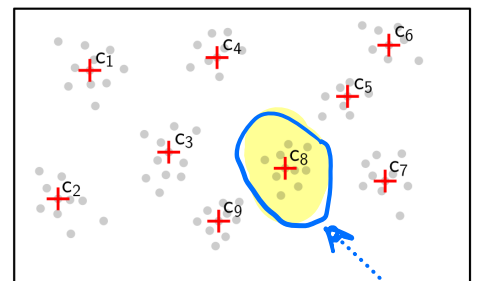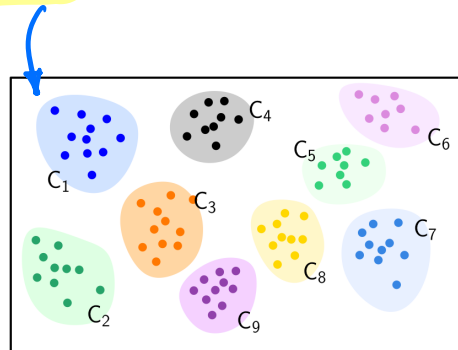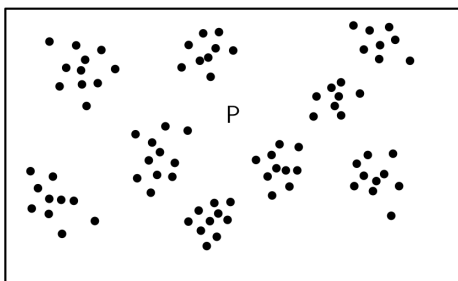CMSC 451 - Algorithm Design
Lecture 6 - k-Center Clustering & Gonzalez's Algorithm

Greedy algorithms often used to approximate NP-hard problems

Clustering - Given a set of points $P$ & distance function, partition it into similar groups, called clusters $\{C_1, ..., C_k\}$



Center-based clustering -
Compute a set of cluster centers $\{c_1, ..., c_k\}$ and clusters are implicitly defined by distance
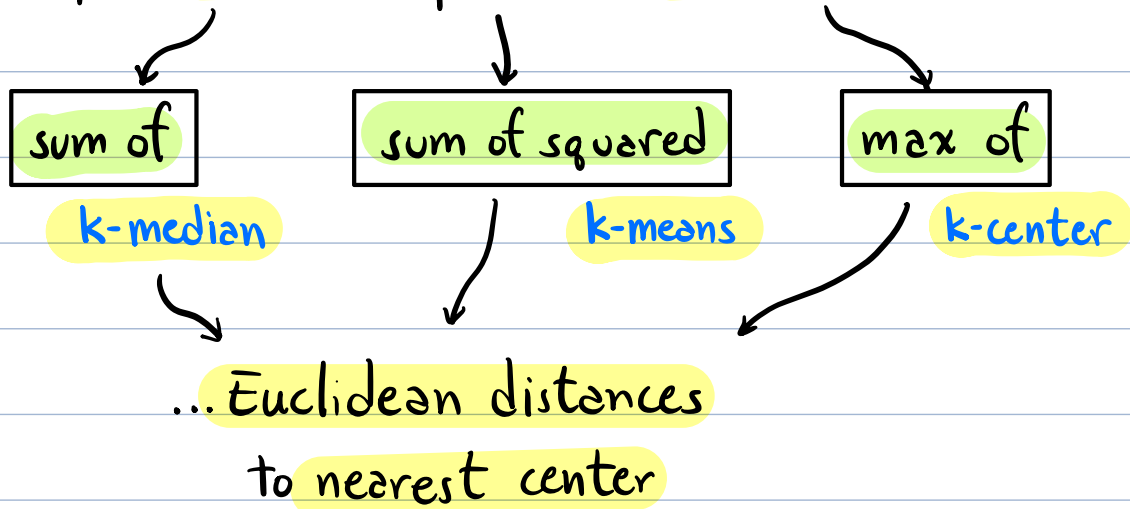
$N(c_i)$ = subset of $P$ closest to $c_i$

