



Coresets and their Applications

Dan Feldman

Elections Queries

- Input: $P = \{$  $\}$

Elections Queries

- Input: $P = \{$  $\}$
- Query: $q \in \{$  ,  $\}$

Elections Queries

- Input: $P = \{$  $\}$
- Query: $q \in \{$  ,  $\}$
- Output: $f(P, q) =$
fraction of P that vote for q

ε -Coreset for Elections

Answer queries in sub-linear time.

ε -Coreset for Elections

Answer queries in sub-linear time.

Key Idea: Random Sampling $C \subseteq P$

ε -Coreset for Elections

Answer queries in sub-linear time.

Key Idea: Random Sampling $C \subseteq P$

Chernoff Inequality: If $|C| \geq 100/\varepsilon^2$,

$$\forall q : |f(P, q) - f(C, q)| \leq \varepsilon$$

with high probability.

Definition

C is an ε -coreset for (P, Q, f)

if for every $q \in Q$ we have

$$f(P, q) \sim f(C, q),$$

Definition

C is an ε -coreset for (P, \mathbb{Q}, f)

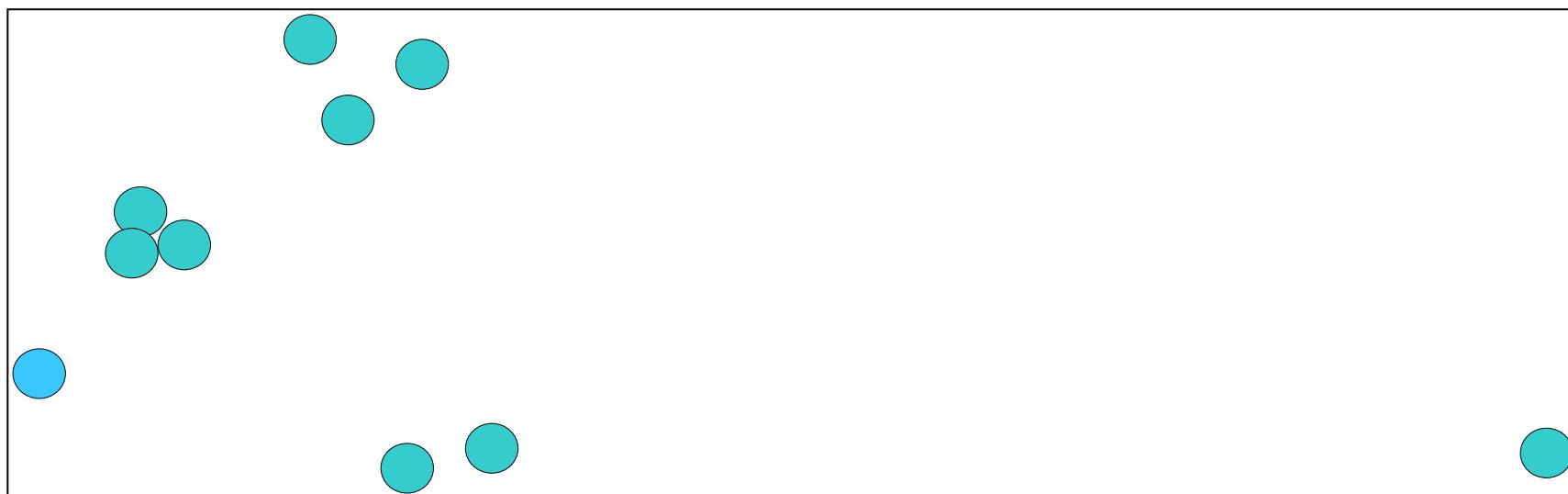
if for every $q \in \mathbb{Q}$ we have w.h.p

$$f(P, q) \sim f(C, q),$$

where the approximation depends on ε .

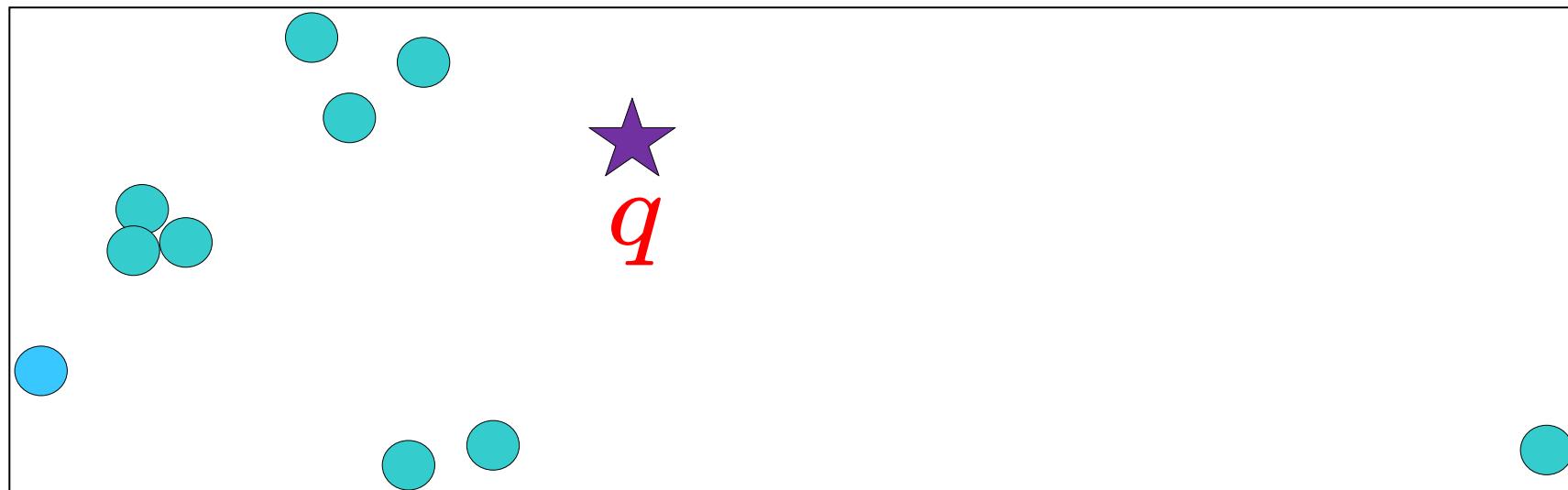
Mean Queries

- Input: P in \mathbb{R}^d



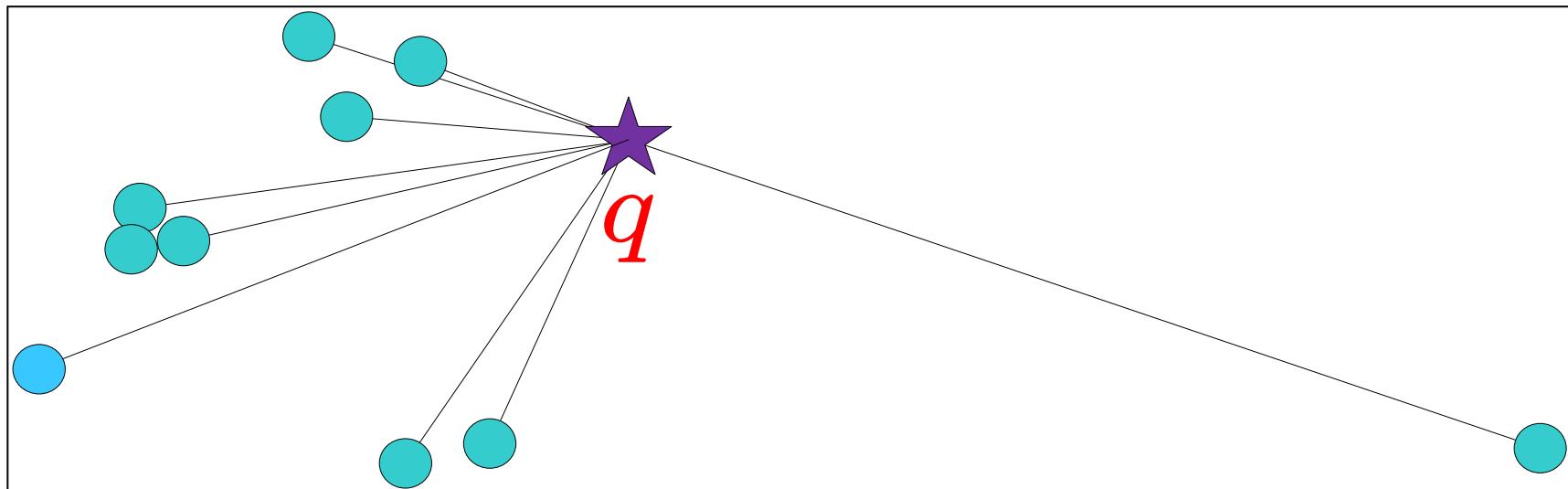
Mean Queries

- Input: P in \mathbb{R}^d
- Query: a point $q \in \mathbb{R}^d$



Mean Queries

- Input: P in \mathbb{R}^d
- Query: a point $q \in \mathbb{R}^d$
- Output: $f(P, q) = \sum_{p \in P} (\text{dist}(p, q))^2$

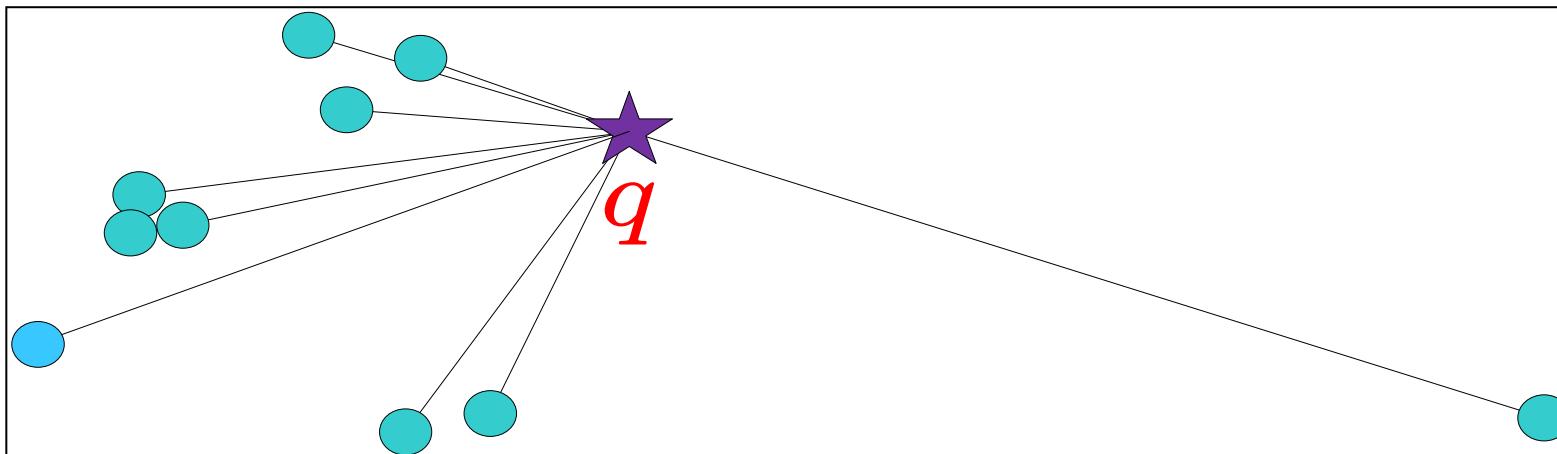


Coreset For Mean Queries

- Assume $\sum_{p \in P} = 0$.

Otherwise, translate P .

- $(\text{dist}(p, q))^2 = \|p - q\|^2$
 $= \|p\|^2 + \|q\|^2 - 2p \cdot q$



Coreset For Mean Queries

- Assume $\sum_{p \in P} = 0$.
Otherwise, translate P .
- $(\text{dist}(p, q))^2 = \|p - q\|^2$
 $= \|p\|^2 + \|q\|^2 - 2p \cdot q$
- $\sum_{p \in P} (\text{dist}(p, q))^2 = \sum_{p \in P} \|p\|^2 + \sum_{p \in P} \|q\|^2 - 2q \cdot \sum_{p \in P} p$
 $= \sum_{p \in P} \|p\|^2 + |P| \|q\|^2$

Coreset For Mean Queries

$$f(P, q) = \sum_{p \in P} \|p\|^2 + |P| \cdot \|q\|^2$$
$$= g(C, q)$$

where $C = \sum_{p \in P} p = 0$,

$$g(C, q) = |P| \cdot \text{dist}(C, q) + \sum_{p \in P} \|p\|^2$$

Properties of Coresets

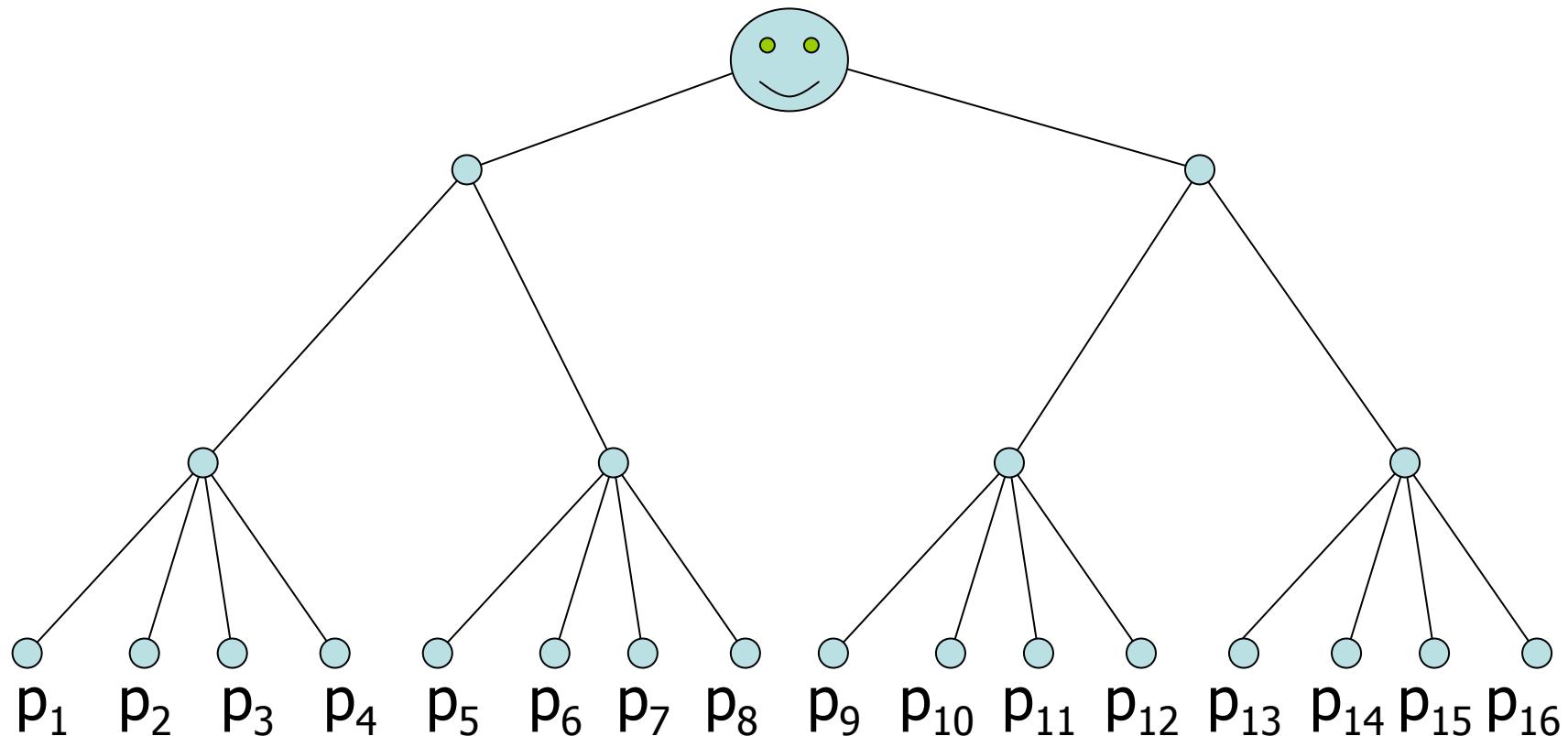
- **Transitivity:** A coreset of a coreset is a coreset
- **Additivity:** A union of coresets is a coreset

Streaming: The Model

- Single pass over the data: p_1, p_2, \dots, p_n
- $O(\log n)$ space
- Fast processing time per element

From a presentation by Piotr Indyk

Tree Computation



From a presentation by Piotr Indyk

Other Applications

- Distributed Computing
- GPUs
- Faster Heuristics
- Optimization with/without constraints

Other Coresets

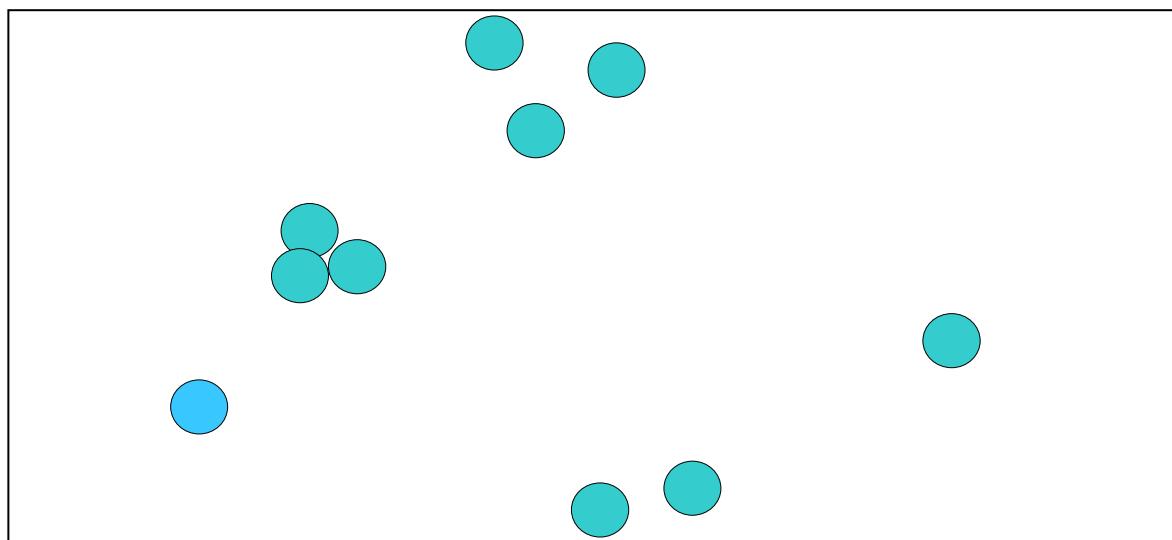
- k -Line median/means [with A. Fiat, M. Sharir]
- Private Coresets [with A. Fiat, H. Kaplan, K. Nissim]
- Regression/Matrix approximations [with C. Sohler]
- Compressed sensing [with B. Hassibi]
- Supported vector machines [with R. Condor]
- Mixture of Gaussians [with Andreas Krause]
- Dynamic coresets [with D. Golonov]
- k -Sparsest-cut [with M. Mahoney]
- A Unified Framework [with L. Schulman, M. Langberg]

Implementations

- Google's PageRank [with R. Bar-Yaets, Y. Koren]
- Image Compression [with M. Feige, N. Sochen]
- Text Mining [MyOgger Ltd.]

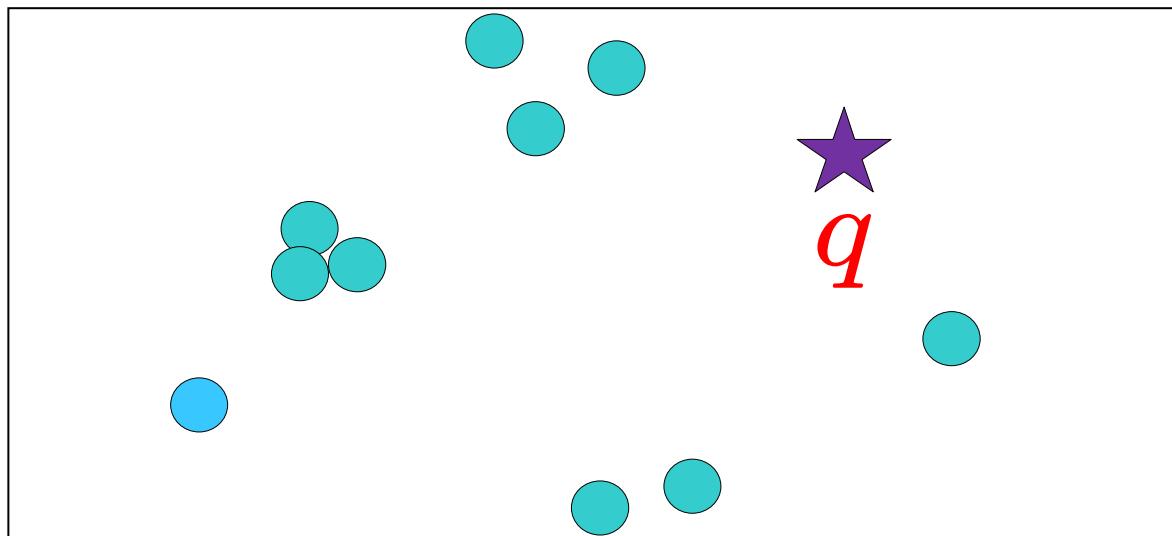
Minimum Enclosing Ball

- Input: P in \mathbb{R}^d



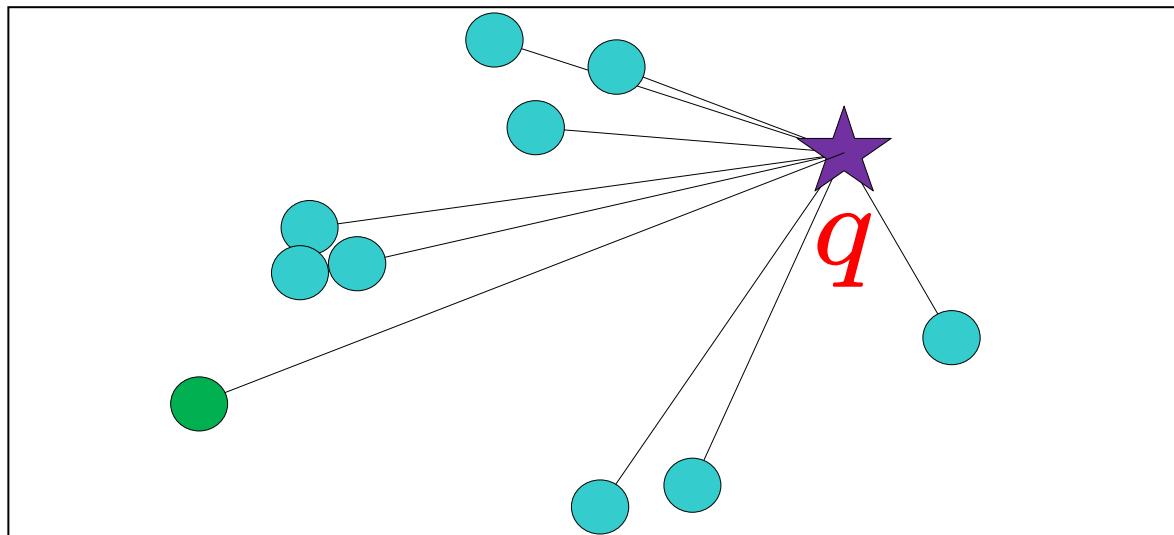
Minimum Enclosing Ball

- Input: P in \mathbb{R}^d
- Query: a point $q \in \mathbb{R}^d$



Minimum Enclosing Ball

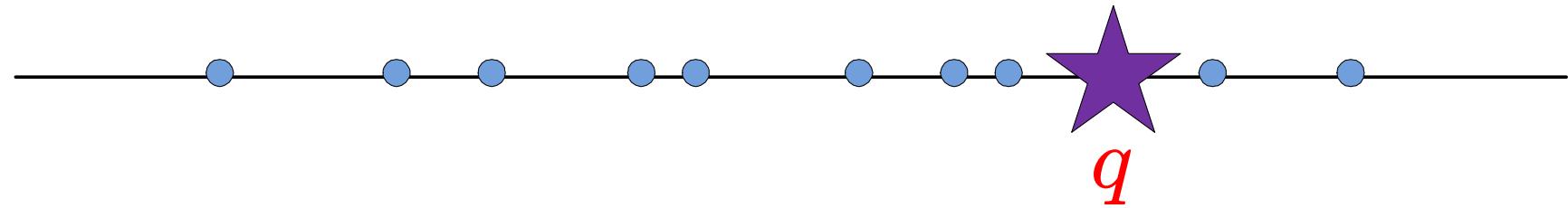
- Input: P in \mathbb{R}^d
- Query: a point $q \in \mathbb{R}^d$
- Output: $\text{far}(P, q) = \max_{p \in P} \text{dist}(p, q)$



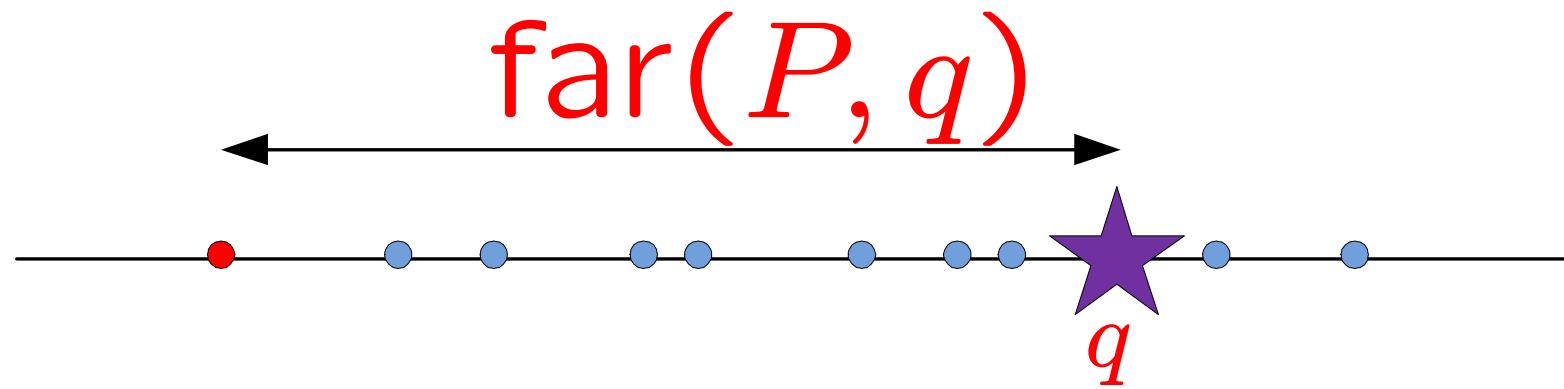
Coreset for Enclosing Balls $P \subseteq \mathbb{R}$



Coreset for Enclosing Balls $P \subseteq \mathbb{R}$



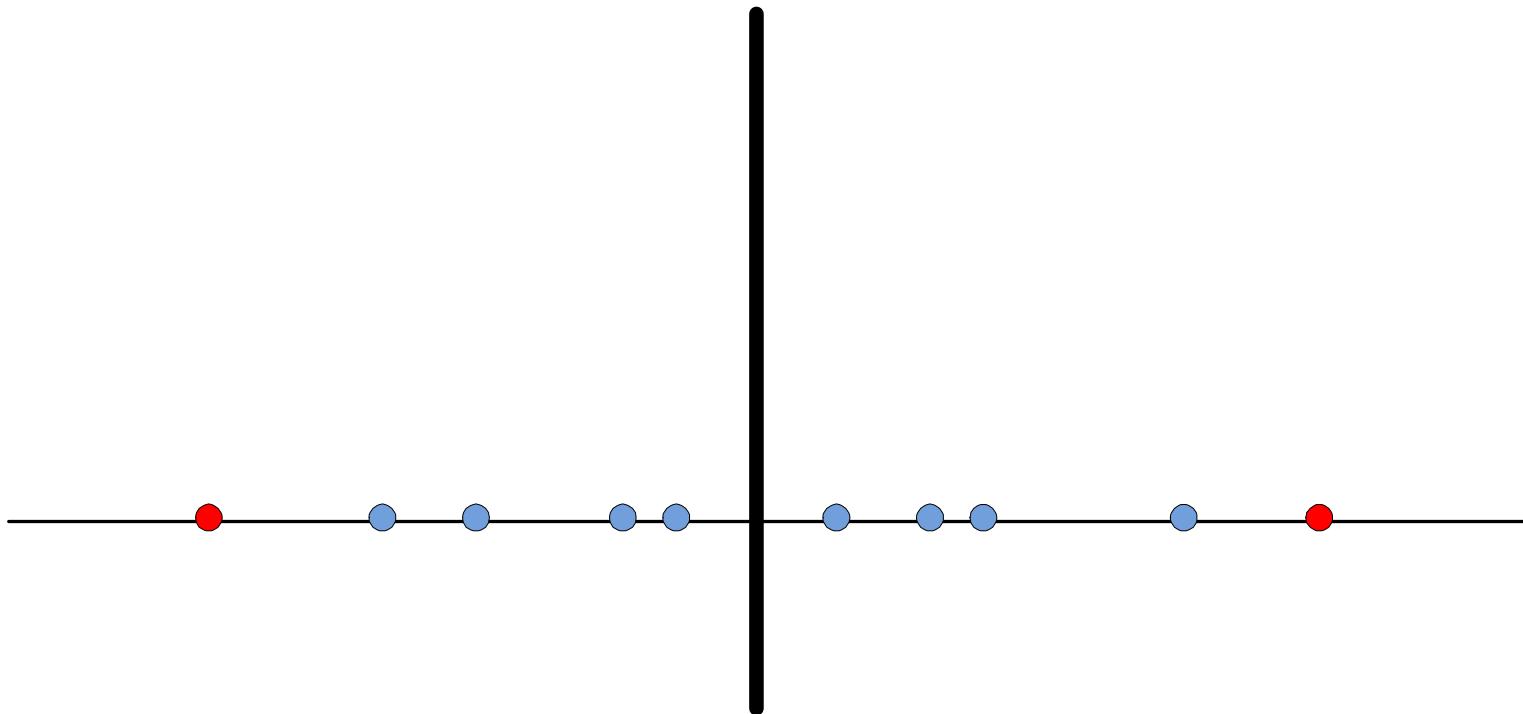
Coreset for Enclosing Balls $P \subseteq \mathbb{R}$



