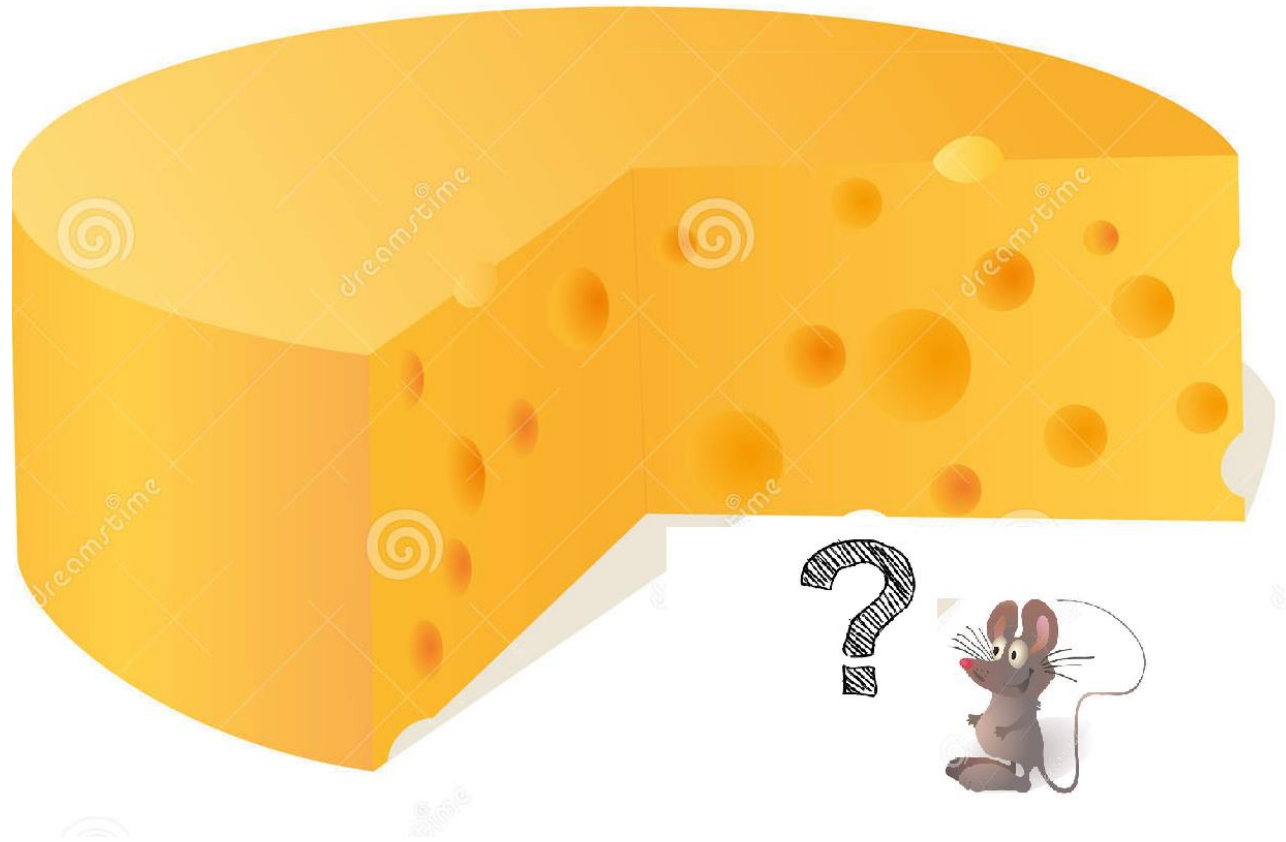


# Big Data Class



---

LECTURER: DAN FELDMAN

TEACHING ASSISTANTS:

IBRAHIM JUBRAN

ALAA MAALOUF

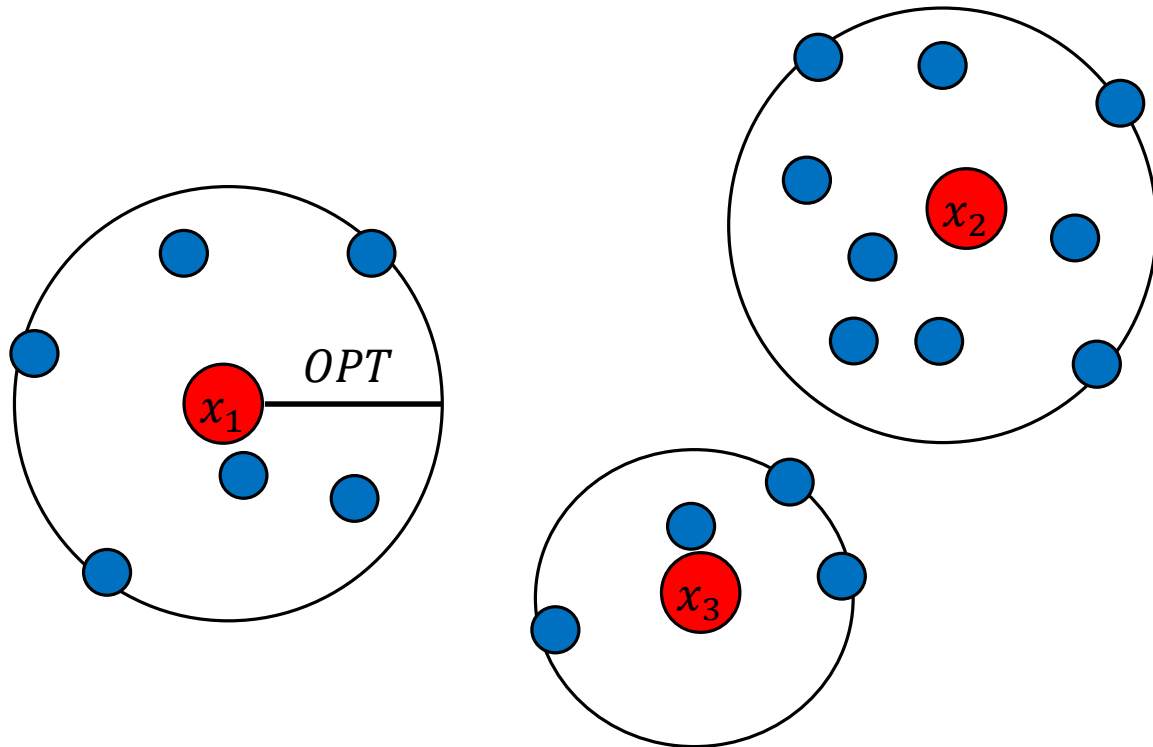


# $k$ -Center / $k$ -Minimum Enclosing Balls

- Given a set of  $n$  points  $P$  in  $R^d$  and an integer  $k$ , find a set of  $k$  points (centers)  $X = \{x_1, \dots, x_k\} \subseteq R^d$  that minimizes:

$$far_k(P, X) = \max_{p \in P} \|p - X(p)\|$$

Closest point in  $X$  to  $p$   
 $X(p) \in \arg \min_{x \in X} \|x - p\|$



# $k$ -Center / $k$ -Minimum Enclosing Balls

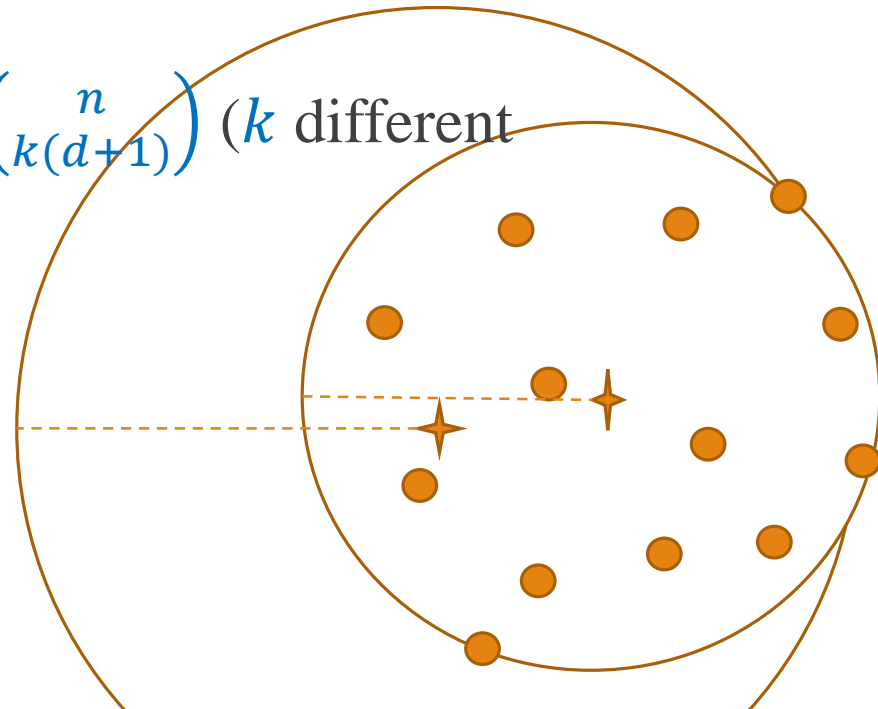
## Optimal solution in $R^d$ :

**Claim 1:** A sphere in  $R^d$  is determined by  $d + 1$  points.

**Claim 2:** A sphere with minimal radius enclosing a set of points in  $R^d$  passes through  $d + 1$  points from the set.