Two varieties -
- Centers must be chosen from $P$ (discrete clustering)
- Centers can be any point in space

# Three Common Center-Based Clusterings:
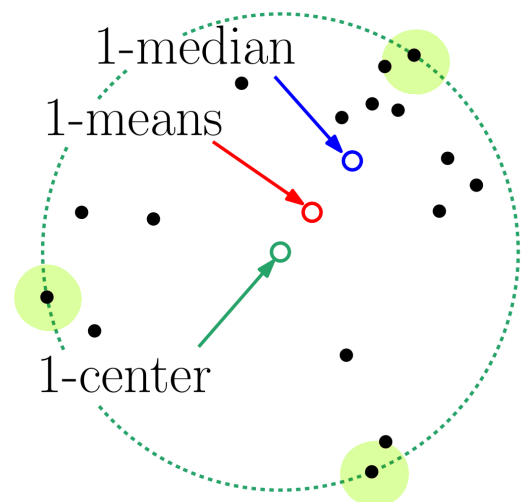
Compute **k center points to minimize...**

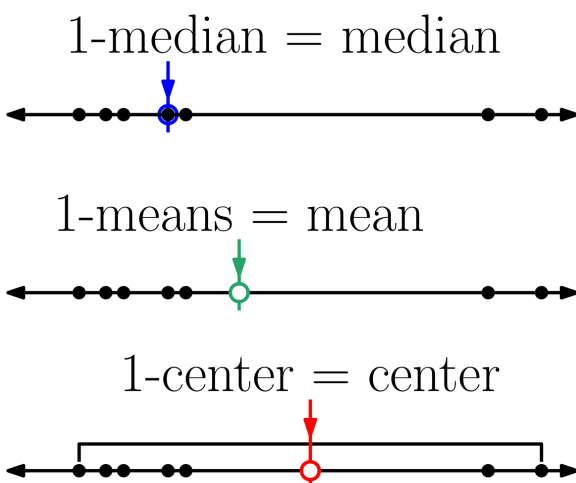| sum of | sum of squared | max of |
|--------|----------------|--------|
| k-median | k-means | k-center |

...**Euclidean distances to nearest center**

**Which is best?** Depends on **application** (e.g., sensitivity to **outliers**)

Helps to understand the **single-cluster case (k=1)**

**1-median** - **1-D** → median
  **d-D** → hard! Fermat-Weber problem

1-median = median

1-means = mean

1-center = center

1-median
1-means
1-center

1-means — 1-D → mean

d-D → centroid (center of mass)

Easy to compute in any dimension!

Take mean coord. value in each dim.

k-Means is very popular - Lloyd's Algorithm

1-center — 1-D → midpoint of min & max

d-D → center of min enclosing ball

(Can compute in $O(n)$ time,

but tricky algorithm - take CMSC 754)

## Metric Space:

Distance function $\delta : \underline{P} \times \underline{P} \to \mathbb{R}^{\geq 0}$

- $\delta(p,q) \geq 0$ & $\delta(p,p) = 0$ — Positive
- $\delta(p,q) = \delta(q,p)$ — Symmetric
- $\delta(p,r) \leq \delta(p,q) + \delta(q,r)$ — $\Delta$-Inequality

## k-Center Problem:

Given point set $P$ in a metric space and $k \geq 1$, compute $C \subseteq P$ of size $k$ to minimize max distance to closest center in C.

Note: Centers must be drawn from $P$

– Given $C \subseteq P$, define objective fn.