Coreset for Enclosing Balls $P \subseteq \mathbb{R}$

The farthest point from every query $q \in \mathbb{R}$

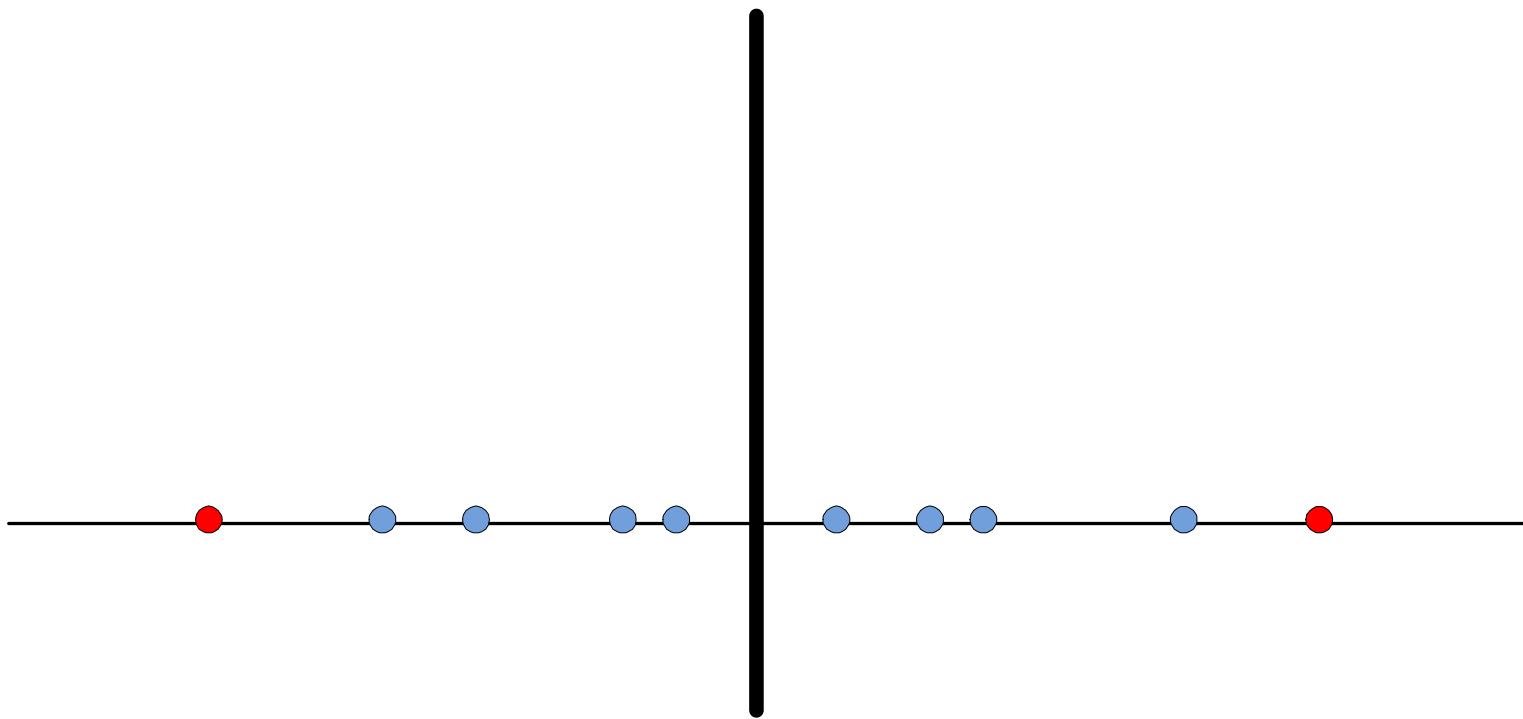
is a red point



Coreset for Enclosing Balls $P \subseteq \mathbb{R}^d$

The fattest point from every query $q \in \mathbb{R}^d$

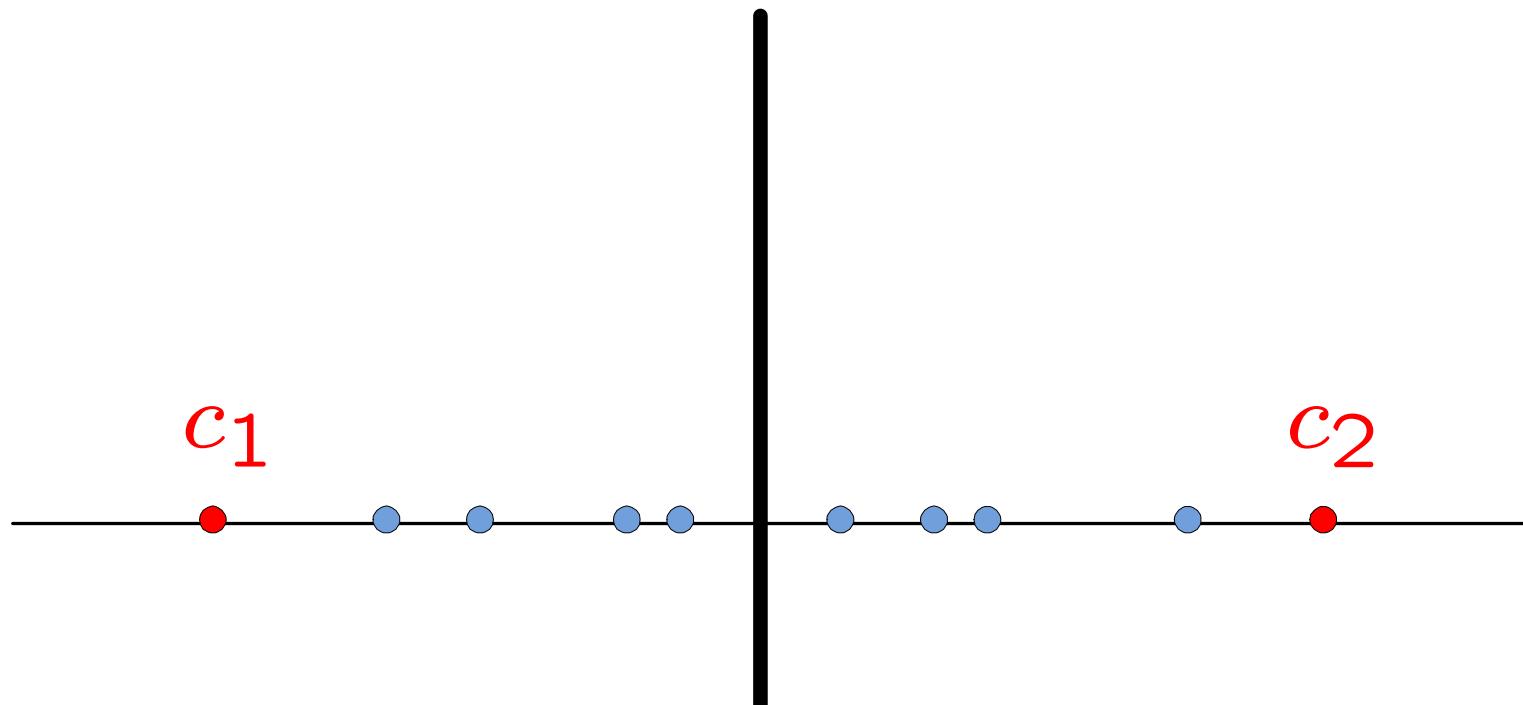
is a red point



Coreset for Enclosing Balls $P \subseteq \mathbb{R}^d$

The fattest point from every query $q \in \mathbb{R}^d$

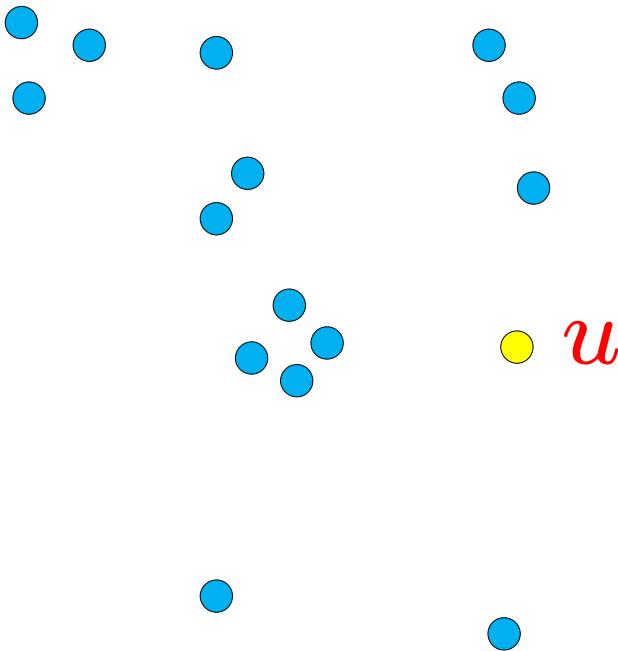
is a red point



$C := \{c_1, c_2\}$ is a coresset for P

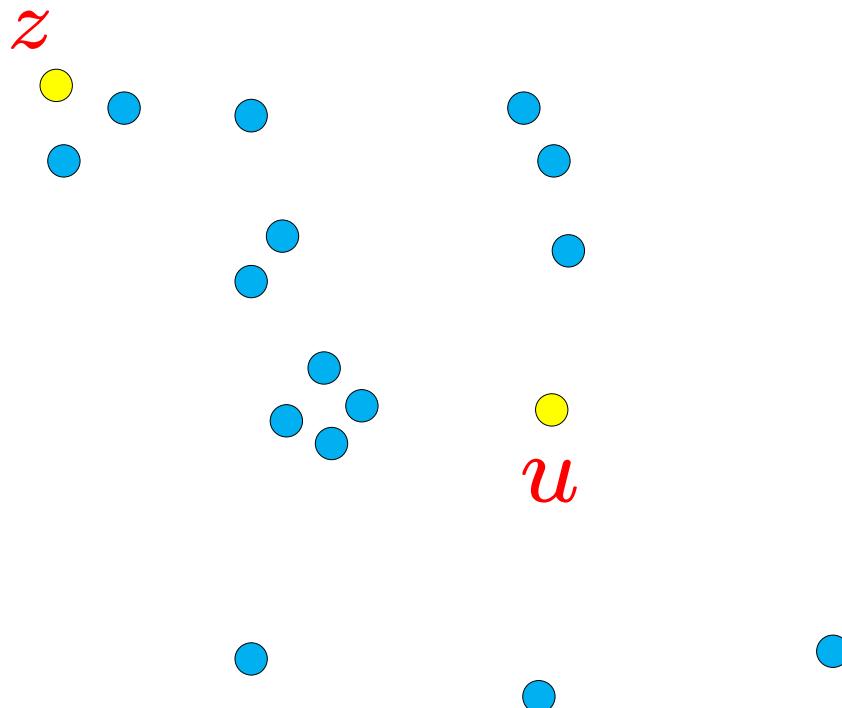
Coreset for Enclosing Balls

1) Choose an arbitrary point $u \in P$



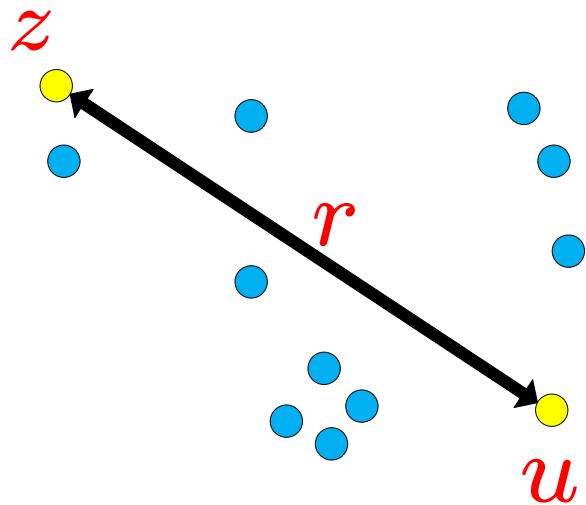
Coreset for Enclosing Balls

2) Find the farthest point $z \in P$ from u



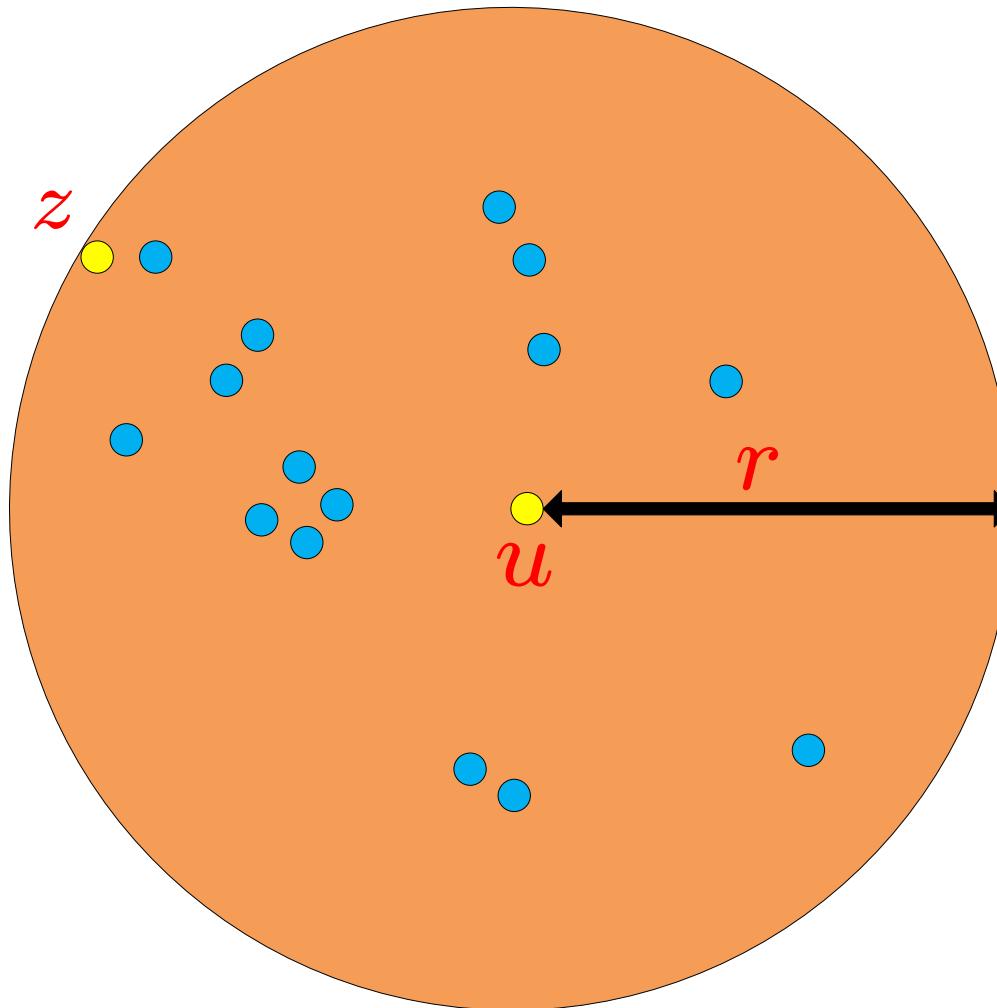
Coreset for Enclosing Balls

$$r := \text{dist}(u, z)$$



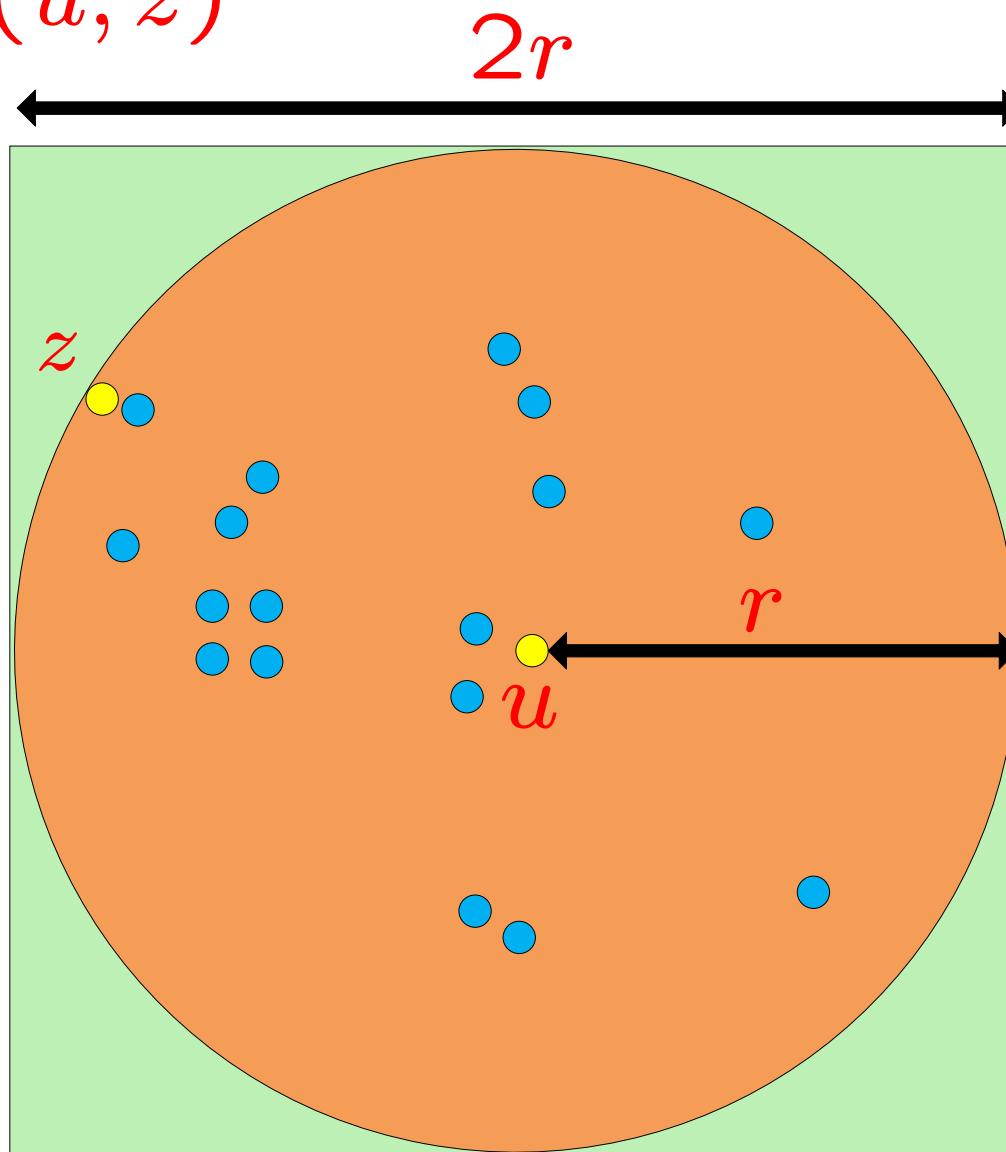
Coreset for Enclosing Balls

$$r := \text{dist}(u, z)$$



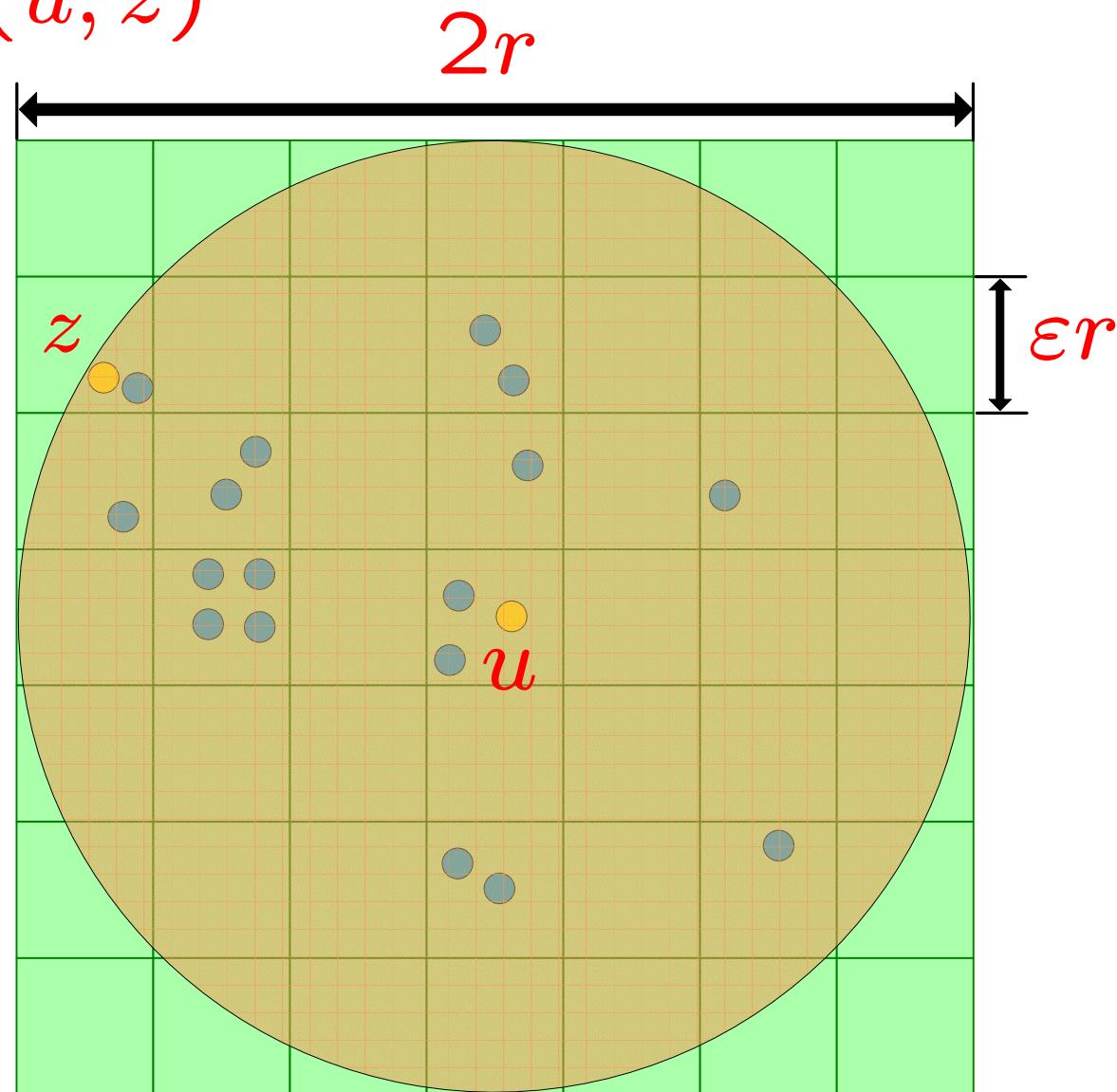
Coreset for Enclosing Balls

$$r := \text{dist}(u, z)$$



Coreset for Enclosing Balls

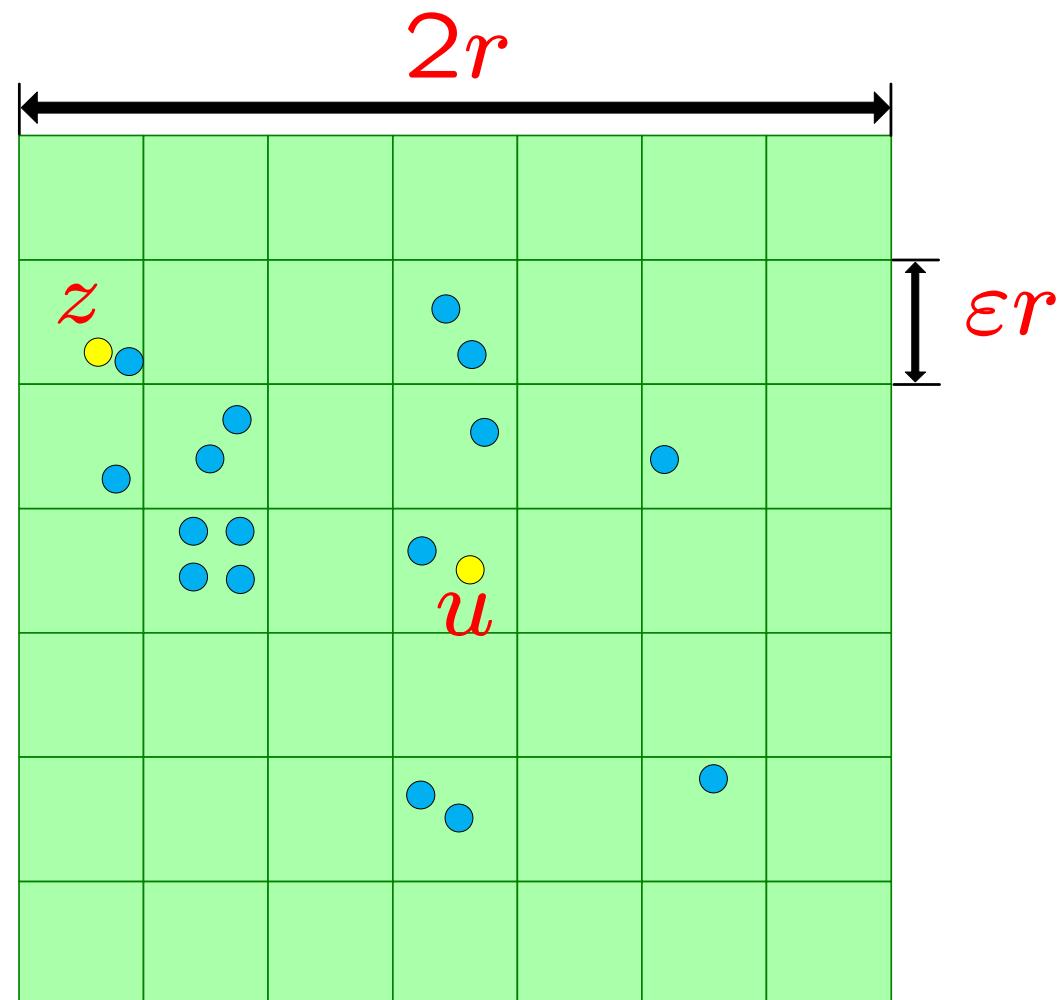
$$r := \text{dist}(u, z)$$



Coreset for Enclosing Balls

- 3) Construct a grid of $\frac{2}{\varepsilon^2}$ cells of size $\varepsilon r \times \varepsilon r$, centered at u

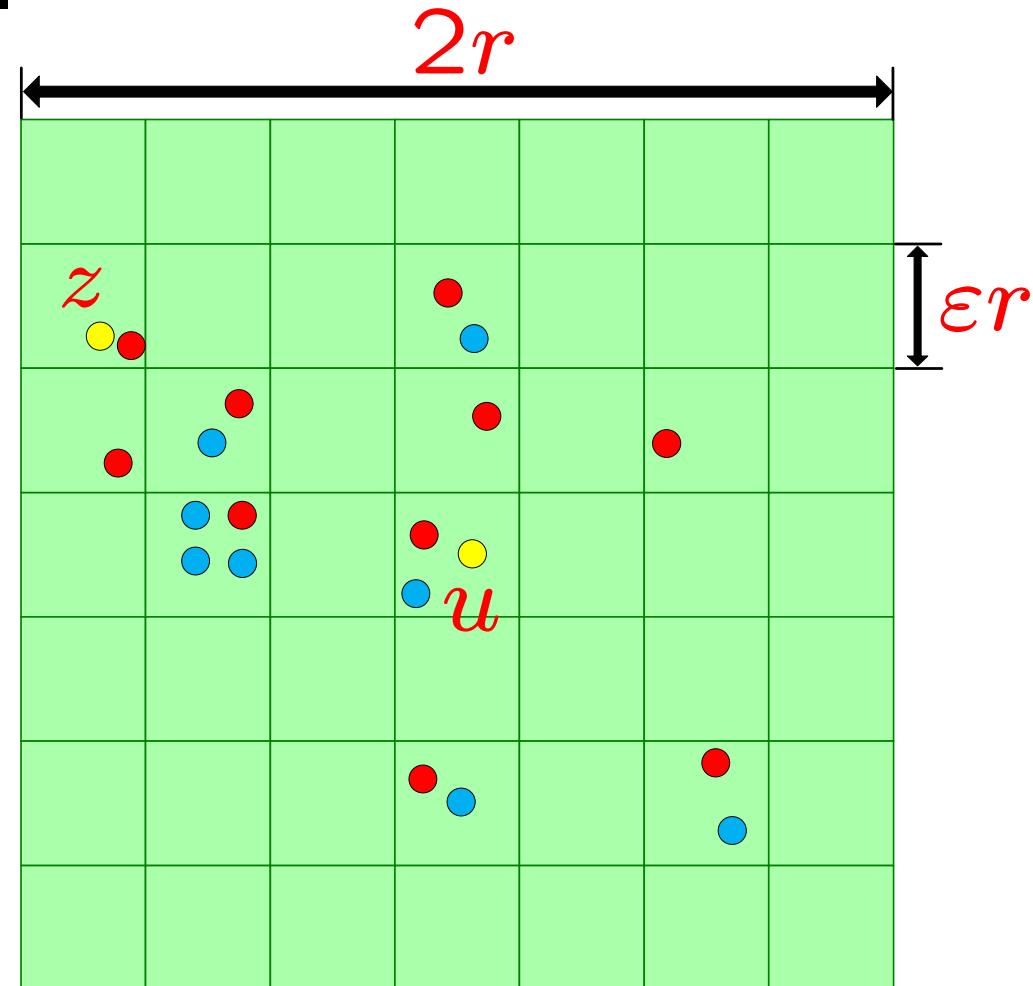
$$r := \text{dist}(u, z)$$



Coreset for Enclosing Balls

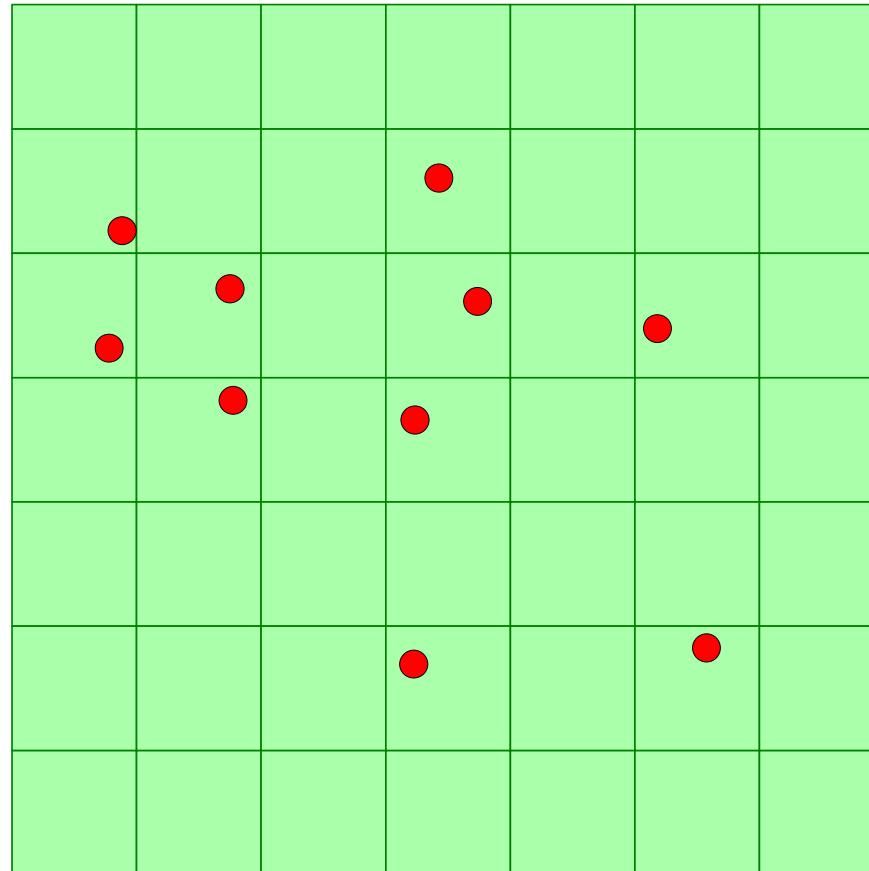
- 4) Pick a representative point from each non-empty cell

$$r := \text{dist}(u, z)$$



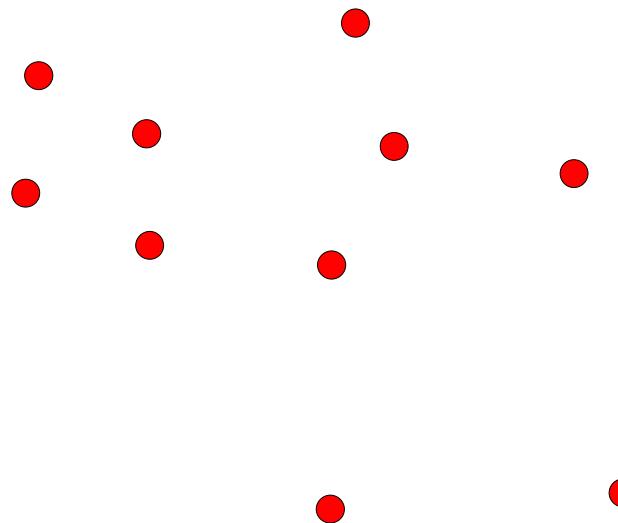
Coreset for Enclosing Balls

5) $C :=$ the set of the $O\left(\frac{1}{\varepsilon^2}\right)$ representatives



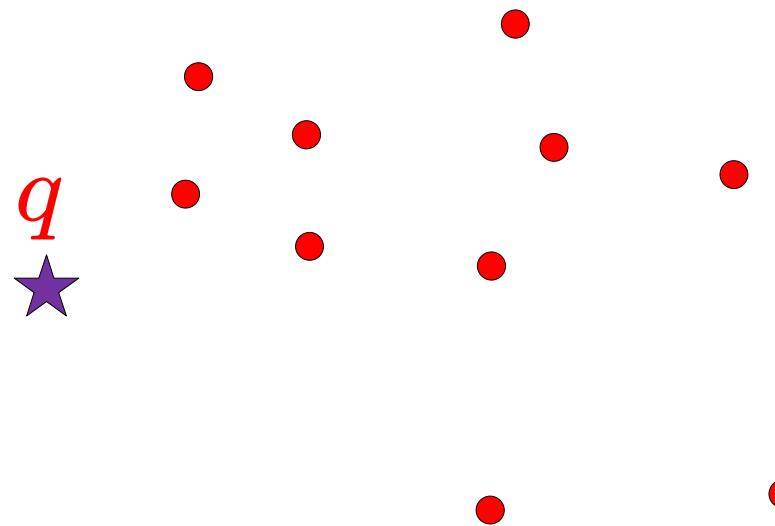
Coreset for Enclosing Balls

6) Return C



Proof of Correctness

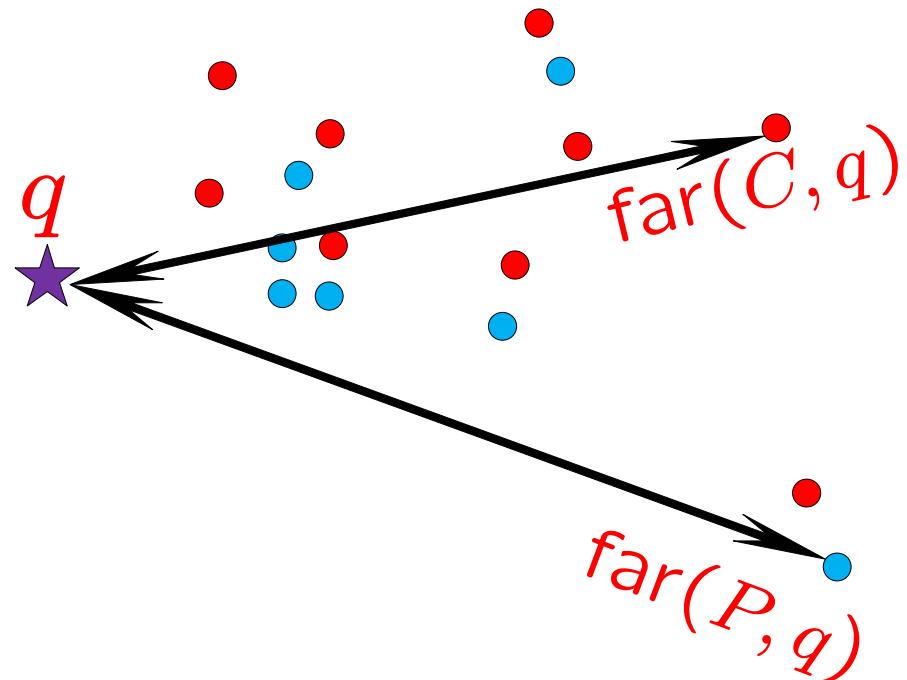
$q :=$ an arbitrary query point



Proof of Correctness

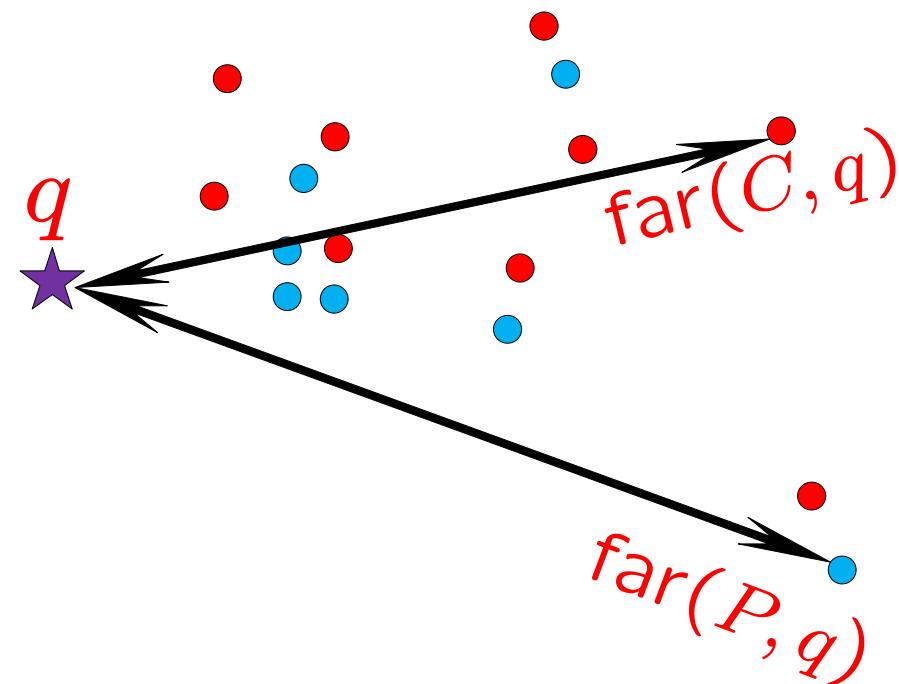
$$\text{far}(P, q) = \max_{p \in P} \text{dist}(p, q)$$

$$\text{far}(C, q) = \max_{c \in C} \text{dist}(c, q)$$



Proof of Correctness

$$C \subseteq P \quad \longrightarrow \quad \text{far}(C, q) \leq \text{far}(P, q)$$

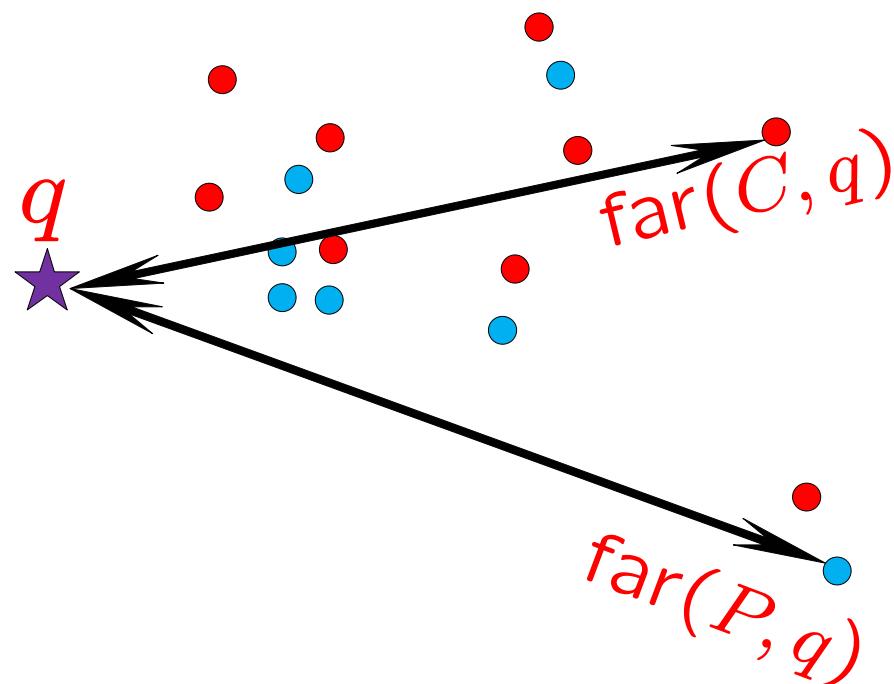


Proof of Correctness

$$C \subseteq P \quad \longrightarrow \quad \text{far}(C, q) \leq \text{far}(P, q)$$