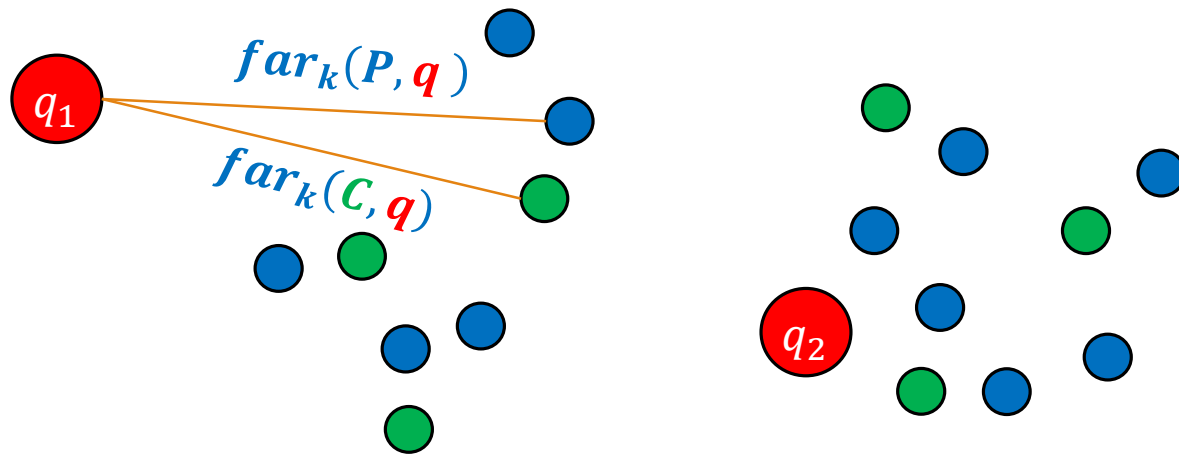
**Algorithm:** Exhaustive search over all possible tuples  $\binom{n}{k(d+1)}$  ( $k$  different circles, each determined by  $d + 1$  points).

**Running time:**  $n^{O(dk)}$ .

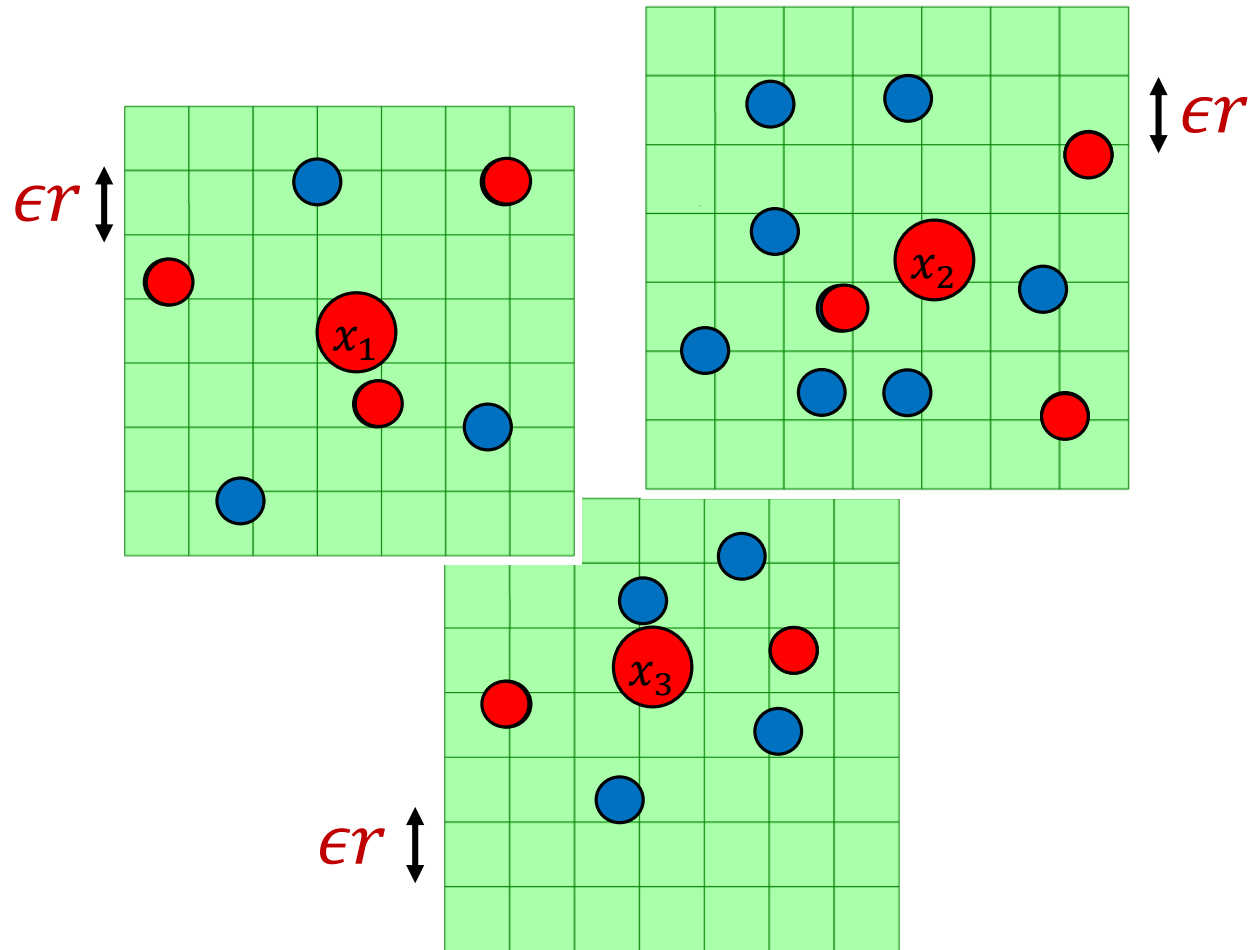


# Coreset for $k$ -Center

- Input:  $(P, k, Q)$  where  $P \subseteq \mathbb{R}^d$ ,  $k$  is an integer and  $Q \subseteq (\mathbb{R}^d)^k$ .
- Output:  $C \subseteq P, |C| = k \left(\frac{1}{\epsilon}\right)^{O(d)}$  s. t. for every  $q \in Q$ :  
 $far_k(P, q) - far_k(C, q) \leq O(\epsilon) \cdot far_k(P, q)$



# Coreset for $k$ -Center



## $(1 + \epsilon)$ -Coreset

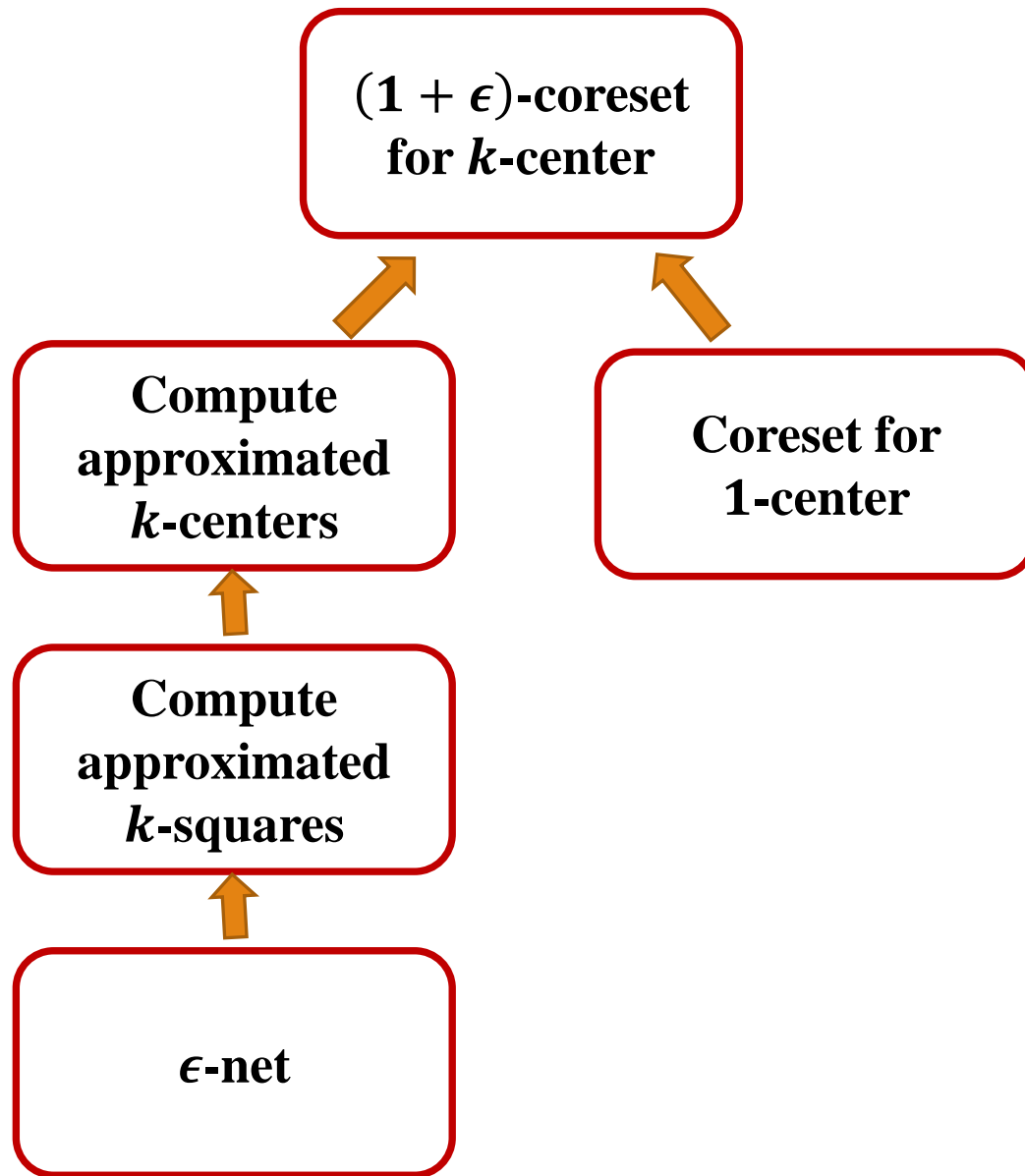
### Algorithm:

- 1) Find optimal  $k$ -centers  
(To find  $k$  “clusters” and the optimal radius  $OPT$ ).
- 2) Compute 1-center coreset for each cluster where  $r = OPT$ .

