Data mining
clustering

---

* it is mainly use for data visualization.
* clustering means groups, grouping, similar properties, similar characteristics. to find homogeneous group
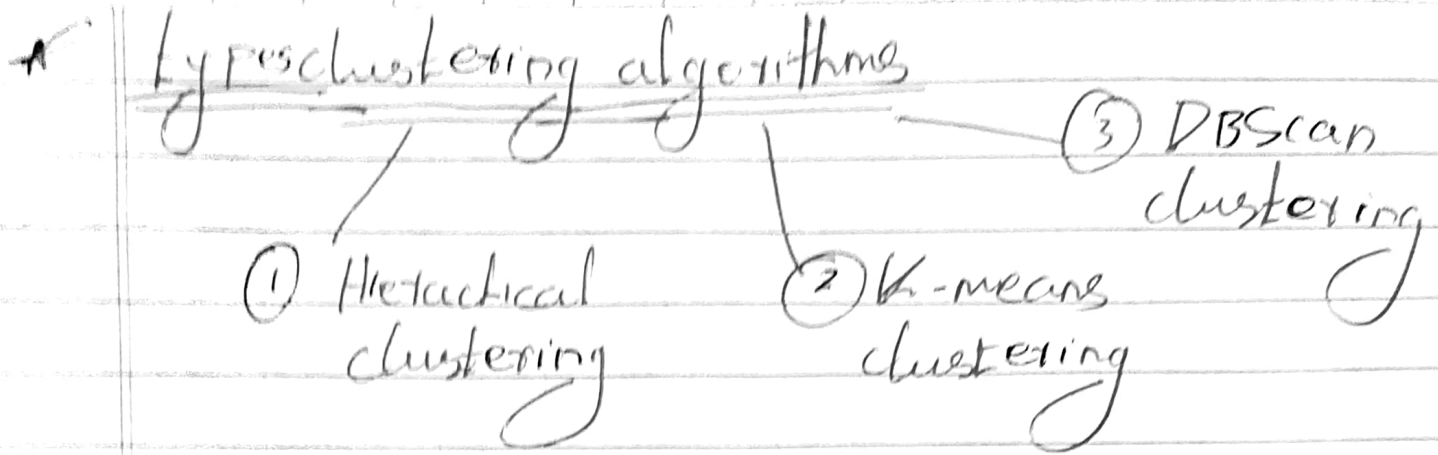
example Suppose
if I have dataframe

| 1 | 12 | 20 |
|---|----|----|
| 2 | 15 | 23 |
| 3 | 16 | 43 |
| 4 | 17 | 60 |

Now I have to group the datapoints
So how we will do this by techique
called clustering

---
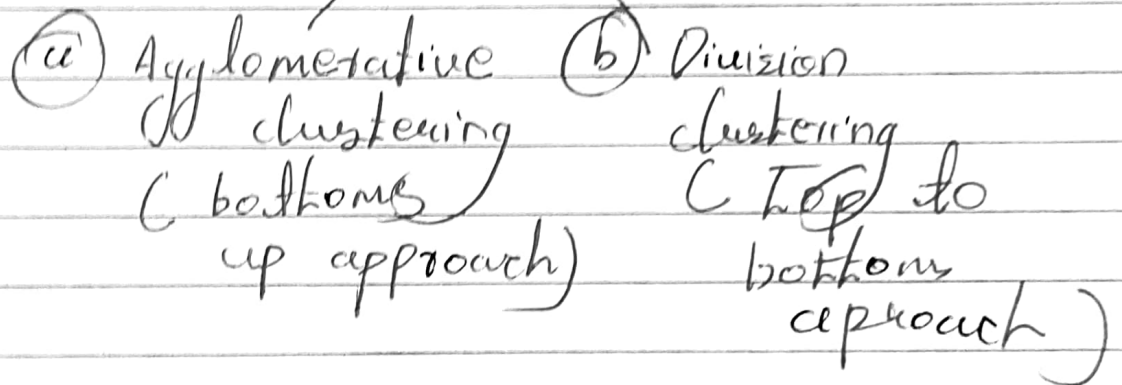
* how to find a good cluster

① Similar records should belong to the same cluster

② Dissimilar records should belong to different cluster

# Types clustering algorithms

③ DBScan clustering

① Hierarchical clustering

② K-means clustering

---

① Hierarchical clustering Algorithms

(a) Agglomerative clustering (bottom up approach)
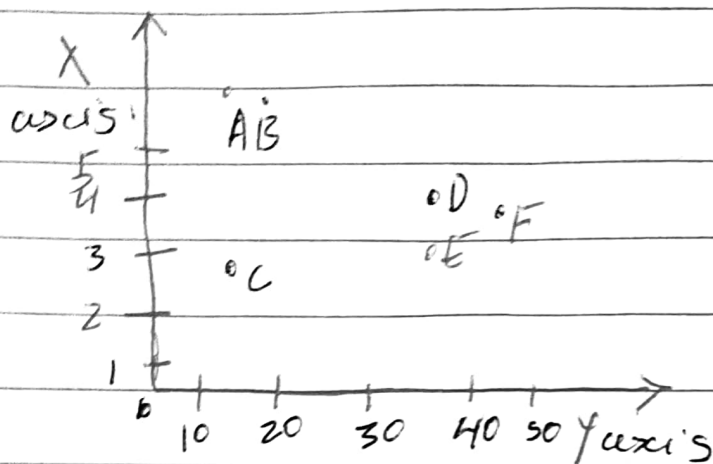
(b) Division clustering (Top to bottom approach)