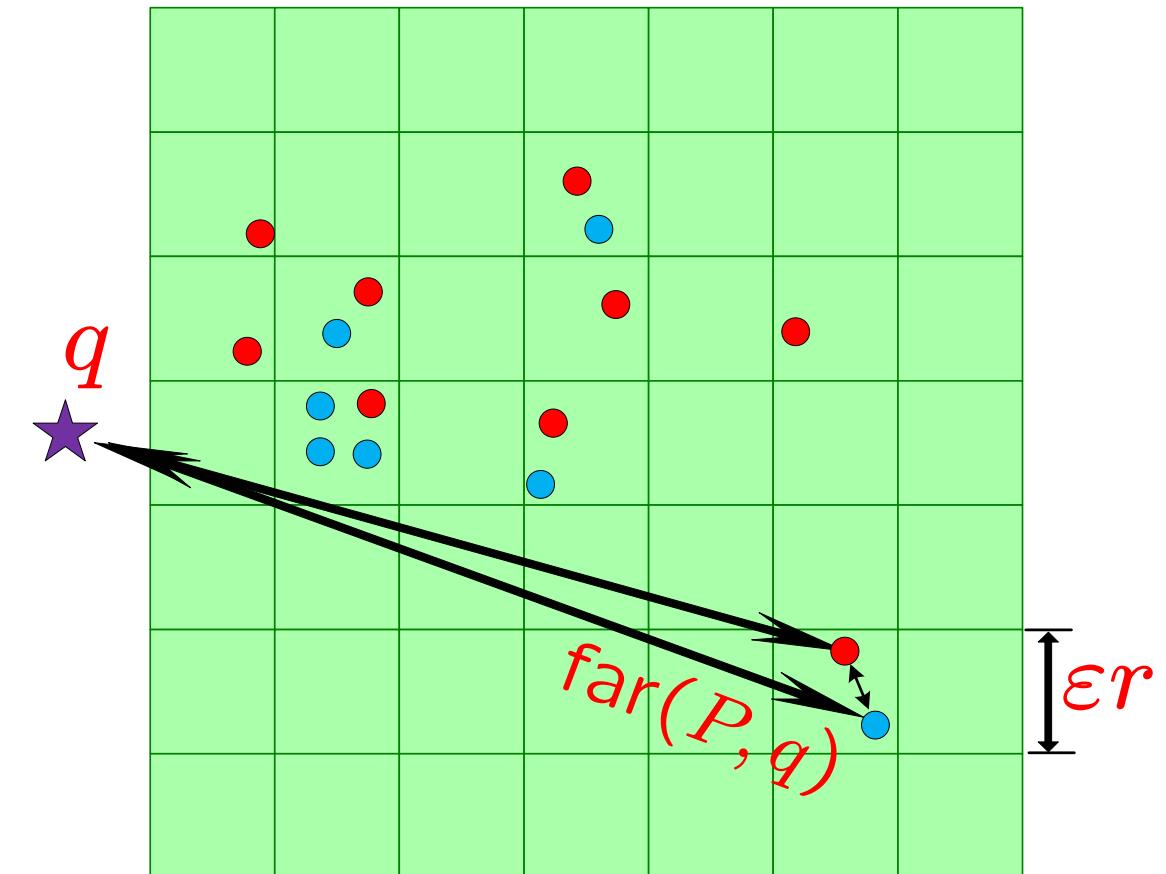
Need to prove:

$$\text{far}(P, q) - \text{far}(C, q) \leq O(\varepsilon) \text{far}(P, q)$$



Proof of Correctness

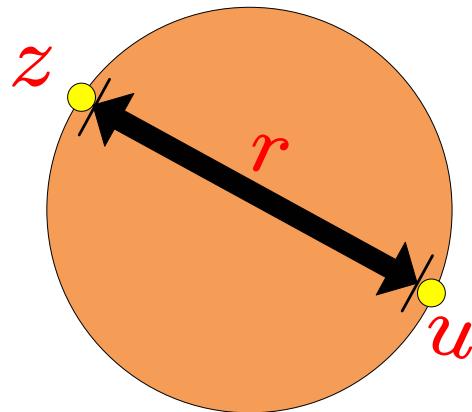
$$\text{far}(P, q) \leq \text{far}(C, q) + O(\varepsilon r)$$



Main Observation

Every ball that covers u and z ,

has a diameter of at least r :



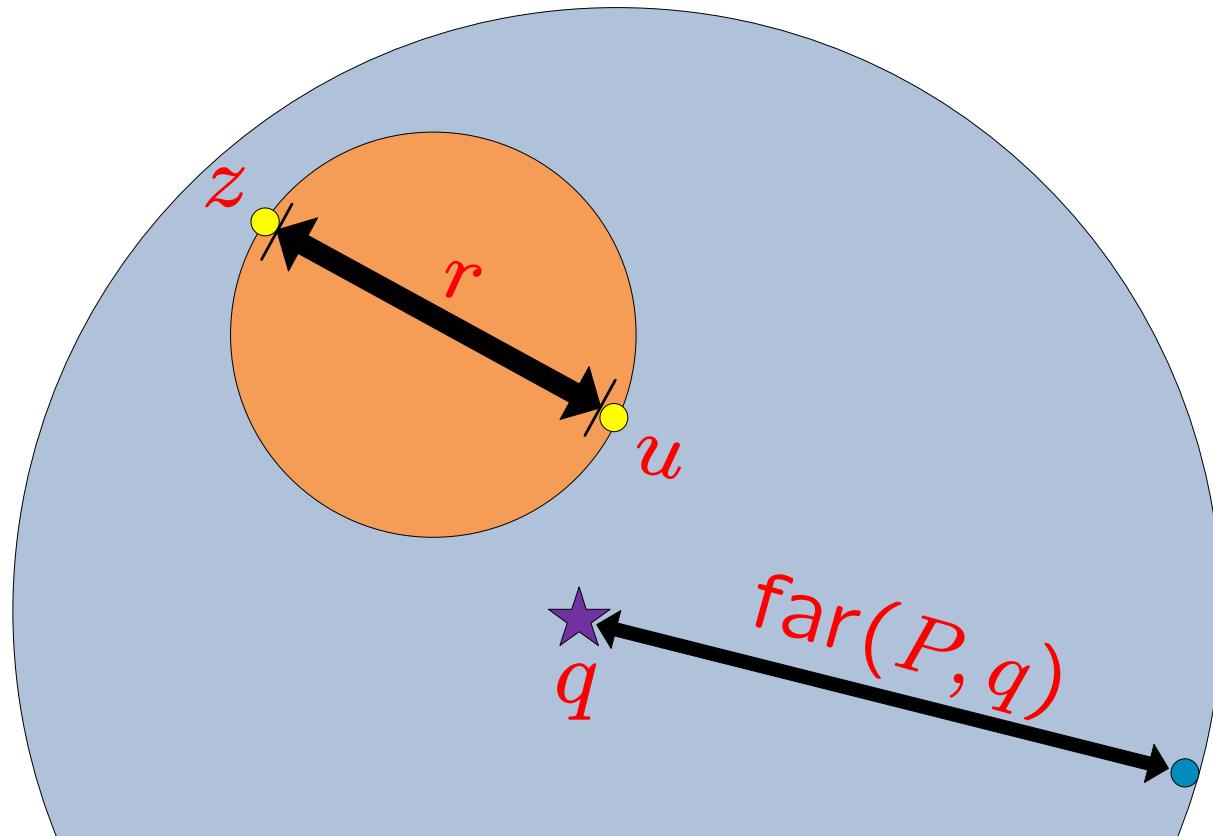
|

Minimum enclosing ball of $\{u, z\}$

Main Observation

Every ball that covers u and z ,

has a diameter of at least r :

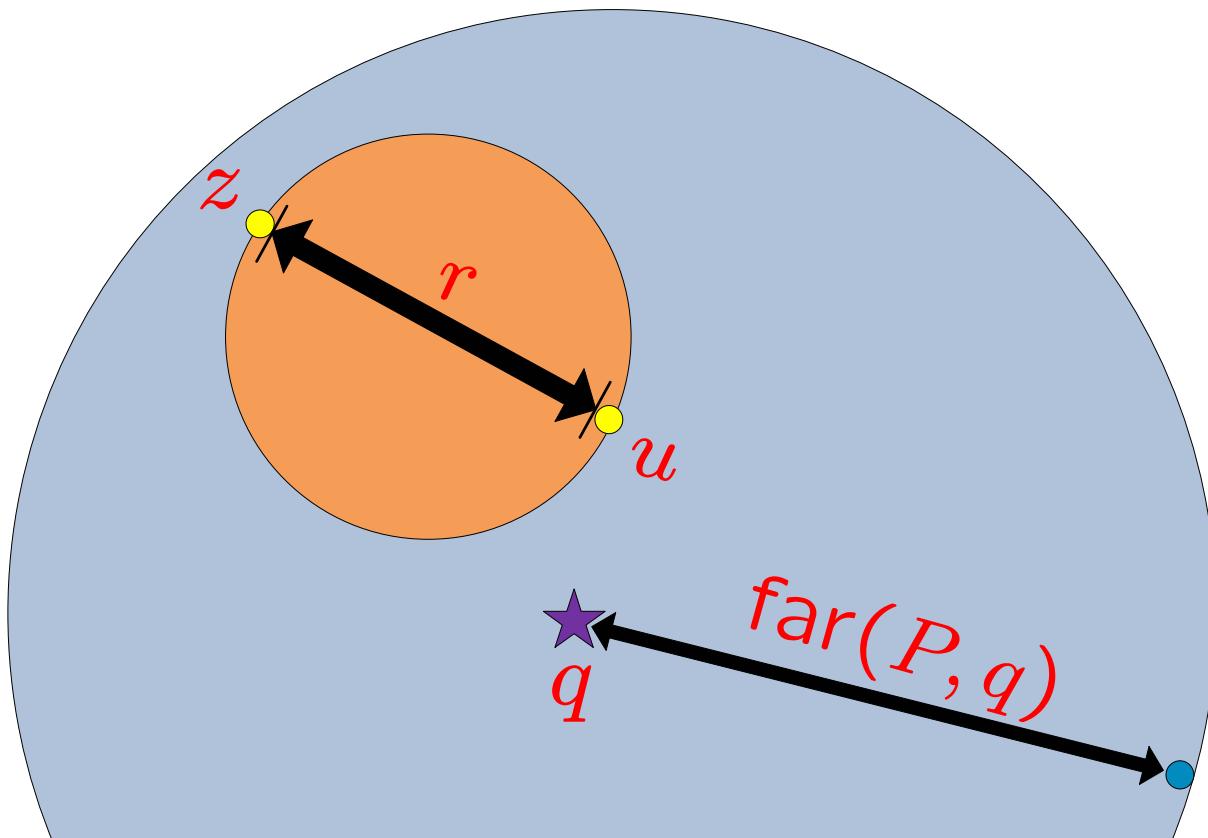


$$r \leq 2\text{far}(P, q)$$

Main Observation

Every ball that covers u and z ,

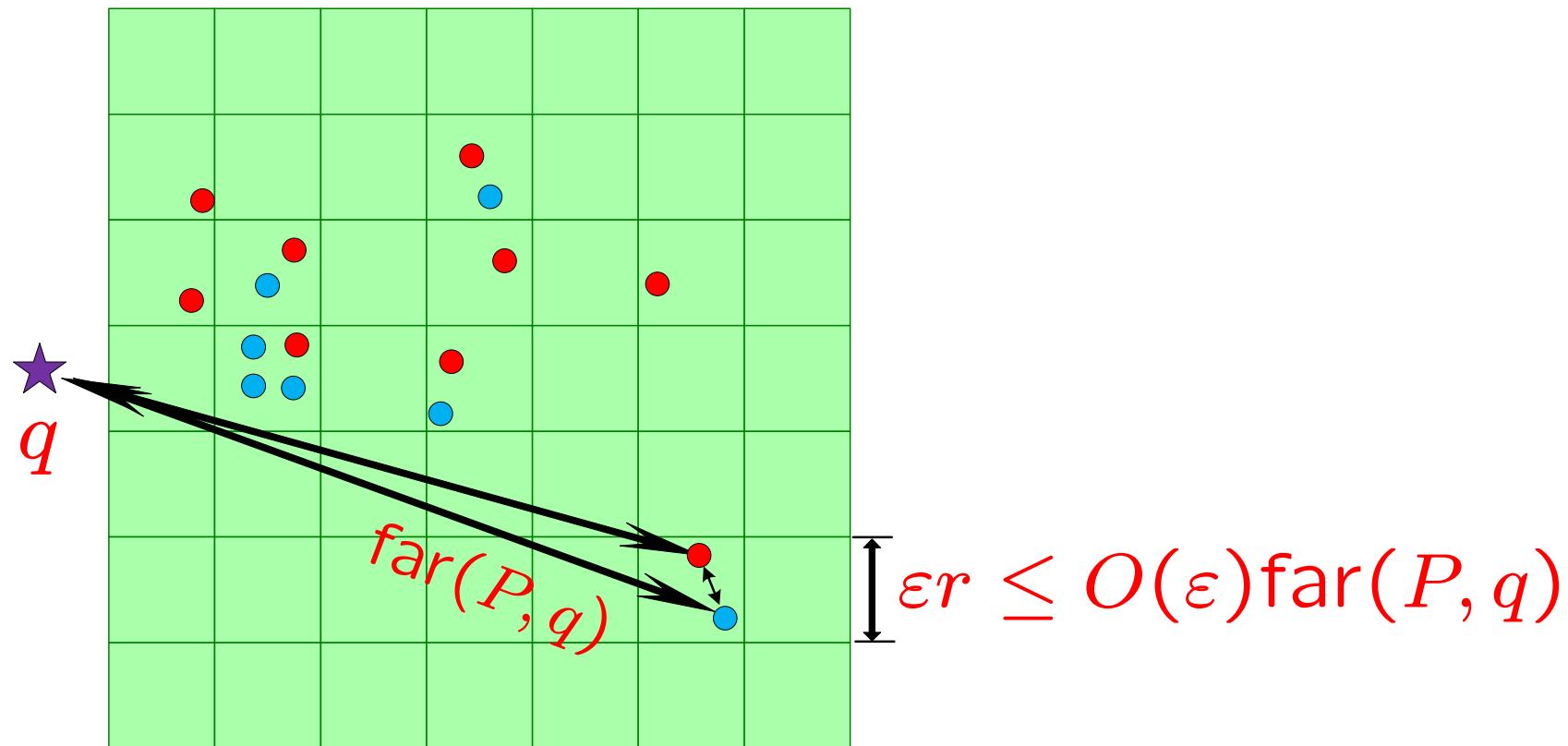
has a diameter of at least r :



$$r \leq 2\text{far}(P, q) \quad \Rightarrow \quad O(\varepsilon r) \leq O(\varepsilon)\text{far}(P, q)$$

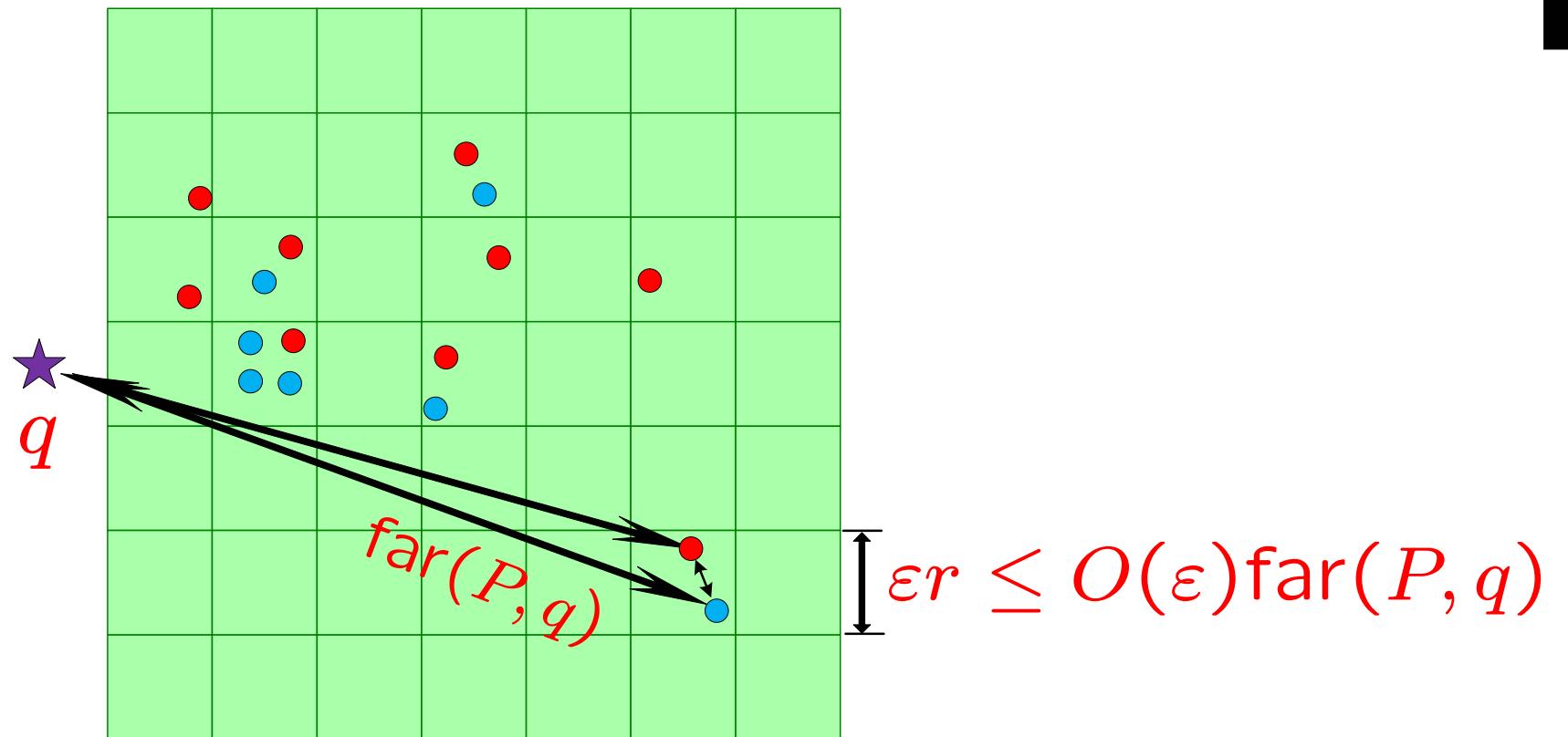
Proof of Correctness

$$\begin{aligned}\text{far}(P, q) &\leq \text{far}(C, q) + O(\varepsilon r) \\ &\leq \text{far}(C, q) + O(\varepsilon) \text{far}(P, q)\end{aligned}$$



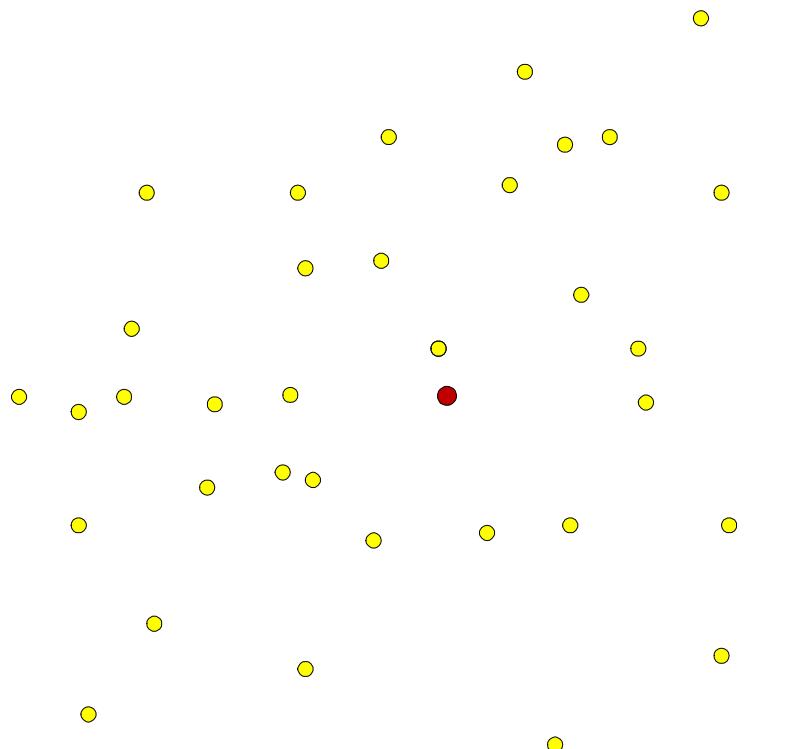
Proof of Correctness

$$\begin{aligned}\text{far}(P, q) &\leq \text{far}(C, q) + O(\varepsilon r) \\ &\leq \text{far}(C, q) + O(\varepsilon) \text{far}(P, q)\end{aligned}$$



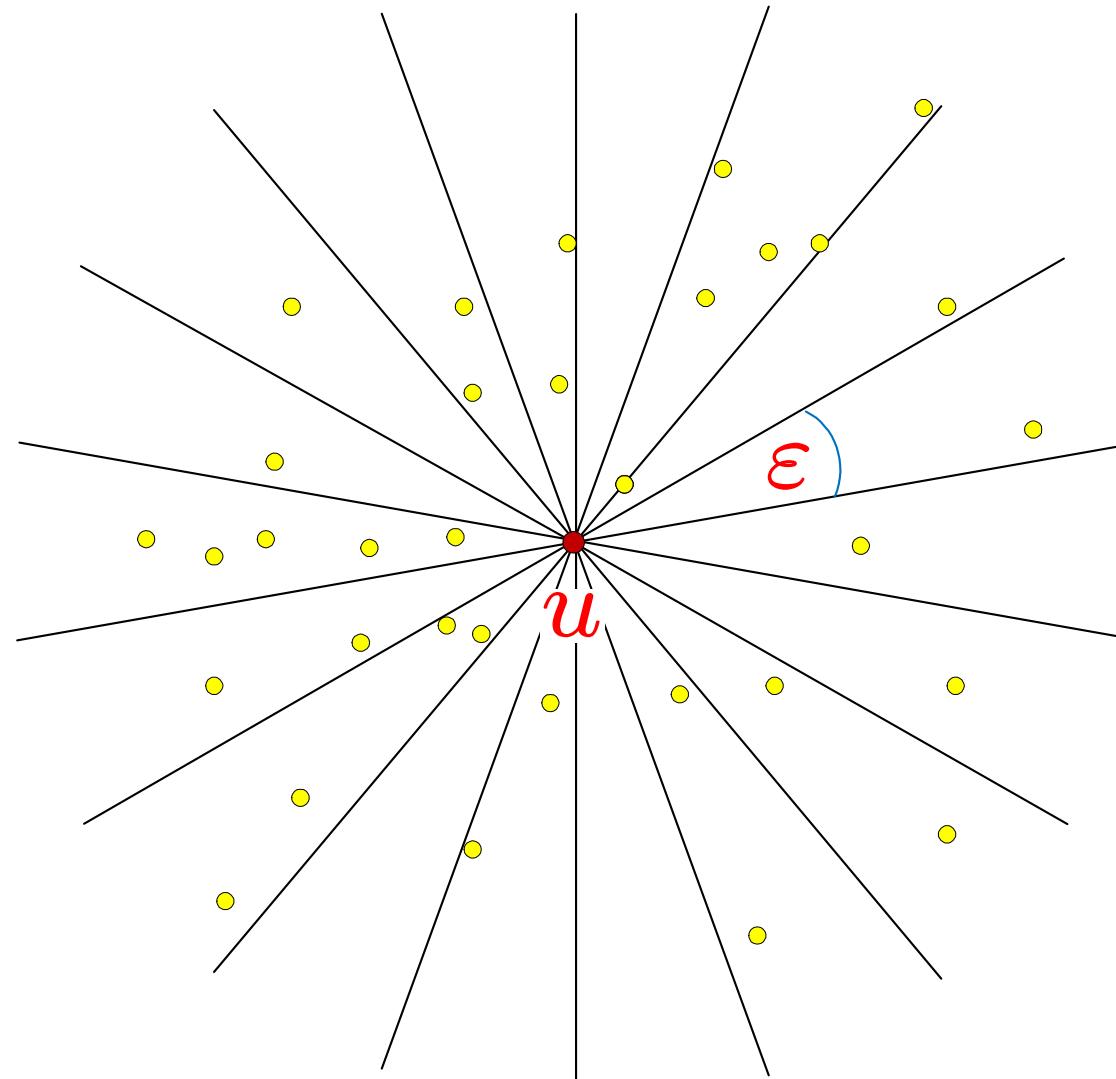
Smaller Coreset

1) Choose an arbitrary point $u \in P$



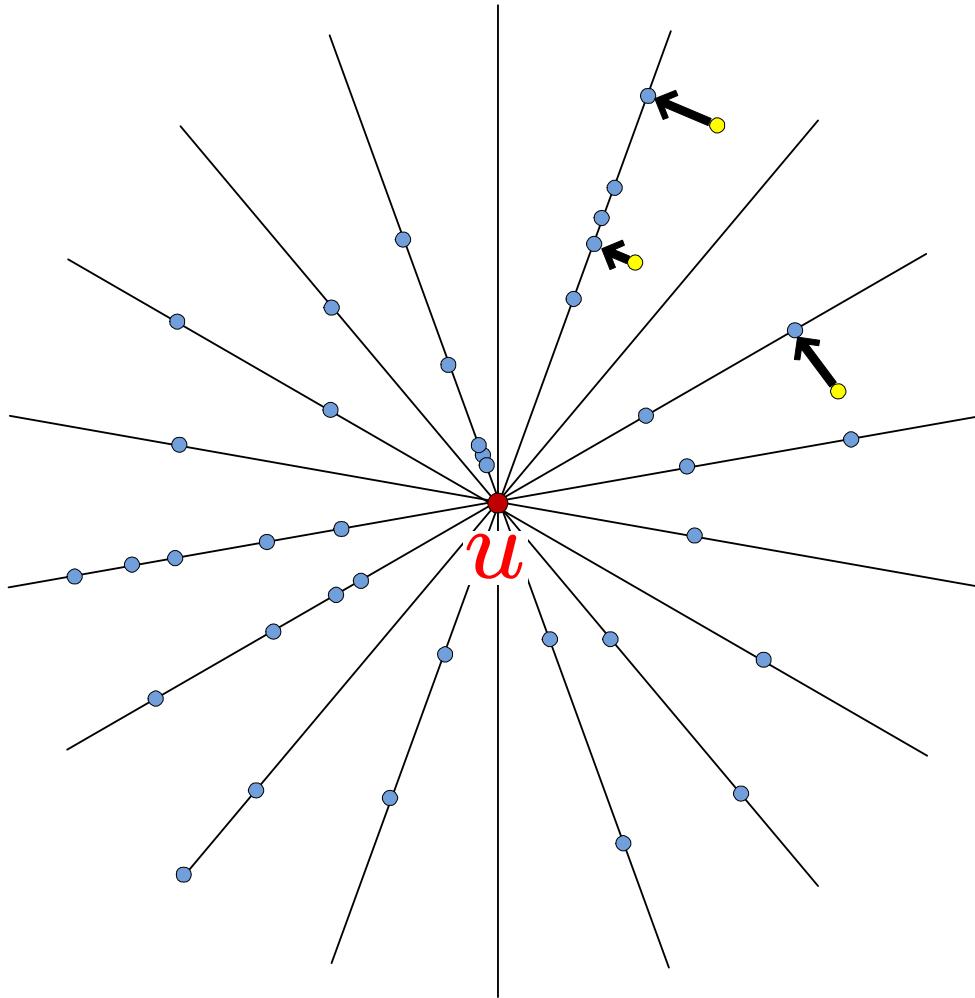
Smaller Coreset

2) Draw a "star" of $\frac{2\pi}{\varepsilon}$ lines around u



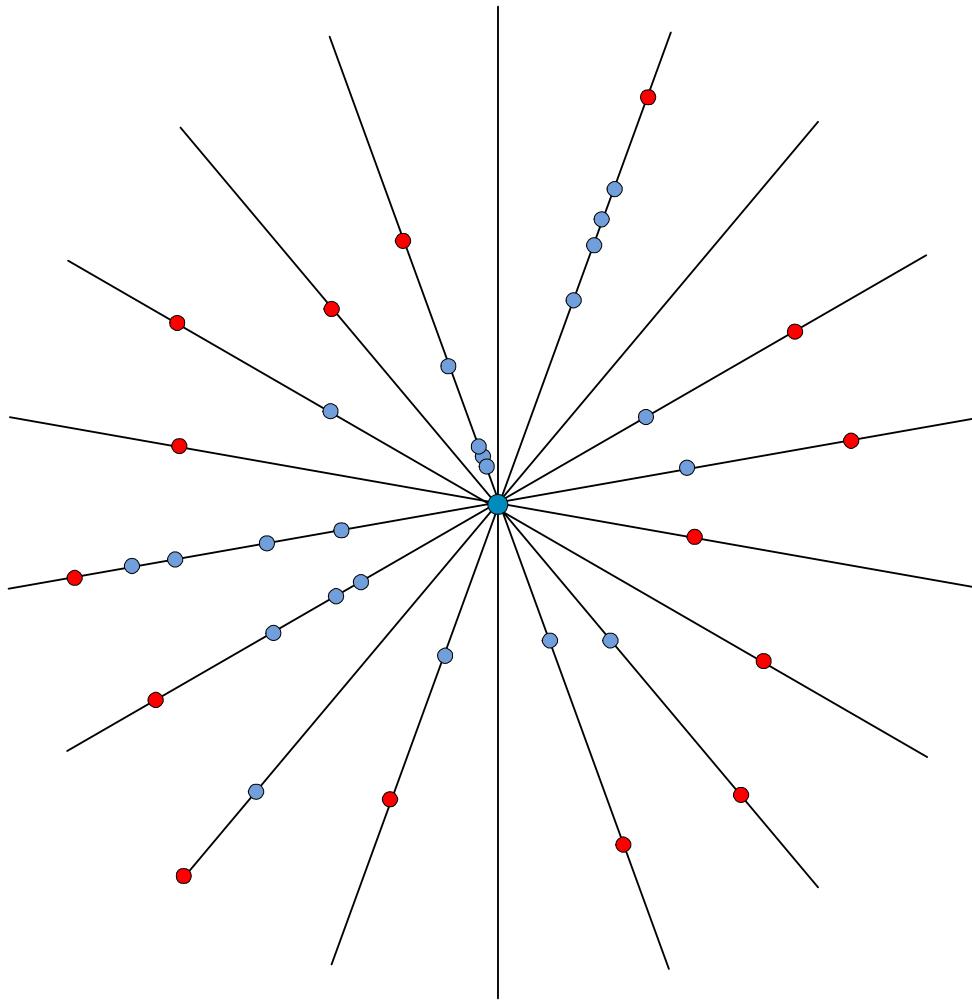
Smaller Coreset

3) $P' :=$ Projection of P onto the lines



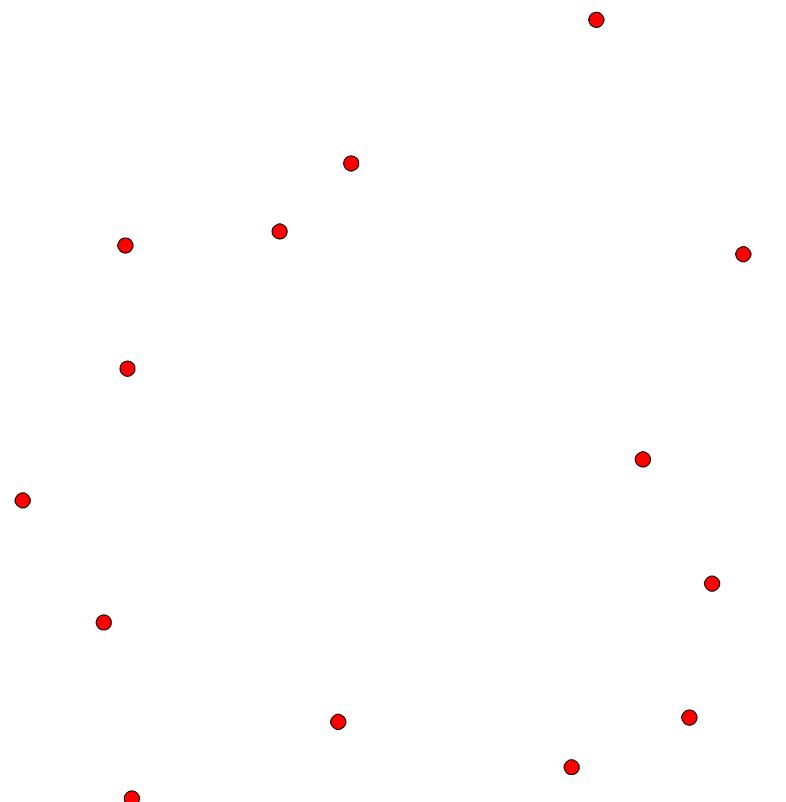
Smaller Coreset

4) $C :=$ union of endpoints on the lines



Smaller Coreset

5) Return C



Proof: Overview

C is a coresnet for P'

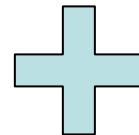
Proof: Overview

C is a coresnet for P'

P' is a coresnet for P

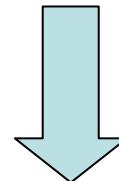
Proof: Overview

C is a coresset for P'



P' is a coresset for P

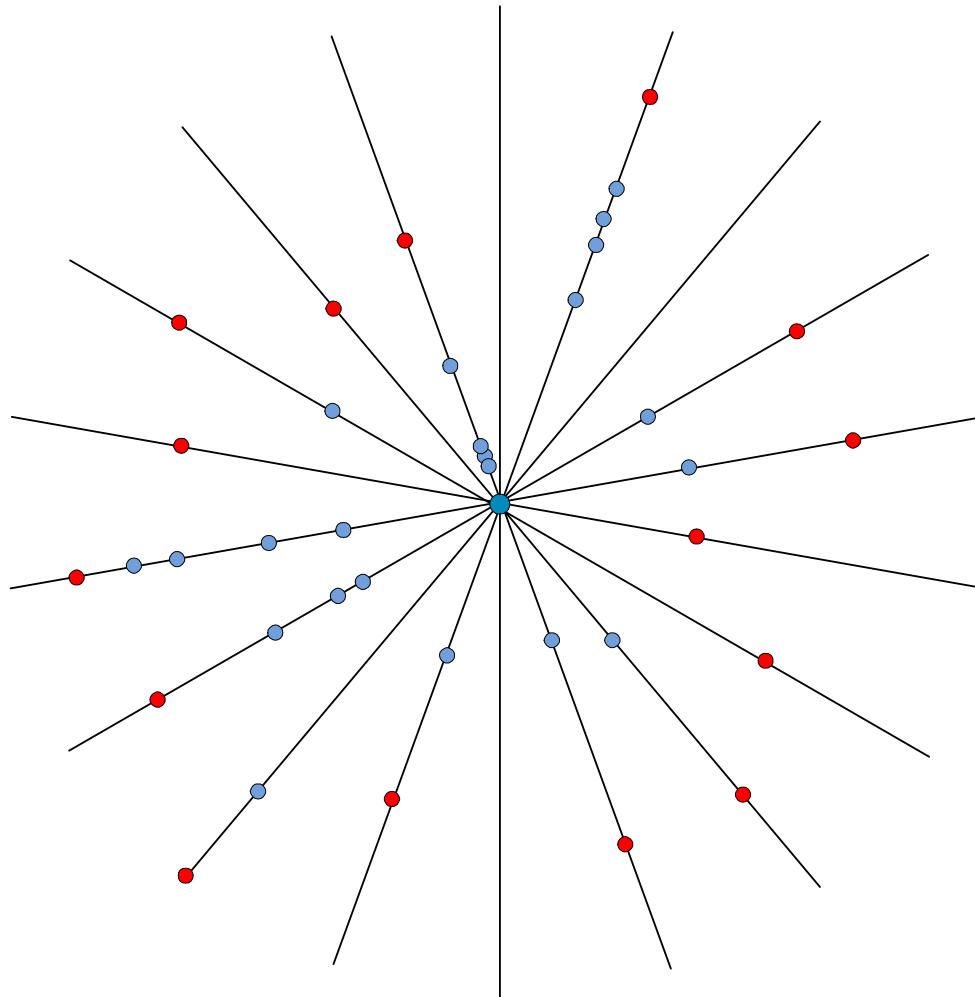
(Large coresset but on few lines)



