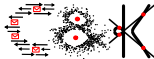
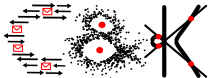


# diameter, radius, discrete radius



$D : M \times M \rightarrow \mathbb{R}$  distance function,  $S \subset M, |S| < \infty$

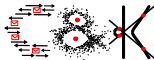
■  $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )





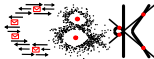
$D : M \times M \rightarrow \mathbb{R}$  distance function,  $S \subset M, |S| < \infty$

- $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )
- $\text{rad}^D(S) := \min_{m \in M} \max_{x \in S} D(x,m)$  (radius of  $S$ )



$D : M \times M \rightarrow \mathbb{R}$  distance function,  $S \subset M, |S| < \infty$

- $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )
- $\text{rad}^D(S) := \min_{m \in M} \max_{x \in S} D(x,m)$  (radius of  $S$ )
- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x,m)$  (discrete radius of  $S$ )

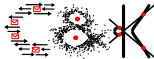


$D : M \times M \rightarrow \mathbb{R}$  distance function,  $S \subset M, |S| < \infty$

- $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )
- $\text{rad}^D(S) := \min_{m \in M} \max_{x \in S} D(x,m)$  (radius of  $S$ )
- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x,m)$  (discrete radius of  $S$ )

$P \subset M, |P| < \infty, \mathcal{C} = \{C_1, \dots, C_k\}$  partition of  $P$

- $\text{cost}_{\text{diam}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{diam}^D(C_i)$  (diameter cost)

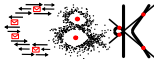


$D : M \times M \rightarrow \mathbb{R}$  distance function,  $S \subset M, |S| < \infty$

- $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )
- $\text{rad}^D(S) := \min_{m \in M} \max_{x \in S} D(x,m)$  (radius of  $S$ )
- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x,m)$  (discrete radius of  $S$ )

$P \subset M, |P| < \infty, \mathcal{C} = \{C_1, \dots, C_k\}$  partition of  $P$

- $\text{cost}_{\text{diam}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{diam}^D(C_i)$  (diameter cost)
- $\text{cost}_{\text{rad}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{rad}^D(C_i)$  (radius cost)



$D : M \times M \rightarrow \mathbb{R}$  distance function,  $S \subset M, |S| < \infty$

- $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )
- $\text{rad}^D(S) := \min_{m \in M} \max_{x \in S} D(x,m)$  (radius of  $S$ )
- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x,m)$  (discrete radius of  $S$ )

$P \subset M, |P| < \infty, \mathcal{C} = \{C_1, \dots, C_k\}$  partition of  $P$

- $\text{cost}_{\text{diam}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{diam}^D(C_i)$  (diameter cost)
- $\text{cost}_{\text{rad}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{rad}^D(C_i)$  (radius cost)
- $\text{cost}_{\text{drad}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{drad}^D(C_i)$  (discrete radius cost)



### Problem 6.1 (diameter $k$ -clustering)

Given a set  $P$ ,  $|P| < \infty$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{diam}}^D(\mathcal{C})$ .



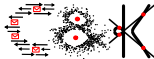
### Problem 6.1 (diameter $k$ -clustering)

Given a set  $P$ ,  $|P| < \infty$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{diam}}^D(\mathcal{C})$ .

### Problem 6.2 (radius $k$ -clustering)

Given a set  $P$ ,  $|P| < \infty$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{rad}}^D(\mathcal{C})$ .





### Problem 6.1 (diameter $k$ -clustering)

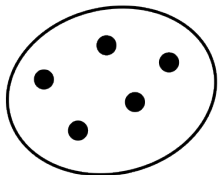
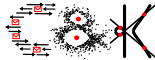
Given a set  $P$ ,  $|P| < \infty$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{diam}}^D(\mathcal{C})$ .

### Problem 6.2 (radius $k$ -clustering)

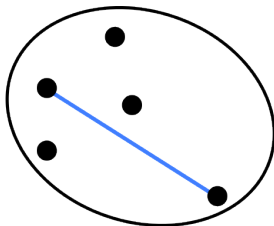
Given a set  $P$ ,  $|P| < \infty$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{rad}}^D(\mathcal{C})$ .

### Problem 6.3 (discrete radius $k$ -clustering)

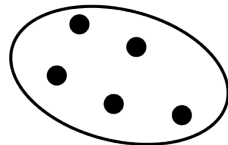
Given a set  $P$ ,  $|P| < \infty$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{drad}}^D(\mathcal{C})$ .