$$\Delta_p(C) = \max_{p \in P} \min_{c \in C} d(p,c)$$
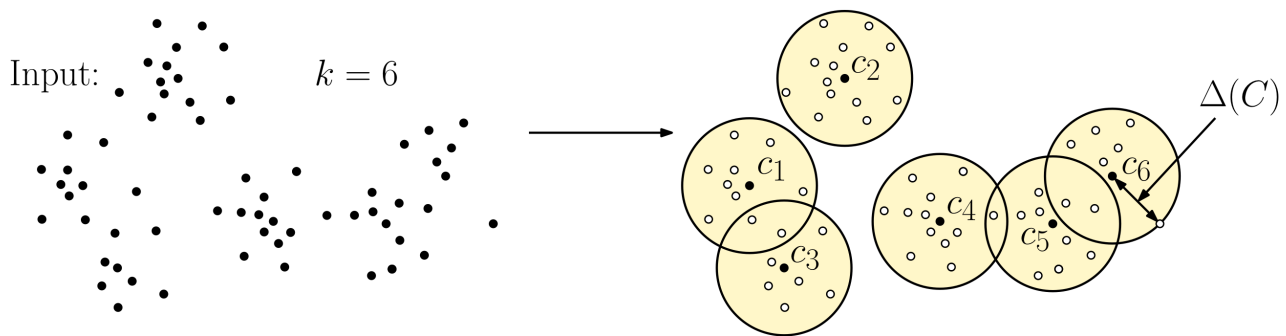
Problem – Compute a k-element set to minimize this:

$$\min_{\substack{C \subseteq P \\ |C| = k}} \Delta_c(P)$$

Geometric interpretation –
Cover all pts of $P$
  – k balls (centered at pts of $P$)
  – minimum radius $r = \Delta_p(C)$

Input: $k = 6$

$\Delta(C)$

$c_1$ $c_2$ $c_3$ $c_4$ $c_5$ $c_6$

Gonzalez's Algorithm –
  – Greedy + very simple
  – 2x-approx. to k-center
  – $O(k \cdot n)$ time

Intuitive Explanation:
   Repeatedly add the point that is
   farthest from its closest center

gonzalez (P, k)        // Gonzalez's k-center
  G ← ∅
  for each (p ∈ P) d[p] ← +∞        // init. dists
  for (i ← 1 to k)
      p ← pt of P s.t. d[p] is max        // farthest pt
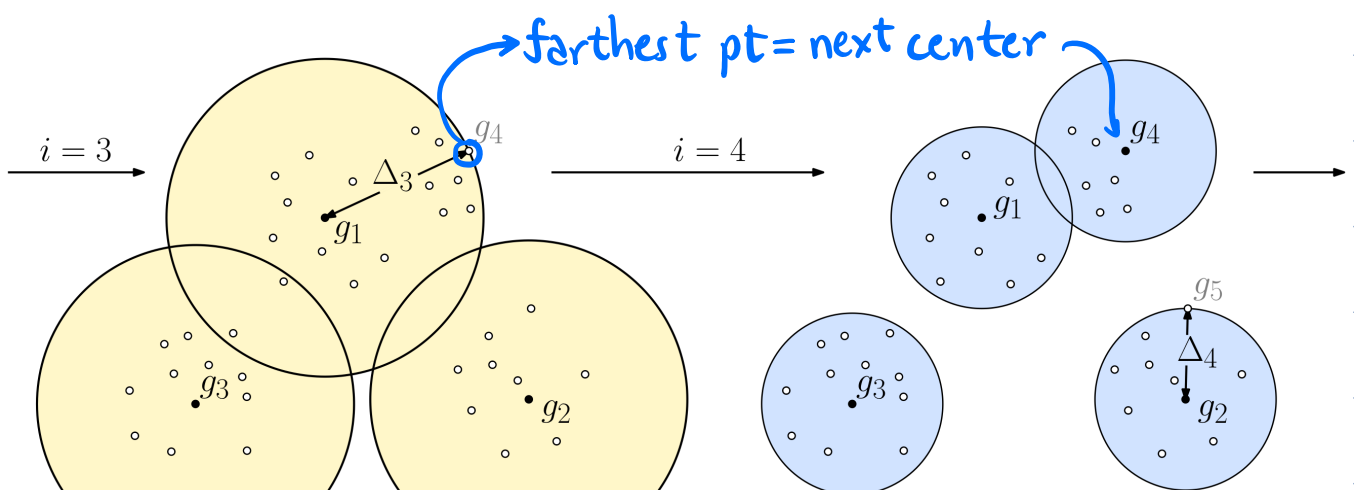      add p to G        // ... is next center
      for each (q ∈ P)        // update dist to closest
          d[q] ← min(d[q], dist(p, q))
  return G        // final centers