---

(a) Agglomerative clustering

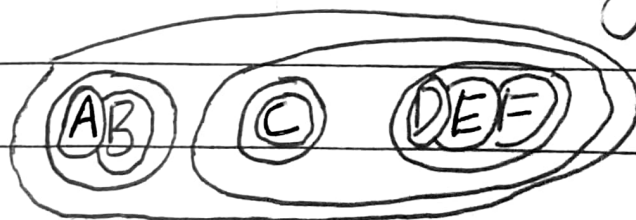To explain Agglomeratic clustering we will take an example.
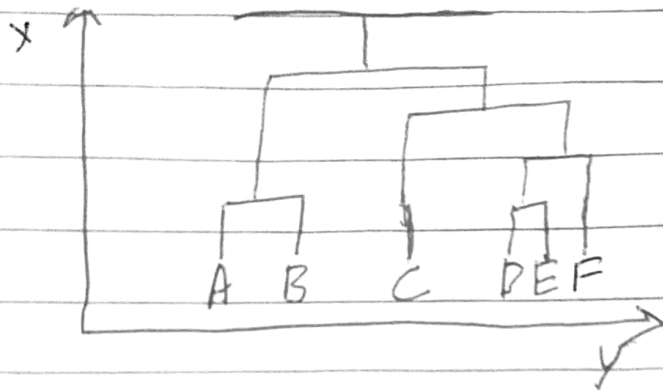
let say we have datapoints called

A, B, C, D, E, F
Now we will plot the graph.

Aglomerative is an iterative process in first iteration everydata point will be treated as signal cluster and then 2nd iteration whoever is near of datapoint will be it will group them suppose in 1st iteration ABCDEF will be signal cluster then in 2nd iteration it found AB and DEF is close to each other so they will group them and lastly c is to the cluster of DEF so it make clucter with them and lastly the cluster of AB and DEFC will merge together will form one cluster to represent this cluster for better data visucalisation we will use always dendogram.
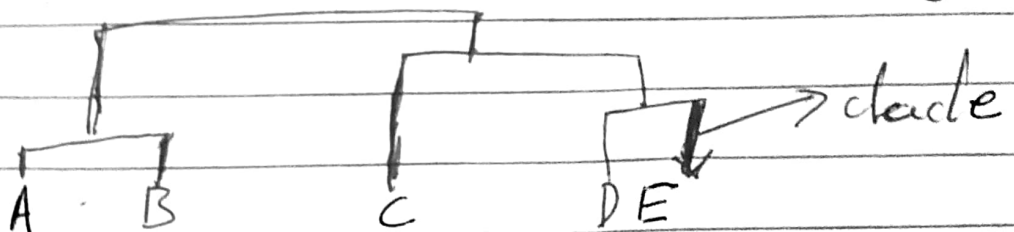
## Datogram visualization.



## cutting the detogram

before going to this topic we will
learn 2 New word called
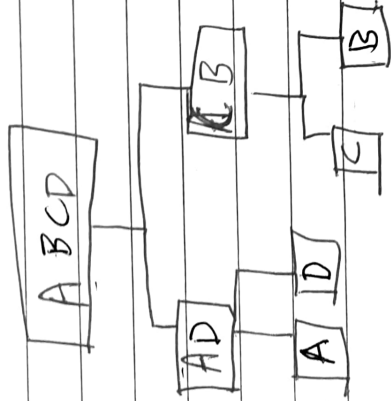1. clade ( intra)
2. intracluster distance.

① clade — its branch in dextogram which
represents the similarity between the
datapoint it tells you the similarity



when the distance is more then there
will be disimilarity between datapoint
So which is the best cluster is AB because
the clade measure is less.

# ⓑ Division clustering

Its same as Agglomerative clustering but instead of going bottom up approach it use top down approach it means firstly it calculate the cluster as whole then it will divide the cluster which has longest distand and iterative process goes on till it reaches singal clusterpoint

```
        +--------+
        | A B C D|
        +--------+
         /      \
     +----+     +----+
     | A D|     | C B|
     +----+     +----+
      /  \       /  \
   +--+ +--+  +--+ +--+
   | A| | D|  | C| | B|
   +--+ +--+  +--+ +--+
```

① "intracluster distance"
→ its distance between datapoint
of cluster and different cluster
for example



intracluster distance

So intracluster distance should always
grows simultaneously when its going
bottoms up otherwise its really forced
dedogram.