C is a coresset for P

(Transitive Property)

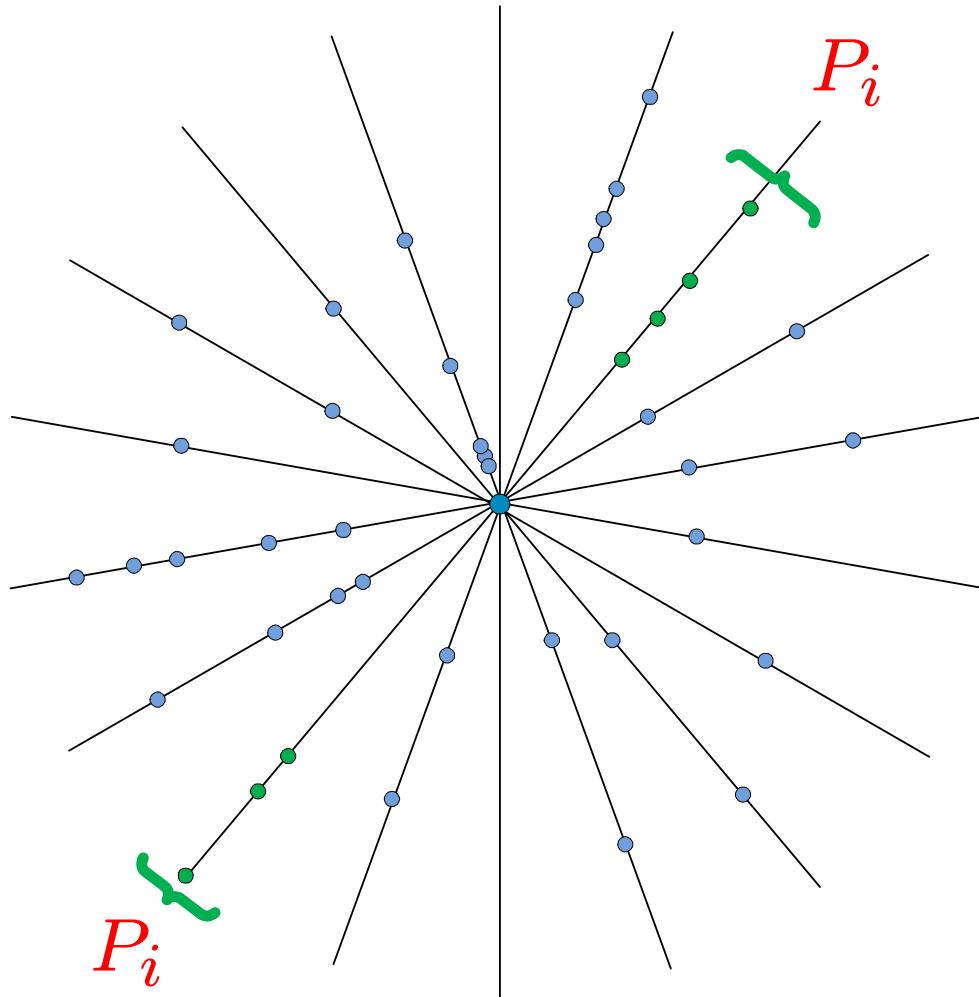
Claim: C is a coresnet for P'



$$P' := \bullet \cup \bullet$$
$$C := \bullet$$

Claim: C is a coresnet for P'

$P_i :=$ intersection of P' with the i th line

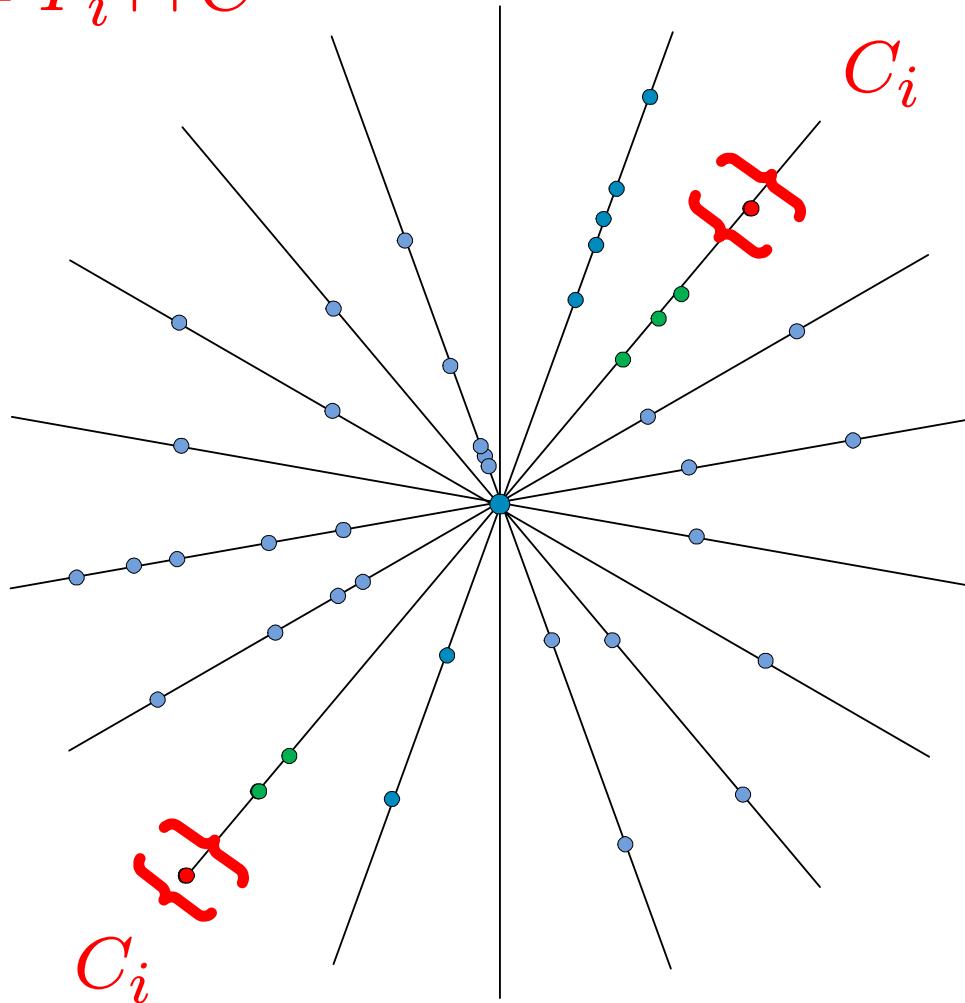


$$P' := \bullet \cup \bullet$$
$$P_i := \bullet$$

Claim: C is a coresset for P'

$P_i :=$ intersection of P' with the i th line

$C_i := P_i \cap C$

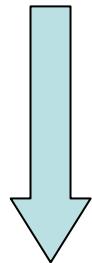


$$P_i := \bullet \cup \textcolor{red}{\bullet}$$
$$C_i := \textcolor{red}{\bullet}$$

Claim: C is a coresset for P'

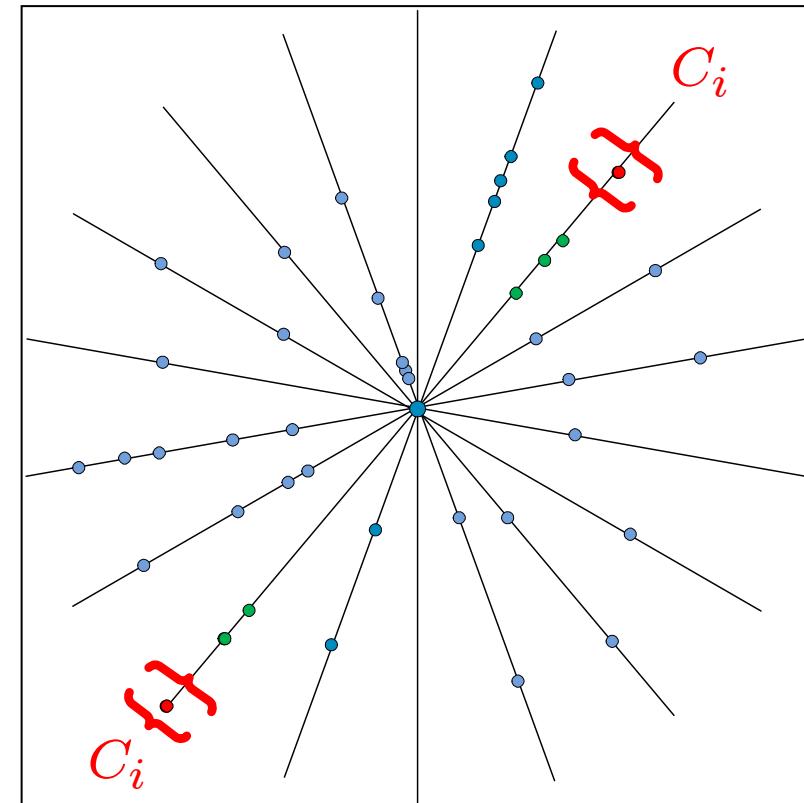
$P_i :=$ intersection of P' with the i th line

$C_i := P_i \cap C$



By proof for $d = 1$

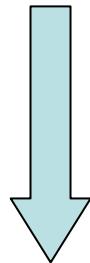
C_i is a coresset for P_i



Claim: C is a coresset for P'

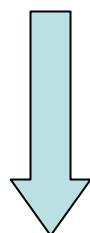
$P_i :=$ intersection of P' with the i th line

$C_i := P_i \cap C$



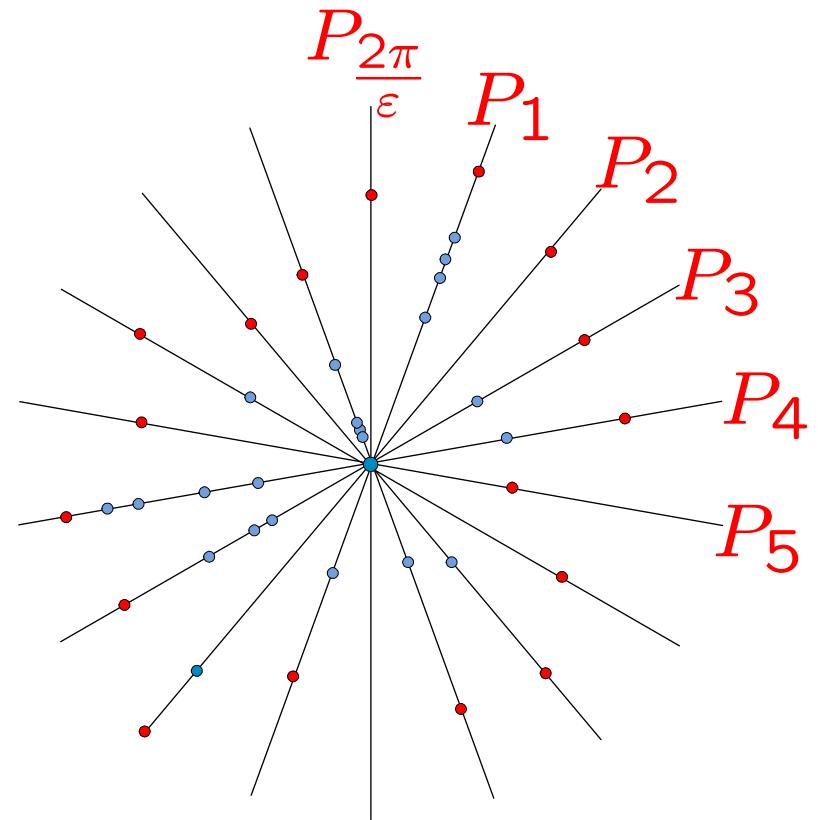
By proof for $d = 1$

C_i is a coresset for P_i



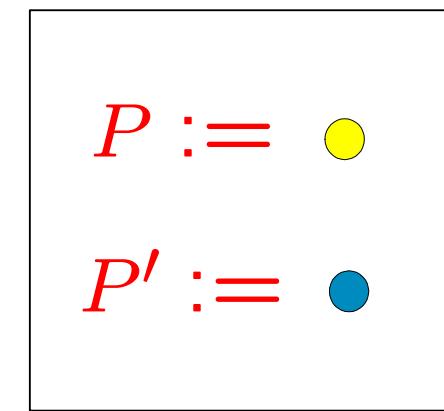
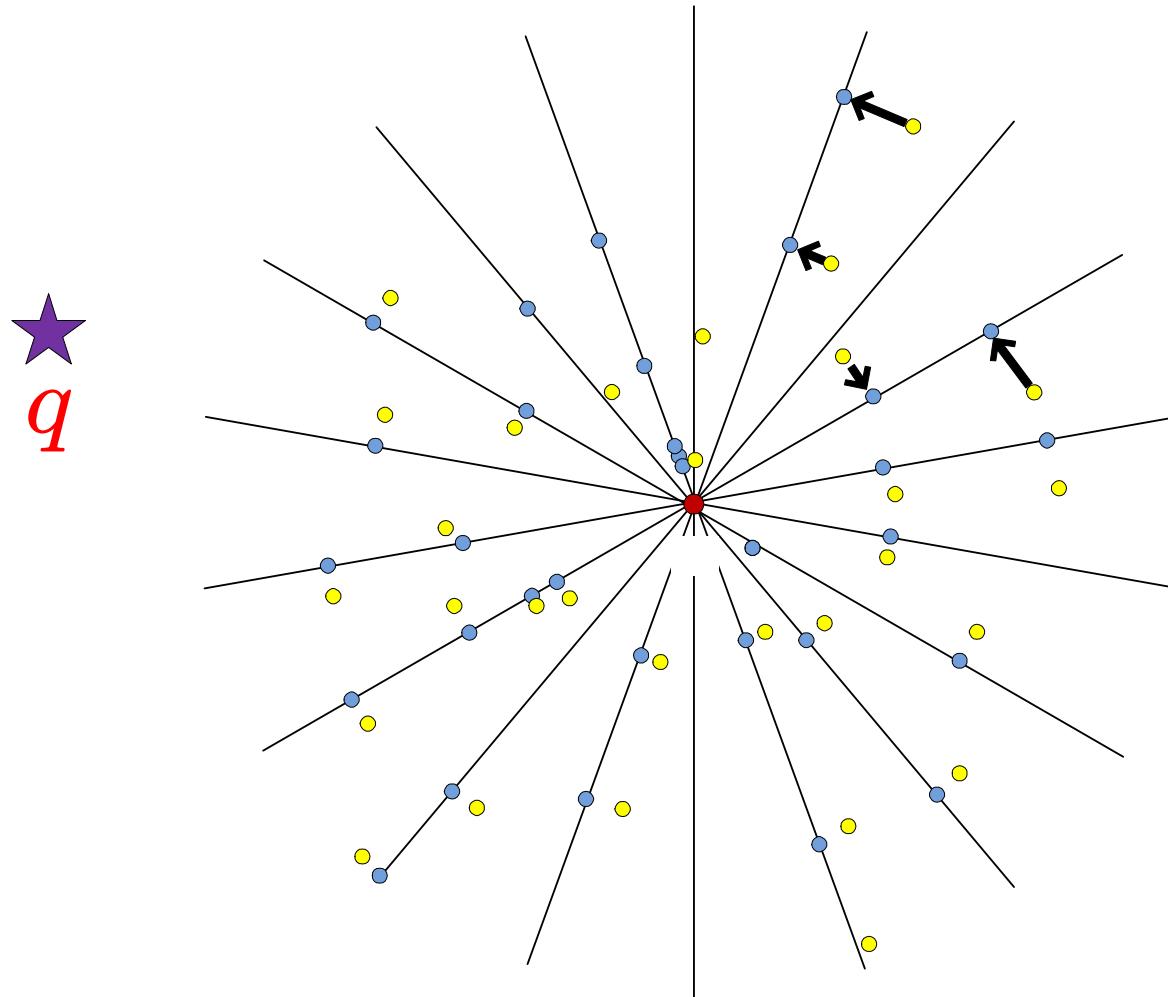
Union Rule

$C = \bigcup_i C_i$ is a coresset for $P' = \bigcup_i P_i$



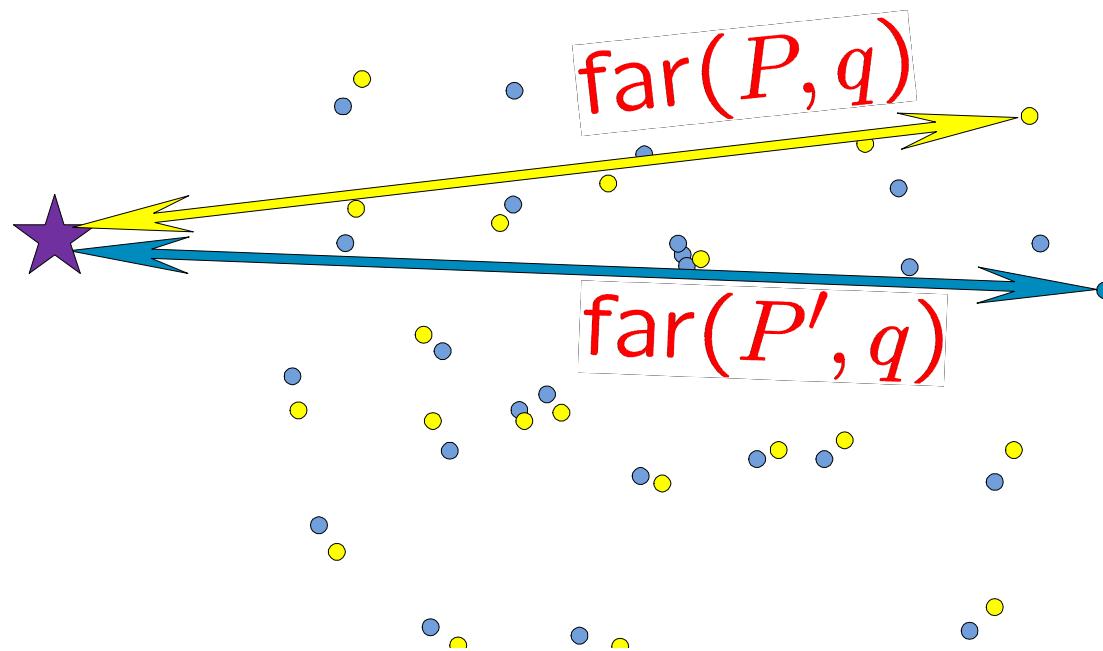
Claim: P' is a coresset for P

$q :=$ an arbitrary query point



Claim: P' is a coresset for P

$q :=$ an arbitrary query point



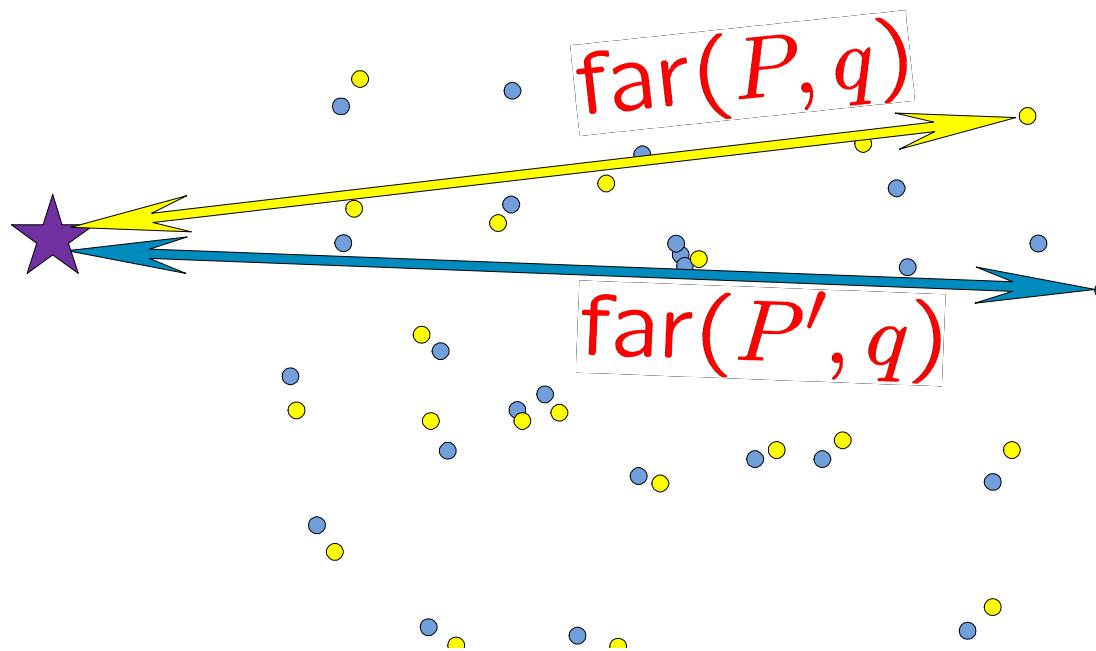
$P :=$

$P' :=$

Claim: P' is a coresset for P

$q :=$ an arbitrary query point

$$\cancel{P' \subseteq P} \rightarrow \cancel{\text{far}(P', q) \leq \text{far}(P, q)}$$

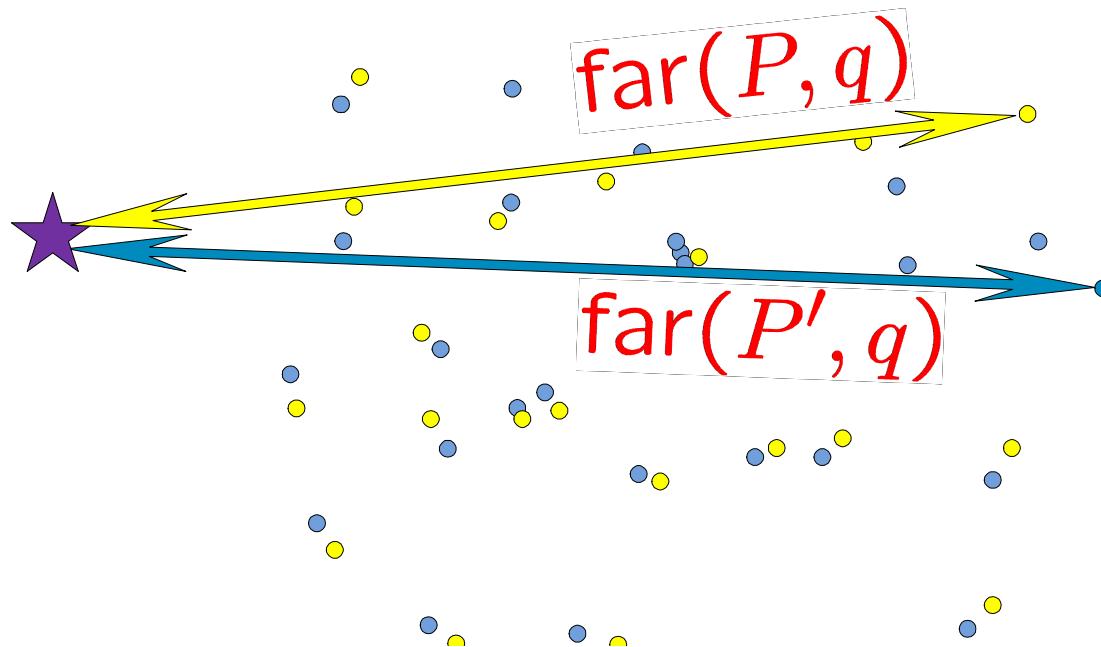


$$\boxed{\begin{aligned} P &:= \text{yellow dot} \\ P' &:= \text{blue dot} \end{aligned}}$$

Claim: P' is a coresset for P

Need to prove:

$$\text{far}(P, q) - \text{far}(P', q) \leq \varepsilon \text{far}(P, q)$$



$P :=$

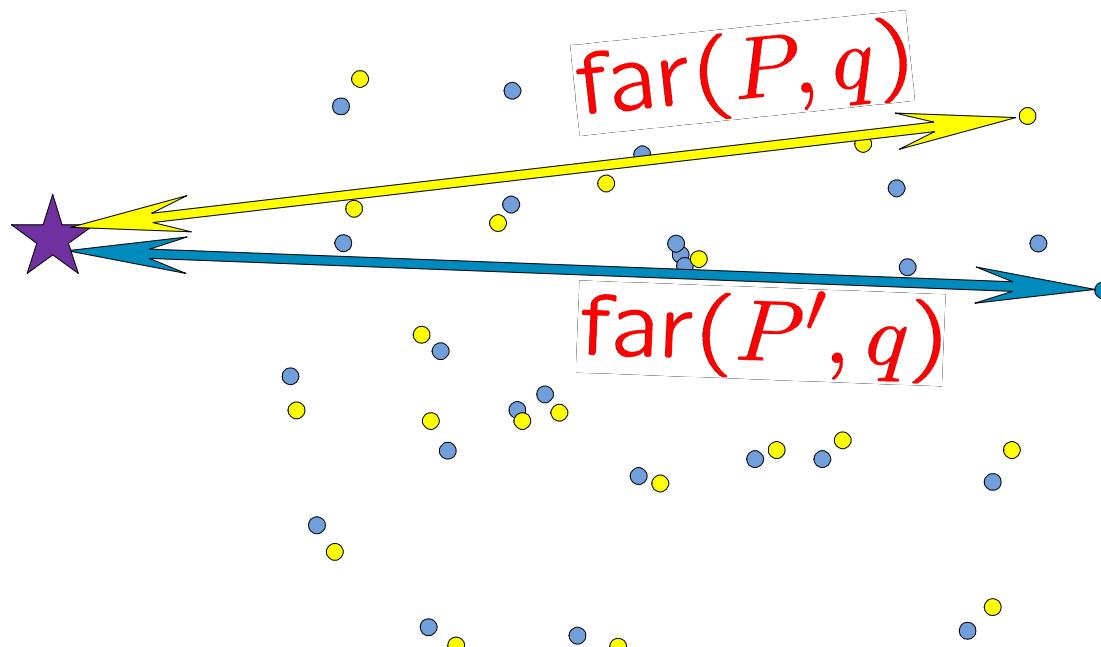
$P' :=$

Claim: P' is a coresset for P

Need to prove:

$$\text{far}(P, q) - \text{far}(P', q) \leq \varepsilon \text{far}(P, q)$$

$$\text{far}(P', q) - \text{far}(P, q) \leq \varepsilon \text{far}(P, q)$$

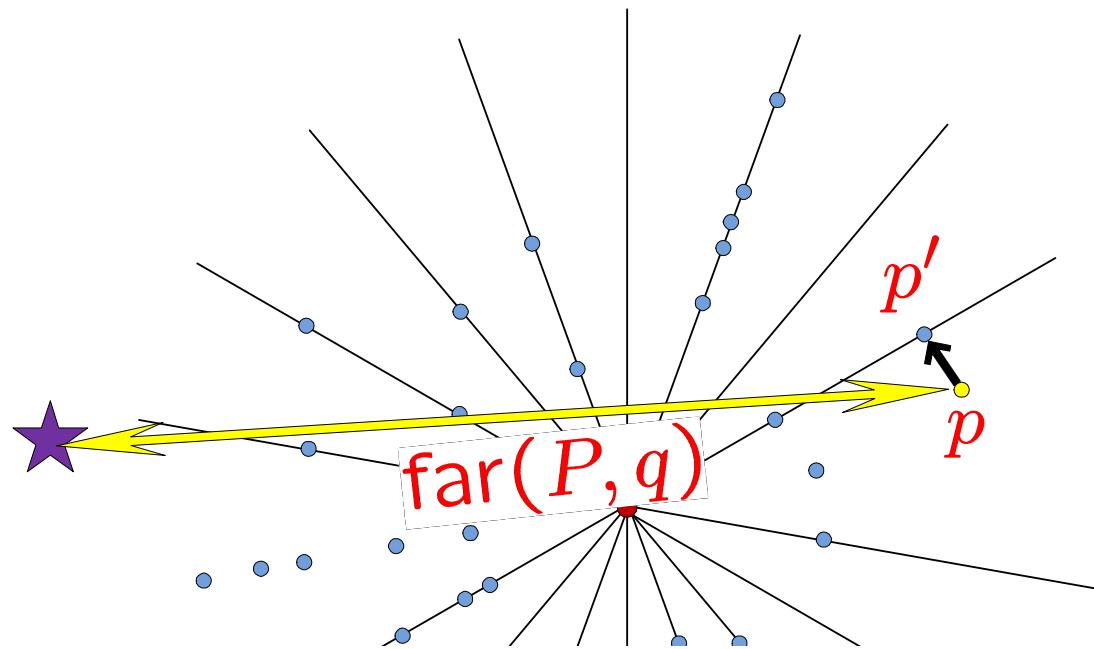


$P :=$	●
$P' :=$	●

Claim: P' is a coresset for P

Let $\text{far}(P, q) = \text{dist}(p, q)$

$p' :=$ the projection of p on the "star"



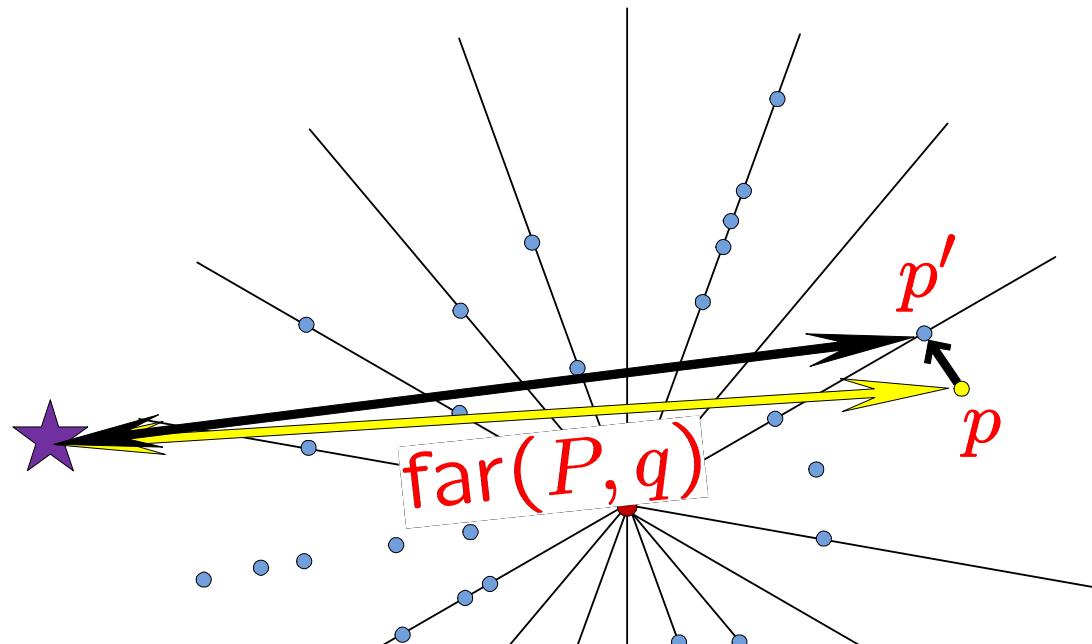
Claim: P' is a coresset for P

Let $\text{far}(P, q) = \text{dist}(p, q)$

$p' :=$ the projection of p on the "star"



$$\text{far}(P, q) - \text{far}(P', q) \leq \text{far}(P, q) - \text{dist}(p', q)$$



Claim: P' is a coresset for P

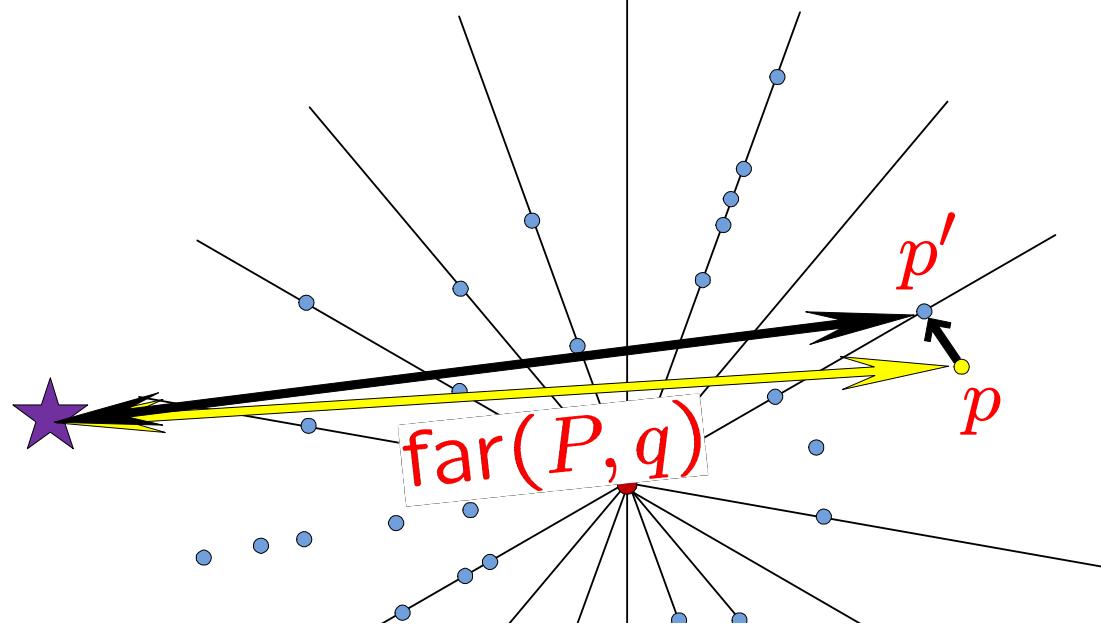
Let $\text{far}(P, q) = \text{dist}(p, q)$

$p' :=$ the projection of p on the "star"



$$\text{far}(P, q) - \text{far}(P', q) \leq \text{far}(P, q) - \text{dist}(p', q)$$