Example − $i = 4$    $\Delta_i$ = max distance to closest center
                        = ball radius



$i = 3$ → $\Delta_3$ $g_1$ $g_4$ farthest pt = next center → $i = 4$ → $g_4$ $g_1$ $g_3$ $g_5$ $\Delta_4$ $g_2$ $g_3$ $g_2$

## Correctness-

**Feasibility** – Clearly the algorith returns a valid clustering (provided $|P| \geq k$)

**(Approx.) Optimality** – Will show that our final radius $\leq 2 \cdot$ opt radius

Given any set $\dot{C} \subseteq P$, recall that obj. fn. is

$$\Delta_p(\dot{C}) = \max_{p \in P} \min_{c \in \dot{C}} \delta(p, c)$$

Let $G$ = output of Gonzalez

$O$ = opt. k-center solution

.

$$\boxed{\text{Thm:} \quad \Delta_p(G) \leq 2 \cdot \Delta_p(O)}$$

→ We'll drop this subscript

At first glance this seems hopeless!

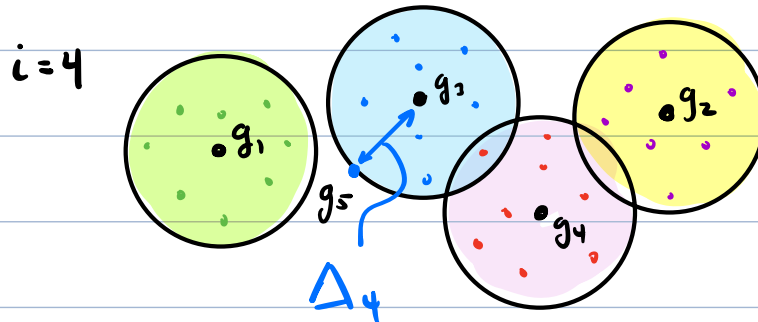- k-center is NP-hard

- We cannot know what $\Delta(O)$ is !

## Strategy -

- Derive an easily computable estimate, $\Delta_{min}$

- Show: $\Delta(O) \geq \Delta_{min}$

- Show: $\Delta(G) \leq 2 \cdot \Delta_{min}$

$\Rightarrow \quad \Delta(G) \leq 2 \cdot \Delta_{min} \leq 2 \cdot \Delta(O)$ ✓

$G_i = \{g_1, ..., g_i\}$ — the **first** $i$ greedy ctrs.

$\Delta_i = \Delta(G_i)$ — **farthest dist** to these ctrs.

$i = 4$



$\Delta_4$

Imagine that we ran one additional iteration to get **k+1 centers** $G_{k+1}$

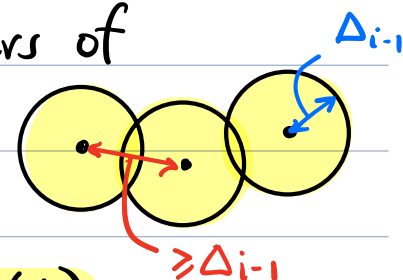**The theorem follows from 3 claims:**

**Claim 1:** (Greedy distances decrease)

For $1 \leq i \leq k+1$, $\Delta_{i+1} \leq \Delta_i$

Pf: As we **add more centers**, the dist to each pt's closest ctr. **can never increase.**

**Claim 2:** (Greedy centers are never close)

For $1 \leq i \leq k+1$, every pair of centers of $G_i$ are at **dist $\geq \Delta_{i-1}$**

$\Delta_{i-1}$



$\geq \Delta_{i-1}$

Corollary: $g, g' \in G_{k+1} \Rightarrow \delta(g,g') \geq \Delta_k = \Delta(G)$