$C_1$

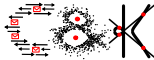


$C_2$



$C_3$

# Agglomerative clustering - setup and idea



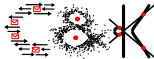
$D : M \times M \rightarrow \mathbb{R}$  distance function,  
 $P \subset M, |P| = n, P = \{p_1, \dots, p_n\}$



$D : M \times M \rightarrow \mathbb{R}$  distance function,  
 $P \subset M, |P| = n, P = \{p_1, \dots, p_n\}$

## Basic idea of agglomerative clustering

- start with  $n$  clusters  $C_i, 1 \leq i \leq n, C_i := \{p_i\}$
- in each step replace two clusters  $C_i, C_j$  that are "closest" by their union  $C_i \cup C_j$
- until single cluster is left.

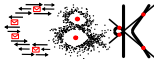


$D : M \times M \rightarrow \mathbb{R}$  distance function,  
 $P \subset M, |P| = n, P = \{p_1, \dots, p_n\}$

## Basic idea of agglomerative clustering

- start with  $n$  clusters  $C_i, 1 \leq i \leq n, C_i := \{p_i\}$
- in each step replace two clusters  $C_i, C_j$  that are "closest" by their union  $C_i \cup C_j$
- until single cluster is left.

**Observation** Computes  $k$ -clustering for  $k = n, \dots, 1$ .

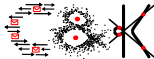


## Definition 6.4

For  $C_1, C_2 \subset M$

$$D_{CL}(C_1, C_2) := \max_{x \in C_1, y \in C_2} D(x, y)$$

*is called the complete linkage cost of  $C_1, C_2$ .*

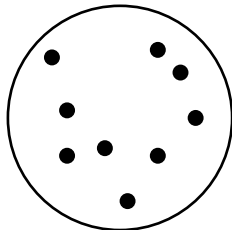
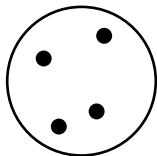


## Definition 6.4

For  $C_1, C_2 \subset M$

$$D_{CL}(C_1, C_2) := \max_{x \in C_1, y \in C_2} D(x, y)$$

is called the complete linkage cost of  $C_1, C_2$ .



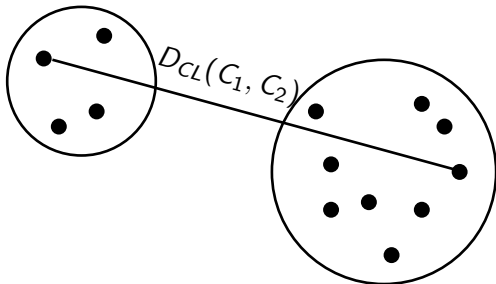


## Definition 6.4

For  $C_1, C_2 \subset M$

$$D_{CL}(C_1, C_2) := \max_{x \in C_1, y \in C_2} D(x, y)$$

is called the complete linkage cost of  $C_1, C_2$ .







---

AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$C_n := \{\{p_i\} \mid p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

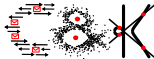
    | find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{CL}(A, B);$

    |  $C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---



---

AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$C_n := \{\{p_i\} | p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

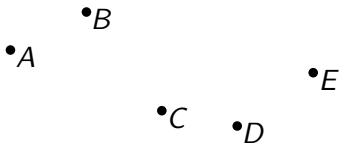
    | find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{CL}(A, B);$

    |  $C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---





---

AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$C_n := \{\{p_i\} | p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

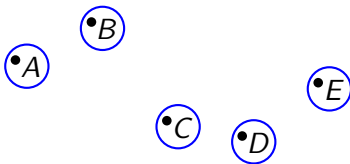
    | find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{CL}(A, B);$

    |  $C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---





---

## AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$C_n := \{\{p_i\} | p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

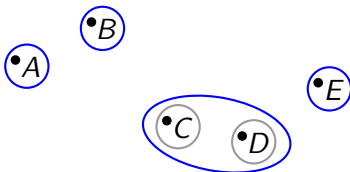
    | find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{CL}(A, B);$

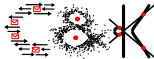
    |  $C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---





---

## AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$C_n := \{\{p_i\} | p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

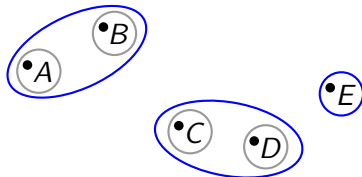
    find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{CL}(A, B);$

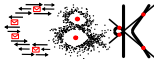
$C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---





---

## AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$C_n := \{\{p_i\} | p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

