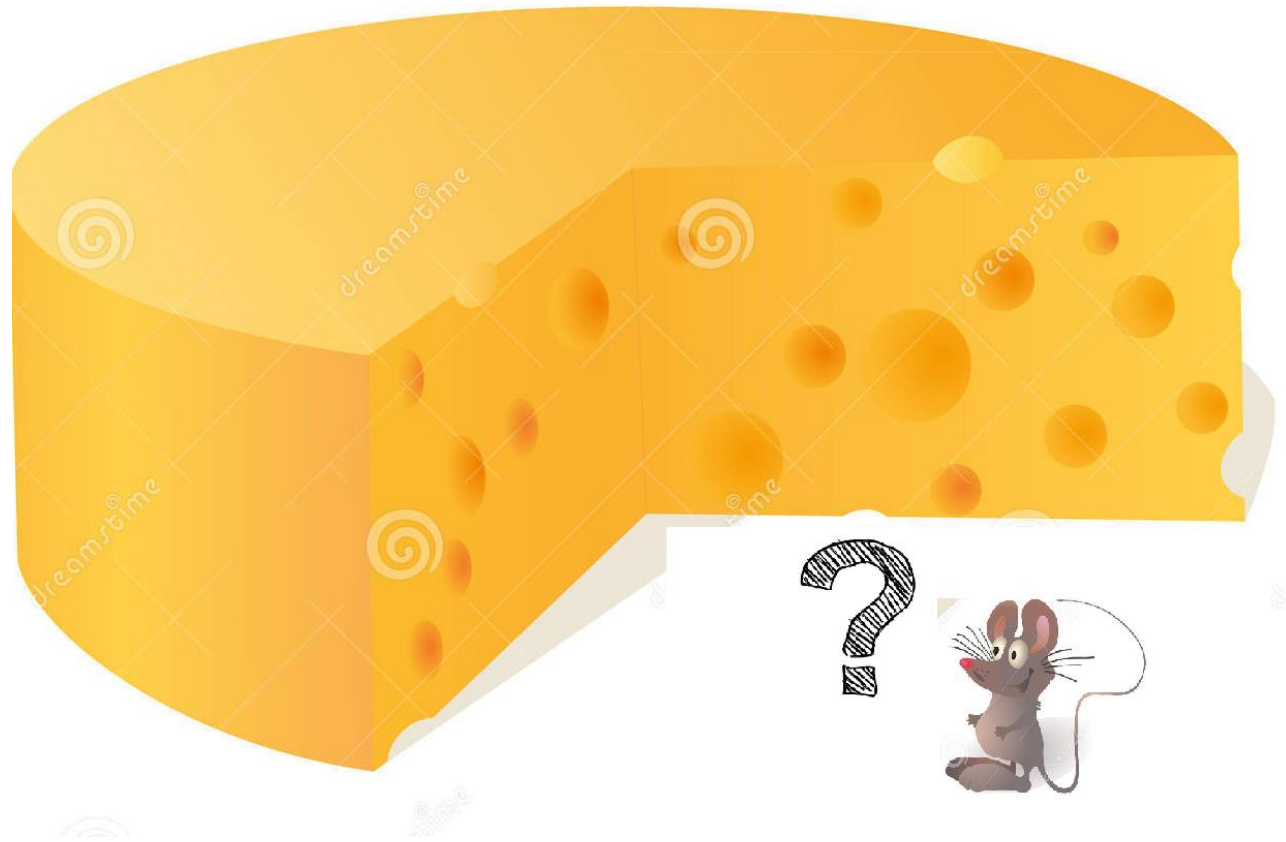


# Big Data Class



---

LECTURER: DAN FELDMAN

TEACHING ASSISTANTS:

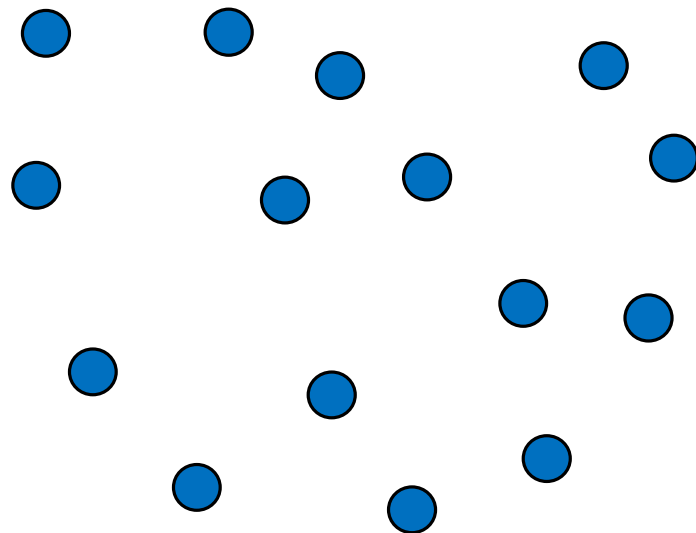
IBRAHIM JUBRAN

ALAA MAALOUF



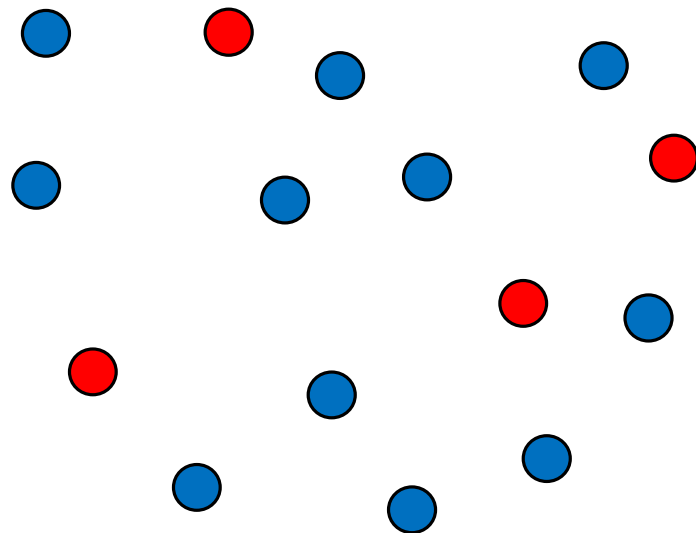
# Coreset for 1-Line in $R^2$

- Input:  $P \subseteq R^2$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^2\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$



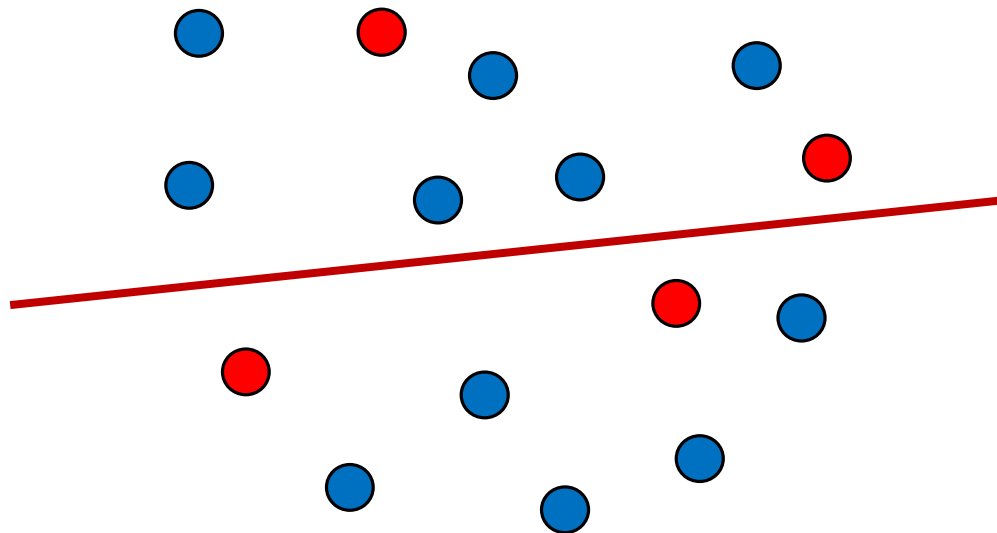
# Coreset for 1-Line in $R^2$

- Input:  $P \subseteq R^2$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^2\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s.t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$