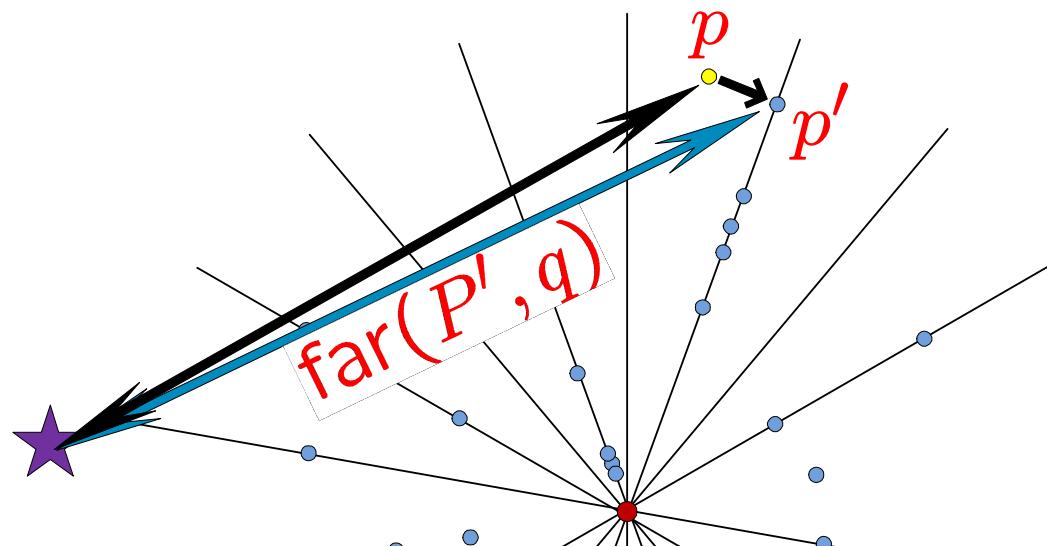
$$\leq \text{dist}(p, p')$$



Claim: P' is a coresset for P

Let $\text{far}(P', q) = \text{dist}(p', q)$

$p :=$ the point whose projection is p'



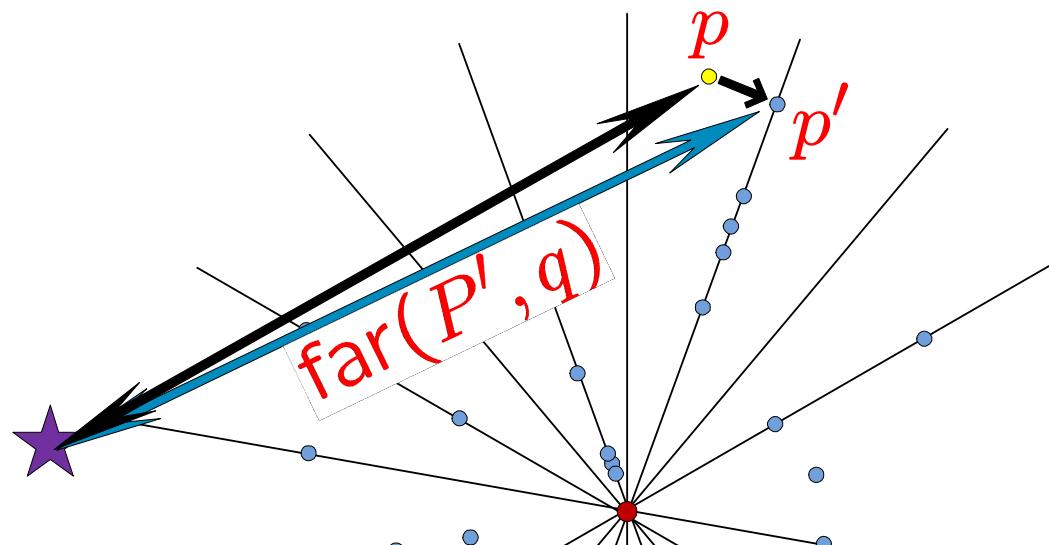
Claim: P' is a coresset for P

Let $\text{far}(P', q) = \text{dist}(p', q)$

$p :=$ the point whose projection is p'



$$\text{far}(P', q) - \text{far}(P, q) \leq \text{far}(P', q) - \text{dist}(p, q)$$



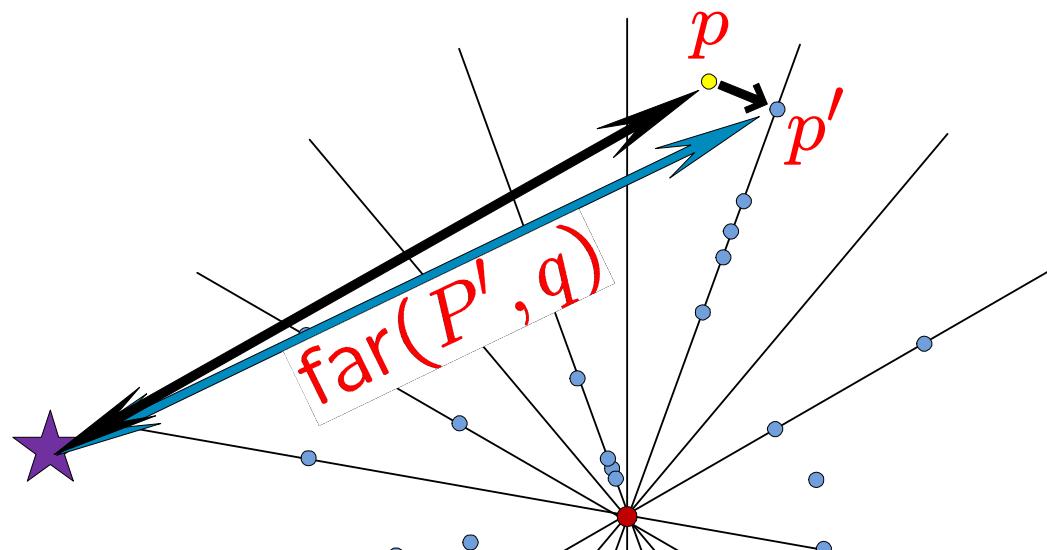
Claim: P' is a coresset for P

Let $\text{far}(P', q) = \text{dist}(p', q)$

$p :=$ the point whose projection is p'

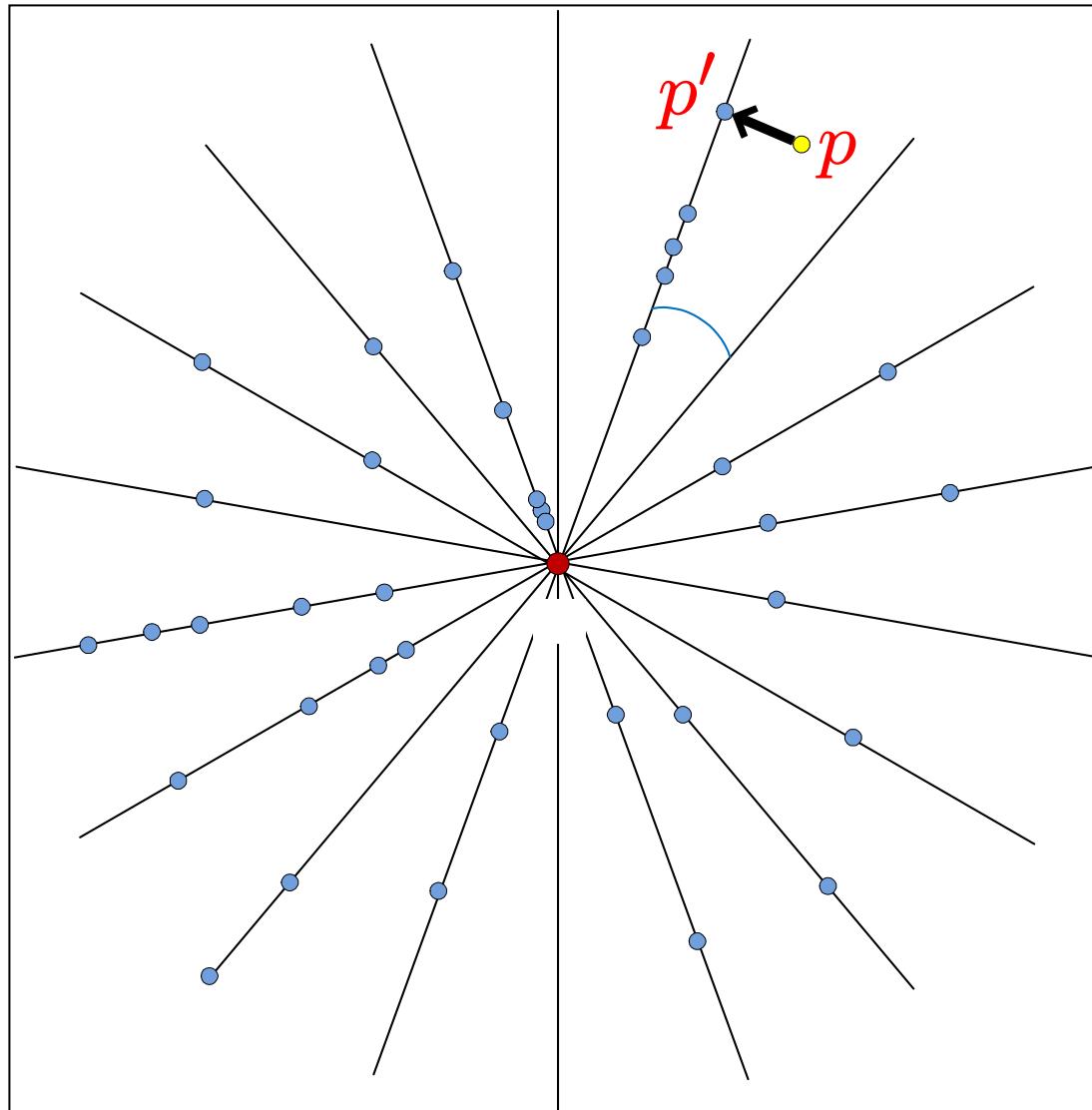


$$\begin{aligned} \text{far}(P', q) - \text{far}(P, q) &\leq \text{far}(P', q) - \text{dist}(p, q) \\ &\leq \text{dist}(p, p') \end{aligned}$$



Bounding $\text{dist}(p, p')$

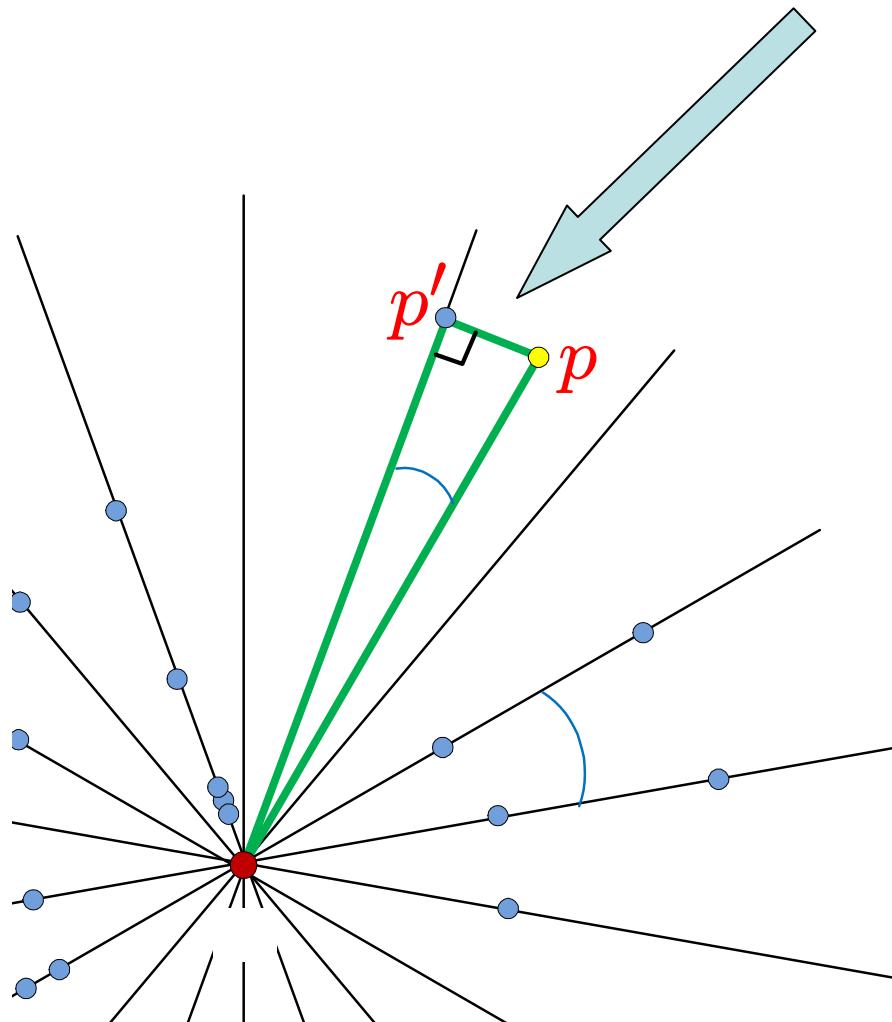
$p' :=$ the projection of $p \in P$



Bounding $\text{dist}(p, p')$

p' := the projection of $p \in P$

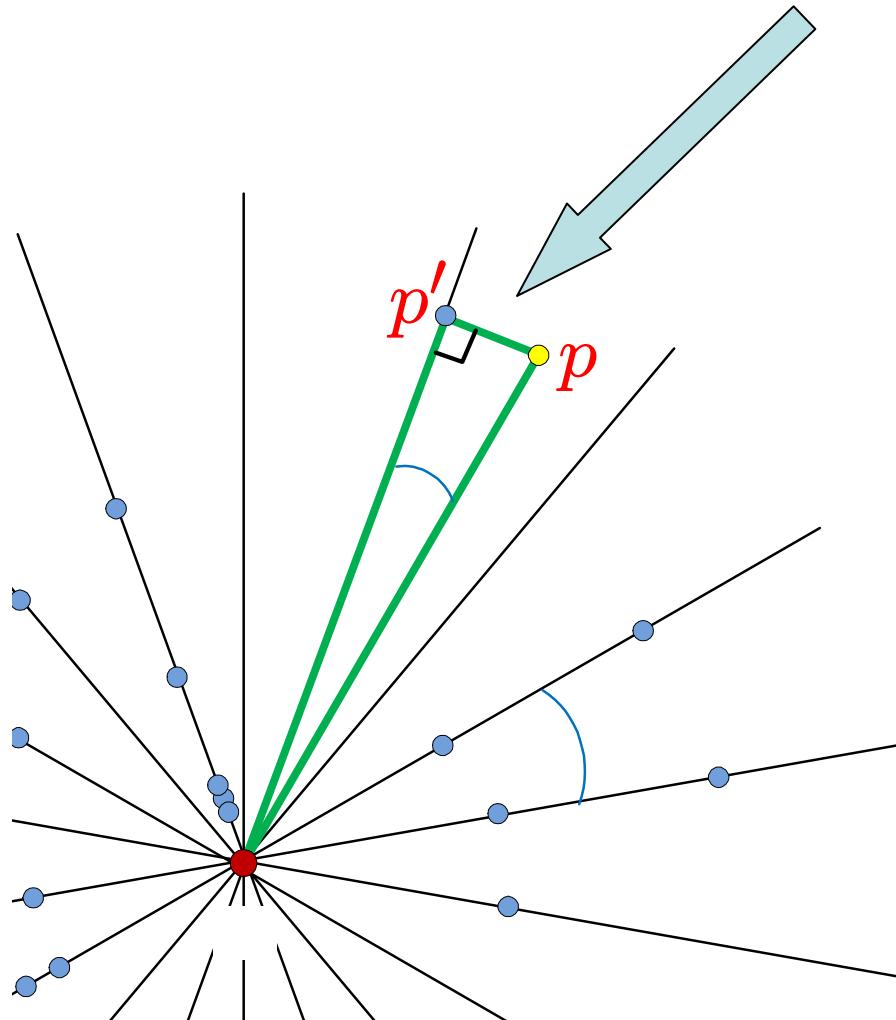
$$\text{dist}(p, p') = \sin \alpha \cdot \text{dist}(u, p)$$



Bounding $\text{dist}(p, p')$

$p' :=$ the projection of $p \in P$

$$\begin{aligned}\text{dist}(p, p') &= \sin \alpha \cdot \text{dist}(u, p) \\ &\leq O(\varepsilon) \cdot \text{dist}(u, p)\end{aligned}$$



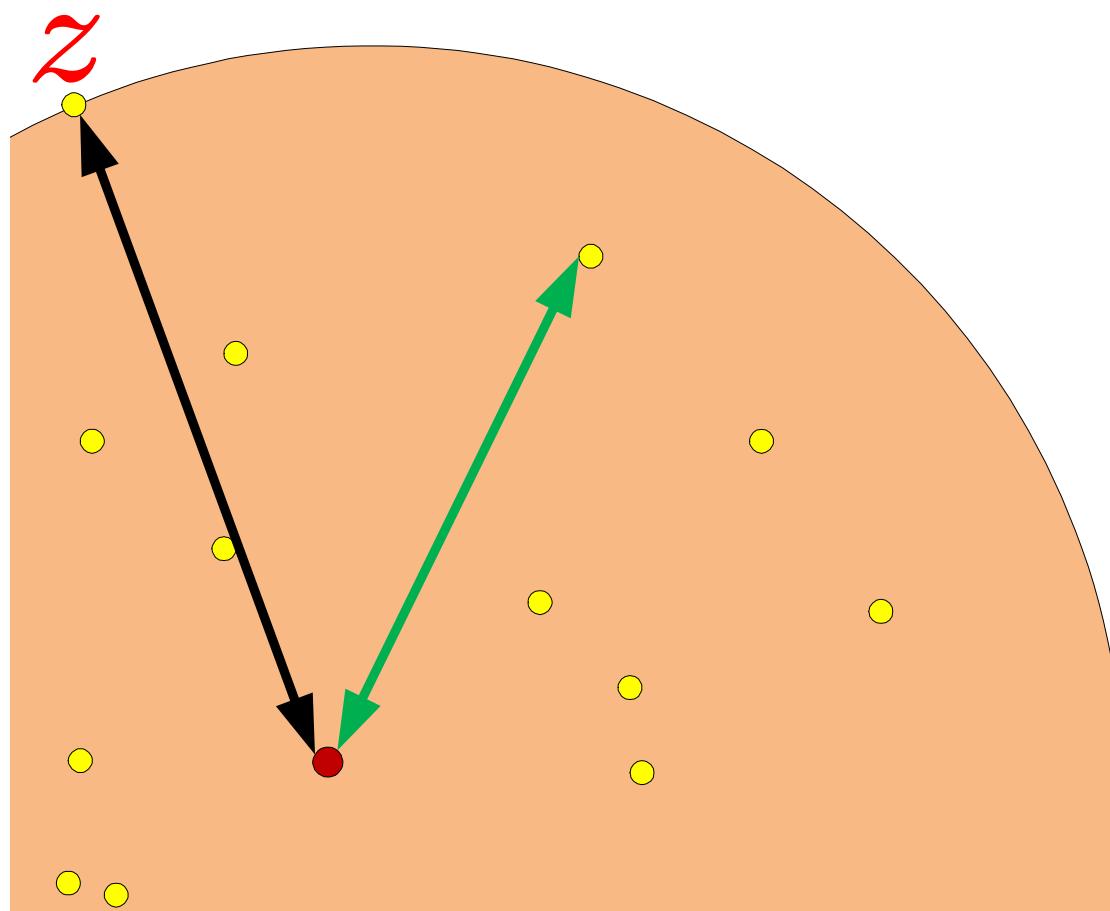
Bounding $\text{dist}(p, p')$

$p' :=$ the projection of $p \in P$

$$\text{dist}(p, p') = \sin \alpha \cdot \text{dist}(u, p)$$

$$\leq O(\varepsilon) \cdot \text{dist}(u, p)$$

$$\leq O(\varepsilon) \cdot r$$

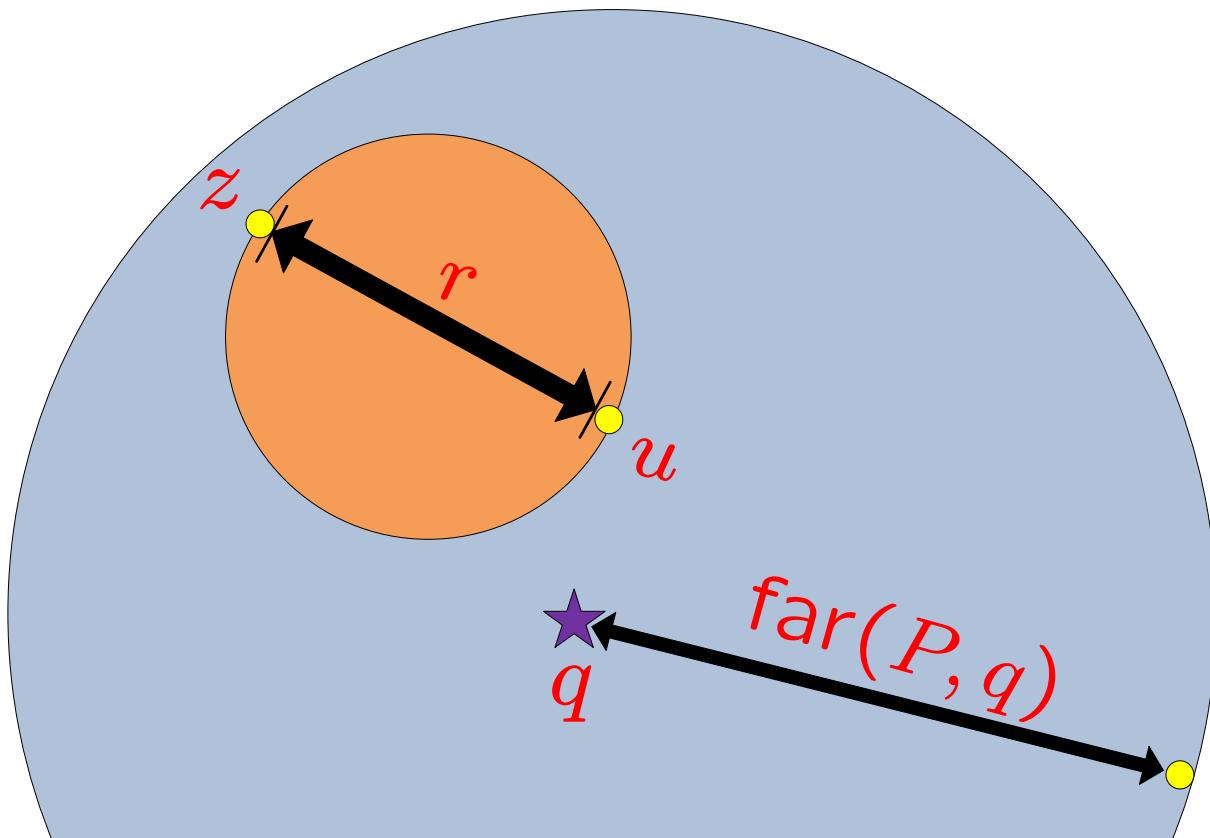


$$r := \max_{p \in P} \text{dist}(u, p)$$

Main Observation

Every ball that covers u and z ,

has a diameter of at least r :



$$r \leq 2\text{far}(P, q) \quad \Rightarrow \quad O(\varepsilon r) \leq O(\varepsilon)\text{far}(P, q)$$

Bounding $\text{dist}(p, p')$

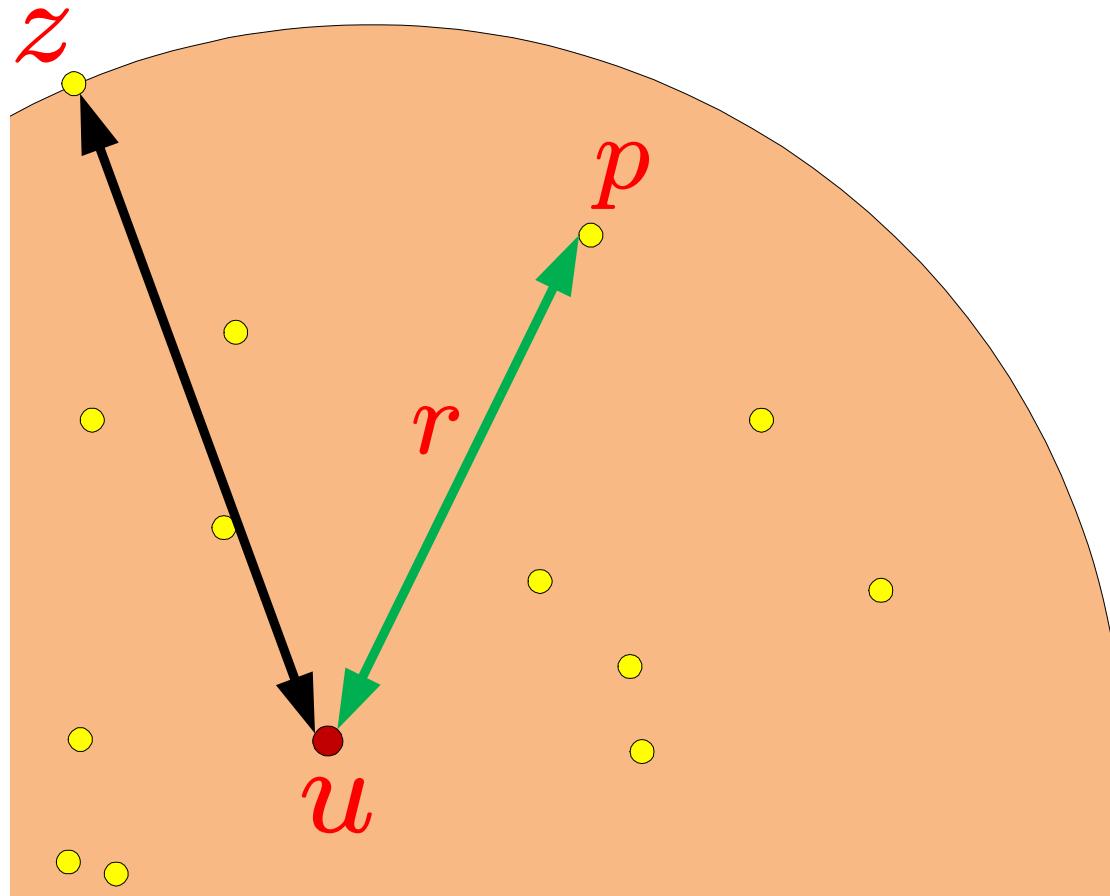
$p' :=$ the projection of $p \in P$

$$\text{dist}(p, p') = \sin \alpha \cdot \text{dist}(u, p)$$

$$\leq O(\varepsilon) \cdot \text{dist}(u, p)$$

$$\leq O(\varepsilon) \cdot r$$

$$\leq O(\varepsilon) \cdot \text{far}(P, q)$$



$$r := \max_{p \in P} \text{dist}(u, p)$$

Bounding $\text{dist}(p, p')$

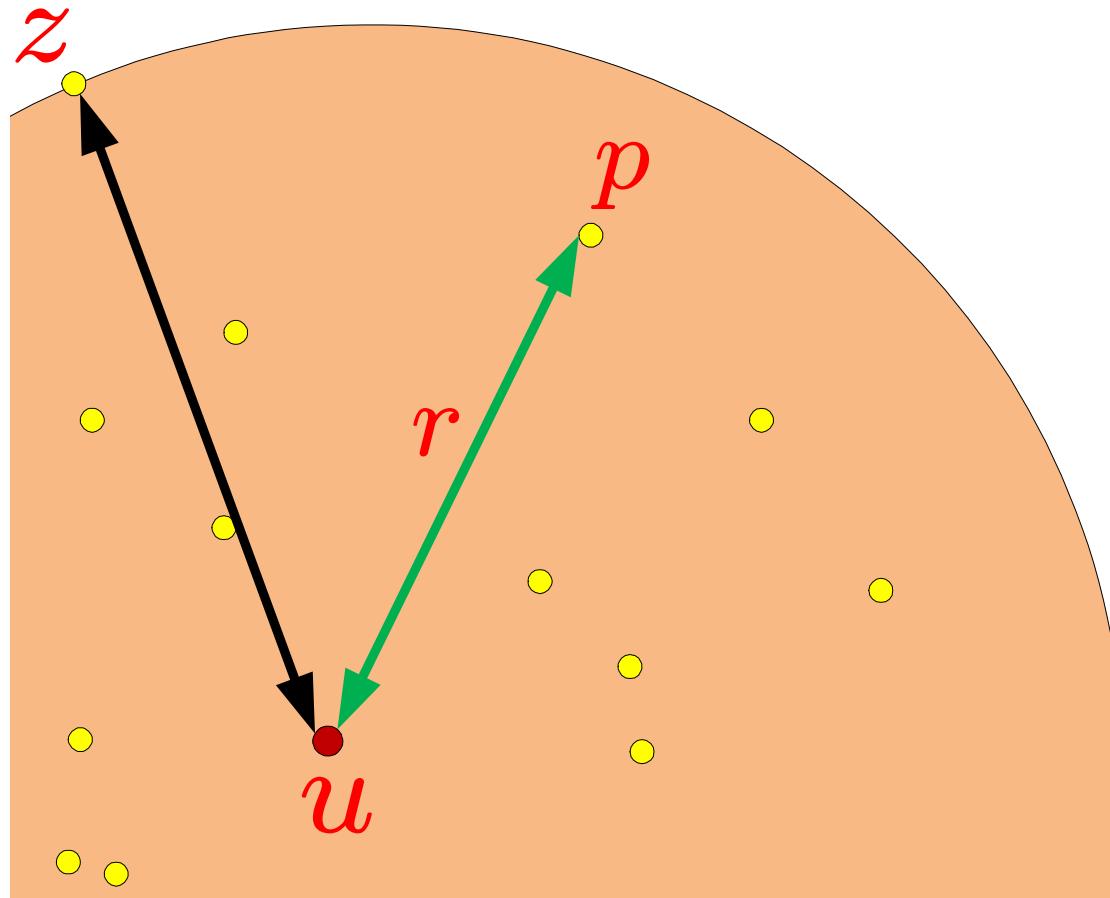
$p' :=$ the projection of $p \in P$

$$\text{dist}(p, p') = \sin \alpha \cdot \text{dist}(u, p)$$

$$\leq O(\varepsilon) \cdot \text{dist}(u, p)$$

$$\leq O(\varepsilon) \cdot r$$

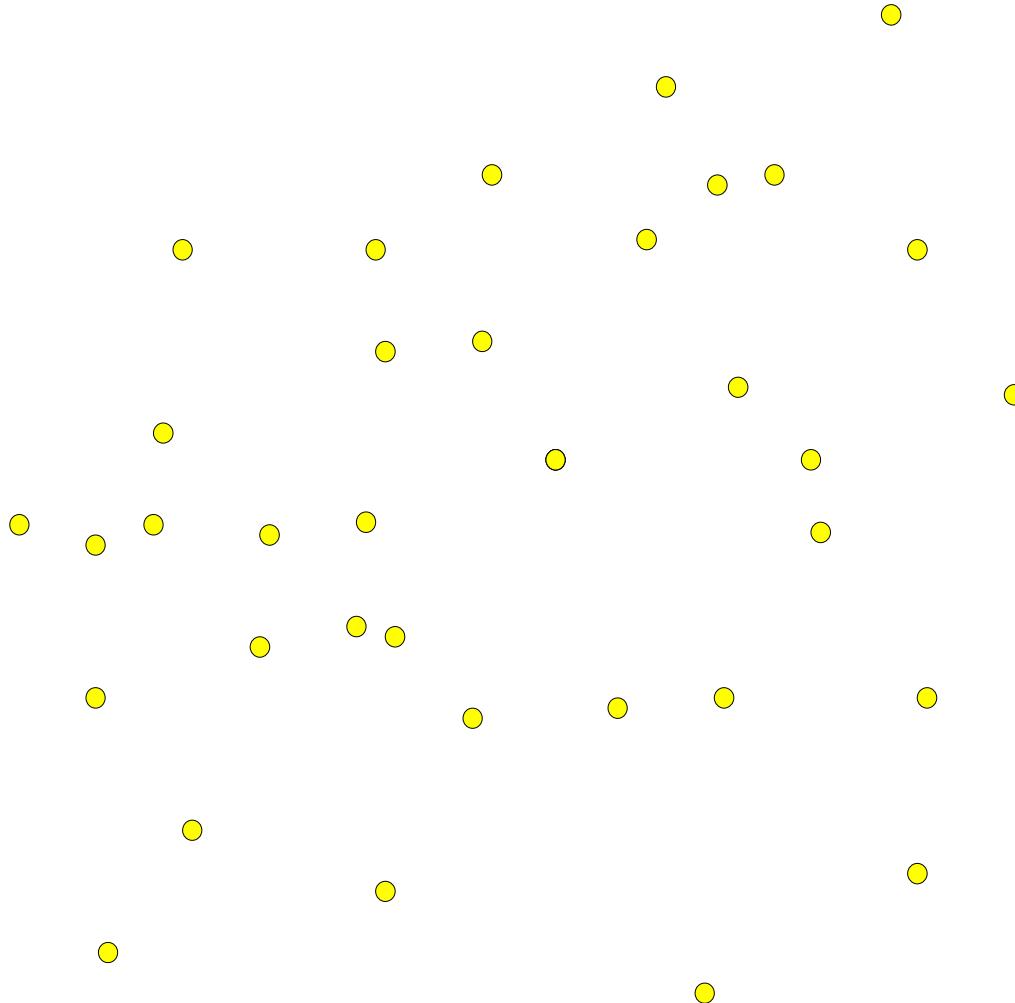
$$\leq O(\varepsilon) \cdot \text{far}(P, q)$$



$$r := \max_{p \in P} \text{dist}(u, p)$$

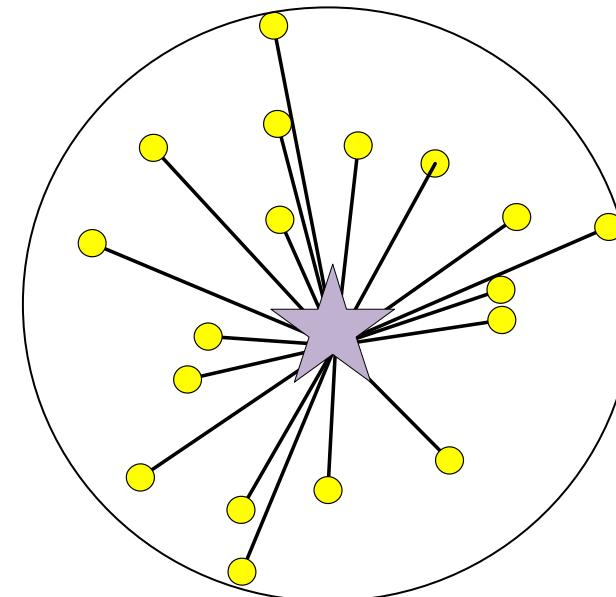
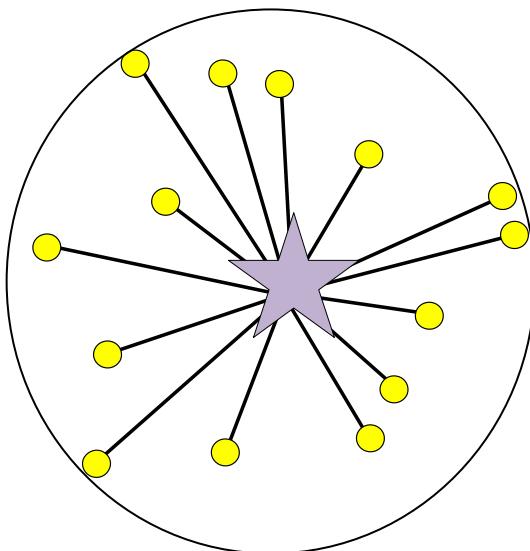