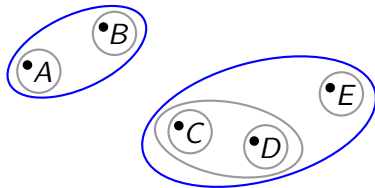
    | find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{CL}(A, B);$

    |  $C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---





---

AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$C_n := \{\{p_i\} | p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

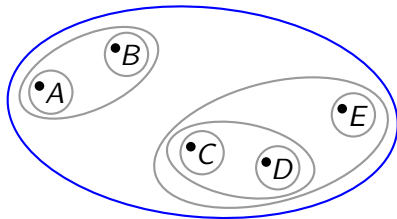
    | find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{CL}(A, B);$

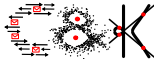
    |  $C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---





---

AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$\mathcal{C}_n := \{\{p_i\} \mid p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

    | find distinct cluster  $A, B \in \mathcal{C}_{i+1}$  minimizing  $D_{CL}(A, B);$

    |  $\mathcal{C}_i := (\mathcal{C}_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $\mathcal{C}_1, \dots, \mathcal{C}_n$  (or single  $\mathcal{C}_k$ )

---





---

AGGLOMERATIVECOMPLETELINKAGE( $P$ )

---

$\mathcal{C}_n := \{\{p_i\} \mid p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

    find distinct cluster  $A, B \in \mathcal{C}_{i+1}$  minimizing  $D_{CL}(A, B);$

$\mathcal{C}_i := (\mathcal{C}_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

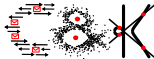
**return**  $\mathcal{C}_1, \dots, \mathcal{C}_n$  (or single  $\mathcal{C}_k$ )

---

## Theorem 6.5

*Algorithm AGGLOMERATIVECOMPLETELINKAGE requires time  $O(n^2 \log n)$  and space  $O(n^2)$ .*

# Approximation guarantees



- $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )
- $\text{cost}_{\text{diam}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{diam}^D(C_i)$  (diameter cost)
- $\text{opt}_k^{\text{diam}}(P) := \min_{|\mathcal{C}|=k} \text{cost}_{\text{diam}}^D(\mathcal{C})$



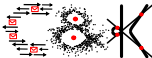
- $\text{diam}^D(S) := \max_{x,y \in S} D(x,y)$  (diameter of  $S$ )
- $\text{cost}_{\text{diam}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{diam}^D(C_i)$  (diameter cost)
- $\text{opt}_k^{\text{diam}}(P) := \min_{|\mathcal{C}|=k} \text{cost}_{\text{diam}}^D(\mathcal{C})$

## Theorem 6.6

Let  $D$  be a distance metric on  $M \subseteq \mathbb{R}^d$ . Then for all sets  $P$  and all  $k \leq |P|$ , Algorithm AGGLOMERATIVECOMPLETELINKAGE computes a  $k$ -clustering  $\mathcal{C}_k$  with

$$\text{cost}_{\text{diam}}^D(\mathcal{C}_k) \leq O\left(\text{opt}_k^{\text{diam}}(P)\right),$$

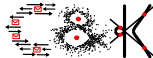
where the constant hidden in the  $O$ -notation is double exponential in  $d$ .



## Theorem 6.7

*There is a point set  $P \subset \mathbb{R}^2$  such that for the metric  $D_{l_\infty}$  algorithm `AGGLOMERATIVECOMPLETELINKAGE` computes a clustering  $\mathcal{C}_k$  with*

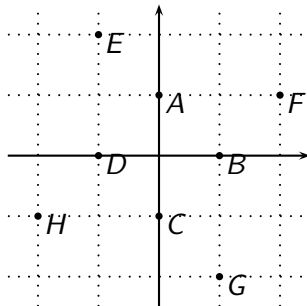
$$\text{cost}_{\text{diam}}^D(\mathcal{C}_k) = 3 \cdot \text{opt}_k^{\text{diam}}(P).$$



## Theorem 6.7

There is a point set  $P \subset \mathbb{R}^2$  such that for the metric  $D_{l_\infty}$  algorithm `AGGLOMERATIVECOMPLETELINKAGE` computes a clustering  $\mathcal{C}_k$  with

$$\text{cost}_{\text{diam}}^D(\mathcal{C}_k) = 3 \cdot \text{opt}_k^{\text{diam}}(P).$$

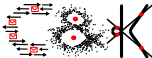




## Theorem 6.8

There is a point set  $P \subset \mathbb{R}^d$ ,  $d = k + \log k$  such that for the metric  $D_{l_1}$  algorithm `AGGLOMERATIVECOMPLETELINKAGE` computes a clustering  $\mathcal{C}_k$  with

$$\text{cost}_{\text{diam}}^{D_{l_1}}(\mathcal{C}_k) \geq \frac{1}{2} \log k \cdot \text{opt}_k^{\text{diam}}(P).$$