Total time:  $n^{O(dk)}$

Coreset size:  $k \cdot \left(\frac{1}{\epsilon}\right)^{O(d)}$

# Coreset for $k$ -Center



## $(1 + \epsilon)$ -Coreset

### Algorithm:

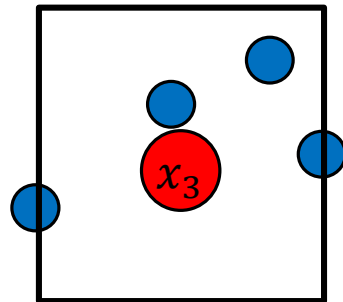
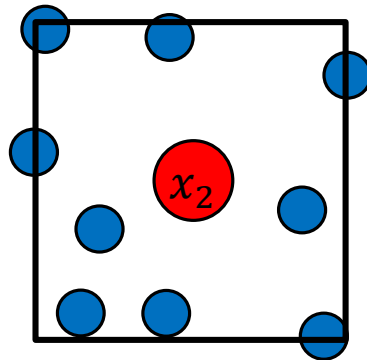
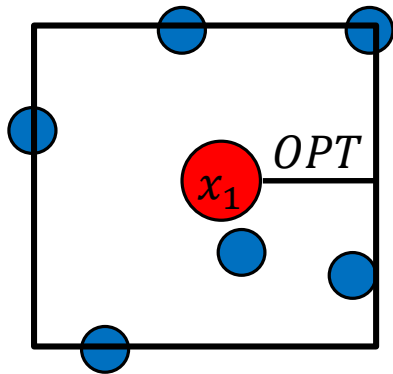
- 1) Find optimal  $k$ -centers (To find  $k$  “clusters” and the optimal radius  $OPT$ ).
- 2) Compute 1-center coreset for each cluster where  $r = OPT$ .

# $K$ -Minimum Enclosing Squares

- Given a set of  $n$  points  $P$  in  $R^d$ , find a set of  $k$  points (centers)  $X = \{x_1, \dots, x_k\} \subseteq R^d$  that minimizes:

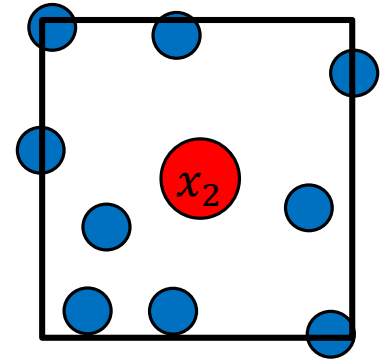
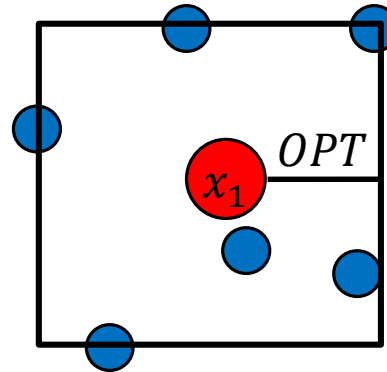
$$far_k(P, X) = \max_{p \in P} \|p - X(p)\|_\infty$$

Closest point in  $X$  to  $p$   
 $X(p) \in \arg \min_{x \in X} \|x - p\|_\infty$



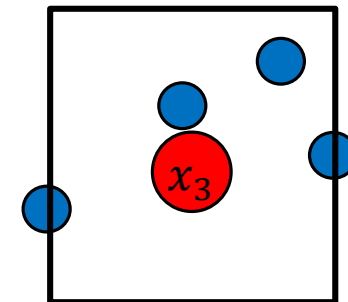
# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

- Given  $K$ -minimum Enclosing Squares  
where  $OPT = \mathit{far}_k(P, X) = \max_{p \in P} \|p - X(p)\|_\infty$



We want to compute:

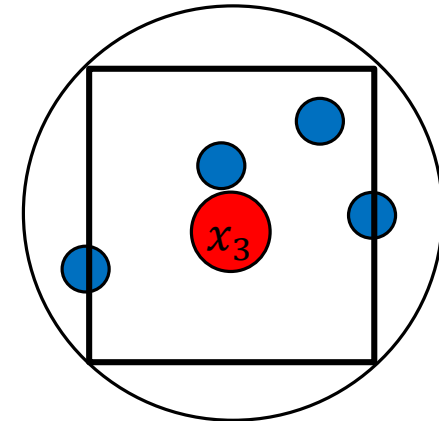
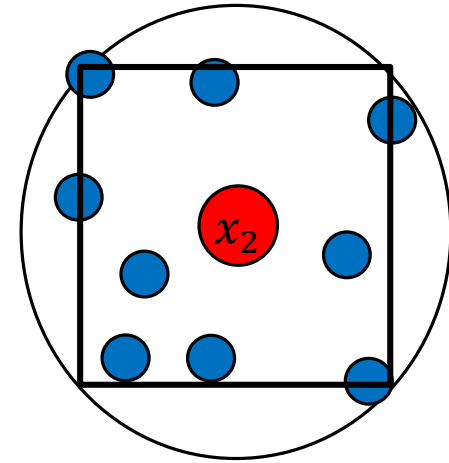
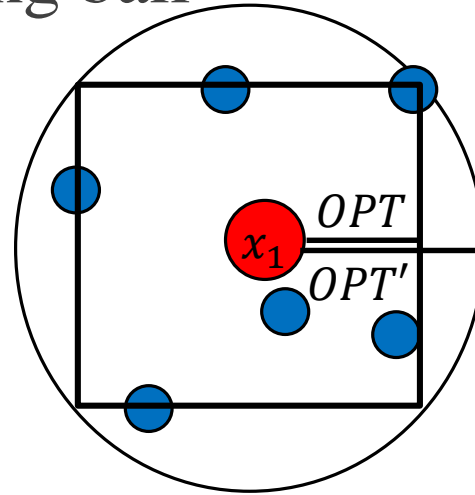
$$\widehat{OPT} = \mathit{far}_k(P, X) = \max_{p \in P} \|p - X(p)\|$$





# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

- Given  $K$ -minimum Enclosing Squares
- For each square, draw an enclosing ball



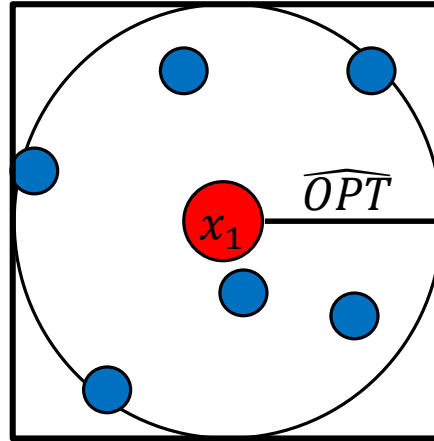
Claim:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

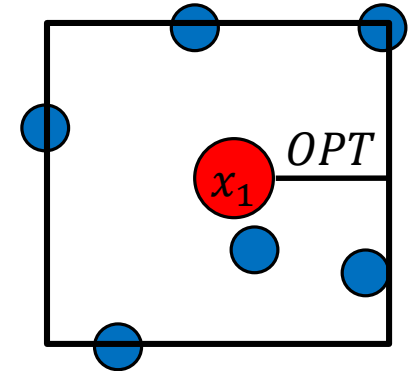
Proof:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

1)  $OPT \leq \widehat{OPT}$  (by definition)

$$\widehat{OPT} = \max_{p \in P} \|p - X(p)\|$$



$$OPT = \max_{p \in P} \|p - X(p)\|_{\infty}$$

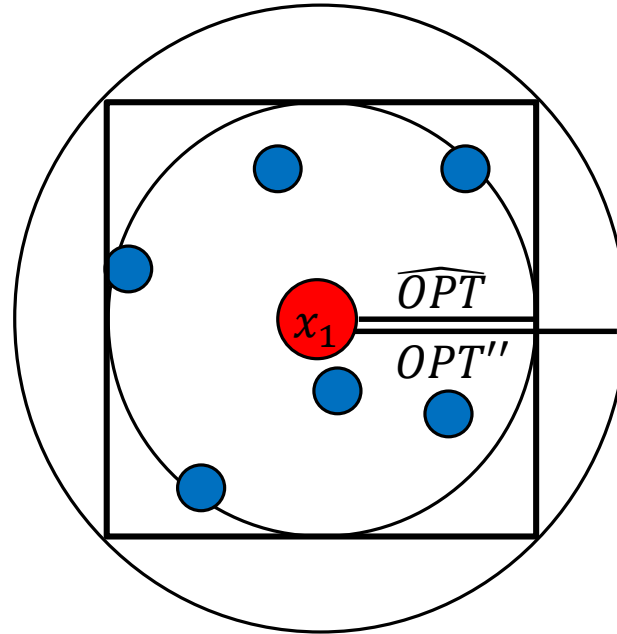


# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

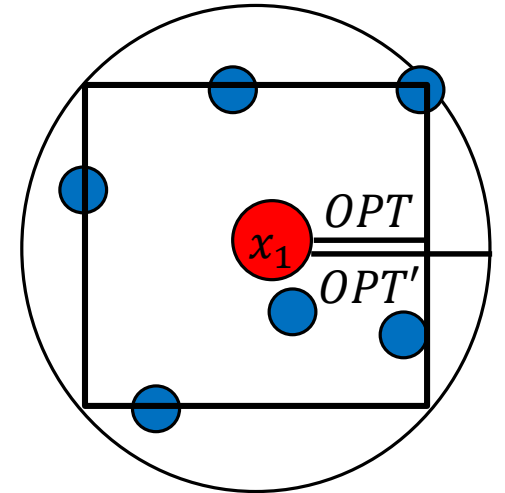
Proof:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

