡

# Same as k-means but different process.

# * Hierarchical clustering

Agglomera
-tive

Divisive



* They are two types clustering.

    * Algomerative ( From bottom to top)

    * Divisive ( From Top to bottom)

Most of
The Time, we
use, This
Techinque.

→ Steps of Hierachical clustering of (Algomerative)

STEP : 1  Make Each data point a single point cluster
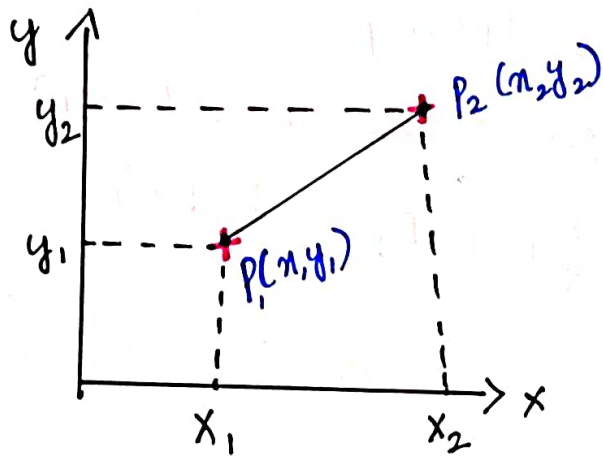(# that forms N clusters)

**Step: 2**    Take the two closest data points and Make them

One cluster   (# that forms N-1 clusters)

**Step: 3**    Take the two closest clusters and make them

One cluster   (# that forms N-2 clusters)

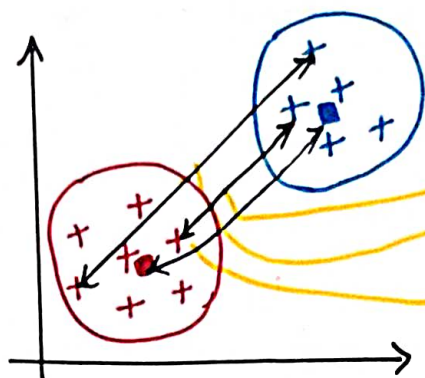**Step: 4**    Repeat Step 3, until there is only one cluster

**Step 5:**    Final model.

\# **Euclidean Distance** :- it is used to identify distance between two points.



$$P_2 (x_2, y_2)$$
$$P(x_1, y_1)$$

*Euclidean distance between $P_1$ and $P_2$

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$
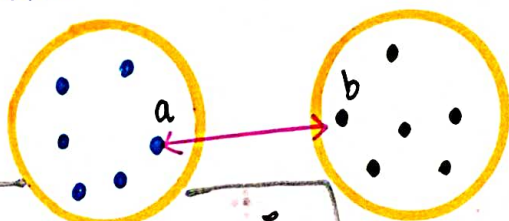
→ Distance between clusters

# Distance b/w two clusters.

- option 1 : closest Points
- option 2 : furthest Points
- option 3 : Average Distance
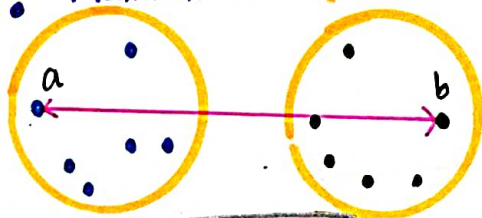- option 4 : Distance between Centroids

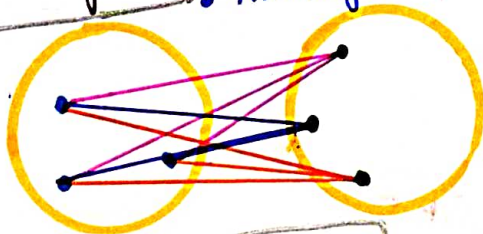* **Single - linkage clustering**
- Minimum distance between Clusters.