## Theorem 6.8

*There is a point set  $P \subset \mathbb{R}^d$ ,  $d = k + \log k$  such that for the metric  $D_{l_1}$  algorithm AGGLOMERATIVECOMPLETELINKAGE computes a clustering  $\mathcal{C}_k$  with*

$$\text{cost}_{\text{diam}}^{D_{l_1}}(\mathcal{C}_k) \geq \frac{1}{2} \log k \cdot \text{opt}_k^{\text{diam}}(P).$$

## Corollary 6.9

*For every  $1 \leq p < \infty$ , there is a point set  $P \subset \mathbb{R}^d$ ,  $d = k + \log k$  such that for the metric  $D_{l_p}$  algorithm AGGLOMERATIVECOMPLETELINKAGE computes a clustering  $\mathcal{C}_k$  with*

$$\text{cost}_{\text{diam}}^{D_{l_p}}(\mathcal{C}_k) \geq \sqrt[p]{\frac{1}{2} \log k} \cdot \text{opt}_k^{\text{diam}}(P).$$

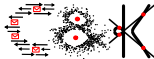


## Theorem 6.10

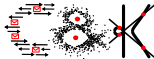
*For the metric  $D_{l_2}$  the diameter  $k$ -clustering problem is **NP-hard**. Moreover, assuming  $\mathbf{P} \neq \mathbf{NP}$ , there is no polynomial time approximation for the diameter  $k$ -clustering with approximation factor  $\leq 1.96$ .*



# Hardness of diameter clustering



- $\Delta \in \mathbb{R}_{\geq 0}^{n \times n}$ ,  $\Delta_{xy} := (x, y)$ -entry in  $\Delta$ ,  $1 \leq x, y \leq n$
- $\mathcal{C} = \{C_1, \dots, C_k\}$  partition of  $\{1, \dots, n\}$
- $\text{cost}_{\text{diam}}^{\Delta} := \max_{1 \leq i \leq k} \max_{x, y \in C_i} \Delta_{xy}$



- $\Delta \in \mathbb{R}_{\geq 0}^{n \times n}$ ,  $\Delta_{xy} := (x, y)$ -entry in  $\Delta$ ,  $1 \leq x, y \leq n$
- $\mathcal{C} = \{C_1, \dots, C_k\}$  partition of  $\{1, \dots, n\}$
- $\text{cost}_{\text{diam}}^{\Delta} := \max_{1 \leq i \leq k} \max_{x, y \in C_i} \Delta_{xy}$

## Problem 6.11 (matrix diameter $k$ -clustering)

Given a matrix  $\Delta \in \mathbb{R}_{\geq 0}^{n \times n}$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $\{1, \dots, n\}$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{diam}}^{\Delta}(\mathcal{C})$ .



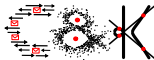
- $\Delta \in \mathbb{R}_{\geq 0}^{n \times n}$ ,  $\Delta_{xy} := (x, y)$ -entry in  $\Delta$ ,  $1 \leq x, y \leq n$
- $\mathcal{C} = \{C_1, \dots, C_k\}$  partition of  $\{1, \dots, n\}$
- $\text{cost}_{\text{diam}}^{\Delta} := \max_{1 \leq i \leq k} \max_{x, y \in C_i} \Delta_{xy}$

## Problem 6.11 (matrix diameter $k$ -clustering)

Given a matrix  $\Delta \in \mathbb{R}_{\geq 0}^{n \times n}$ ,  $k \in \mathbb{N}$ , find a partition  $\mathcal{C}$  of  $\{1, \dots, n\}$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{diam}}^{\Delta}(\mathcal{C})$ .

## Theorem 6.12

The matrix diameter  $k$ -clustering problem is **NP-hard**. Moreover, assuming **P**  $\neq$  **NP**, there is no polynomial time approximation for the diameter  $k$ -clustering with approximation factor  $\alpha \geq 1$  arbitrary.