1)  $OPT \leq \widehat{OPT}$  (by definition)

$$\widehat{OPT} = \max_{p \in P} \|p - X(p)\|$$



$$OPT = \max_{p \in P} \|p - X(p)\|_{\infty}$$



# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

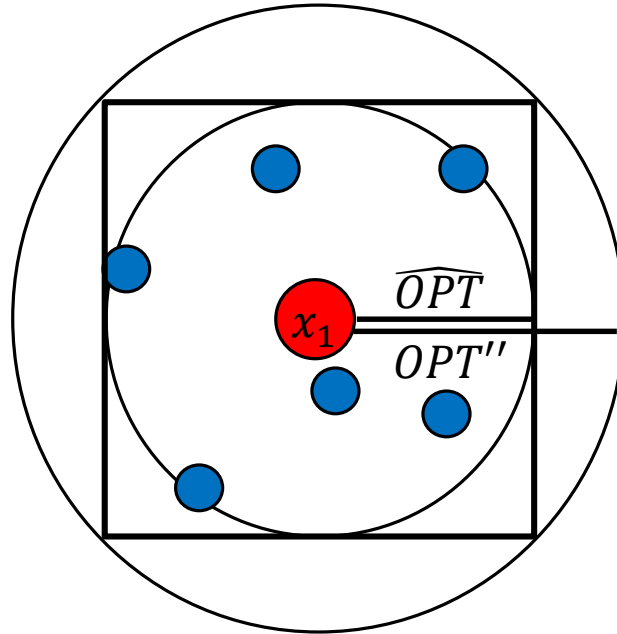
Proof:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

1)  $OPT \leq \widehat{OPT}$  (by definition)

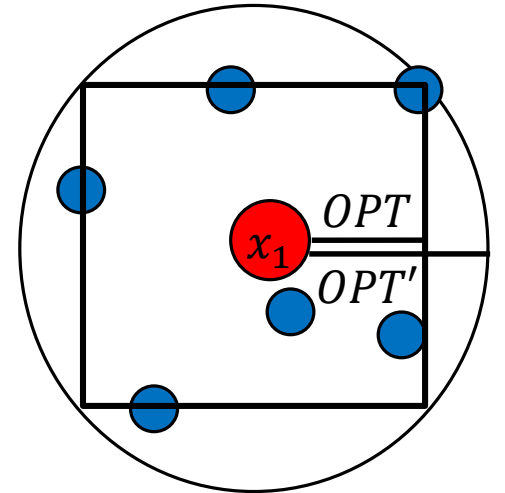
2)  $OPT' \leq OPT''$

Claim:  $OPT'' \leq \sqrt{d} \cdot \widehat{OPT}$

$$\widehat{OPT} = \max_{p \in P} \|p - X(p)\|$$



$$OPT = \max_{p \in P} \|p - X(p)\|_{\infty}$$



# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

Proof:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

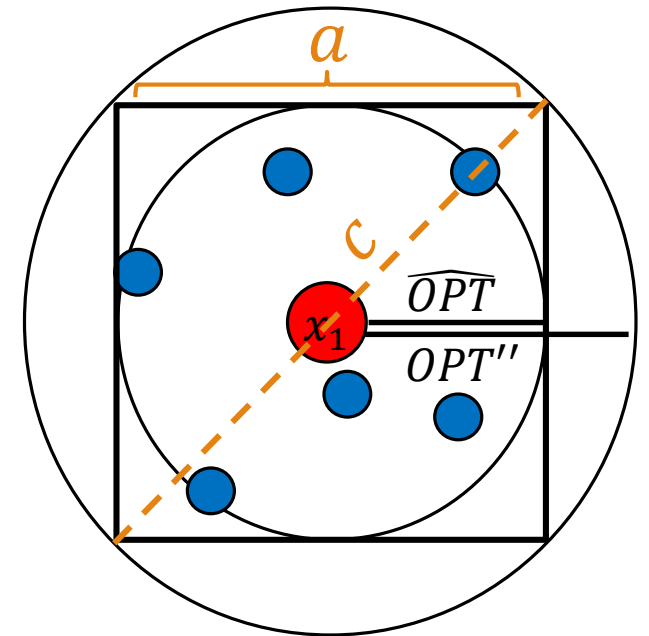
Proof:  $OPT'' \leq \sqrt{d} \cdot \widehat{OPT}$

$$a = 2 \cdot \widehat{OPT}$$

$$\begin{aligned} c &= \sqrt{a^2 + \dots + a^2} = \sqrt{da^2} \\ &= \sqrt{d} a = 2\sqrt{d} \cdot \widehat{OPT} \end{aligned}$$

$$OPT'' = \frac{c}{2} = \sqrt{d} \cdot \widehat{OPT}$$

$$\widehat{OPT} = \max_{p \in P} \|p - X(p)\|$$



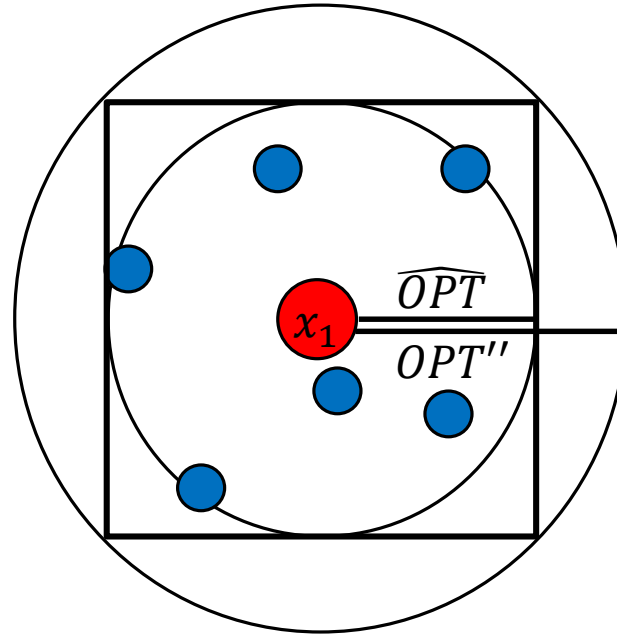
# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

Proof:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

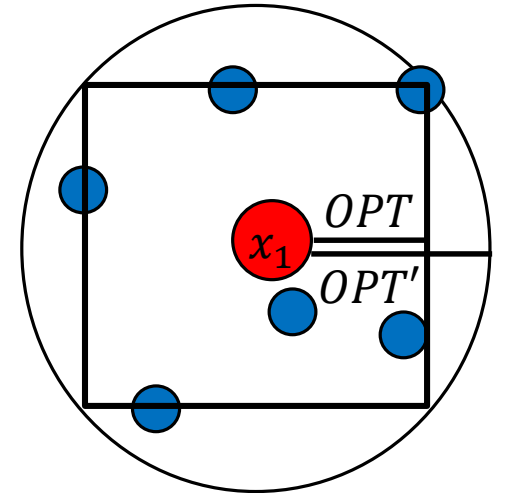
1)  $OPT \leq \widehat{OPT}$  (by definition)

2)  $OPT' \leq OPT''$

$$\widehat{OPT} = \max_{p \in P} \|p - X(p)\|$$



$$OPT = \max_{p \in P} \|p - X(p)\|_{\infty}$$



# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

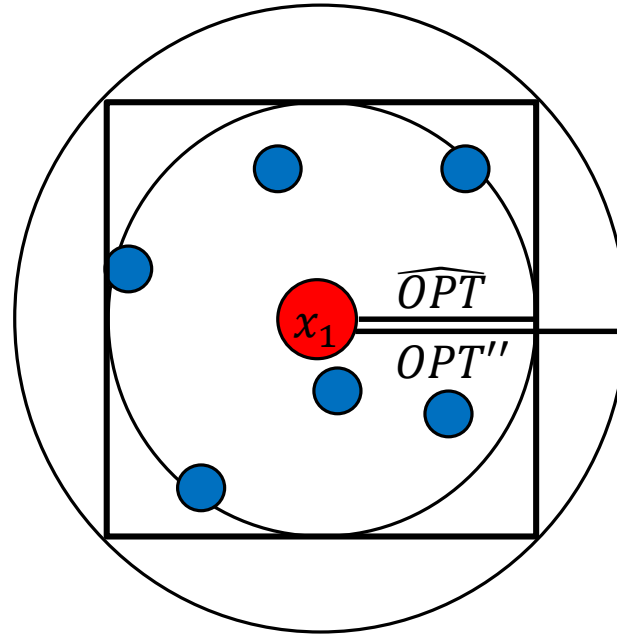
Proof:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

1)  $OPT \leq \widehat{OPT}$  (by definition)