Generalization for k -Center



The k -Center of P

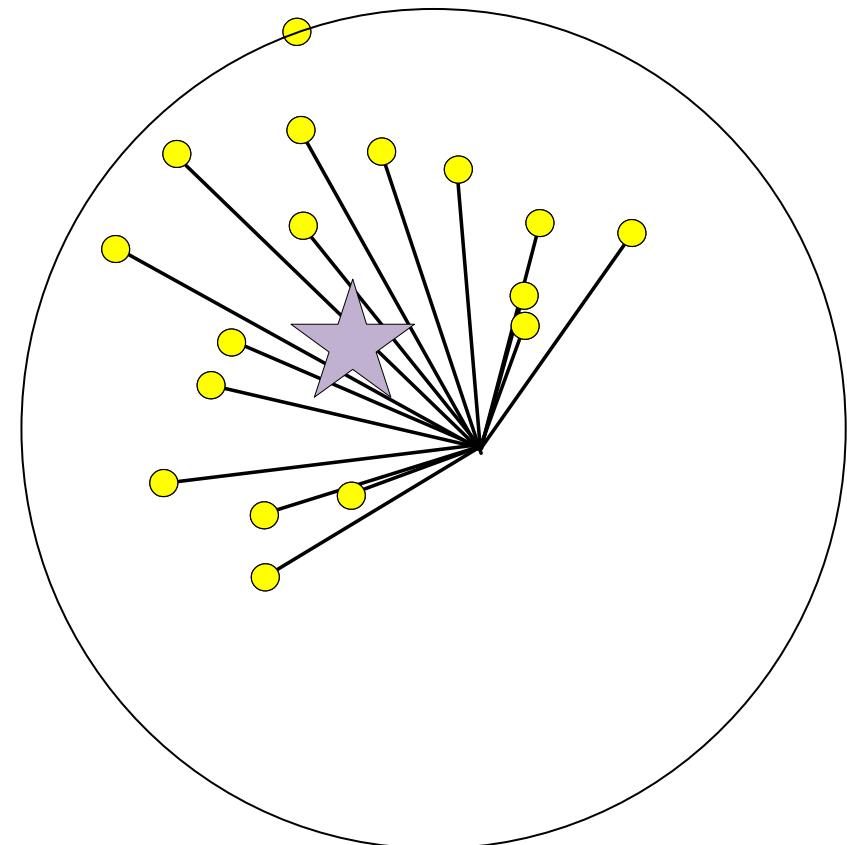
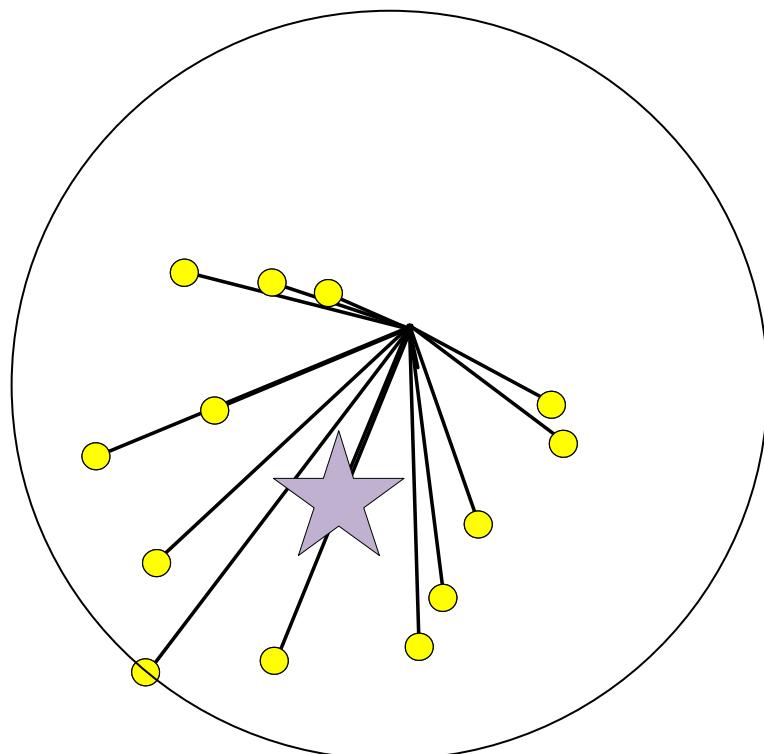
$$\text{opt} = \min_{|\text{OPT}|=k} \max_{p \in P} \text{dist}(p, \text{OPT})$$



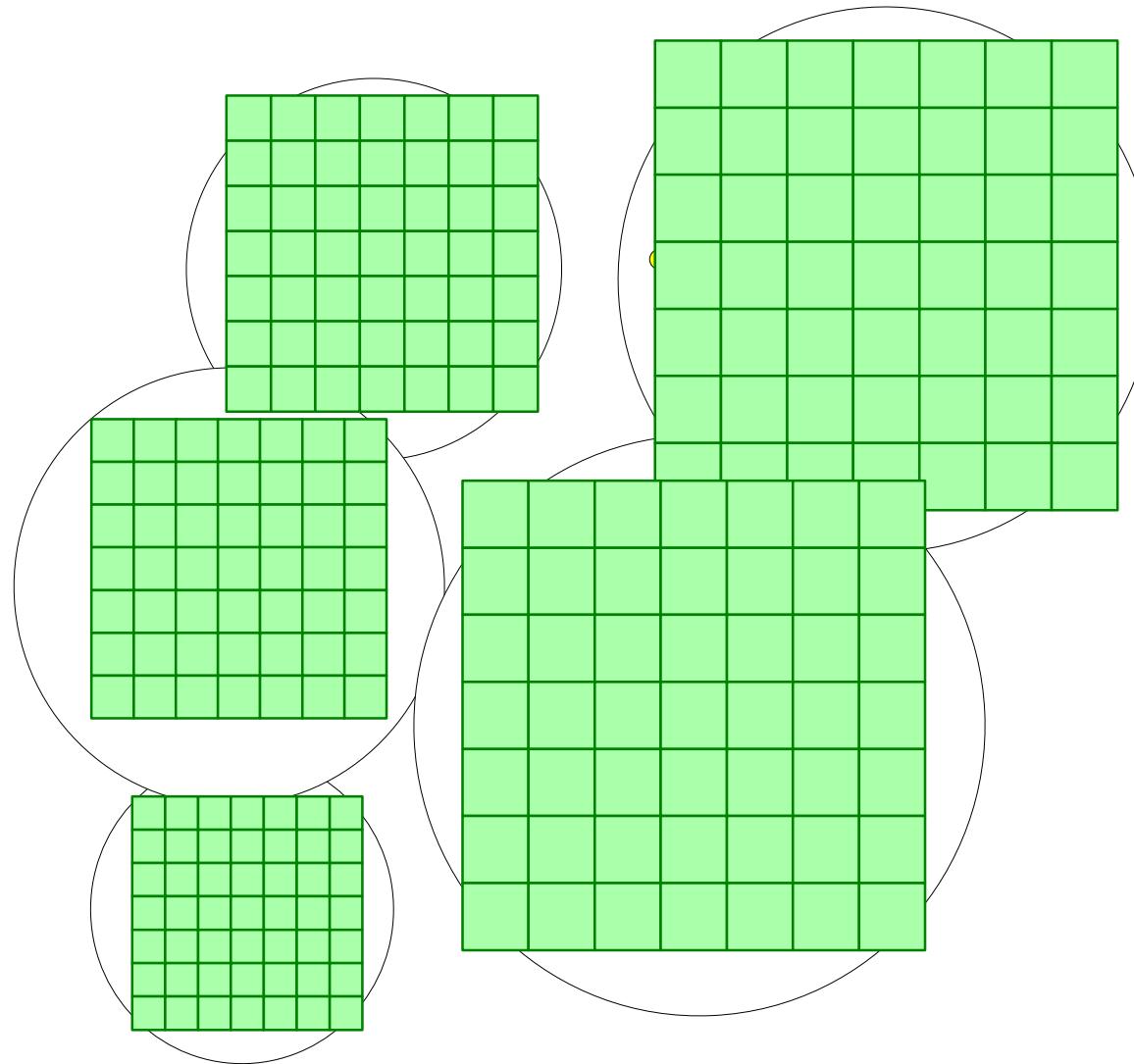
Constant Approximation

$$|\widetilde{\text{OPT}}| = k$$

$$\max_{p \in P} \text{dist}(p, \widetilde{\text{OPT}}) \leq O(1) \cdot \text{opt}$$



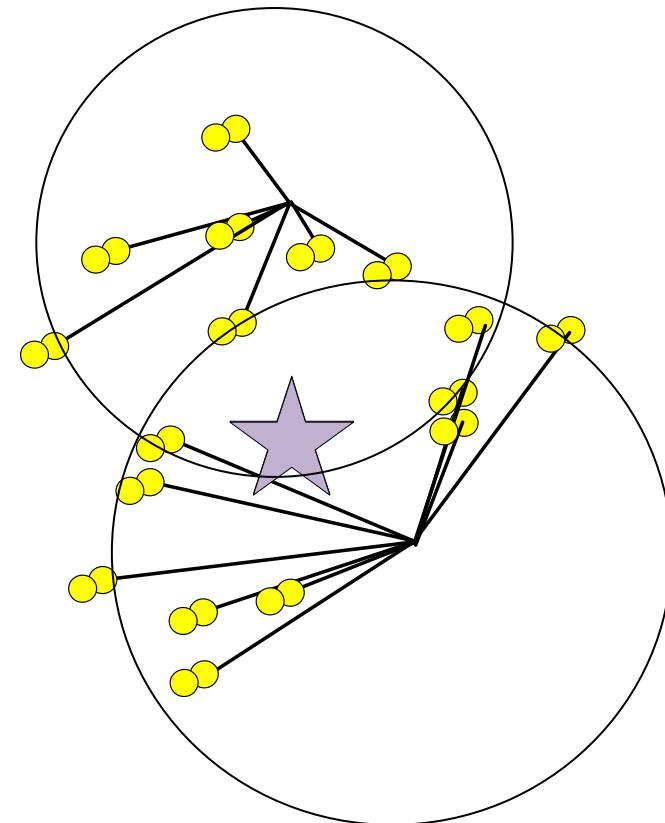
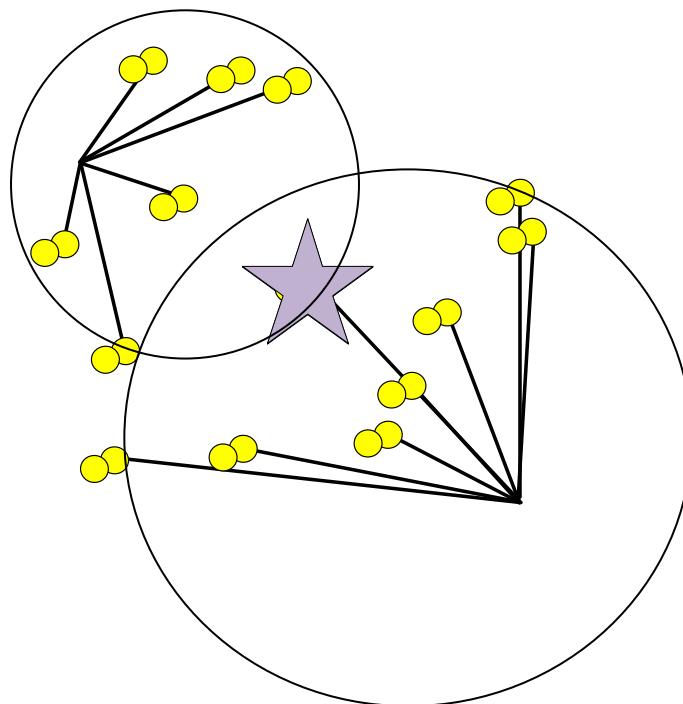
For each cluster:
Apply construction for $k = 1$



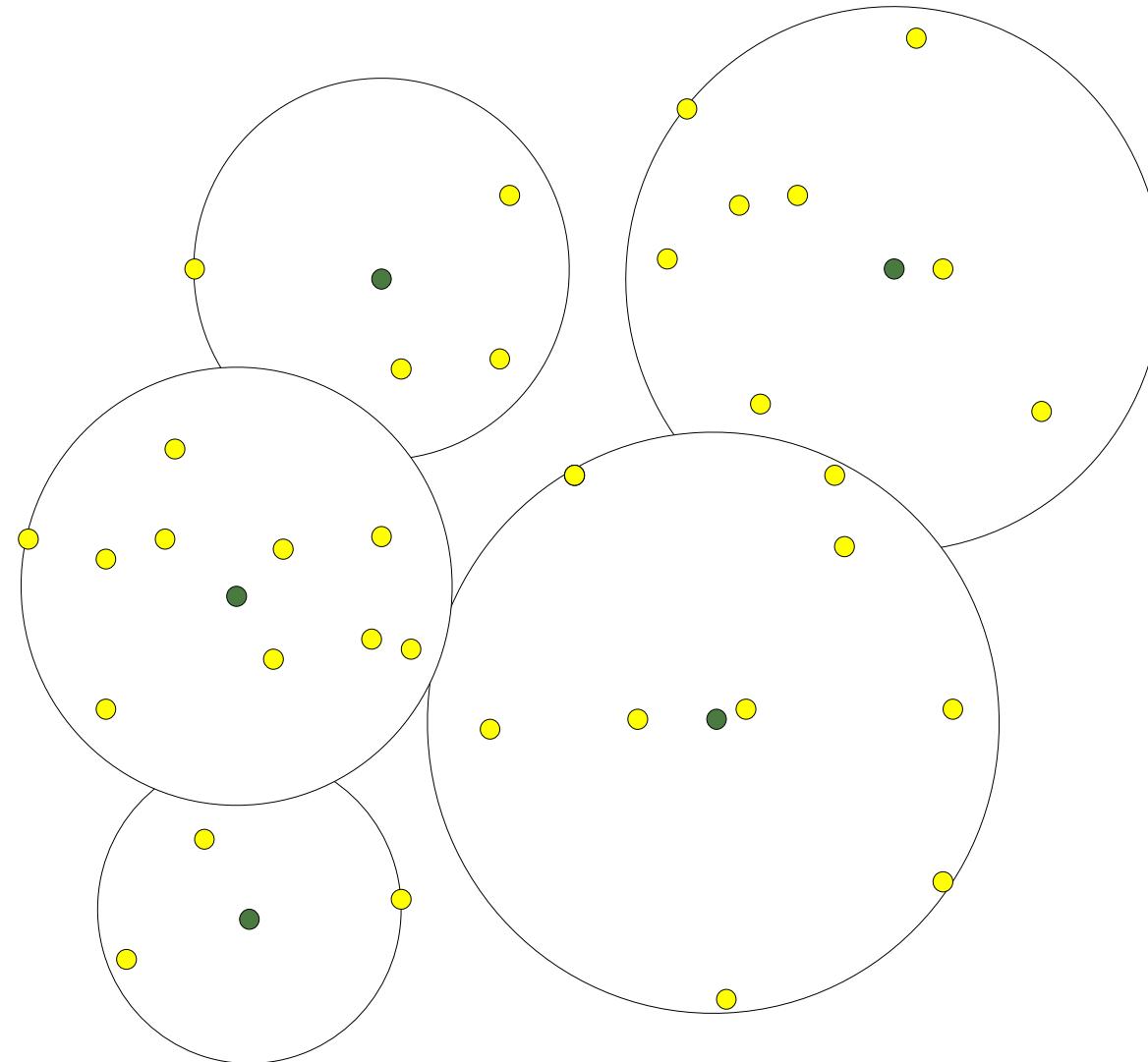
Bi-criteria Approximation

$$|B| = O(k \log n)$$

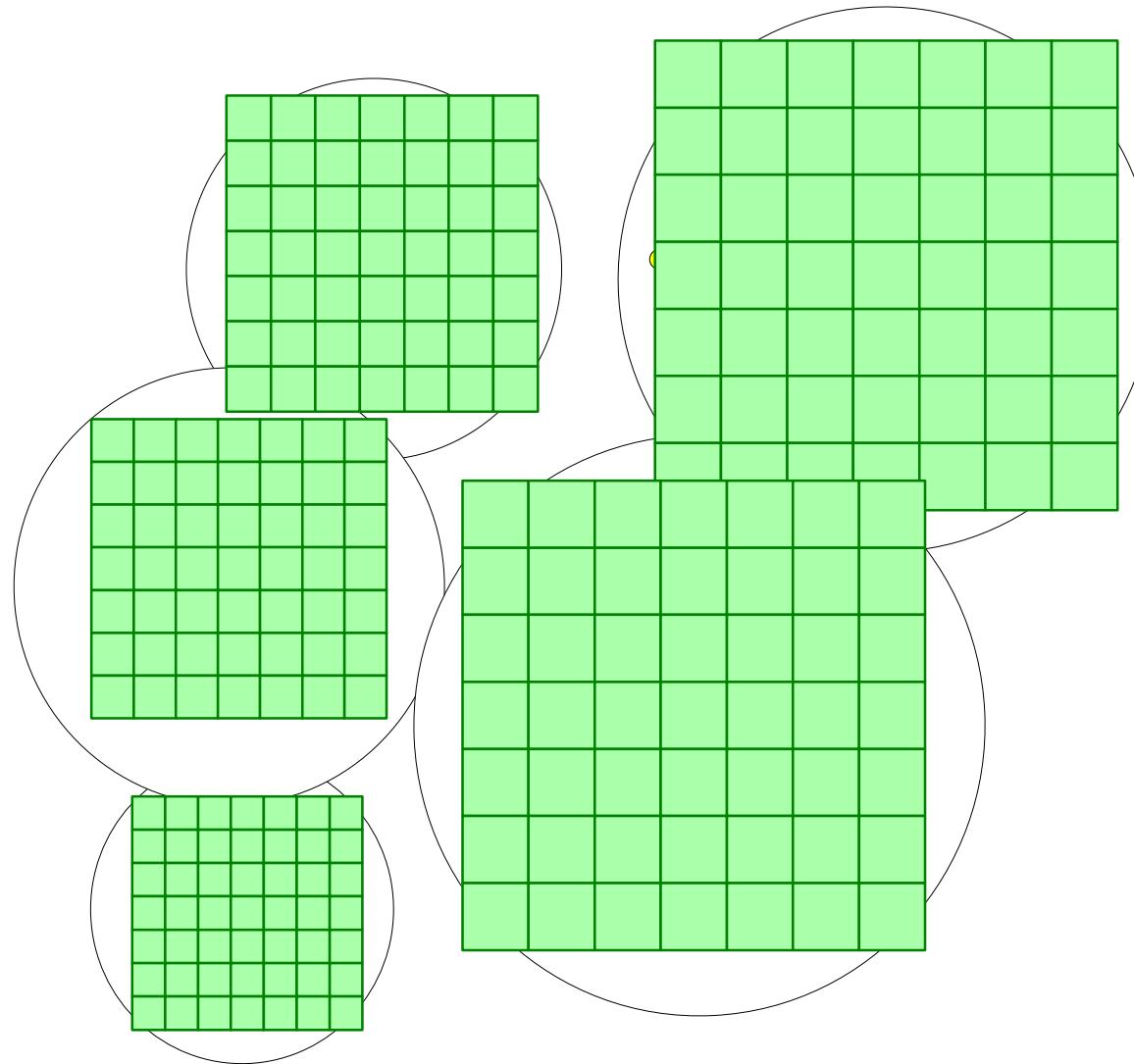
$$\max_{p \in P} \text{dist}(p, B) \leq O(1) \cdot \text{opt}$$



Compute Bi-criteria Approximation Using ε -Approximations

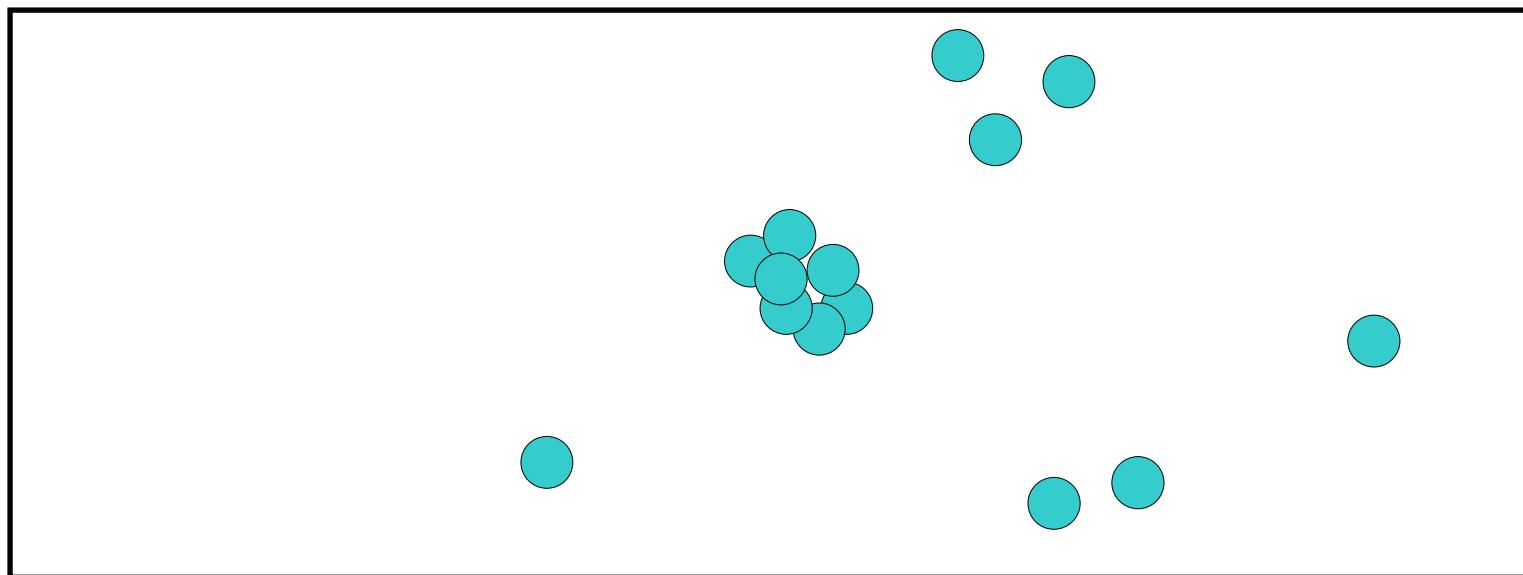


For each cluster:
Apply construction for $k = 1$



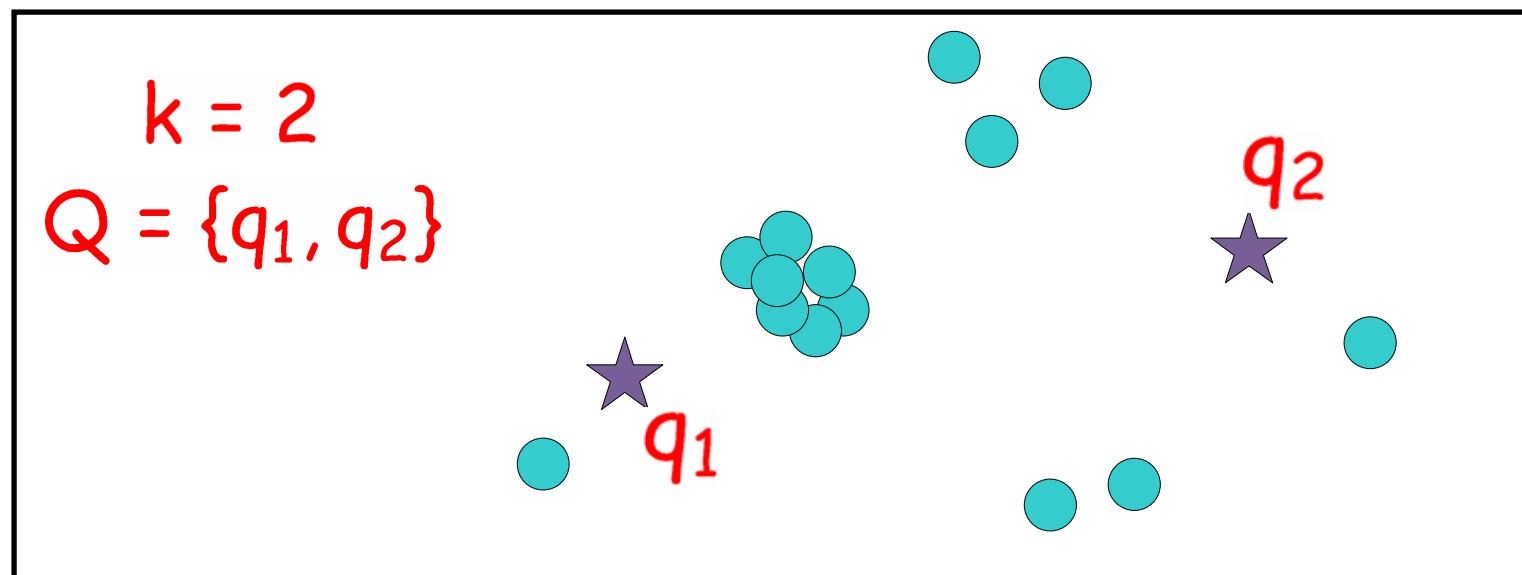
k -Median Queries

- Input: $P \subseteq \mathbb{R}^d$



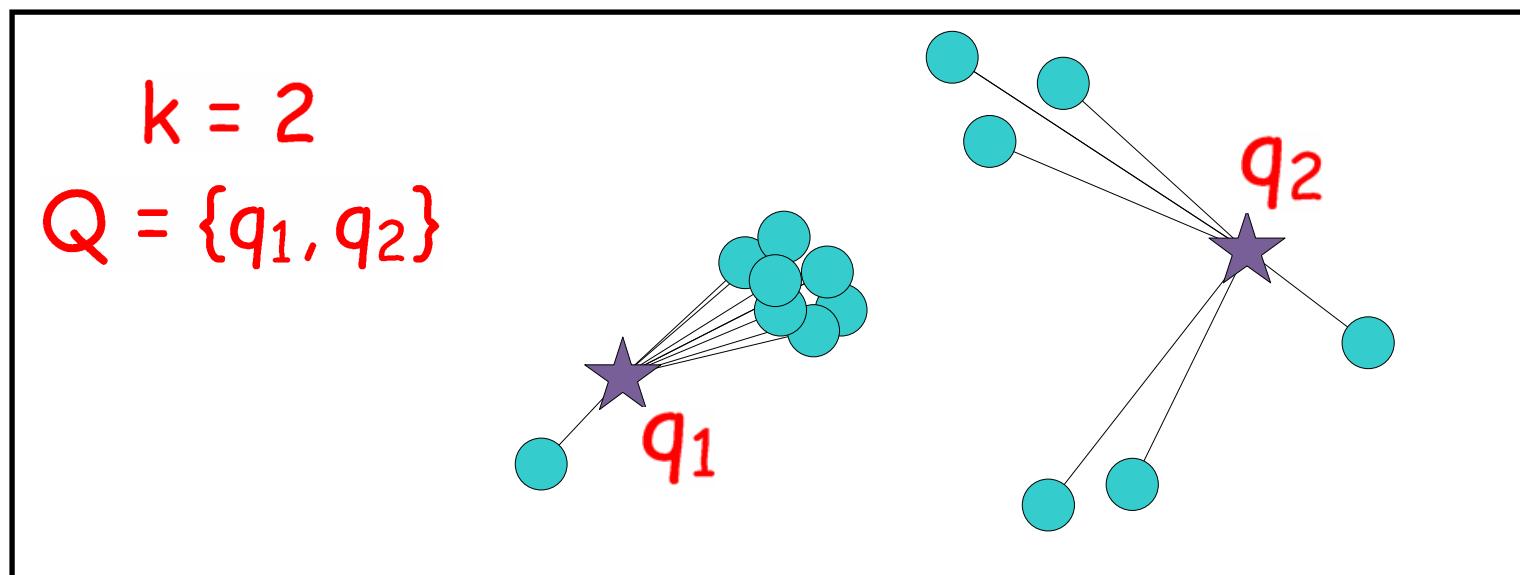
k -Median Queries

- Input: $P \subseteq \mathbb{R}^d$
- Query: A set Q of k points



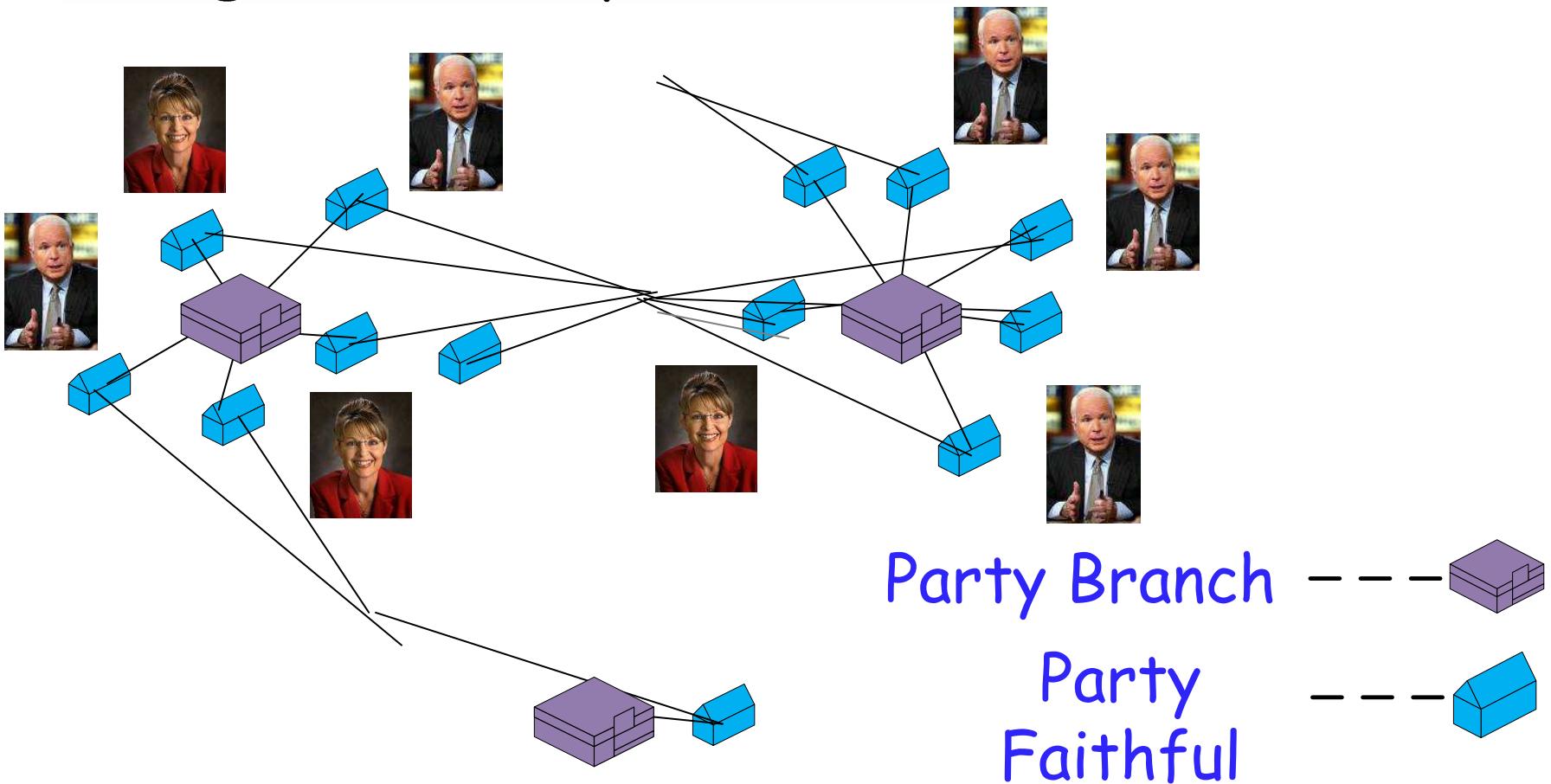
k -Median Queries

- Input: $P \subseteq \mathbb{R}^d$
- Query: A set Q of k points
- Output: $\sum_{p \in P} \text{dist}(p, Q) = \sum_{p \in P} \min_{q \in Q} \|p - q\|$



Motivation

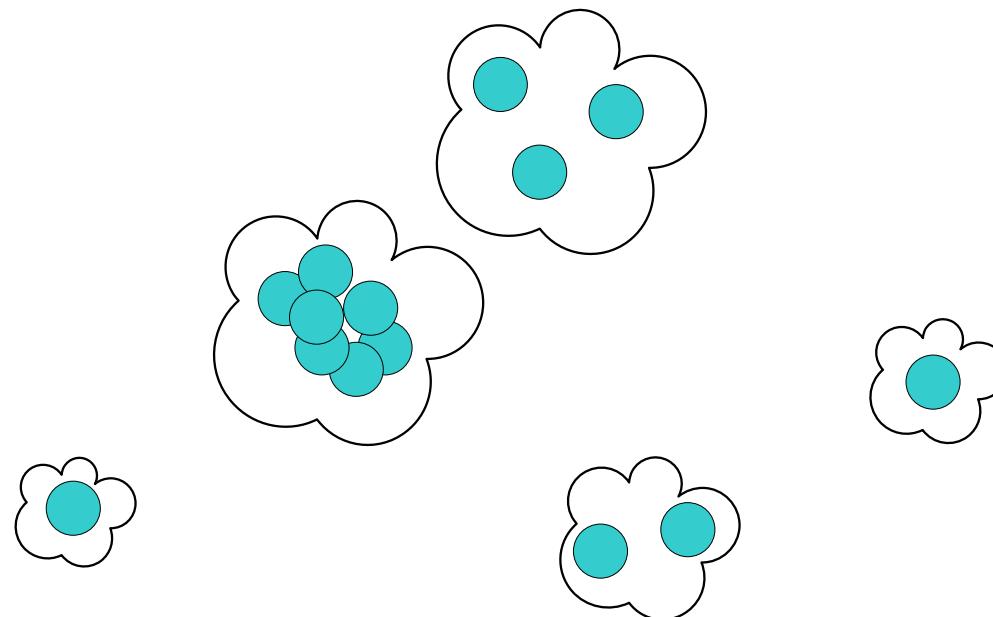
Comparing alternatives:
How good is this placement?