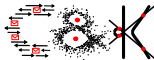


## Problem 6.13 (maximum distance $k$ -clustering)

Given distance measure  $D : M \times M \rightarrow \mathbb{R}$ ,  $k \in \mathbb{N}$ , and  $P \subset M$ , find a partition  $\mathcal{C} = \{C_1, \dots, C_k\}$  of  $P$  into  $k$  clusters that maximizes

$$\min_{x \in C_i, y \in C_j, i \neq j} D(x, y),$$

*i.e. a partition that maximizes the minimum distance between points in different clusters.*



## Problem 6.13 (maximum distance $k$ -clustering)

Given distance measure  $D : M \times M \rightarrow \mathbb{R}$ ,  $k \in \mathbb{N}$ , and  $P \subset M$ , find a partition  $\mathcal{C} = \{C_1, \dots, C_k\}$  of  $P$  into  $k$  clusters that maximizes

$$\min_{x \in C_i, y \in C_j, i \neq j} D(x, y),$$

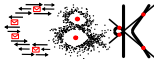
*i.e. a partition that maximizes the minimum distance between points in different clusters.*

## Definition 6.14

For  $C_1, C_2 \subset M$

$$D_{SL}(C_1, C_2) := \min_{x \in C_1, y \in C_2} D(x, y)$$

*is called the single linkage cost of  $C_1, C_2$ .*



---

AGGLOMERATIVE SINGLE LINKAGE( $P$ )

---

$\mathcal{C}_n := \{\{p_i\} \mid p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

    find distinct cluster  $A, B \in \mathcal{C}_{i+1}$  minimizing  $D_{SL}(A, B);$

$\mathcal{C}_i := (\mathcal{C}_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $\mathcal{C}_1, \dots, \mathcal{C}_n$  (or single  $\mathcal{C}_k$ )

---



---

AGGLOMERATIVE SINGLE LINKAGE( $P$ )

---

$C_n := \{\{p_i\} \mid p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

    find distinct cluster  $A, B \in C_{i+1}$  minimizing  $D_{SL}(A, B);$

$C_i := (C_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

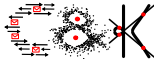
**end**

**return**  $C_1, \dots, C_n$  (or single  $C_k$ )

---

## Theorem 6.15

*Algorithm AGGLOMERATIVE SINGLE LINKAGE optimally solves the maximum distance  $k$ -clustering problem.*



- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x, m)$  (discrete radius of  $S$ )
- $\text{cost}_{\text{drad}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{drad}^D(C_i)$  (discrete radius cost)
- find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{drad}}^D(\mathcal{C})$  or  $\text{cost}_{\text{rad}}^D(\mathcal{C})$ .





- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x, m)$  (discrete radius of  $S$ )
- $\text{cost}_{\text{drad}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{drad}^D(C_i)$  (discrete radius cost)
- find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{drad}}^D(\mathcal{C})$  or  $\text{cost}_{\text{rad}}^D(\mathcal{C})$ .

### Theorem 6.16

Let  $D : M \times M \rightarrow \mathbb{R}$  be a metric,  $P \subset M$  and  $\mathcal{C} = \{C_1, \dots, C_k\}$  a partition of  $P$ . Then

- 1  $\text{cost}_{\text{drad}}(\mathcal{C}) \leq \text{cost}_{\text{diam}}(\mathcal{C}) \leq 2 \cdot \text{cost}_{\text{drad}}(\mathcal{C})$
- 2  $\frac{1}{2} \cdot \text{cost}_{\text{drad}}(\mathcal{C}) \leq \text{cost}_{\text{rad}}(\mathcal{C}) \leq \text{cost}_{\text{drad}}(\mathcal{C})$



### Corollary 6.17

Let  $D : M \times M \rightarrow \mathbb{R}$  be a metric,  $k \in \mathbb{N}$ , and  $P \subset M$ . Then

- 1  $opt_k^{drad}(P) \leq opt_k^{diam}(P) \leq 2 \cdot opt_k^{drad}(P)$
- 2  $\frac{1}{2} \cdot opt_k^{drad}(P) \leq opt_k^{rad}(P) \leq opt_k^{drad}(P)$



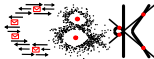
### Corollary 6.17

Let  $D : M \times M \rightarrow \mathbb{R}$  be a metric,  $k \in \mathbb{N}$ , and  $P \subset M$ . Then

- 1  $opt_k^{drad}(P) \leq opt_k^{diam}(P) \leq 2 \cdot opt_k^{drad}(P)$
- 2  $\frac{1}{2} \cdot opt_k^{drad}(P) \leq opt_k^{rad}(P) \leq opt_k^{drad}(P)$

### Corollary 6.18

Assume there is a polynomial time  $c$ -approximation algorithm for the discrete radius  $k$ -clustering problem. Then there is a polynomial time  $2c$ -approximation algorithm for the diameter  $k$ -clustering problem.