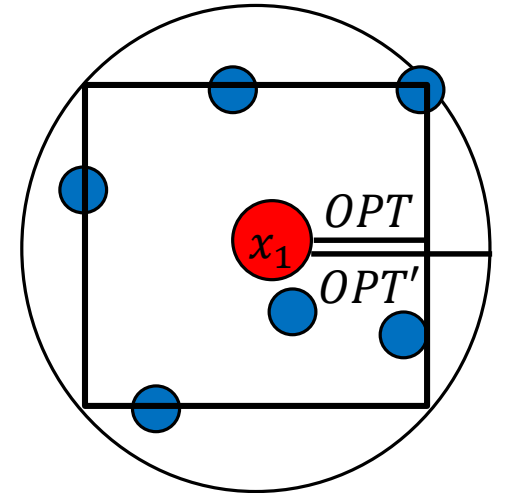
2)  $OPT' \leq OPT''$

3)  $OPT'' \leq \sqrt{d} \cdot \widehat{OPT}$

$$\widehat{OPT} = \max_{p \in P} \|p - X(p)\|$$



$$OPT = \max_{p \in P} \|p - X(p)\|_{\infty}$$



# $K$ -Minimum Enclosing Squares $\rightarrow$ $K$ -Center

Proof:  $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$

1)  $OPT \leq \widehat{OPT}$  (by definition)

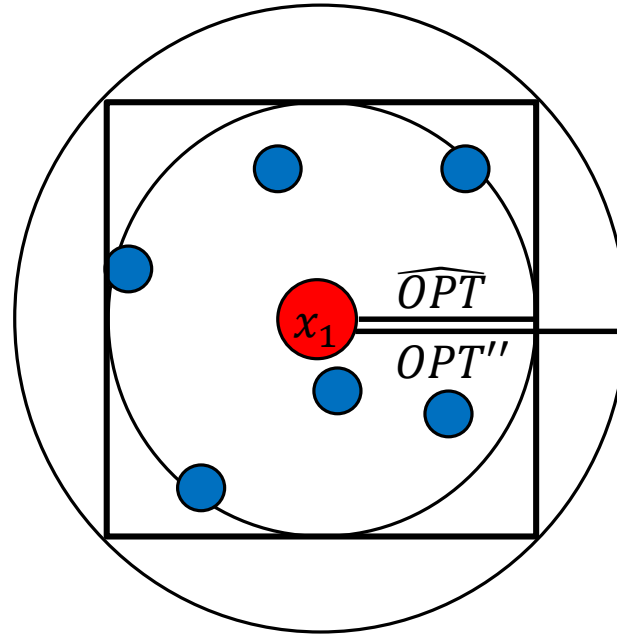
2)  $OPT' \leq OPT''$

3)  $OPT'' \leq \sqrt{d} \cdot \widehat{OPT}$

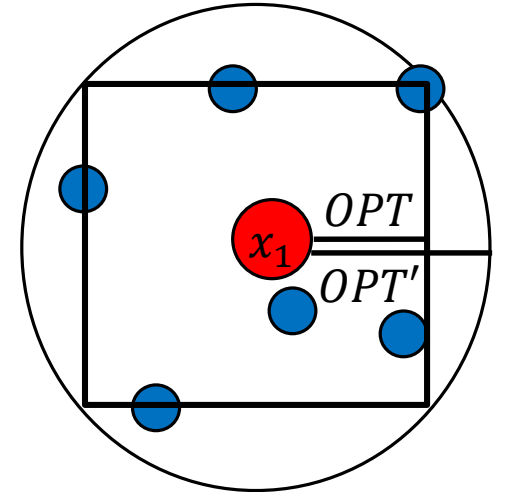
$\rightarrow$   $OPT' \leq \sqrt{d} \cdot \widehat{OPT}$



$$\widehat{OPT} = \max_{p \in P} \|p - X(p)\|$$



$$OPT = \max_{p \in P} \|p - X(p)\|_{\infty}$$





# More Formal: $k$ -Squares $\rightarrow k$ -Centers

## Claim 1:

$$far_{\infty}(P, q) \leq far_2(P, q) \leq \sqrt{d} \cdot far_{\infty}(P, q)$$

## Proof of claim 1:

$$\begin{aligned} far_{\infty}(P, q) &= \max_{p \in P} \|p - q\|_{\infty} \\ &\leq \max_{p \in P} \|p - q\|_2 = far_2(P, q) \\ &= \max_{p \in P} \sqrt{(p(1) - q(1))^2 + \dots + (p(d) - q(d))^2} \\ &\leq \max_{p \in P} \sqrt{d \cdot \max_i (p(i) - q(i))^2} \\ &= \sqrt{d} \cdot \max_{p \in P} (\max_i |p(i) - q(i)|) \\ &= \sqrt{d} \cdot \max_{p \in P} (\|p - q\|_{\infty}) = \sqrt{d} \cdot far_{\infty}(P, q) \end{aligned}$$

## Definitions:

$$far_2(P, q) = \max_{p \in P} \|p - q\|_2$$

$$OPT_2 = \operatorname{argmin}_{q \in Q} far_2(P, q)$$

$$far_{\infty}(P, q) = \max_{p \in P} \|p - q\|_{\infty}$$

$$OPT_{\infty} = \operatorname{argmin}_{q \in Q} far_{\infty}(P, q)$$

# More Formal: $k$ -Squares $\rightarrow$ $k$ -Centers

## Claim 2:

An  $\alpha$ -approximation for  $k$ -squares is an  $O(\alpha \cdot \sqrt{d})$ -approximation for  $k$ -centers.

## Proof of claim 2:

Let  $a_\infty$  be the  $\alpha$ -approximation for  $k$ -squares  $\rightarrow \text{far}_\infty(P, a_\infty) \leq \alpha \cdot \text{far}_\infty(P, OPT_\infty)$ .

$$\rightarrow \text{far}_2(P, a_\infty) \leq \sqrt{d} \cdot \text{far}_\infty(P, a_\infty) \quad (\text{Right side of Claim 1})$$

$$\leq \alpha \sqrt{d} \cdot \text{far}_\infty(P, OPT_\infty) \quad (\text{Definition of } a_\infty)$$

$$\leq \alpha \sqrt{d} \cdot \text{far}_\infty(P, OPT_2) \quad (\text{Definition of } OPT_\infty)$$

$$\leq \alpha \sqrt{d} \cdot \text{far}_2(P, OPT_2) \quad (\text{Left side of Claim 1})$$

# $\epsilon$ -net (d=1)

• Input:

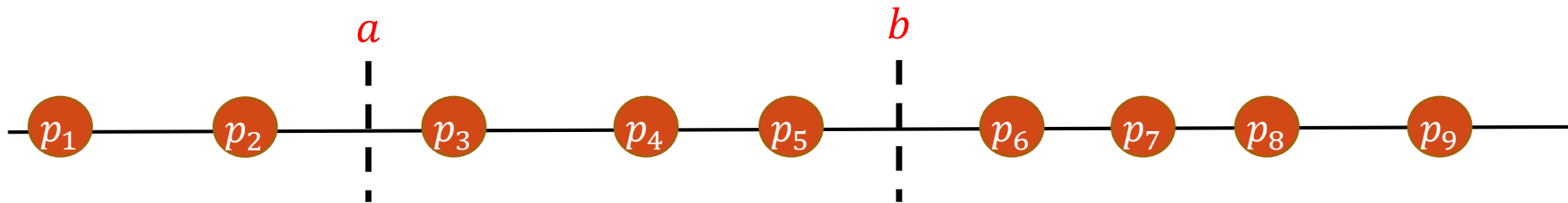
$$P \subseteq R, Q = \{[a, b] \mid a, b \in R, a \leq b\}$$

• Output:

$$C \subseteq P, |C| = \frac{1}{\epsilon} \text{ s. t. for every } q \in Q:$$

$$\left| \frac{|P \cap q|}{|P|} - \frac{|C \cap q|}{|C|} \right| \leq \epsilon$$

$q = [a, b]$



$$\frac{|P \cap q|}{|P|} = \frac{3}{9} = \frac{1}{3}$$

How should we pick the subset  $C$ ?

We need  $C$  to have the same percentage of points in  $[a, b]$  as  $P$ .

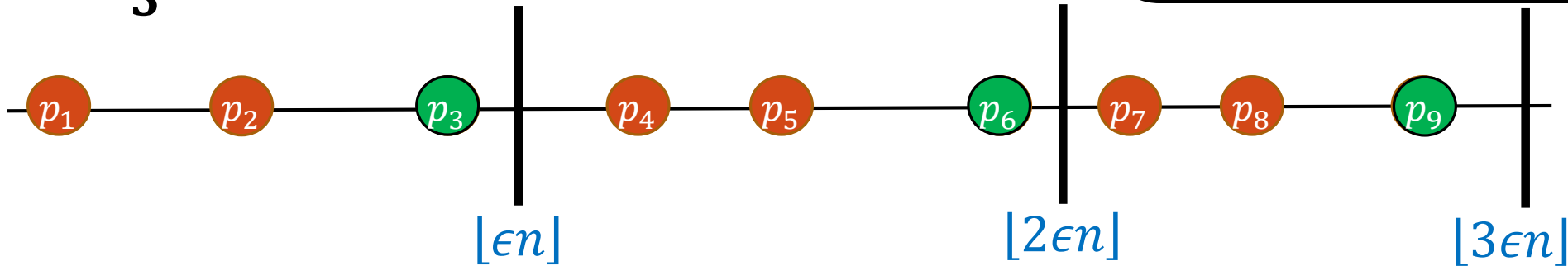
→ we want to “sample points” uniformly from  $P$ .