* **Complete - linkage clustering**
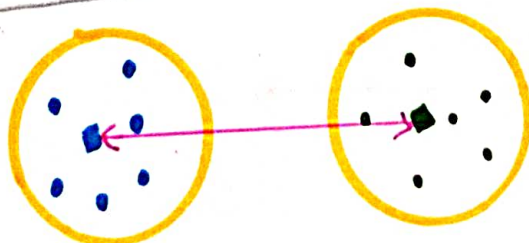- Maximum distance between Clusters

* **Average Linkage clustering**
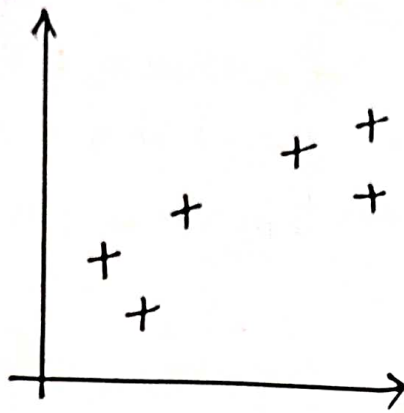- Average distance between clusters

* **Centroid Linkage clustering**
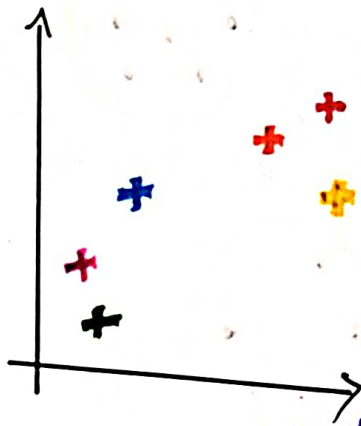- distance between cluster Centroids

# Agglomerative HC

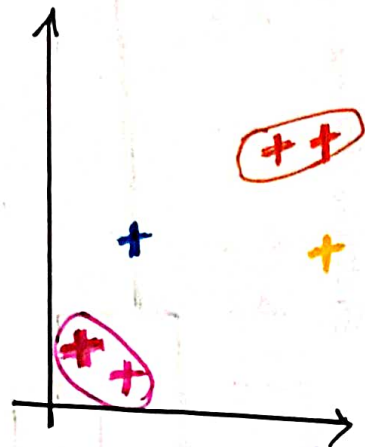Consider the following dataset of N=6 data points



STEP : 1 ( Make Each data point a single point cluster)
that forms 6 clusters)



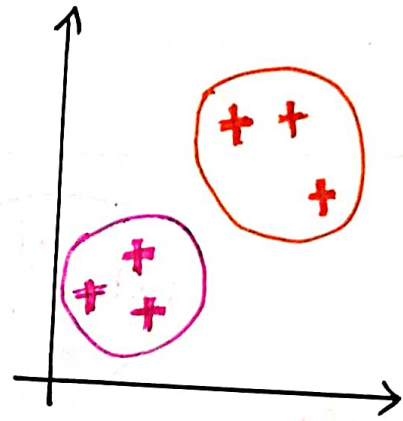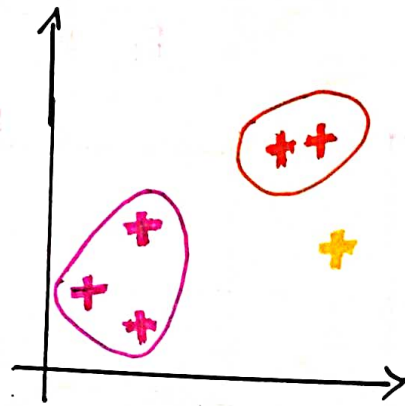STEP : 2 ( Take the two closest data points and make them
one cluster → that forms 5 clusters)
(n-1)

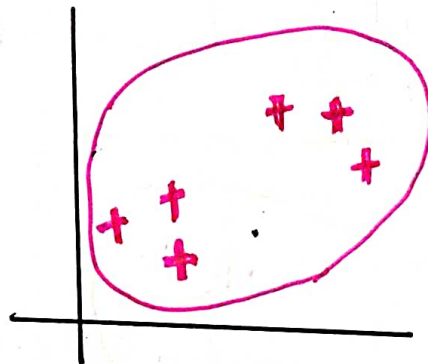: Take the two closest clusters and make them one
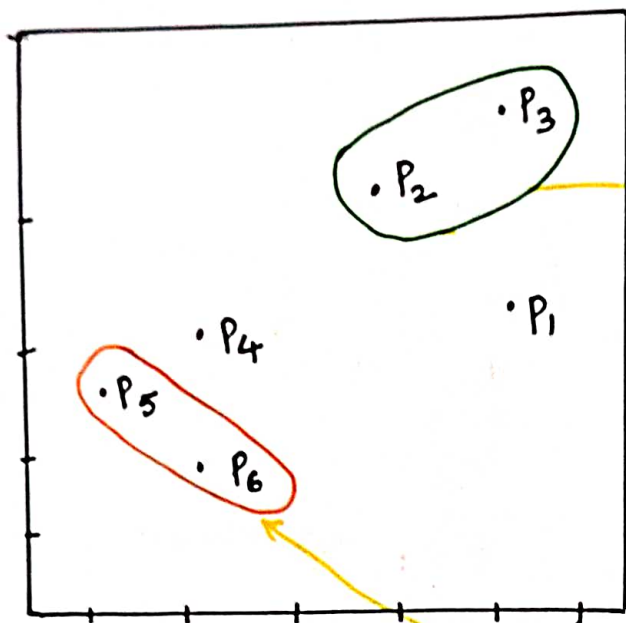
cluster → That forms 4 clusters $(n-2)$