---

GONZALESALGORITHM( $P, k$ )

---

$C := \{p\}$  for  $p \in P$  arbitrary;

**for**  $i = 1, \dots, k - 1$  **do**

$q := \operatorname{argmax}_{y \in P} D(y, C);$

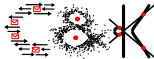
$C := C \cup \{q\};$

**end**

compute partition  $\mathcal{C} = \{C_1, \dots, C_k\}$  corresponding to  $C$ ;

**return**  $\mathcal{C}$  and  $C$

---



---

**GONZALESALGORITHM**( $P, k$ )

---

$C := \{p\}$  for  $p \in P$  arbitrary;

**for**  $i = 1, \dots, k - 1$  **do**

$q := \operatorname{argmax}_{y \in P} D(y, C)$ ;

$C := C \cup \{q\}$ ;

**end**

compute partition  $\mathcal{C} = \{C_1, \dots, C_k\}$  corresponding to  $C$ ;

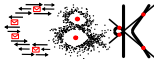
**return**  $\mathcal{C}$  and  $C$

---

## Theorem 6.19

*For any metric  $D$ , Algorithm GONZALESALGORITHM is a 2-approximation algorithm for the diameter, radius, and discrete radius  $k$ -clustering problem.*

## Proof of Theorem 6.19 (for diameter)



- $C := \{c_1, \dots, c_k\}$  set of points chosen by the algorithm
- chosen in order  $c_1, \dots, c_k$
- $G_l := \{c_1, \dots, c_l\}$ , i.e.  $G_l$  set of first  $l$  points chosen by the algorithm,  $C = G_k$ .
- $c_{k+1} := \operatorname{argmax}_{q \in P} D(q, G_k)$ , i.e.  $c_{k+1}$  is the point that would be chosen by the algorithm in an additional iteration

Show that

- 1  $\forall l \leq k - 1 : D(c_{l+1}, G_l) \geq D(c_{k+1}, G_k)$
- 2  $\operatorname{opt}_k^{\operatorname{diam}}(P) \geq D(c_{k+1}, G_k)$
- 3  $\operatorname{cost}_{\operatorname{diam}}^D(C) \leq 2 \cdot D(c_{k+1}, G_k)$ .

## Proof of Theorem 6.19 (for diameter)



- 2 and 3 imply the theorem
- 1 is used to prove 2

### Proof of 1

- Assume there is an index  $l, 1 \leq l \leq k - 1$ , such that  $D(c_{l+1}, G_l) < D(c_{k+1}, G_k)$ .
- Then  $D(c_{l+1}, G_l) < D(c_{k+1}, G_k) \leq D(c_{k+1}, G_l)$ , since  $G_l \subset G_k$ .
- This contradicts the choice of  $c_{l+1}$  as  $\operatorname{argmax}_{q \in P} D(q, G_l)$ .



## Proof of 2

- It suffices to prove that for any  $k$ -clustering  $\hat{C}$  of  $P$  we have  $\text{cost}_{\text{diam}}^D(\hat{C}) \geq D(c_{k+1}, G_k)$ .
- For any  $k$ -clustering  $\hat{C}$  of  $P$  there are two elements  $c_i, c_j$  of  $C \cup \{c_{k+1}\}$  that belong to the same cluster. Assume without loss of generality that  $i < j$ .
- Then  $\text{cost}_{\text{diam}}^D(\hat{C}) \geq D(c_j, c_i) \geq D(c_j, G_{j-1}) \geq D(c_{k+1}, G_k)$ , where the second inequality follows from  $c_i \in G_{j-1}$  and the last inequality follows from 1.
- This contradicts the choice of  $c_{l+1}$  as  $\text{argmax}_{q \in P} D(q, G_l)$ .





## Proof of 3

- Let  $C_l$  be a cluster in the clustering  $\mathcal{C}$  computed by Gonzales' algorithm and let  $u, v \in C_l$ . It suffices to show that  $D(u, v) \leq 2 \cdot D(c_{k+1}, G_k)$ .
- Let  $c_l$  be the element of  $C$  such that  $C_l$  consists of all elements that are closer to  $c_l$  than to any other element in  $C$ .
- Then  $D(u, c_l) = D(u, G_k) \leq D(c_{k+1}, G_k)$ , where the inequality follows from the definition of  $c_{k+1}$  as  $\operatorname{argmax}_{q \in P} D(q, G_k)$ . Similarly, we get  $D(v, c_l) \leq D(c_{k+1}, G_k)$ .
- By the triangle inequality  $D(u, v) \leq D(u, c_l) + D(v, c_l) \leq 2 \cdot D(c_{k+1}, G_k)$ .