# $\epsilon$ -net (d=1)

Example:

$$\epsilon = \frac{1}{3}, n = 9$$

$$|C| = \frac{1}{\epsilon}$$



## Algorithm:

Assume  $p_1 \leq p_2 \leq \dots \leq p_n$

$$I = \{\lfloor \epsilon n \rfloor, \lfloor 2\epsilon n \rfloor, \dots, \lfloor n \rfloor\}$$

$$C = \{p_i \mid i \in I\}.$$

Claim 1: Assume  $|P \cap q| = k$ . By the construction of  $C$  we

get that:  $\lfloor \frac{k}{\epsilon n} \rfloor \leq |C \cap q| \leq \lceil \frac{k}{\epsilon n} \rceil$ .

Claim 2:  $\frac{|C \cap q|}{|C|} \leq \frac{\lceil \frac{k}{\epsilon n} \rceil}{\lfloor \frac{k}{\epsilon n} \rfloor} = \epsilon \frac{k}{\epsilon n} \leq \frac{k}{n} + \epsilon$

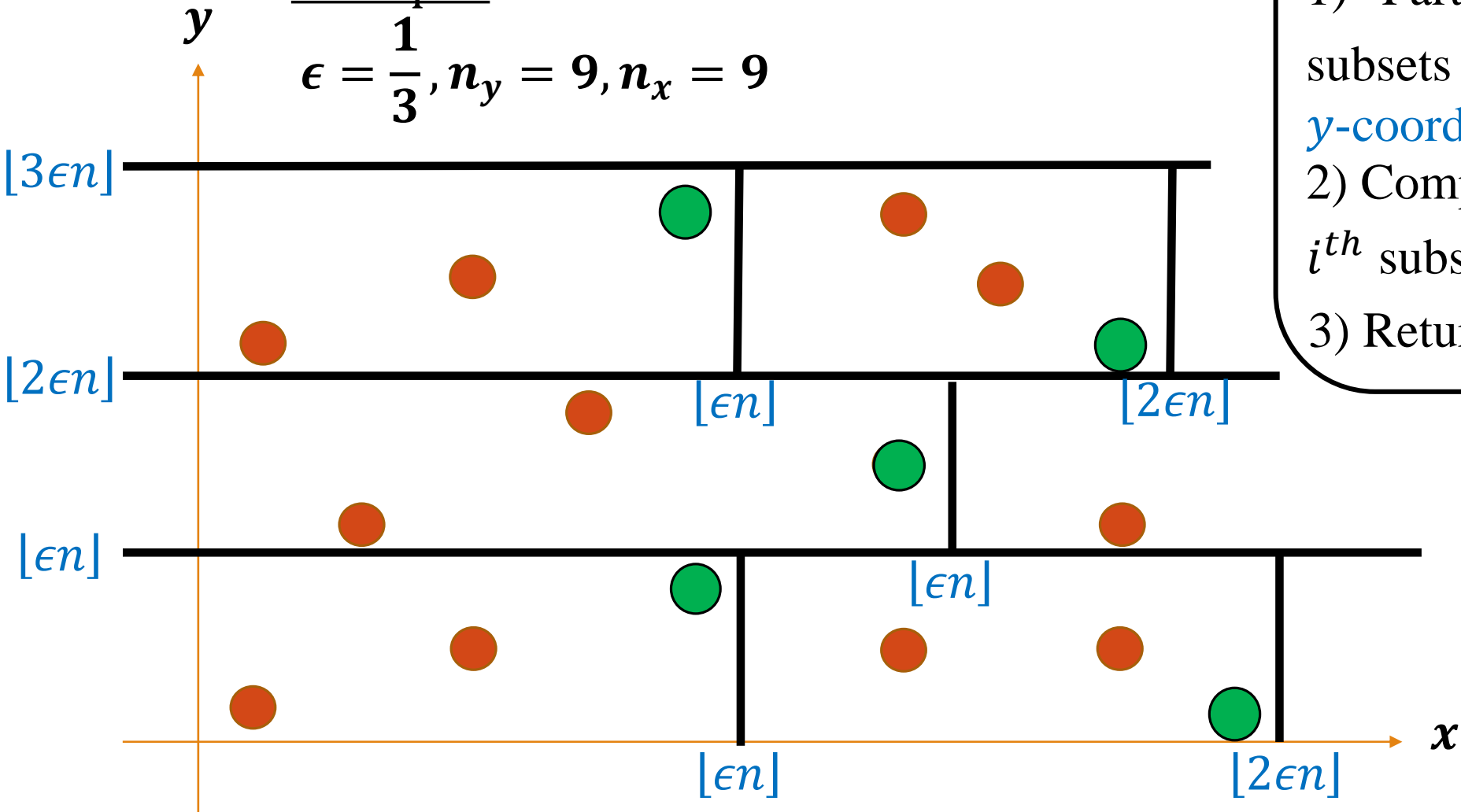
Claim 3:  $\frac{|C \cap q|}{|C|} \geq \frac{\lfloor \frac{k}{\epsilon n} \rfloor}{\lceil \frac{k}{\epsilon n} \rceil} = \epsilon \frac{k}{\epsilon n} \geq \frac{k}{n} - \epsilon$

$$\left| \underbrace{\frac{|P \cap q|}{|P|}}_{\frac{k}{n}} - \frac{|C \cap q|}{|C|} \right| \leq \epsilon$$

# $\epsilon$ -net (d=2)

Example:

$$\epsilon = \frac{1}{3}, n_y = 9, n_x = 9$$



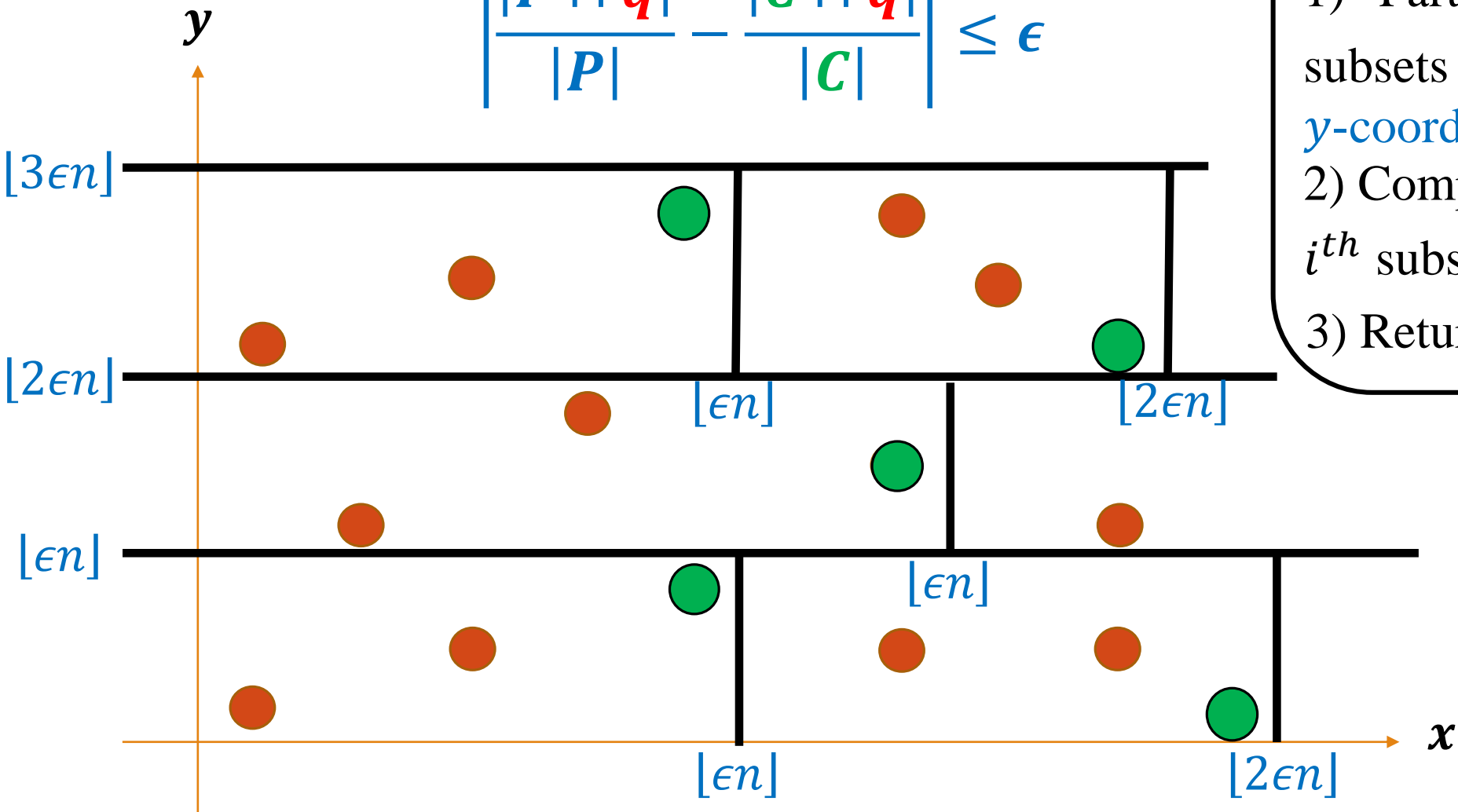
## Algorithm:

