

# Heirarchical Clustering &

## DBSCAN

### [Clustering]

→ 2 new ideas !!

→ Customer segmentation case study.

## 1st big idea:

- ↳ Use geometric intuition to design a metric that can quantify clusters.
- Dunn, WCSS, Silhouette ✓
- ↳ Use an optimisation setup to opt above metric ✗

## 2nd big idea

- ↳ Randomly initialise centers
- assign clusters

→ update centers  
→ repeat!  
→ Kmeans ↪ Lloyd's algo.

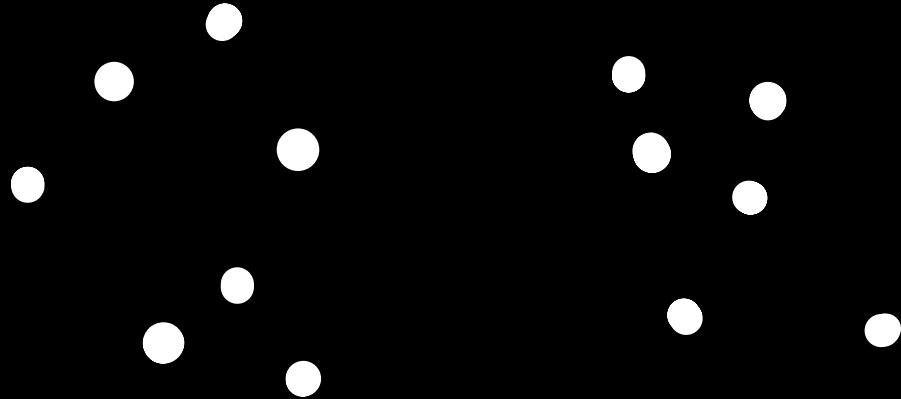
### Extension

↪ smart initialisation  
→ Kmeans ++  
→ code

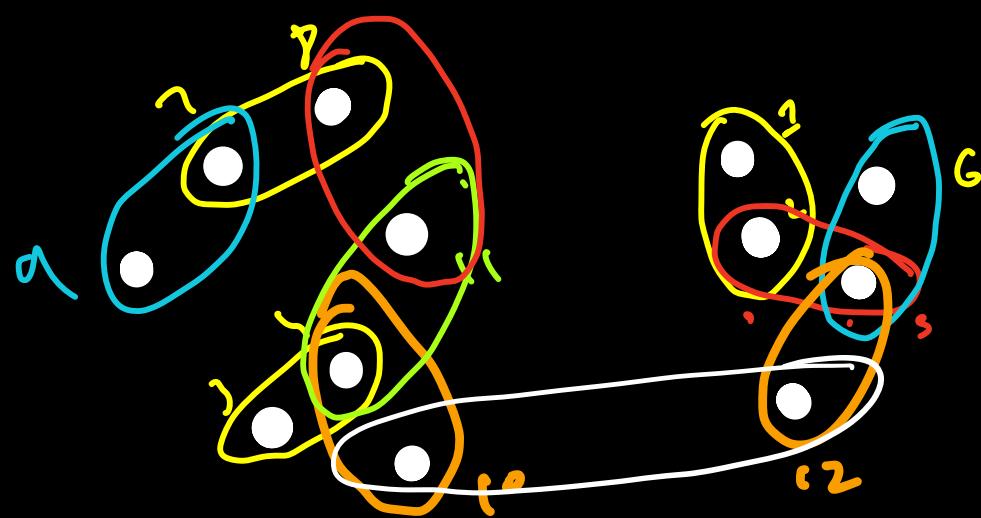
### Next ideas:

→ Hierarchical clustering  
→ DBSCAN  
→ GMM

## Hierarchical Clustering

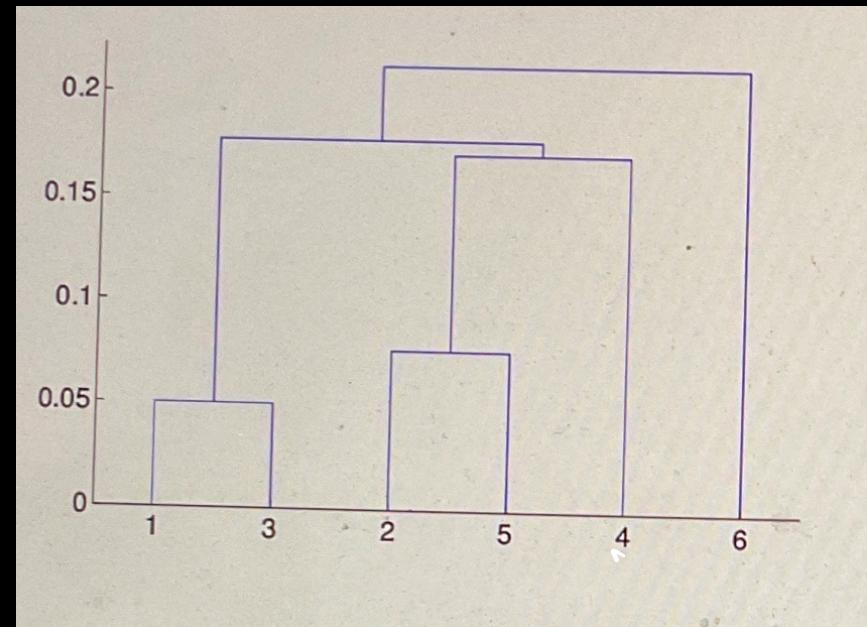
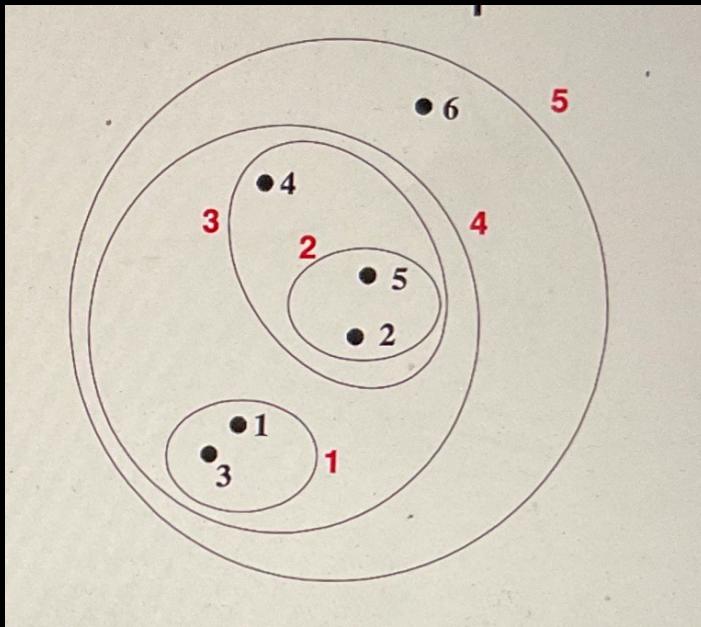


- Compute some form of distance b/w all points.
- combine 2 pts with min distance
- treat the sub-clusters as pts and repeat till only one cluster remains



(another name)

→ Agglomerative clustering:



points

dendrogram

distance b/w clusters:  
→ min

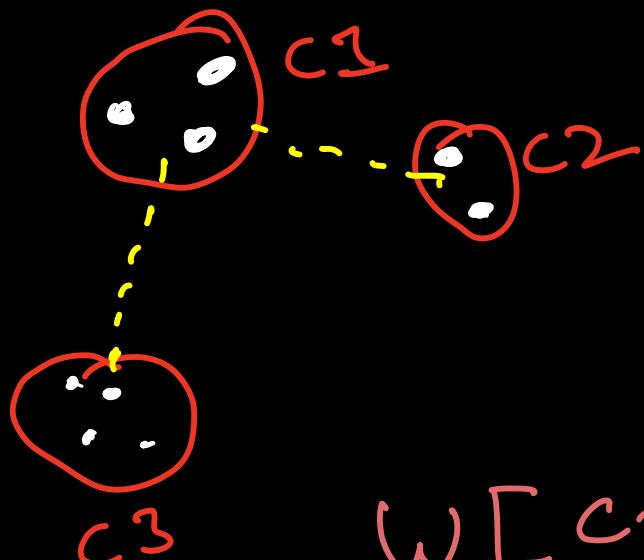
→ animation

→ max  
→ avg

} pros and cons in post method.

→ any other func<sup>n</sup>, eg: squared, cosine etc..  
→ words method

## Word's distance / Word's linkage



Using avg dist

$$\underline{\underline{c_1 + c_2}}, c_3$$

Word's distance

$$w[c_1, c_2] =$$

$$wcss(c_1 + c_2) - wcss(c_1) + wcss(c_2)$$

How much will  $\text{ESS} / \text{WCSS}$  increase  
↓ ↓  
when I merge 2 clusters.

- we will merge clusters where there is min increase in  $\text{ESS} / \text{WCSS}$
- very effective metric in practice.

Pros:

- no need to pre-decide 'K'
- hierarchical may actually be present in real world

Cons:

- Sensitive to choice of linkage form.
- Offline only. ← Can only be used for analysis

Term:

Proximity matrix:

→ code

distance  $(c_i, c_j)$

	$c_1$	$c_2$	$c_3$	$\dots$	$c_n$
$c_1$	0	0.25	1.3	.	.
$c_2$	0.31	0	.	.	.
$c_3$	.	.	.	.	0
$c_n$	.	.	.	.	0