Pf. By induction on i.
- At stage i-1, by induction, all old ctrs.
  are sep. by dist $\geq \Delta_{i-2}$.
  - by Claim 1, $\Delta_{i-2} \geq \Delta_{i-1}$ ✓
- New center is at dist $\Delta_{i-1}$ from
  its closest center $\Rightarrow$ its
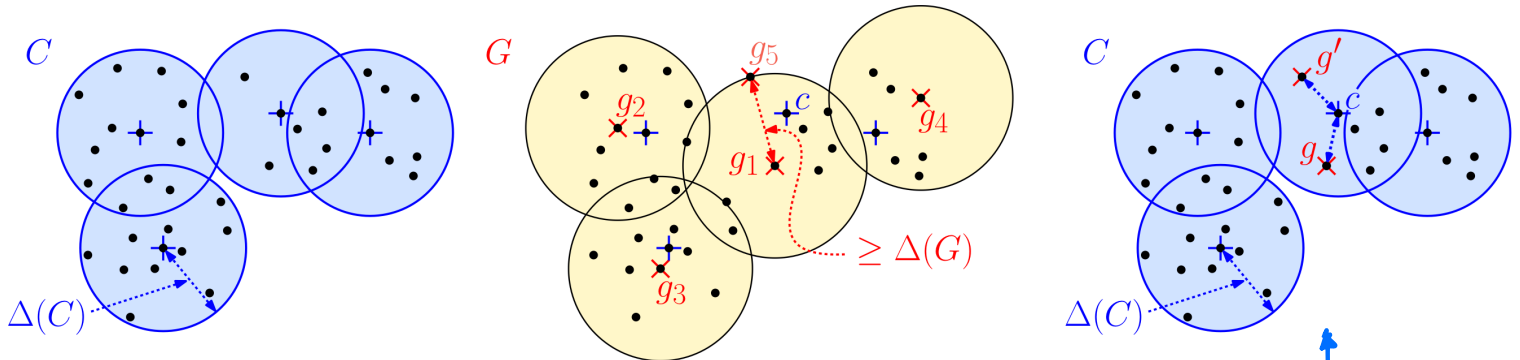  at dist $\geq \Delta_{i-1}$ from all centers ✓

Define: $\Delta_{min} = \Delta(G)/2$

Claim 3: ( $\Delta_{min}$ is a lower bound )
For any set $C \subseteq P$ of size k, $\Delta(G) \geq \Delta_{min}$

Pf: By def, every pt of P lies within
  dist $\Delta(C)$ of some pt of C.
- Since $G_{k+1} \subseteq P$, every pt. of $G_{k+1}$
  is within dist $\Delta(C)$ of some pt of C.

- Since $|G_{k+1}| = k+1$ + $|C| = k$, at least two pts of $G_{k+1}$ are within dist $\Delta(C)$ of same pt of C.

"Pigeonhole principle"

$\Rightarrow \exists g, g' \in G_{k+1}$ $c \in C$ s.t.

$$\delta(g,c) \leq \Delta(C) \quad + \quad \delta(g',c) \leq \Delta(C)^{\text{(b)}}$$

- We have:

$$\Delta(G) \overset{\text{(a)}}{\leq} \delta(g,g')$$
$$\leq \delta(g,c) + \delta(c,g') \quad (\Delta\text{-inequality})$$
$$\leq \delta(g,c) + \delta(g',c) \quad (\text{symmetry})$$
$$\overset{\text{(b)}}{\leq} \Delta(C) + \Delta(C)$$
$$\leq 2\Delta(C)$$

by Def of $\Delta_{min}$

$$\Rightarrow \Delta(C) \geq \Delta(G)/2 = \Delta_{min} \quad \checkmark$$

In conclusion: Applying Claim 3 to opt, $O$,

$$\Delta(G) = 2 \cdot \Delta_{min} \leq 2 \cdot \Delta(O)$$

$\therefore$ Greedy is within factor 2 of opt. $\square$

Summary - k-center - NP-hard clustering problem
- Gonzalez - Greedy alg. for k-center
- Factor 2 approx. to optimum