

L21: Clustering

Voronoi Diagrams → Assignment-based Clustering

Mar 30, 2026

FODA

Jeff M. Phillips



What is Clustering

$$X \subset \Omega$$

• n objects

$$\text{a set } X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^d$$

• distance function $d: \Omega \times \Omega \rightarrow \mathbb{R}_{\geq 0}$

trust distance
as good modeling.

$$d(x_i, x_j) = \|x_i - x_j\|$$

given as input

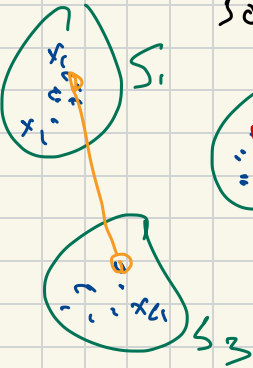
Goal: group objects into k sets

ex: $k=3$

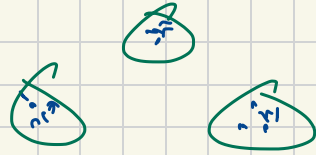
so:

(1) objects in same set \Rightarrow $d(x_i, x_j)$ small.
 x_i, x_j

(2) objects in different sets \Rightarrow $d(x_i, x_j)$ large.
 x_i, x_j



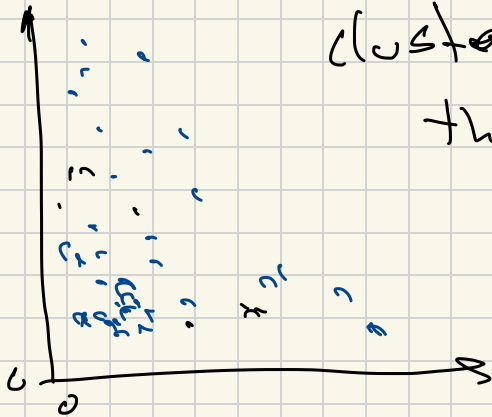
Clusterability



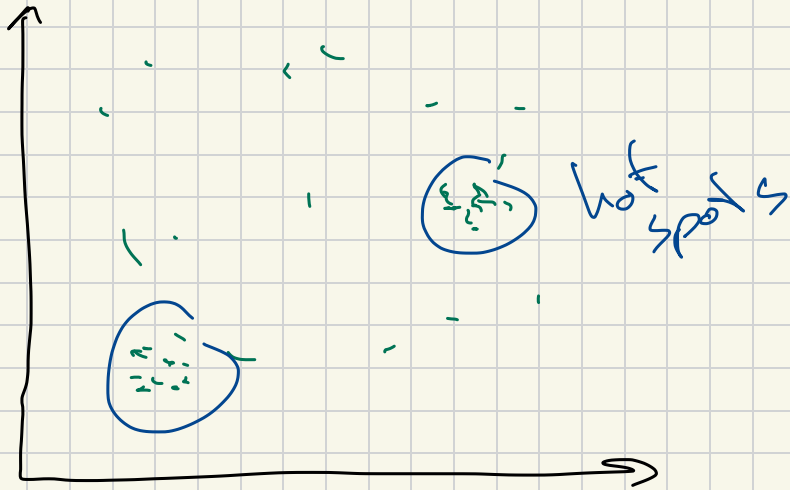
- When data is easy/naturally clusterable, most clustering algos. work well.

- When data is not easy/naturally clusterable, the no magical algorithm that will find a good clustering.