Repeat step 3. until There is only One cluster.



# Final :

# How Do Dendograms Work ?



Euclidean distance

25
2.0
1.5
1.0
0.5
0.0

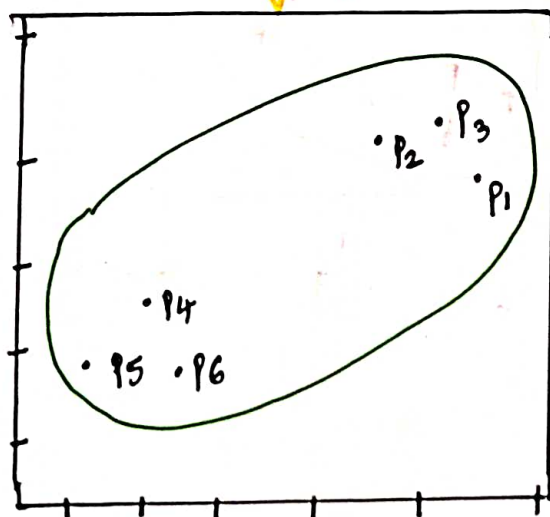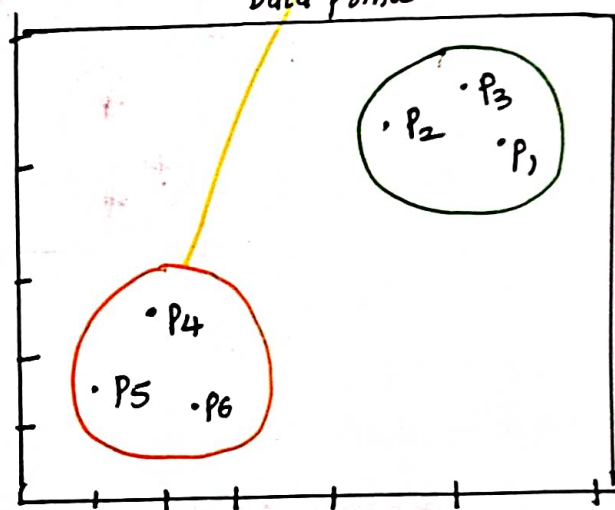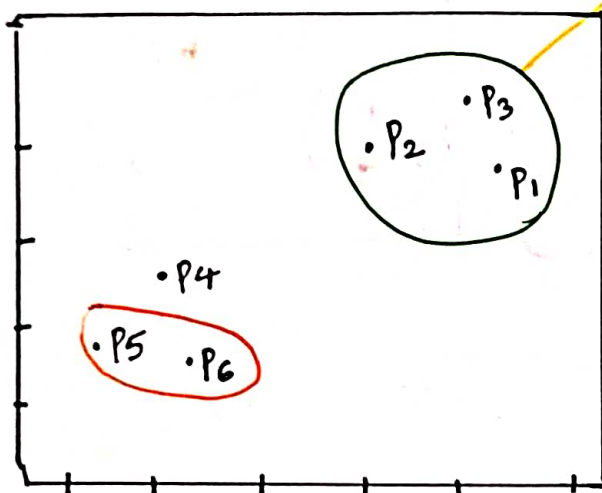P1  P2  P3  P4  P5  P6