A Now comes the part of cutting the
dedogram in order to get best
cluster (Note: we have to do this
visually)

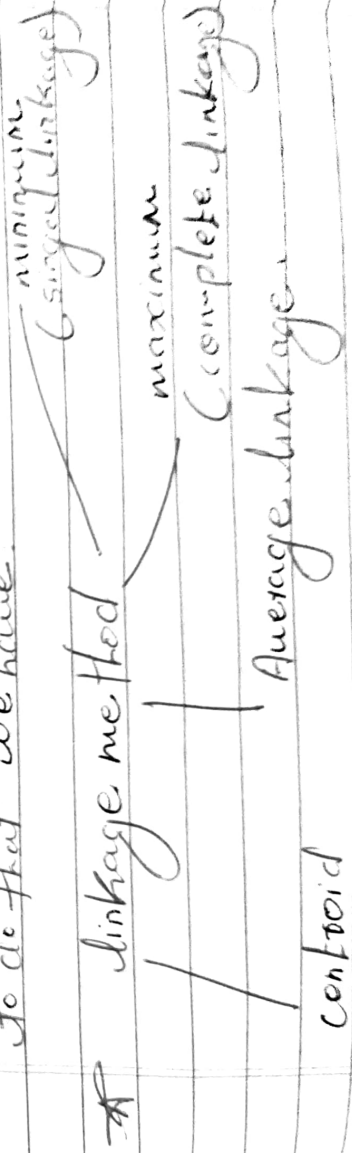lets take example of same
previous dataset.

cutting
the
dendo.
gram

A B c D E F

* Euclidean Distance $\left(d_{ij} = \sqrt{(x_{ip} - x_{jp})^2}\right)$
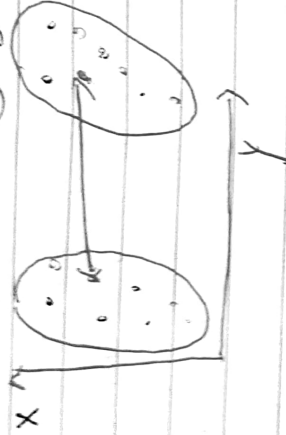
$d_{ij}$ = distance between observction
of $i$ and $j$

* it length of a line segment between two
points.

So how will we measure the length
who how should be the method. inorder
to do that we have.

* linkage method

minimum
(single linkage)

maximum
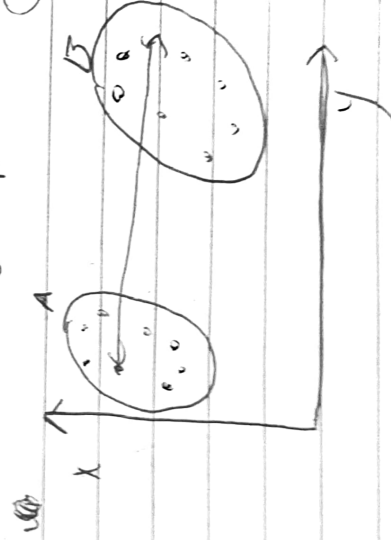(complete linkage)

Average linkage

contvoid
linkage

① minimum (singel linkage)



in singel linkage suppose there are two cluster called A and B and they have a datapoint which is nearest to the other cluster of datapoint so that segment we can use in Euclidean distance.
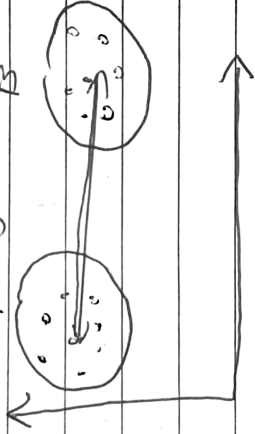
② maximum (complete linkage)



the cluster is same as the above example. Just difference instean instead of mini mum datapoints as segment we took maximum datapoint distance.

③ Average linkage

we will take same graph example
as previous Now just the different is
we will calculate Average distance between
them and then we will use in euclidian
distance

④ centroid linkage



So every cluster has centre point we will
take that point segment of each cluster
and take it as a segment for two calcul
ate euclidian distance

Point to remember

A) Standardization and Normalization
in terms of clustering no while
calculating Euclidean distance.

So whenever we have dataset so the
unit of each coloum will be different
So this will effect the measurement
of Euclidean distance.

So solution is

1) when dataset is normally distributed
we will use standardization (ut
wibe positive or negative values).

2) when dataset is not Normally distributed
So we will use Normalization in our
dataset (ut will always be 0's and 1s
values).

# Distance for Binary data

Suppose what if we have categories data/get columns and we have to find ate euclidian distance

Lets take an example.

| | Married | Smoker. | marriages |
|---|---|---|---|
| S | Yes | yes | yes |
| A | No | yes | No |
| B | Yes | No | yes |

So we have x variable called S,A,B.let assume they are person. and these column of married Smoker marriages they they have got yes or no. Now we will put this above dataset in 2X2 matrix (columns).



So above matrices is especial for B person.
So he have got 1 Nos and 2 yes

So the measure will be

① Binary Euclidean Distance

$$(b+c)/(a+b+c+d)$$

② Simple matching coefficient

$$(a+d)/(a+b+c+d)$$

⑤ Jaquard coefficient.

$$d/(b+c+d).$$

★ For ( Numerical + categorical ) data

Step ① Normalize the data to [0,1]

use

Step ② Gower's General Dissimilarity coefficient.