"smear"



"clustering"
in spatial data
on a map



different

Decompose X
in k sets
each $x_i \in X$
in exactly 1 S_j

Clustering cost function (K-means)

Input data X , $d(x, x') = \|x - x'\|$, value $k > 1$

$$\text{cost}_2(X, S) = \sum_{x_i \in X} \|x_i - \phi_S(x_i)\|^2$$

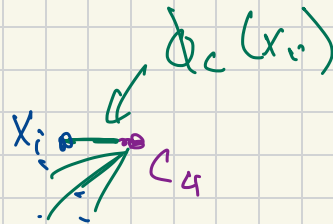
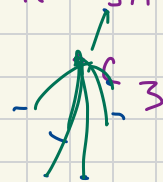
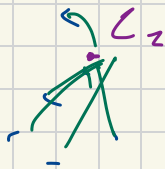
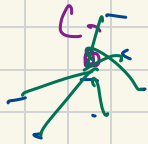
$$S = \{s_1, s_2, \dots, s_k\}$$
$$s_j \subset X$$

↑ similar to TF
(closest point on
subspace F)

Sites $C = \{c_1, c_2, \dots, c_k\} \subset \mathbb{R}^d$

$$\phi_C(x) = \underset{c_j \in C}{\text{argmin}} \|x - c_j\|$$

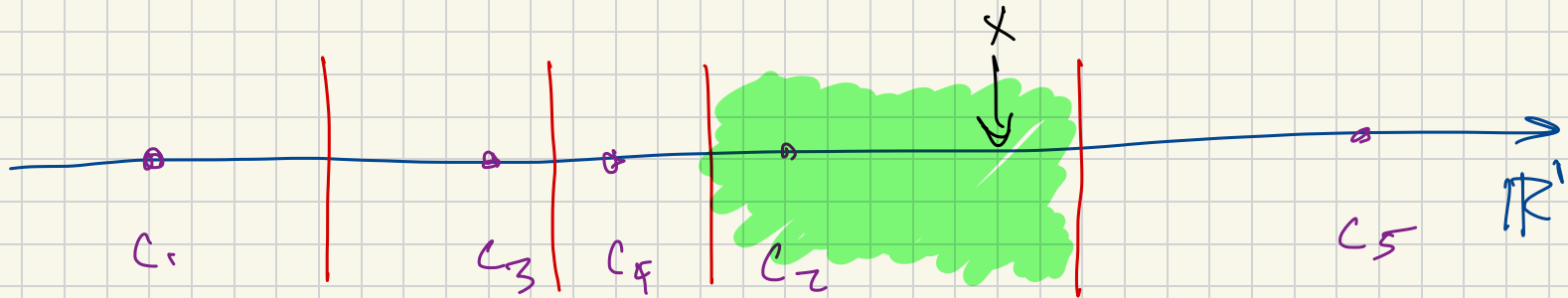
← sites, centers



Voronoi Diagram

for set $C \subset \mathbb{R}^d$ $|C|=k$

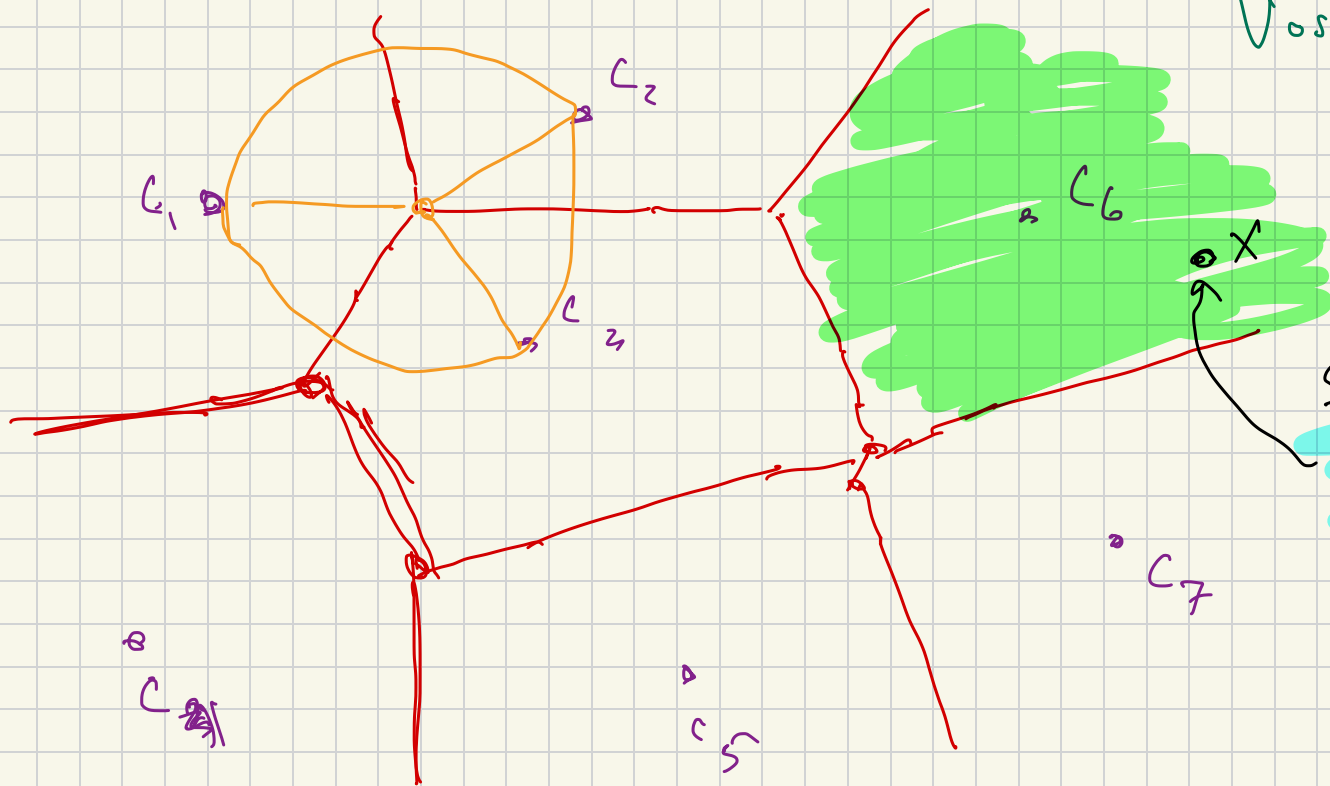
Structure $\Phi_C(x) = \operatorname{argmin}_{c_j \in C} \|x - c_j\|$



$$\{\Phi_C(x) = c_2\}$$

Voronoi cell
of c_2

Voronoi Diagram in \mathbb{R}^2



Voronoi cell