Data points

# Example

|      | TO  | OT   | VA   | MO   | WI   | ED   |
|------|-----|------|------|------|------|------|
| TO   | 0   | 351  | 3363 | 505  | 1510 | 2699 |
| OT   |     | 0    | 3543 | 167  | 1676 | 2840 |
| VA   |     |      | 0    | 3690 | 1867 | 819  |
| MO   |     |      |      | 0    | 1824 | 2976 |
| WI   |     |      |      |      | 0    | 1195 |
| ED   |     |      |      |      |      | 0    |

TO  OT  MO  VA  ED  WI

# OT and MO (combined)
# Distances are recalculated again.

|        | TO  | OT/MO | VA   | WI   | ED   |
|--------|-----|-------|------|------|------|
| TO     |     | 351   | 3363 | 1510 | 2699 |
| OT/MO  |     |       | 3543 | 1676 | 2840 |
| VA     |     |       |      | 1867 | 819  |
| WI     |     |       |      |      | 1195 |
| ED     |     |       |      |      |      |

TO OT MO VA ED WI

# OT/MO and TO (combined)
# Distances are recalculated again, to make as one cluster.

| | TO/OT/MO | VA | WI | ED |
|---|---|---|---|---|
| TO/OT/MO | | 3543 | 1676 | 2840 |
| VA | | | 1867 | 819 |
| WI | | | | 1195 |
| ED | | | | |

# ED and VA (combined)

# Distances are recalculated again



TO OT MO    VA ED WI

| | TO/OT/MO | VA/ED | WI |
|---|---|---|---|
| TO/OT/MO | | 2840 | 1676 |
| VA/ED | | | 1667 # value wrong. |
| | | | |

# WI and VA/ED (combined)

# Distances are recalculated again



TO OT MO    VA ED WI

| | TO/OT/MO | VA/ED/WI |
|---|---|---|
| TO/OT/MO | | 1676 |
| VA/ED/WI | | |

# Dendogram :-



TO   OT   MO   VA   ED   WI

Ex:-

World Countries

U.S

(Nato)

Russia
(China)

# if U.S attack russia. it is going to identify which country is going to Support russia.

# which country is directly participate in war.

# which are very close to Each other. it is going to identify.

# this is Hierarchical clustering.

**Code :-** Same data & Same steps till modelling

Using the dendogram to find optimal no. of clusters.

import Scipy.cluster.hierarchy as sch

# dendrogram = Sch. dendrogram (sch. linkage (x, method = "ward"))
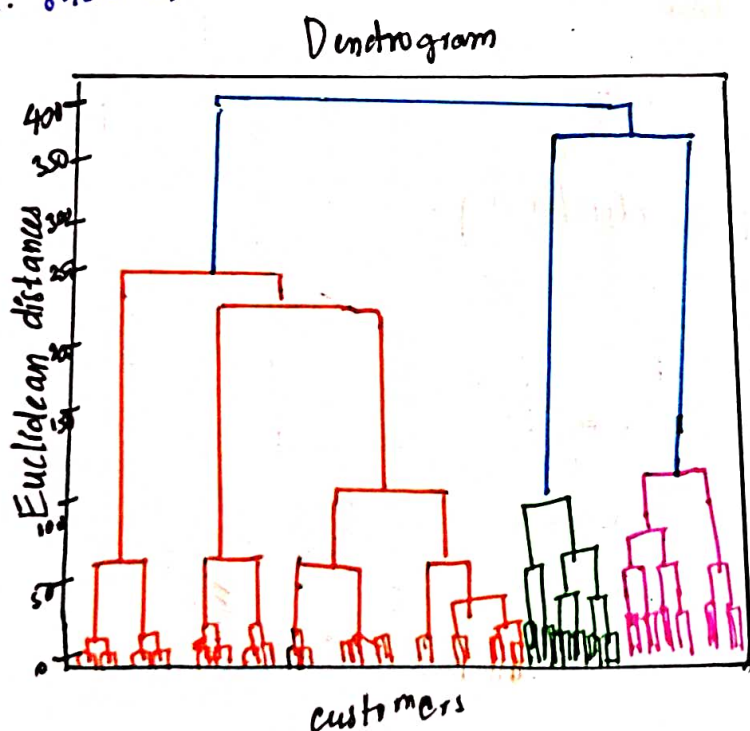
↳ Based on Centroids

# plt. title ("Dendrogram")
# plt. x label (" Customers")
# plt. y label (" Euclidean distances")
# plt. show ()



Dendrogram

Q: How to identify optimal clusters durring dendrogram ?