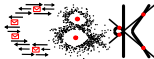
- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x, m)$  (discrete radius of  $S$ )
- $\text{cost}_{\text{drad}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{drad}^D(C_i)$  (discrete radius cost)
- find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{drad}}^D(\mathcal{C})$ .



- $\text{drad}^D(S) := \min_{m \in S} \max_{x \in S} D(x, m)$  (discrete radius of  $S$ )
- $\text{cost}_{\text{drad}}^D(\mathcal{C}) := \max_{1 \leq i \leq k} \text{drad}^D(C_i)$  (discrete radius cost)
- find a partition  $\mathcal{C}$  of  $P$  into  $k$  clusters  $C_1, \dots, C_k$  that minimizes  $\text{cost}_{\text{drad}}^D(\mathcal{C})$ .

## Discrete radius measure

$$D_{\text{drad}}(C_1, C_2) = \text{drad}(C_1 \cup C_2)$$



---

AGGLOMERATIVE DISCRETE RADIUS ( $P$ )

---

$\mathcal{C}_n := \{\{p_i\} \mid p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

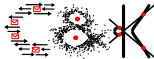
    | find distinct clusters  $A, B \in \mathcal{C}_{i+1}$  minimizing  $D_{\text{drad}}(A, B);$

    |  $\mathcal{C}_i := (\mathcal{C}_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

**return**  $\mathcal{C}_1, \dots, \mathcal{C}_n$  (or single  $\mathcal{C}_k$ )

---



---

AGGLOMERATIVEDISCRETERADIUS( $P$ )

---

$\mathcal{C}_n := \{\{p_i\} \mid p_i \in P\};$

**for**  $i = n - 1, \dots, 1$  **do**

    | find distinct clusters  $A, B \in \mathcal{C}_{i+1}$  minimizing  $D_{\text{drad}}(A, B);$

    |  $\mathcal{C}_i := (\mathcal{C}_{i+1} \setminus \{A, B\}) \cup \{A \cup B\};$

**end**

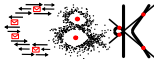
**return**  $\mathcal{C}_1, \dots, \mathcal{C}_n$  (or single  $\mathcal{C}_k$ )

---

## Theorem 6.20

Let  $D$  be a distance metric on  $M \subseteq \mathbb{R}^d$ . Then for all sets  $P \subset M$  and all  $k \leq |P|$ , Algorithm AGGLOMERATIVEDISCRETERADIUS computes a  $k$ -clustering  $\mathcal{C}_k$  with

$$\text{cost}_k^{\text{drad}}(\mathcal{C}_k) < O(d) \cdot \text{opt}_k.$$

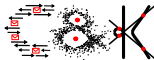


**Hierarchical clustering** Given distance measure

$D : M \times M \rightarrow \mathbb{R}$ ,  $k \in \mathbb{N}$ , and  $P \subset M$ ,  $|P| = n$ , a sequence of clusterings  $\mathcal{C}_n, \dots, \mathcal{C}_1$  with  $|\mathcal{C}_k| = k$  is called *hierarchical clustering of  $P$*  if for all  $A \in \mathcal{C}_k$

**1**  $A \in \mathcal{C}_{k+1}$  or

**2**  $\exists B, C \in \mathcal{C}_{k+1} : A = B \cup C$  and  $\mathcal{C}_k = \mathcal{C}_{k+1} \setminus \{B, C\} \cup \{A\}$ .



**Hierarchical clustering** Given distance measure

$D : M \times M \rightarrow \mathbb{R}$ ,  $k \in \mathbb{N}$ , and  $P \subset M$ ,  $|P| = n$ , a sequence of clusterings  $\mathcal{C}_n, \dots, \mathcal{C}_1$  with  $|\mathcal{C}_k| = k$  is called *hierarchical clustering of  $P$*  if for all  $A \in \mathcal{C}_k$

1  $A \in \mathcal{C}_{k+1}$  or

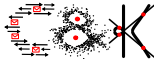
2  $\exists B, C \in \mathcal{C}_{k+1} : A = B \cup C$  and  $\mathcal{C}_k = \mathcal{C}_{k+1} \setminus \{B, C\} \cup \{A\}$ .

**Dendrograms** A dendrogram on  $n$  nodes is a rooted binary tree

$T = (V, E)$  with an index function

$\chi : V \setminus \{\text{leaves of } T\} \rightarrow \{1, \dots, n\}$  such that

- $\forall v \neq w : \chi(v) \neq \chi(w)$
- $\chi(\text{root}) = n$
- $\forall u, v$ : if  $v$  parent of  $u$ , then  $\chi(v) > \chi(u)$ .



$\mathcal{C}_n, \dots, \mathcal{C}_1$  hierarchical clustering of  $P$ .