- 1) Partition the points into  $\lfloor \frac{1}{\epsilon} \rfloor$  subsets by their sorted *y-coordinate*
- 2) Compute an  $\epsilon$ -net  $C_i$  for the  $i^{th}$  subset for every  $i \in \lfloor \frac{1}{\epsilon} \rfloor$ .
- 3) Return  $C = \cup_i C_i$

$$|C| = O\left(\left(\frac{1}{\epsilon}\right)^d\right)$$

# $\epsilon$ -net (d=2)

$$\left| \frac{|P \cap q|}{|P|} - \frac{|C \cap q|}{|C|} \right| \leq \epsilon$$



## Algorithm:

- 1) Partition the points into  $\left\lfloor \frac{1}{\epsilon} \right\rfloor$  subsets by their sorted  $y$ -coordinate
- 2) Compute an  $\epsilon$ -net  $C_i$  for the  $i^{\text{th}}$  subset for every  $i \in \left\lfloor \frac{1}{\epsilon} \right\rfloor$ .
- 3) Return  $C = \bigcup_i C_i$

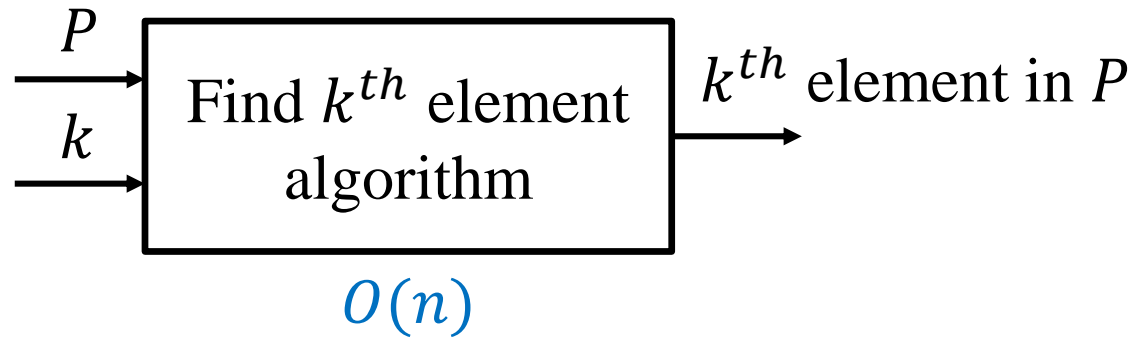
$$|C| = O\left(\left(\frac{1}{\epsilon}\right)^d\right)$$

$$\text{Time} = O\left(\frac{nd}{\epsilon}\right)$$

# $\epsilon$ -net time analysis

## Algorithm for computing the $\epsilon$ -net:

Assuming there is an algorithm for finding the element with rank= $k$  in an unsorted set of  $n$  points:



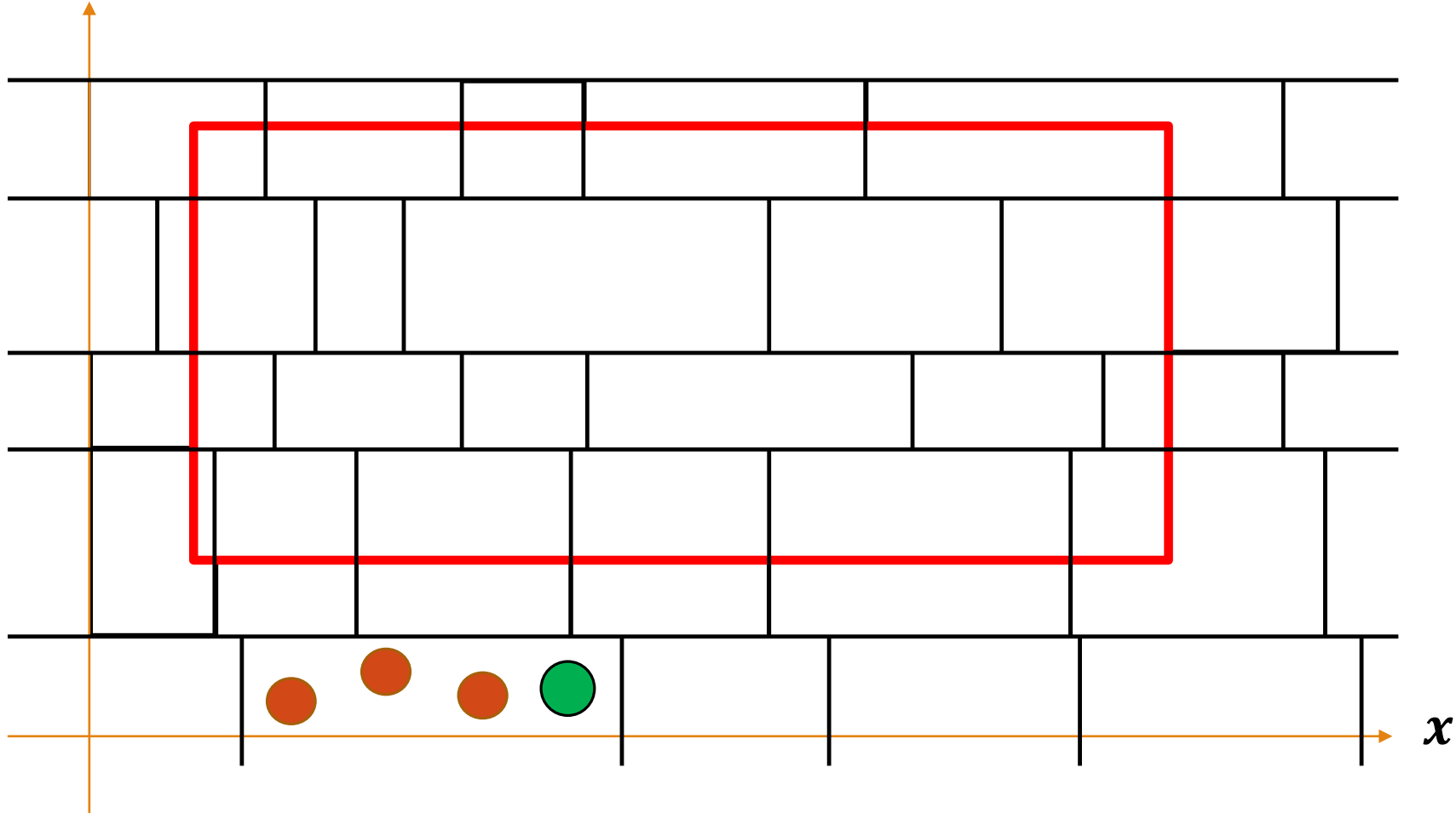
- Find the  $[\epsilon n]^{th}$  point using this algorithm in  $O(n)$  time.
- Repeat  $\left\lfloor \frac{1}{\epsilon} \right\rfloor$  times.
- Repeat for every dimension.

Total time:  $O\left(\frac{nd}{\epsilon}\right)$ .