(k, ϵ) -Median Coreset

Answer k -median queries in sub-linear time

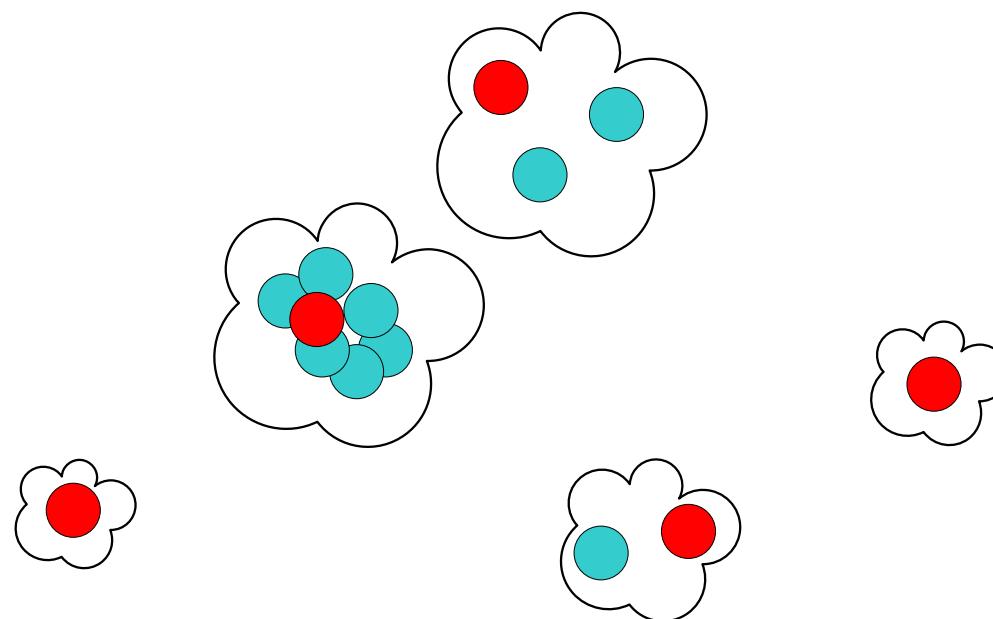
Key Idea: Replace many points by one weighted representative:



(k, ϵ) -Median Coreset

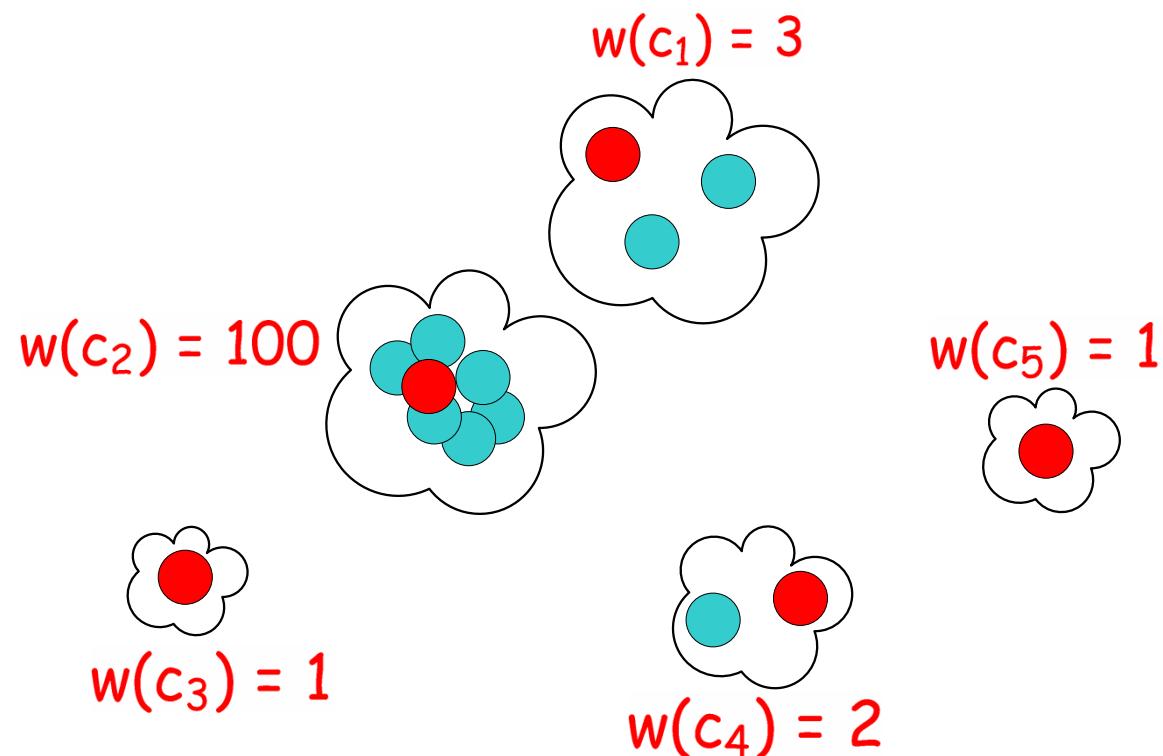
Answer k -median queries in sub-linear time

Key Idea: Replace many points by one weighted representative:



(k, ε) -Median Coreset

Key Idea: Replace many points by one weighted representative:

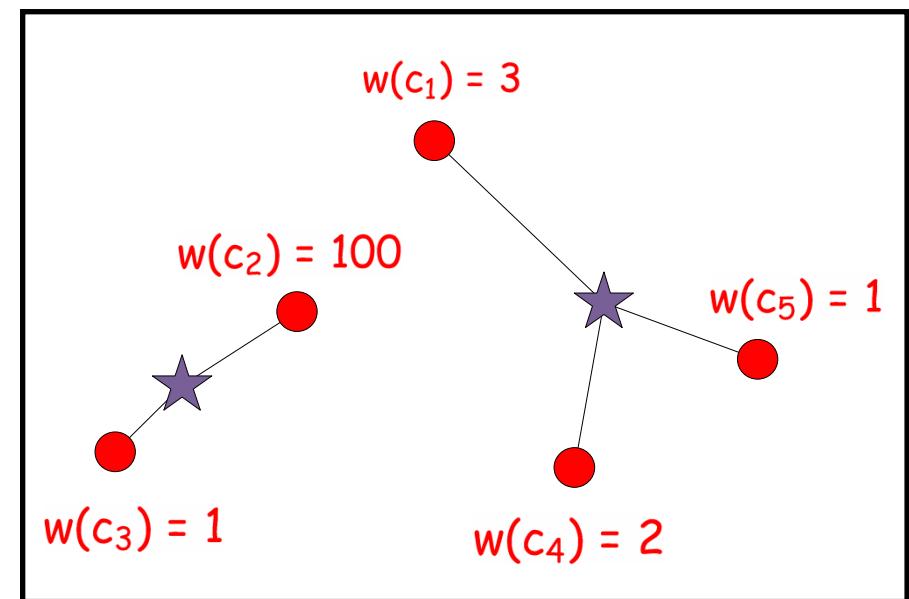


Definition

C is a (k, ε) -coreset for P , if $\forall Q, |Q| = k$:

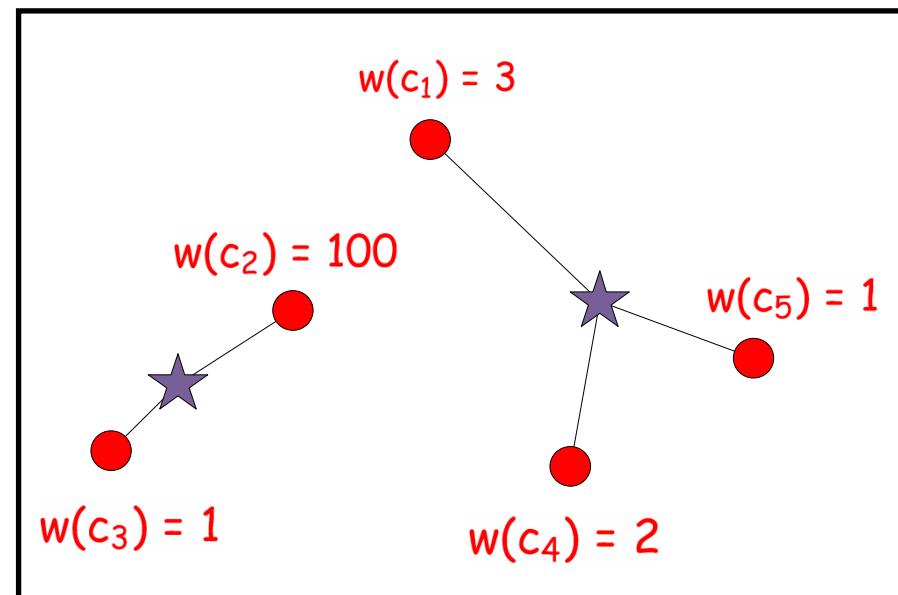
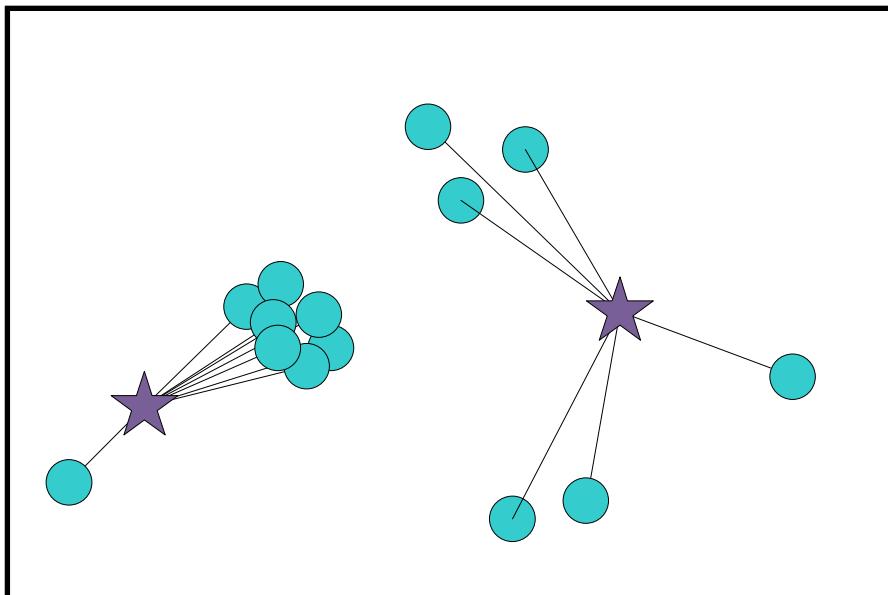
$$\sum_{p \in P} \text{dist}(p, Q) \sim \sum_{c \in C} w(c) \cdot \text{dist}(c, Q)$$

Multiplicative error $\leq 1 + \varepsilon$



(k, ε) -Median Coreset

$$\sum_{p \in P} \text{dist}(p, Q) \sim \sum_{c \in C} w(c) \cdot \text{dist}(c, Q)$$



Related Work

- Strong coresets (also for k -median/mean),

Size: $(k/\varepsilon)^{d/\varepsilon}$, [Har-Peled, Koshar'05]

Size: $\text{poly}(kd \log n / \varepsilon)$, [Ke-Chen'06]

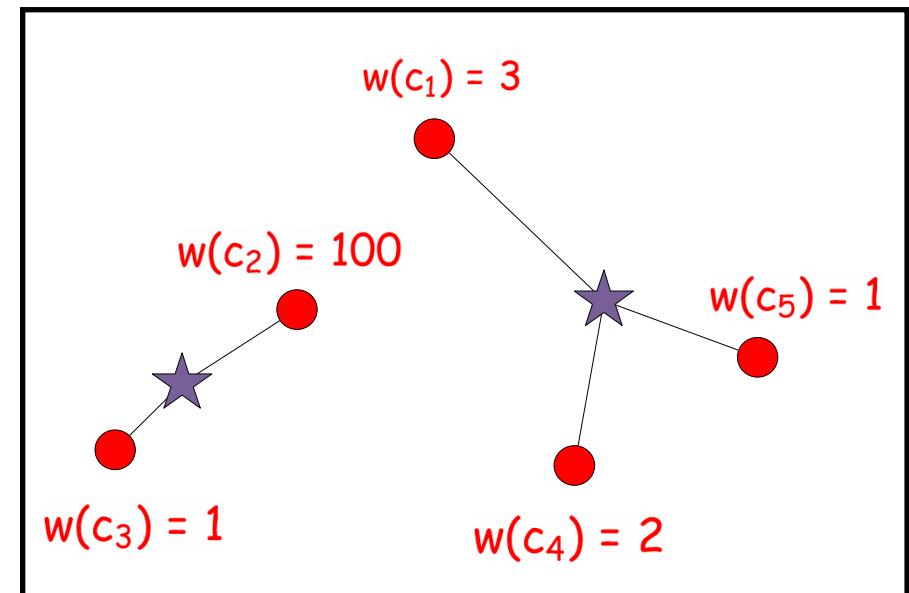
Size: $\text{poly}(kd/\varepsilon)$, [Langberg, Schulman'10],

[F, Langberg '10]

Definition

C is a $\text{weak } (k, \varepsilon)$ -coreset for P if ~~$\forall Q, |Q| = k$~~ :

$$\min_Q \sum_{p \in P} \text{dist}(p, Q) \sim \min_Q \sum_{c \in C} w(c) \cdot \text{dist}(c, Q)$$



Related Work

- Strong coresets (also for k -median/mean),

Size: $(k/\varepsilon)^{d/\varepsilon}$, [Har-Peled, Koshar'05]

Size: $\text{poly}(kd \log n / \varepsilon)$, [Ke-Chen'06]

Size: $\text{poly}(kd/\varepsilon)$, [Langberg, Schulman'10],

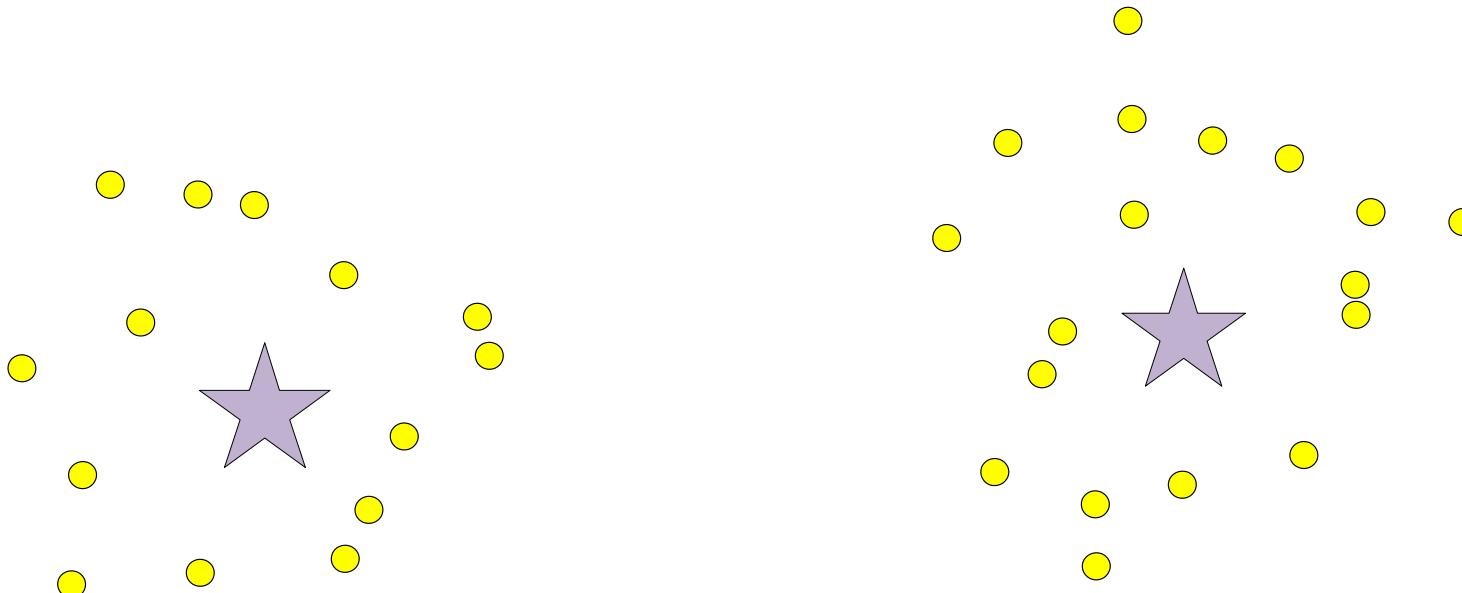
[F, Langberg '10]

- Weak coresets, [F, Sohler, Monemizadeh'07]

Size: $\text{poly}(k/\varepsilon)$,

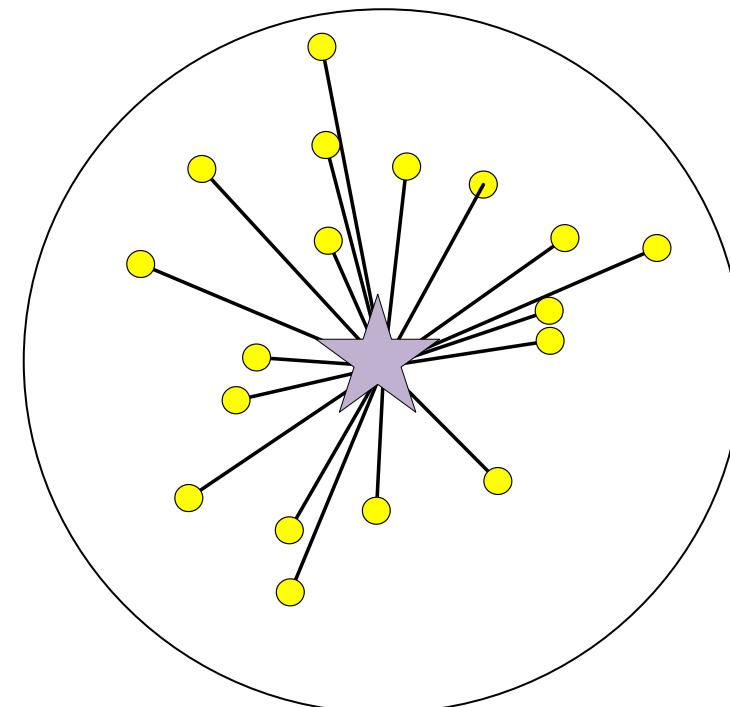
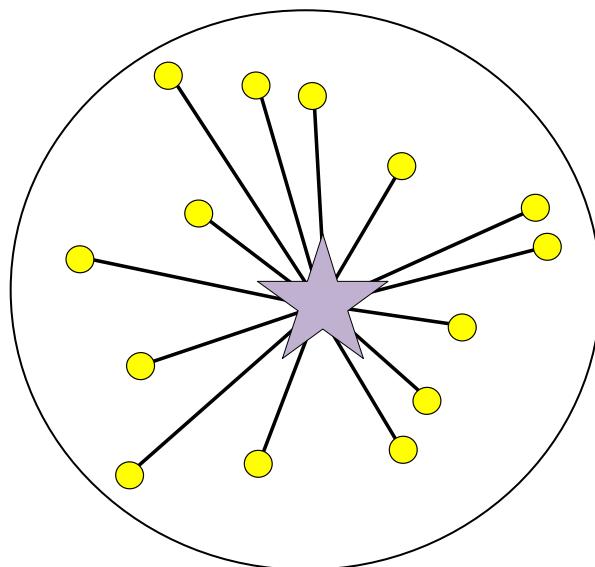
The k-Median of P

$$\text{opt} = \min_{|\text{OPT}|=k} \sum_{p \in P} \text{dist}(p, \text{OPT})$$



The k -Median of P

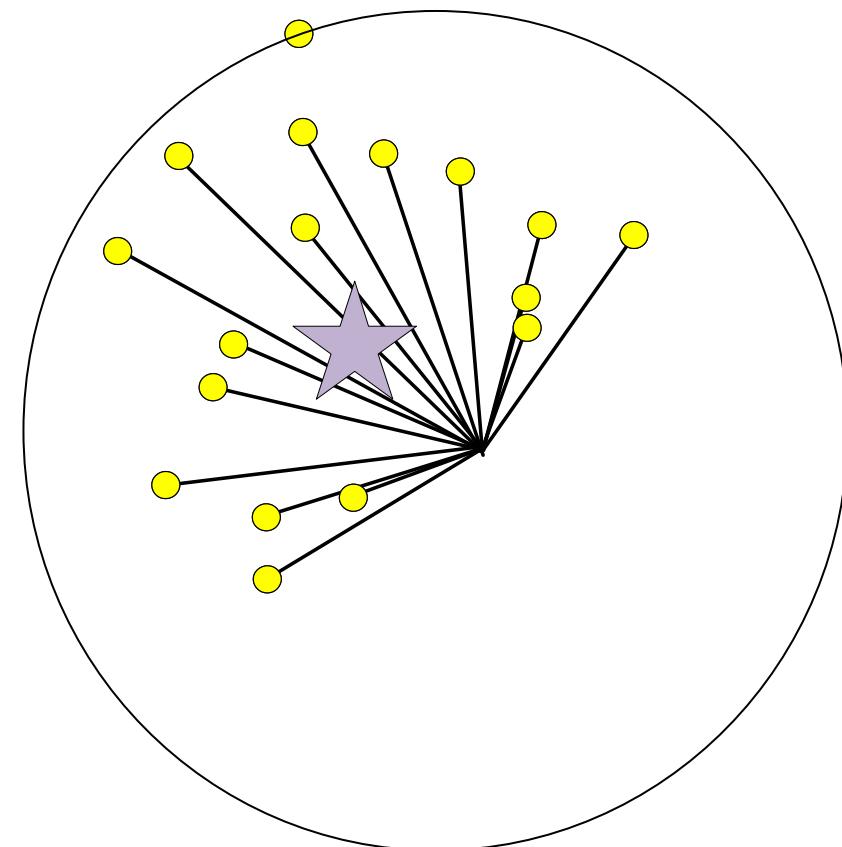
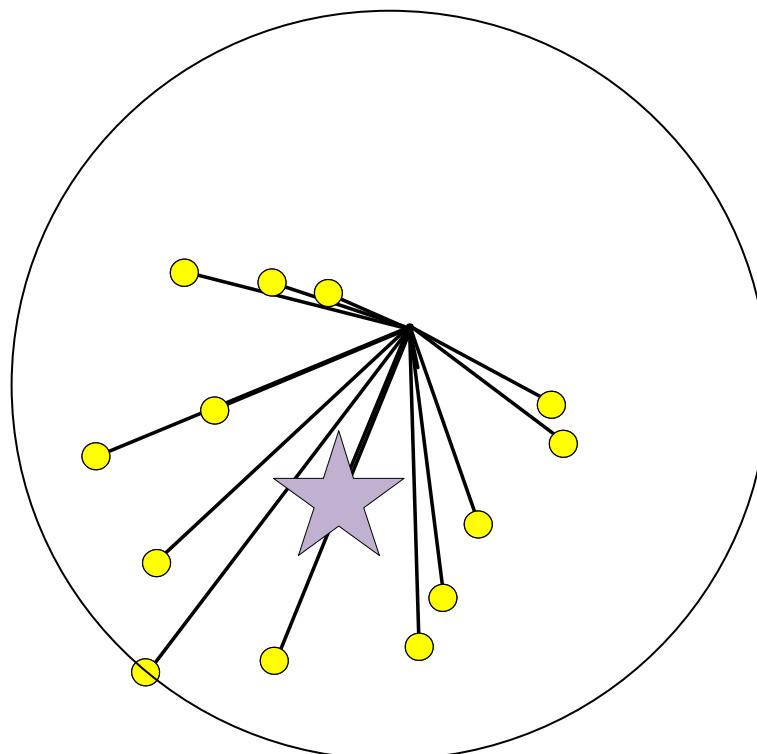
$$\text{opt} = \min_{|\text{OPT}|=k} \sum_{p \in P} \text{dist}(p, \text{OPT})$$



Constant Approximation

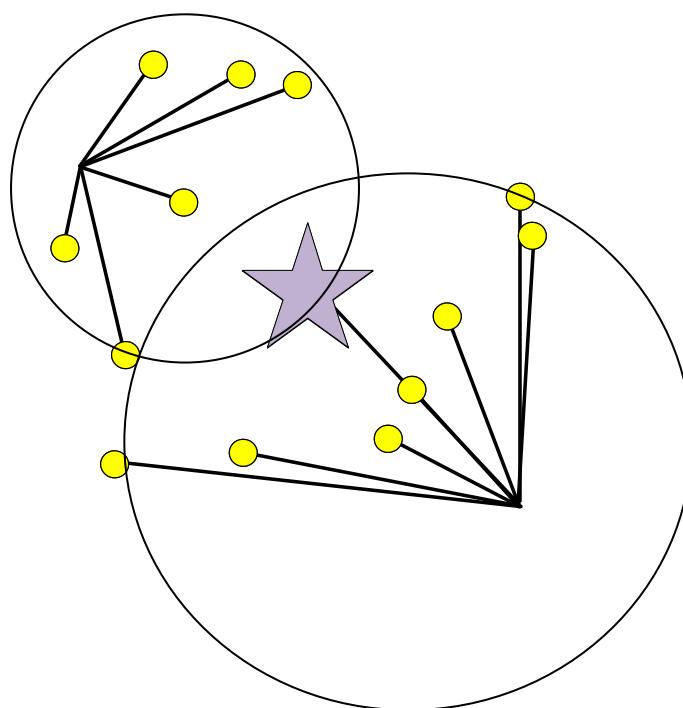
$$|\widetilde{OPT}| = k,$$

$$\sum_{p \in P} \text{dist}(p, \widetilde{OPT}) \leq c \cdot \text{opt}$$

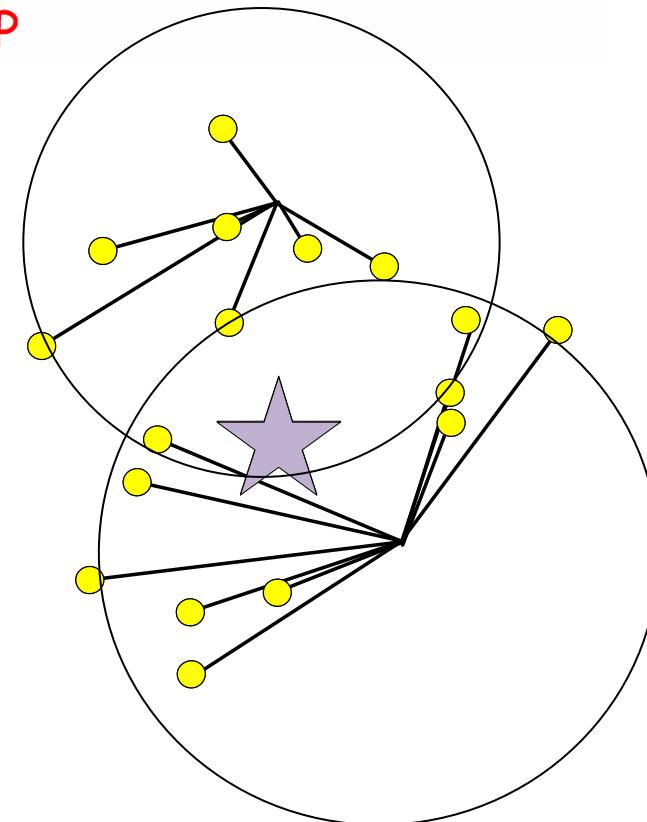


Bi-Criteria Approximation [FFS07]

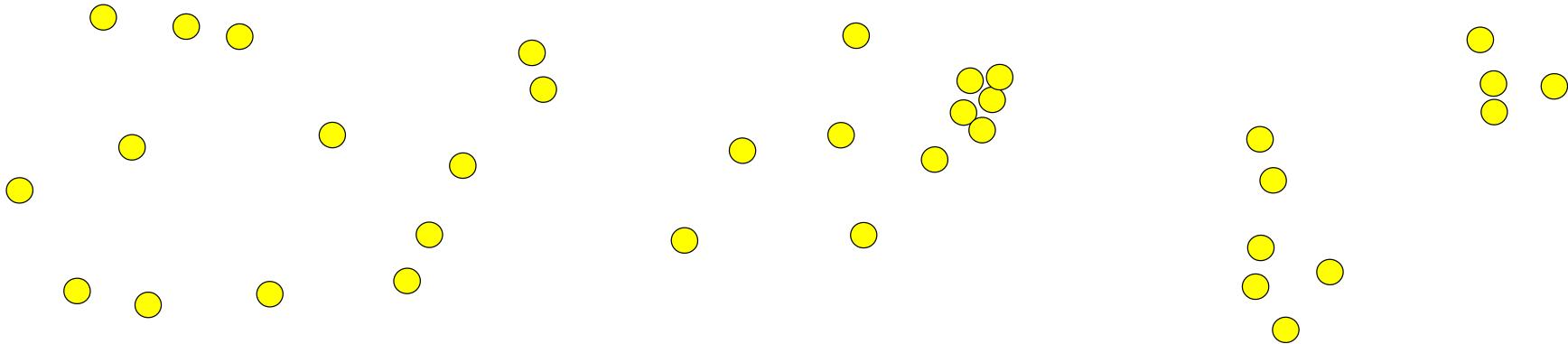
$$|B| = O(k \log n),$$



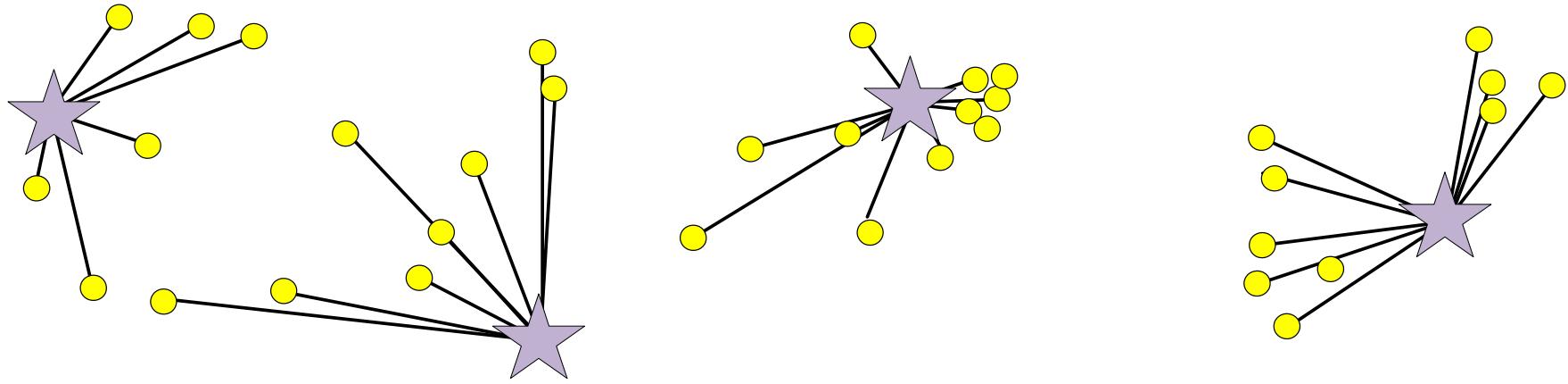
$$\sum_{p \in P} \text{dist}(p, B) \leq c \cdot \text{opt}$$



Coreset for 2-Median [FMS07]



Coreset for 2-Median

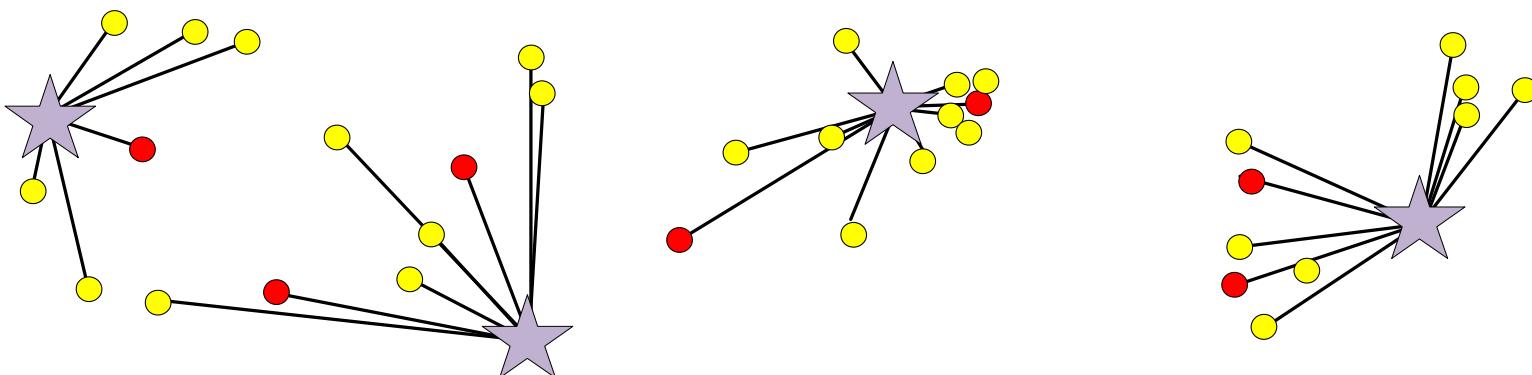


Using a bi-criteria approximation
for 2-Median

Coreset for 2-Median

1. Pick a sample S of $d/\varepsilon^{O(1)}$ points, where

$$\Pr[p] = \frac{\text{dist}(p, B)}{\sum_{p \in P} \text{dist}(p, B)}$$

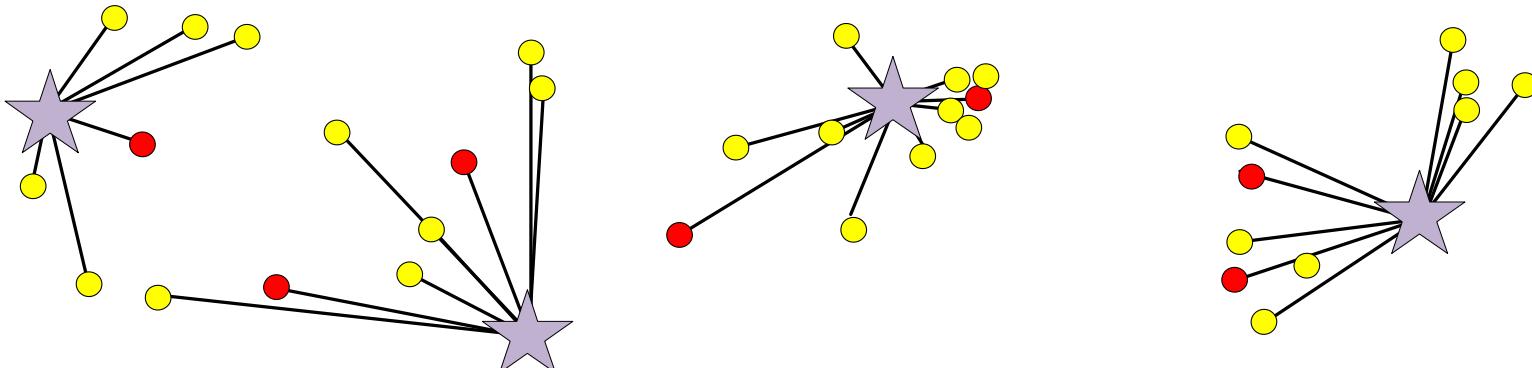


Coreset for 2-Median

1. Pick a sample S of $d/\varepsilon^{O(1)}$ points, where

$$\Pr[p] = \frac{\text{dist}(p, B)}{\sum_{p \in P} \text{dist}(p, B)}$$

2. $\forall p \in S : w(p) \leftarrow \frac{1}{|S| \cdot \Pr[p]}$



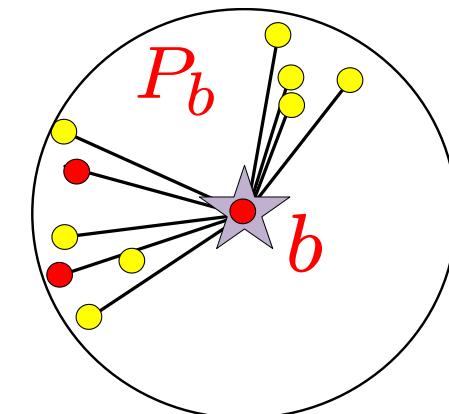
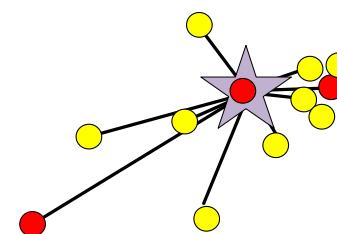
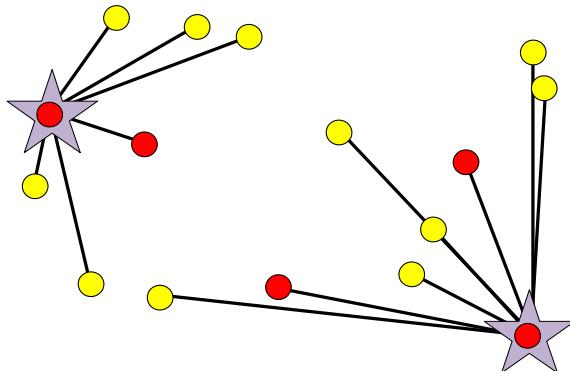
Coreset for 2-Median