of C_1

Solve $\Phi(x)$
in $O(\log k)$
time in
 \mathbb{R}^2

Voronoi Diagrams in \mathbb{R}^d

size $\approx k^{\lceil d/2 \rceil}$

intractable

Course of Dimensionality

Compute $\phi_C(x) = \operatorname{argmin}_{C_j \in C} \|x - C_j\|$

$O(k)$
time

```
 $C^* = C_1$       $j^* = 1$   
for  $j = 2$  to  $k$   
  if  $(\|x - C_j\| < \|x - C_{j^*}\|)$   
     $j^* = j$   
return  $j^*$ 
```

Assignment-based Clustering

$$X \subset \mathbb{R}^d \quad |C| = k$$

→ k-means

most common famous

$$C = \operatorname{arg\,min}_{|C|=k} \sum_{i=1}^n \|x_i - \phi_C(x_i)\|^2$$

↳ soft clusters split decision

→ k-center

Gonzalez's

$$C = \operatorname{arg\,min}_{|C|=k}$$

$$\max_{x_i \in X} \|x_i - \phi_C(x_i)\|$$

MoG

→ k-median

most resistant to outliers

$$C = \operatorname{arg\,min}_{|C|=k}$$

$$\sum_{x_i \in X} \|x_i - \phi_C(x_i)\|$$

→ k-medoid

good attractor more interpretable

$$C = \operatorname{arg\,min}_{|C|=k, C \subset X}$$

$$\sum_{x_i \in X} \|x_i - \phi_C(x_i)\|$$