## Construction of dendrogram

- create leaf for each point  $p \in P$
- interior nodes correspond to union of clusters
- if  $k$ -th cluster is obtained by union of clusters  $B, C$ , create new node with index  $k$  and with children  $B, C$ .



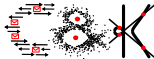


## AGGLOMERATIVE COMPLETE LINKAGE

- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$

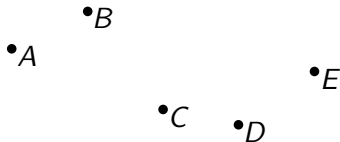


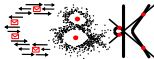
## AGGLOMERATIVE COMPLETE LINKAGE

- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$





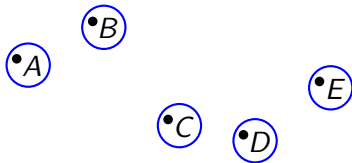
## AGGLOMERATIVE COMPLETE LINKAGE

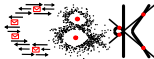
- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$

A B C D E



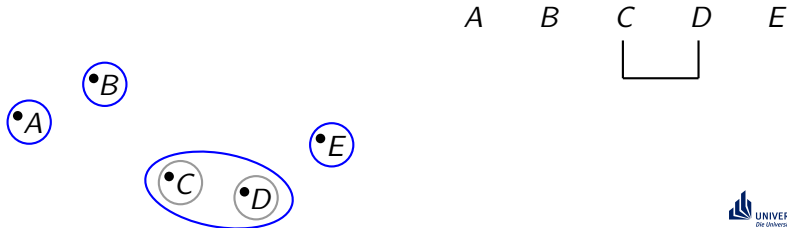


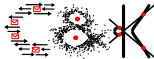
## AGGLOMERATIVE COMPLETE LINKAGE

- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$



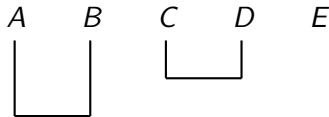
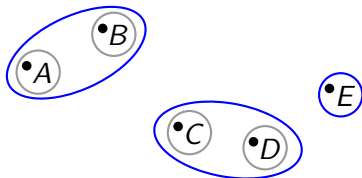


## AGGLOMERATIVE COMPLETE LINKAGE

- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$



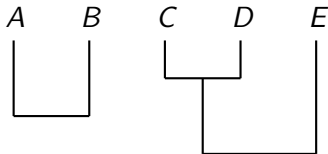
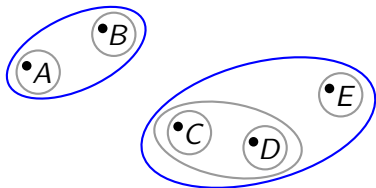


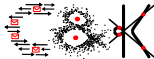
## AGGLOMERATIVE COMPLETE LINKAGE

- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$



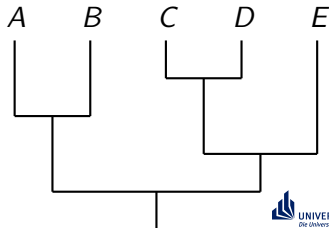
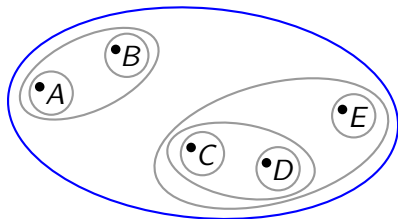


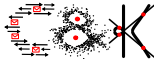
## AGGLOMERATIVE COMPLETE LINKAGE

- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$



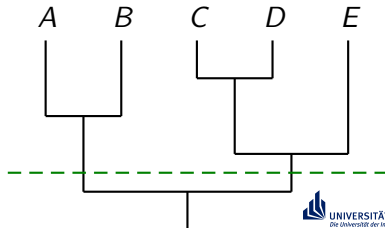
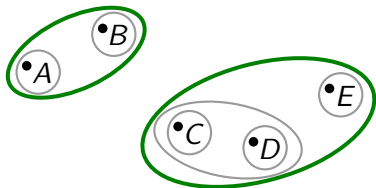


## AGGLOMERATIVE COMPLETE LINKAGE

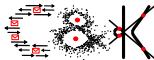
- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$







## AGGLOMERATIVE COMPLETE LINKAGE

- Start with one cluster for each input object.
- Iteratively merge the two closest clusters.

## Complete linkage measure

$$D_{CL}(C_1, C_2) = \max_{x \in C_1, y \in C_2} D(x, y)$$