# $\epsilon$ -net (d=2)

$y$

$k = 1$  squares



#Points in block  $< \epsilon^2 n$

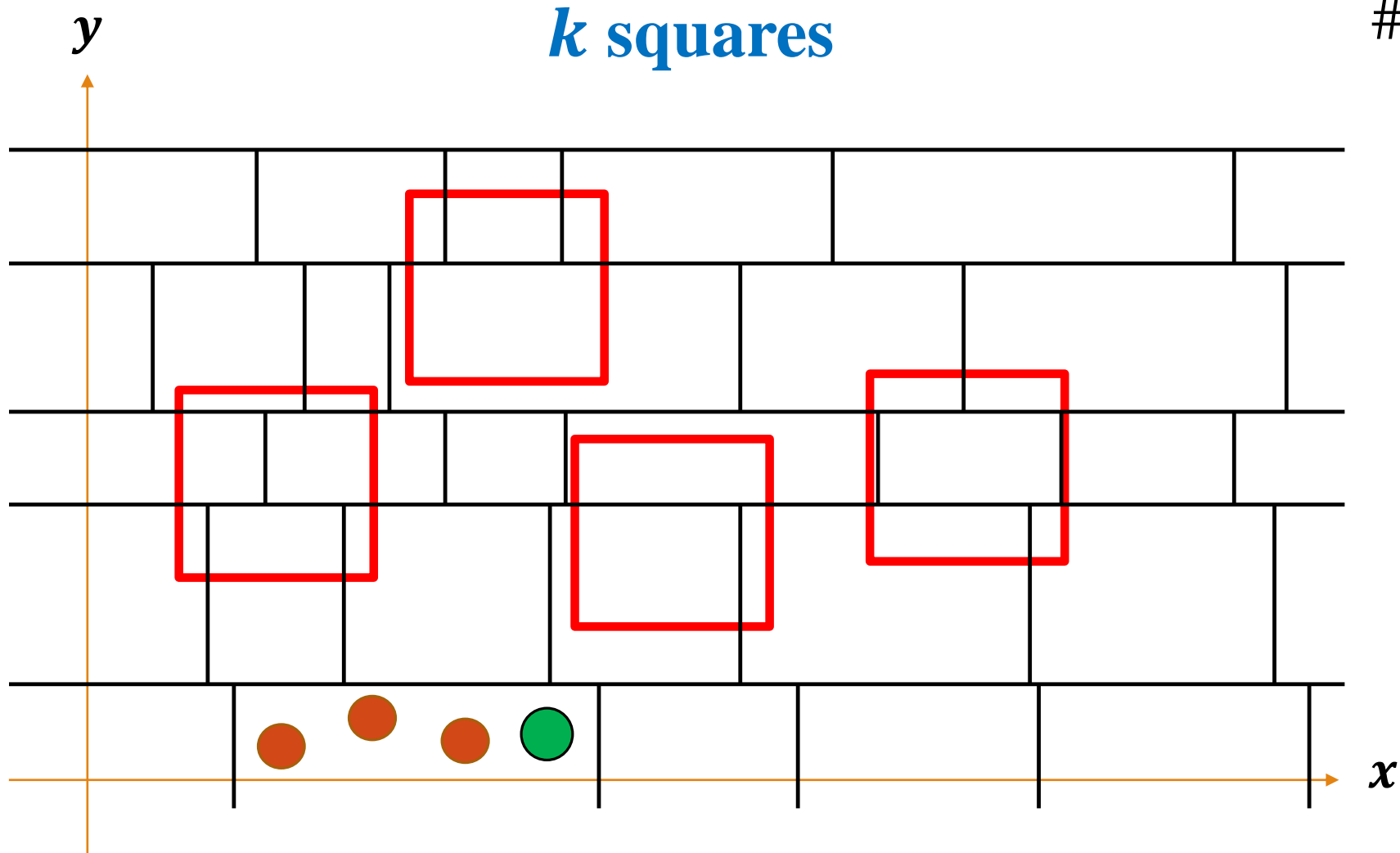
#Bad blocks  $< \frac{4}{\epsilon}$



#Bad points  $< \epsilon n$



# $\epsilon$ -net (d=2)



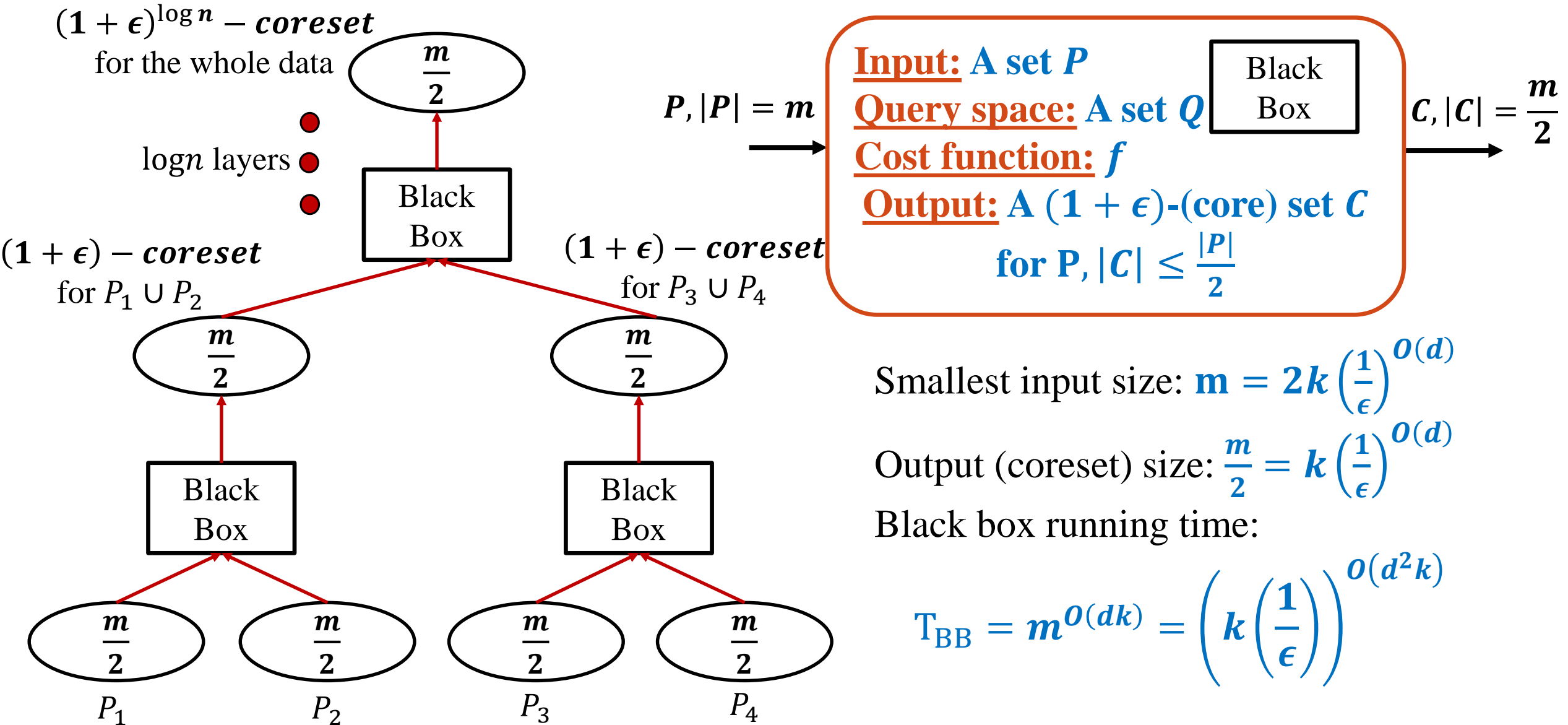
#Points in block  $< \epsilon^2 n$

#Bad blocks  $< k \frac{4}{\epsilon}$

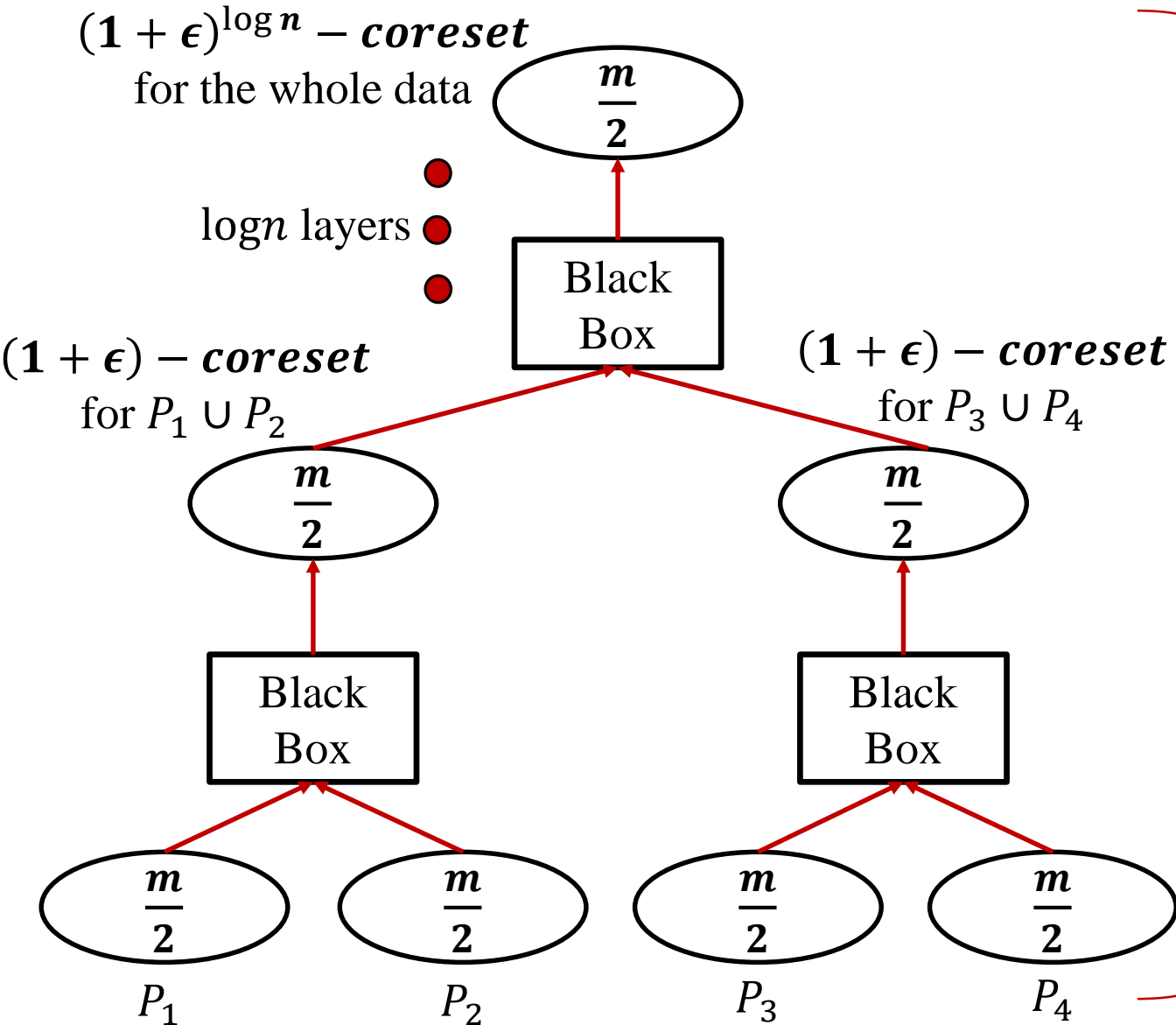


#Bad points  $< k\epsilon n$

# Coreset for $k$ -Center - Streaming



# Coreset for $k$ -Center - Streaming



Final error:

$$(1 + \epsilon)^{\log n} \cong 1 + \epsilon \cdot \log n$$

→ Define  $\epsilon' = \frac{\epsilon}{\log n}$  and run the process with  $\epsilon'$  instead of  $\epsilon$  to get a  $(1 + \epsilon)$ -approximation.

Total running time:

$$\begin{aligned} n \cdot T_{BB} &= n \cdot \left( k \left( \frac{1}{\epsilon'} \right) \right)^{O(d^2 k)} \\ &= n \cdot \left( k \left( \frac{\log n}{\epsilon} \right) \right)^{O(d^2 k)} \end{aligned}$$

# Coreset for $k$ -Center - Streaming

**Problem:** What if  $n$  (the number of input data) is unknown or infinite?

**Solution:** Doubling. (start with a small tree, then double the size).