1. Pick a sample S of $d/\varepsilon^{O(1)}$ points, where

$$\Pr[p] = \frac{\text{dist}(p, B)}{\sum_{p \in P} \text{dist}(p, B)}$$

2. $\forall p \in S : w(p) \leftarrow \frac{1}{|S| \cdot \Pr[p]}$

3. $\forall b \in B : w(b) \leftarrow |P_b| - w(S \cap P_b)$

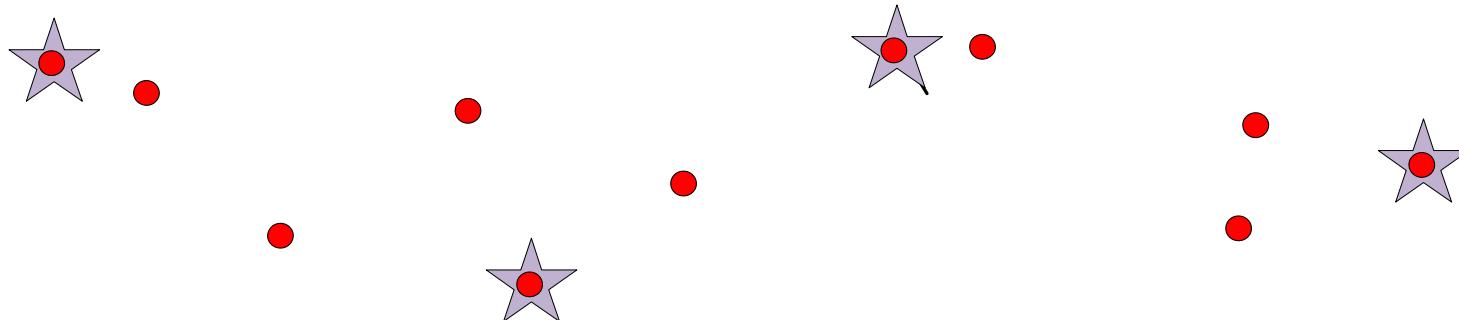


Coreset for 2-Median

1. Pick a sample S of $d/\varepsilon^{O(1)}$ points, where

$$\Pr[p] \leftarrow \frac{\text{dist}(p, B)}{\sum_{p \in P} \text{dist}(p, B)}$$

2. $\forall p \in S : w(p) \leftarrow \frac{1}{|S| \cdot \Pr[p]}$
3. $\forall b \in B : w(b) \leftarrow |P_b| - w(S \cap P_b)$
4. Return $S \cup B$



Correctness

We will prove that, with high probability,

$$\sum_{p \in P} (\text{dist}(p, Q) - \text{dist}(b_p, Q)) \sim \sum_{p \in S} w(p) \cdot (\text{dist}(p, Q) - \text{dist}(b_p, Q))$$

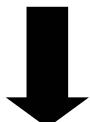


$$\begin{aligned} \sum_{p \in P} \text{dist}(p, Q) &= \sum_{p \in P} \text{dist}(b_p, Q) + \sum_{p \in P_b} (\text{dist}(p, Q) - \text{dist}(b_p, Q)) \\ &\sim \sum_{b \in B} |P_b| \cdot \text{dist}(b, Q) + \sum_{p \in S} w(p) (\text{dist}(p, Q) - \text{dist}(b_p, Q)) \\ &= \sum_{b \in B} (|P_b| - w(S \cap P_b)) \text{dist}(b, Q) + \sum_{p \in S} w(p) \text{dist}(p, Q) . \end{aligned}$$

$$\sum_{p \in S} w(p) \cdot (\text{dist}(p, Q) - \text{dist}(b_p, Q)) \stackrel{?}{\sim} \sum_{p \in P} (\text{dist}(p, Q) - \text{dist}(b_p, Q))$$

Let p be the i th sample in S , and

$$x_i := w(p) \cdot (\text{dist}(p, Q) - \text{dist}(b_p, Q)).$$



$$\begin{aligned} E \left[\sum_{i=1}^{|S|} x_i \right] &= |S| \cdot E[x_1] \\ &= |S| \sum_{p \in P} \Pr[p] \cdot w(p) (\text{dist}(p, Q) - \text{dist}(b_p, Q)) \\ &= \sum_{p \in P} (\text{dist}(p, Q) - \text{dist}(b_p, Q)) \end{aligned}$$

$$\sum_{i=1}^{|S|} x_i \stackrel{?}{\sim} E \left[\sum_{i=1}^{|S|} x_i \right]$$

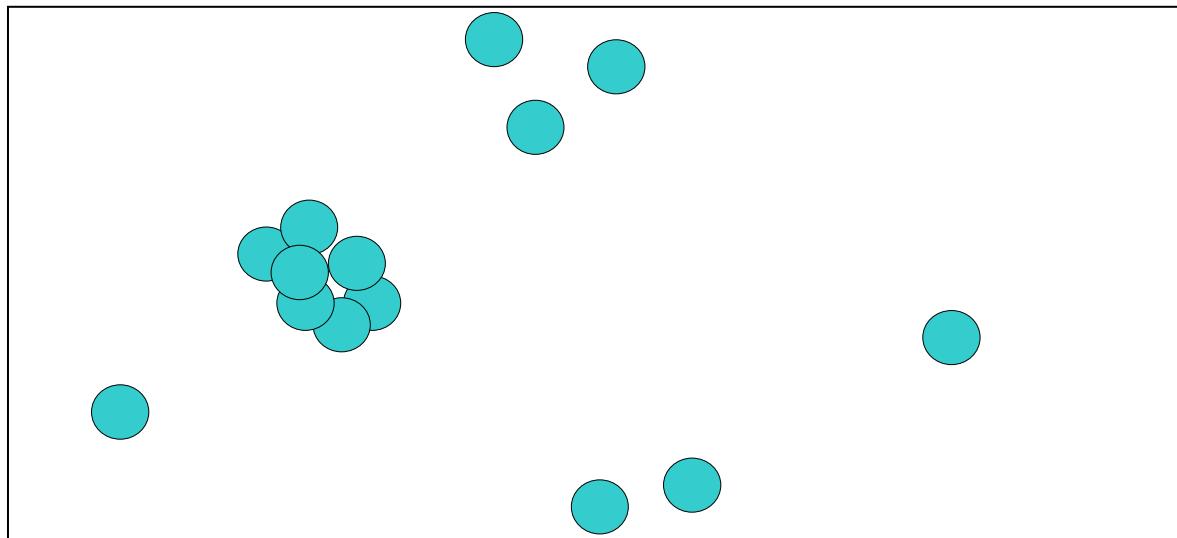
$$\begin{aligned}
|x_i| &\leq w(p) \cdot |\text{dist}(p, Q) - \text{dist}(b_p, Q)| \\
&= \frac{\sum_{p \in P} \text{dist}(p, B)}{|S| \text{dist}(p, b_p)} \cdot |\text{dist}(p, Q) - \text{dist}(b_p, Q)| \\
&\leq \frac{c \cdot \text{opt}}{|S|}
\end{aligned}$$

By Chernoff inequality, with constant probability,

$$\left| \sum_{i=1}^{|S|} x_i - E \left[\sum_{i=1}^{|S|} x_i \right] \right| \leq |S| \cdot \varepsilon \max_i |x_i| \leq \varepsilon c \cdot \text{opt}$$

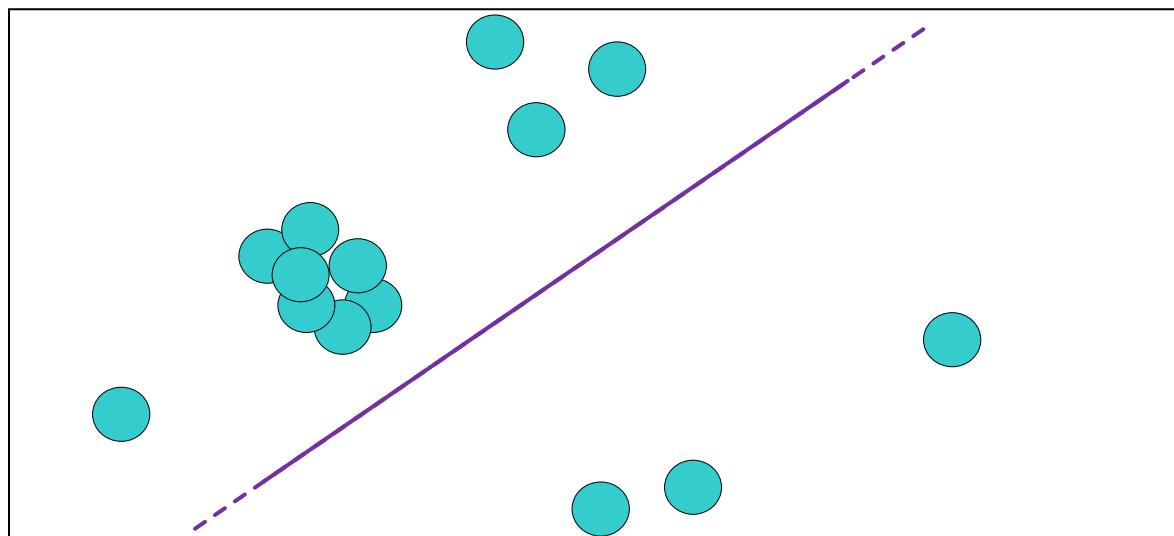
Line Queries

- Input: P in \mathbb{R}^d



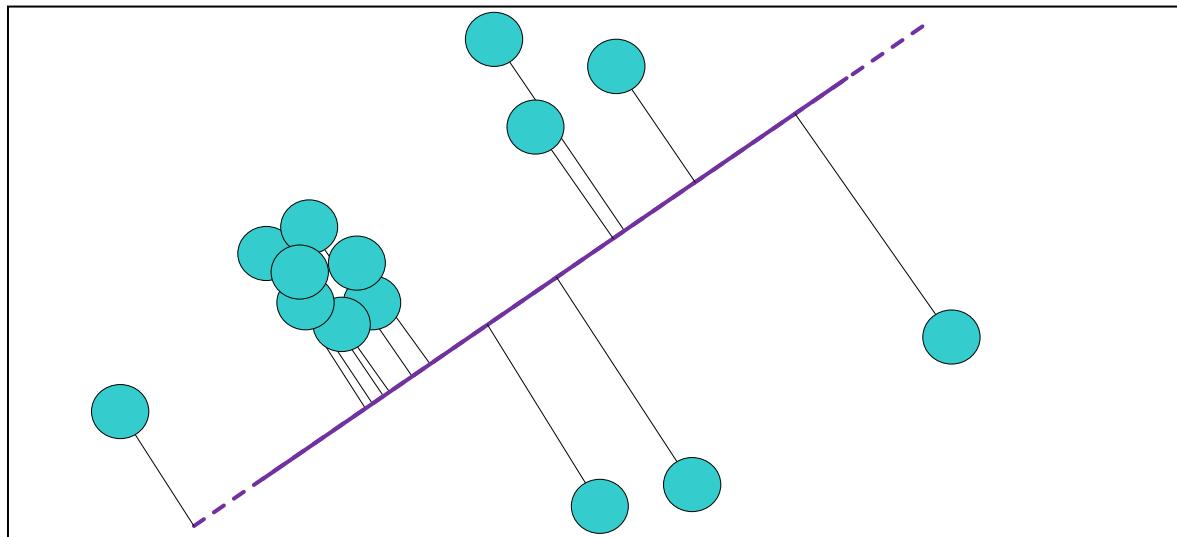
Line Queries

- Input: P in \mathbb{R}^d
- Query: a line Q



Line Queries

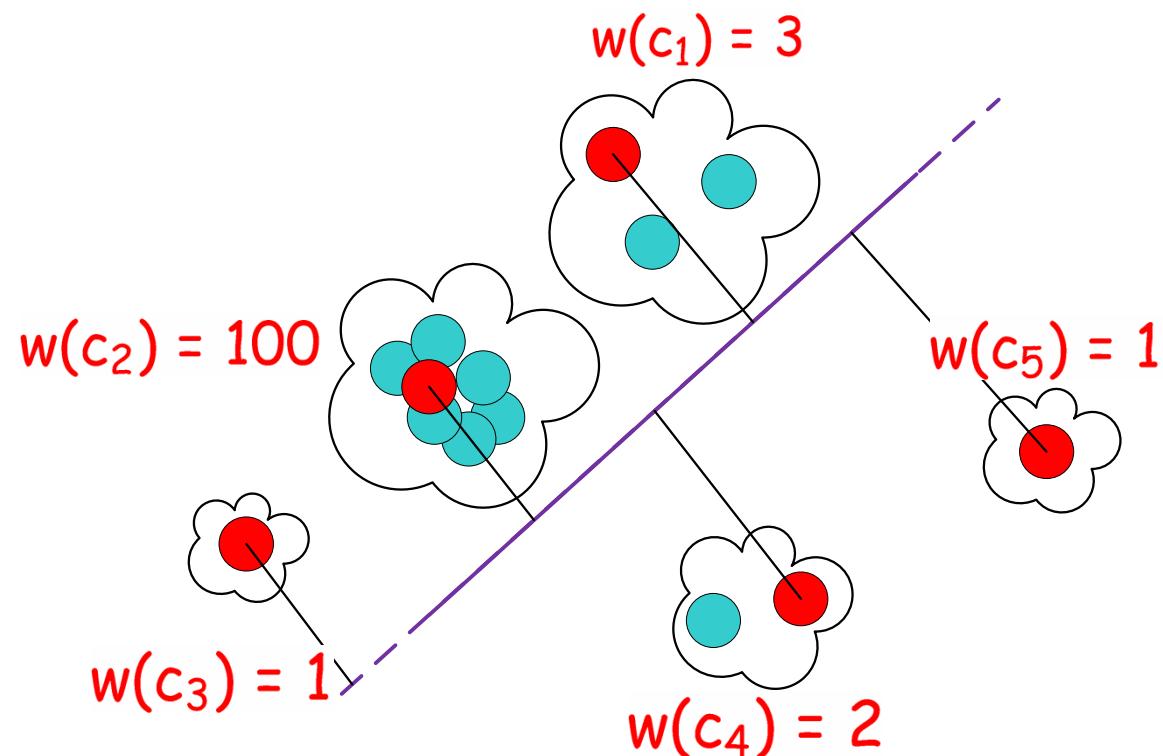
- Input: P in \mathbb{R}^d
- Query: a line Q
- Output: $\sum_{p \in P} \text{dist}(p, Q)$



ε -Line Coreset

Answer line queries in **sub-linear time**

Key Idea: Replace many points by weighted representatives:

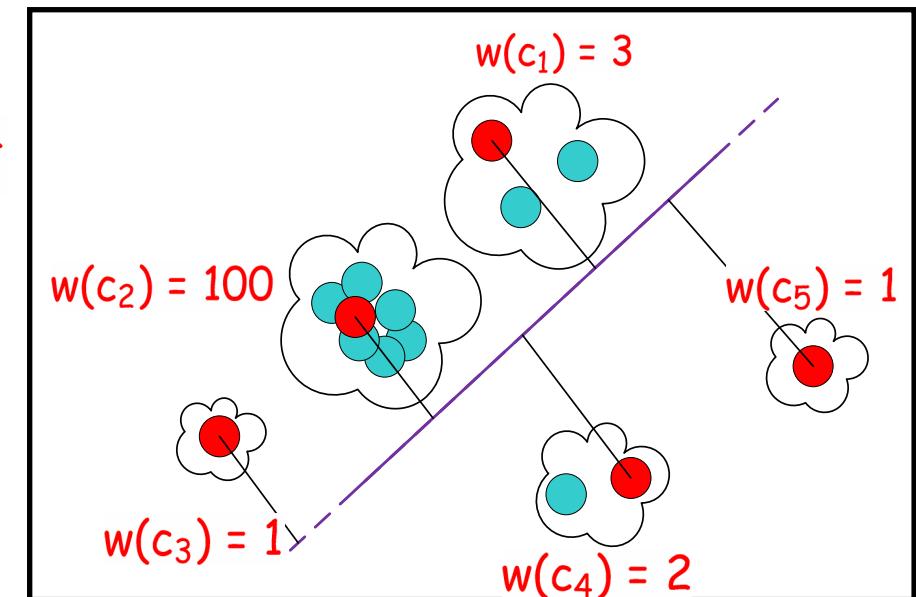


(j, ε) -Subspace Coreset

C is a (j, ε) -subspace coresset for P , if
for every j -dimensional subspace Q ,

$$\sum_{p \in P} \text{dist}(p, Q) \sim \sum_{c \in C} w(c) \cdot \text{dist}(c, Q)$$

Multiplicative error $\leq 1 + \varepsilon$



Approximation Algorithms (PTAS)

- Mean:

$O(nd^2)$ time, Exact (SVD) [Pearson, 1901]

$nd \cdot \text{poly}(j, 1/\epsilon)$ time, PTAS

[Deshpande et al.,] [Sarlos]

[Har-Peled] (2006)

- Median:

$nd \cdot 2^{\text{poly}(j, 1/\epsilon)}$ time, PTAS

[Shyamalkumar & Varadarajan, 2007]

- Center:

$nd \cdot \exp(\text{poly}(2^{(j^2)}, 1/\epsilon))$ time, PTAS

[Har-Peled & Varadarajan, 2004]

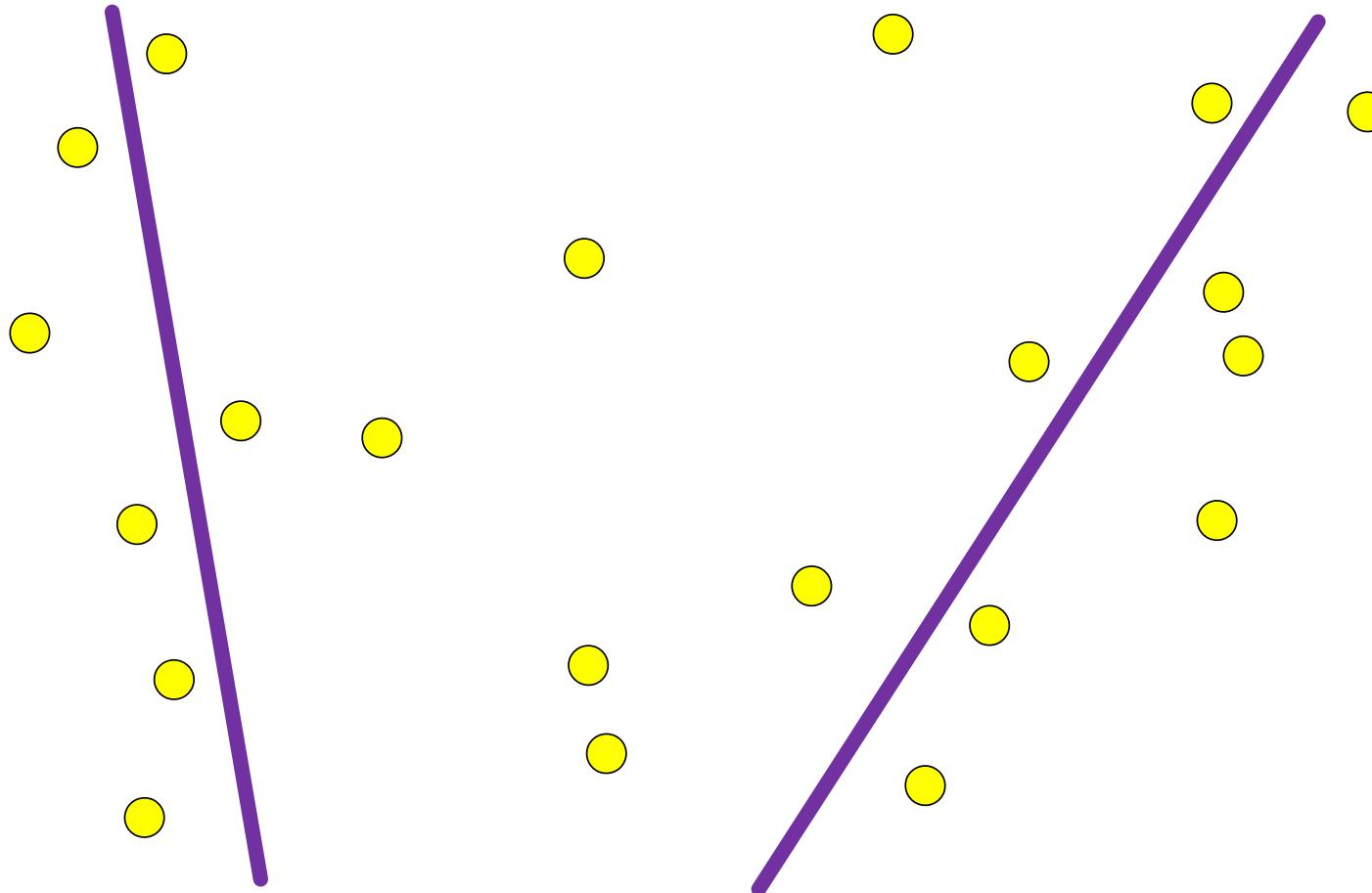
Coresets

- Strong (j, ε) -subspace coresets,
Size: $\exp(d)$, [F, Fiat, Sharir'06]
- Weak coresets for j -subspace Q
which minimizes $\max_{p \in P} \text{dist}(p, Q)$
Size: $\text{poly}(j/\varepsilon)$, [Har-Peled, Varadarajan'04]
- Weak coresets for regression distances,
Size: $\text{poly}(d/\varepsilon)$,
[Dasgupta, Drineas, Harb, Kumar, Mahoney'08]

Recent Results [FMSW10]

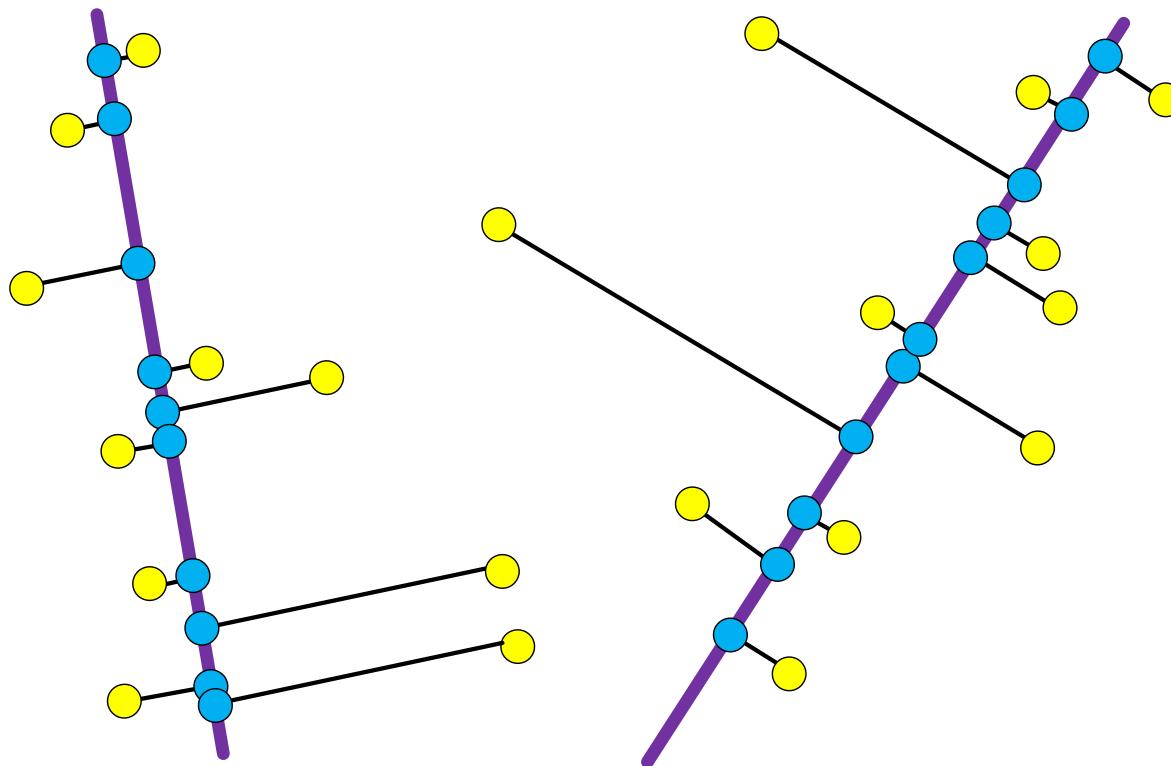
- Weak (j, ε) -coresets,
Size: $\text{poly}(1/\varepsilon)$ for a fixed j
- Strong (j, ε) -coresets,
Size: $\text{poly}(d/\varepsilon)$ for a fixed j
- Small sketches for streaming updates of
coordinates of points
- General dimensional reduction technique
for shape fitting

Coreset for 2-lines Median



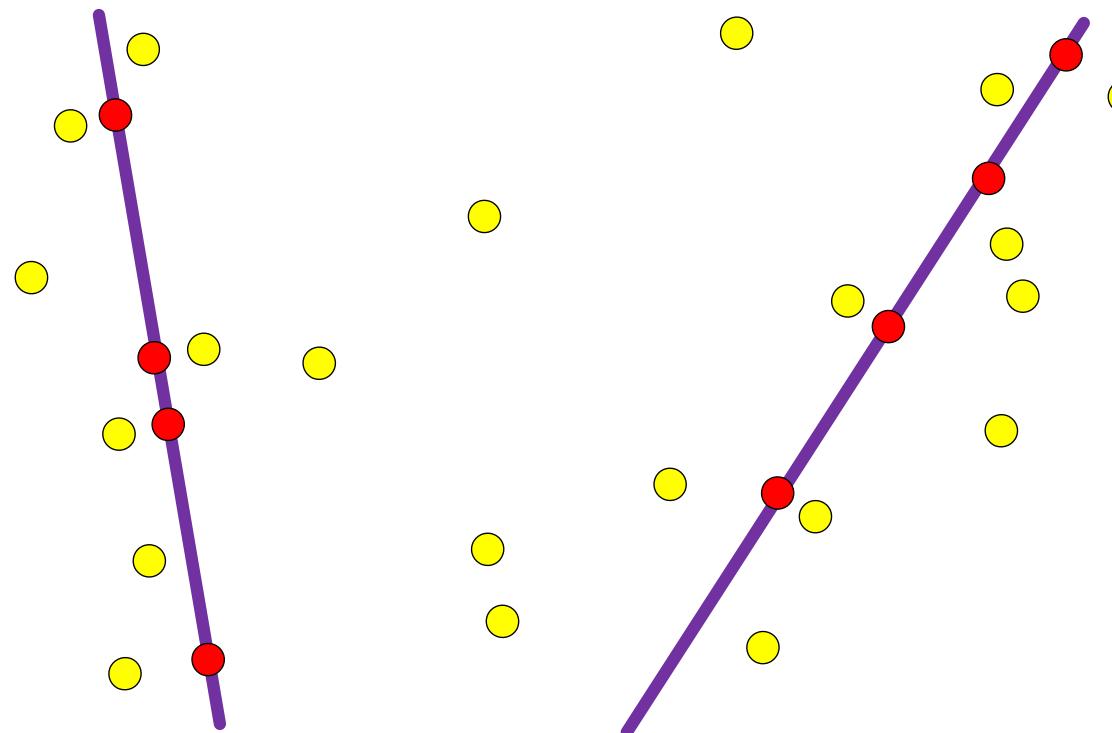
Coreset for 2-lines Median

1. $\forall b \in B : P_b \leftarrow$ projection of P onto b



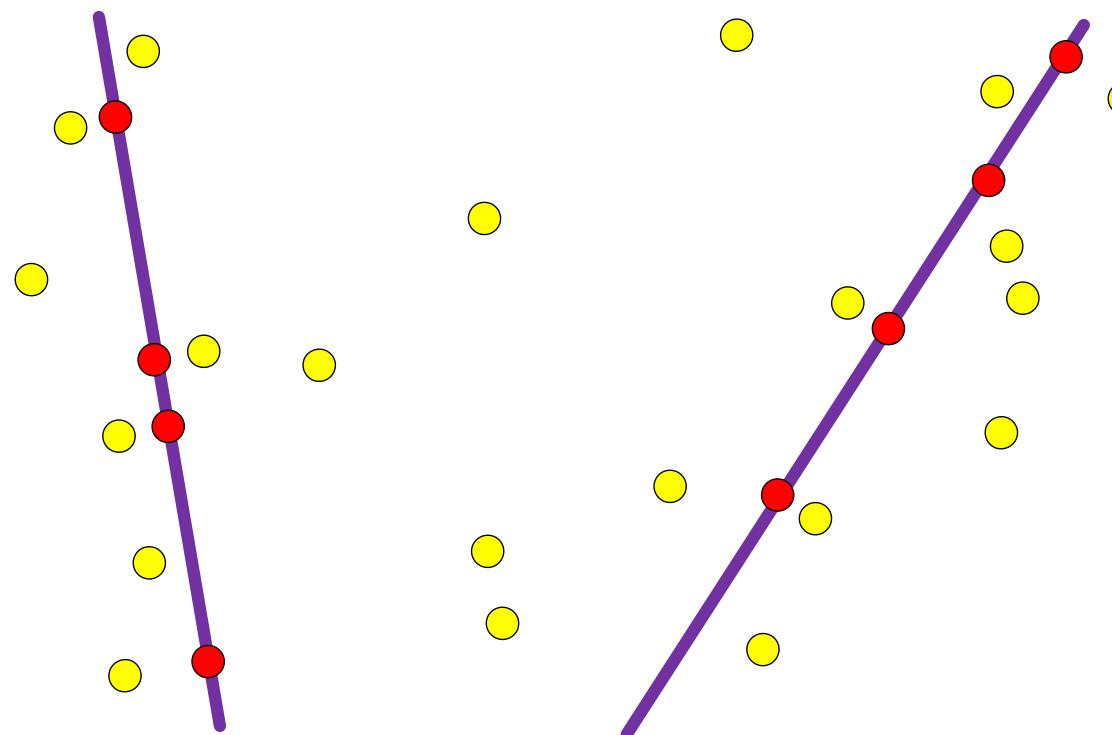
Coreset for 2-lines Median

1. $\forall b \in B : P_b \leftarrow$ projection of P onto b
2. $T_b \leftarrow$ coresset for P_b [FFS06]



Coreset for 2-lines Median

1. $\forall b \in B : P_b \leftarrow$ projection of P onto b
2. $T_b \leftarrow$ coreset for P_b [FFS06]
3. $T \leftarrow \bigcup_{b \in B} T_b$

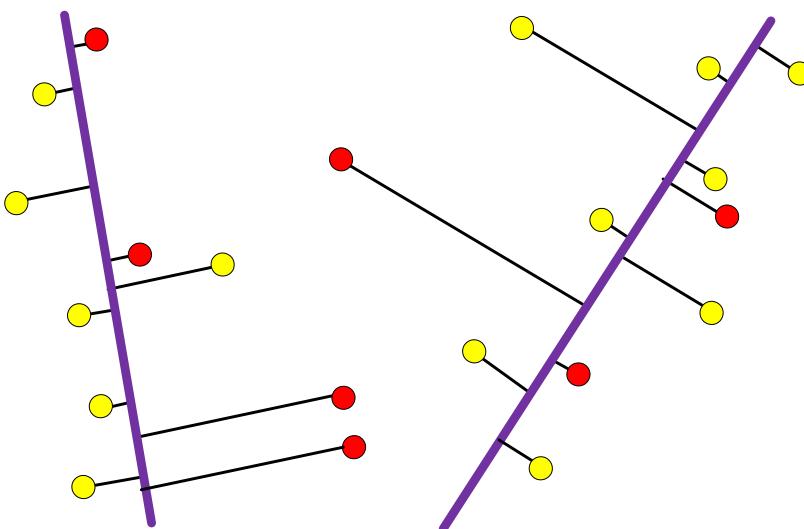


Coreset for 2-lines Median

4. Pick a sample S of $d/\varepsilon^{O(1)}$ points, where

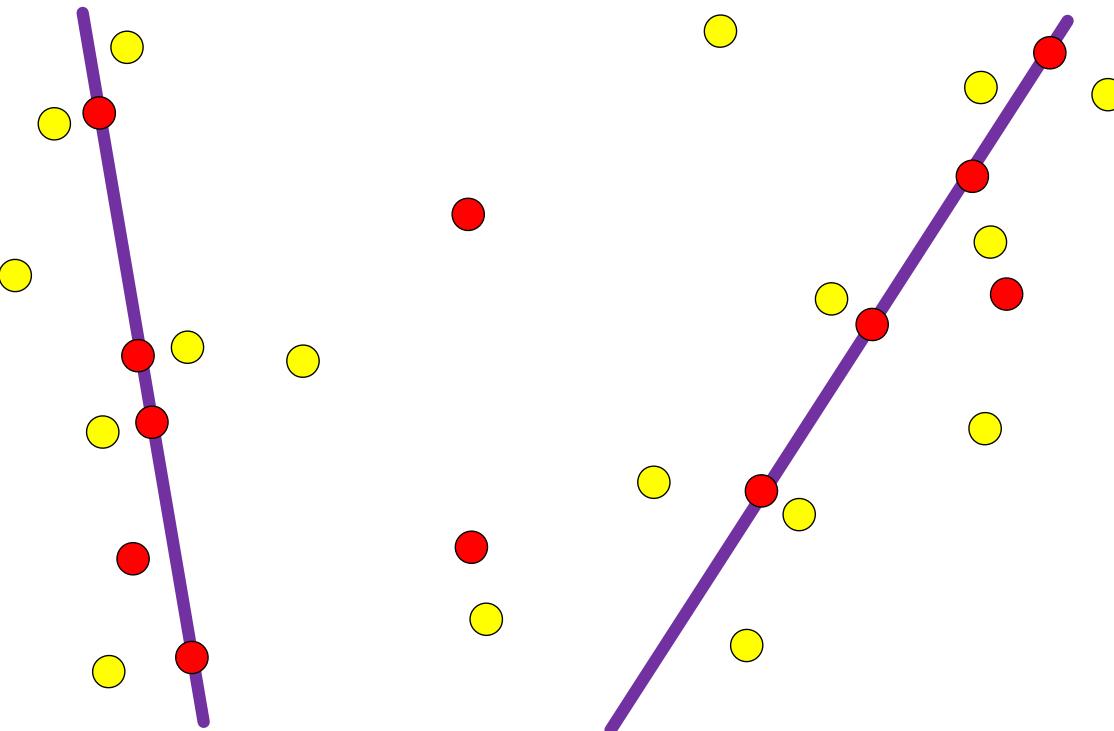
$$\Pr[p] = \frac{\text{dist}(p, B)}{\sum_{p \in P} \text{dist}(p, B)}$$

5. $\forall p \in S : w(p) \leftarrow \frac{1}{|S| \cdot \Pr[p]}$



Coreset for 2-lines Median

6. Return $T \cup S \cup \text{project}(S^-, B)$



Results for $j, k > 1$

PTAS that takes $d \cdot n^{\text{poly}(j,k,1/\epsilon)}$ time.

Mean: [Deshpande et al., 2006]

Median: [Shyamalkumar & Varadarajan, 2007]

Center: [Har-Peled & Varadarajan, 2002]

Thank You!