A :- biggest disadvantage with "HC" is we (dn't identify "optimal no. of clusters". only in K-means we identify.

> ### Hierarchical clustering model :

↗ small

from sklearn. cluster import Agglomerative Clustering.

capital

(no gap)

# hc = Agglomerative Clustering (n_clusters = 5, affinity = "Euclidean", linkage = "ward")

> ### Predict

# y_hc = hc. fit _ predict (x)

# y_hc

Out : array ([4, 3, 4, 3, 4, 3, - - - - - .

1,1,1,1, 1, 1 - - -

0, 2]

# Visualising clusters

```python
# cluster 1
# plt. scatter ( x [y_hc = =0,0], x [y_hc = =0,1], s=100,
                 c= "red", Label = "cluster 1")

# cluster 2
# plt. scatter ( x [y_hc== 1,0] , x [y_hc = = 1,1], s=100,
                 c = "blue", label = "cluster 2")

# cluster 3
# plt. scatter ( x [y_hc ==2,0] , x [y_hc = = 2,1], s=100,
                 c= "green", label = "cluster 3")

# cluster 4
# plt. scatter (n [y_hc = = 3,0], x [y_hc == 3,1], s=100,
                 c= "cyan", Label = "cluster 4")

# cluster 5
# plt. scatter ( x[y_hc = = 4,0] , x [y_hc = = 4,1], s=100,
                 c= "magenta", label= "cluster 5")

# plt. title ( "clusters of customers")
# plt. xlabel ( " Annual Income (k$)")
# plt. ylabel ( " Spending Score (1-100)")
# plt. show ()
        ))
# plt. legend ()
```
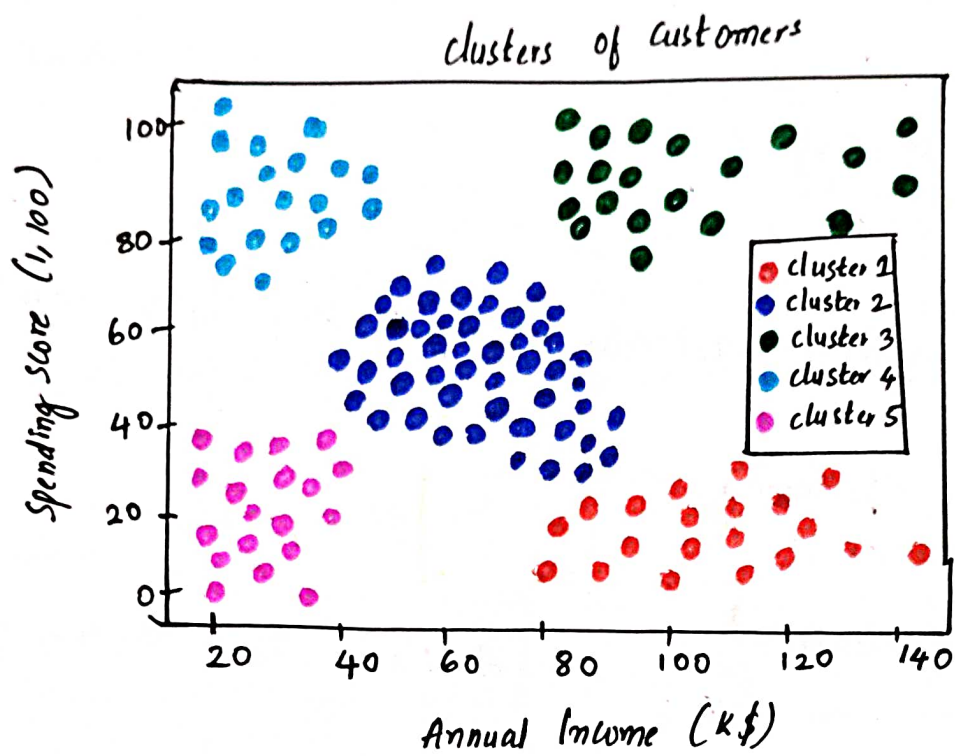
Clusters of Customers

Spending Score (1, 100) vs Annual Income (K$)

Legend:
- cluster 1
- cluster 2
- cluster 3
- cluster 4
- cluster 5

04/05/22

5:00 Am.