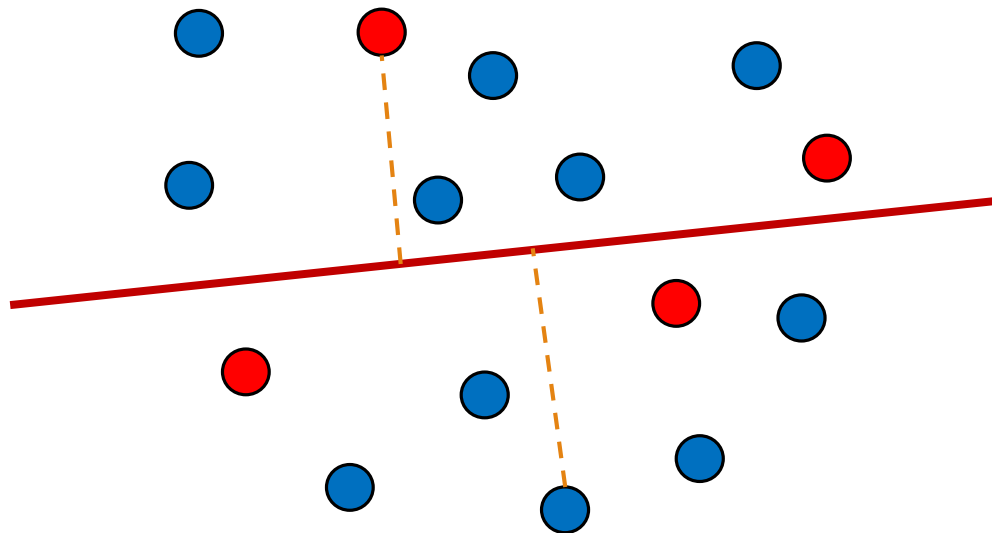
# Coreset for 1-Line in $R^2$

- Input:  $P \subseteq R^2$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^2\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$

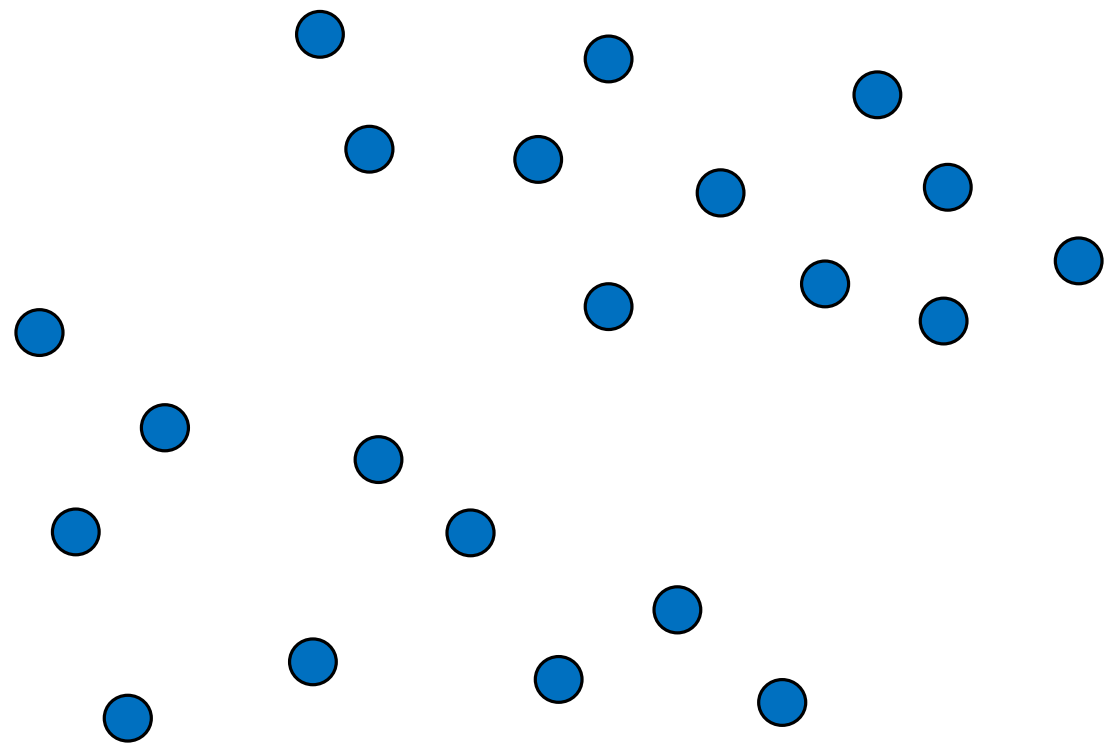


# Coreset for 1-Line in $R^2$

- Input:  $P \subseteq R^2$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^2\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s.t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$

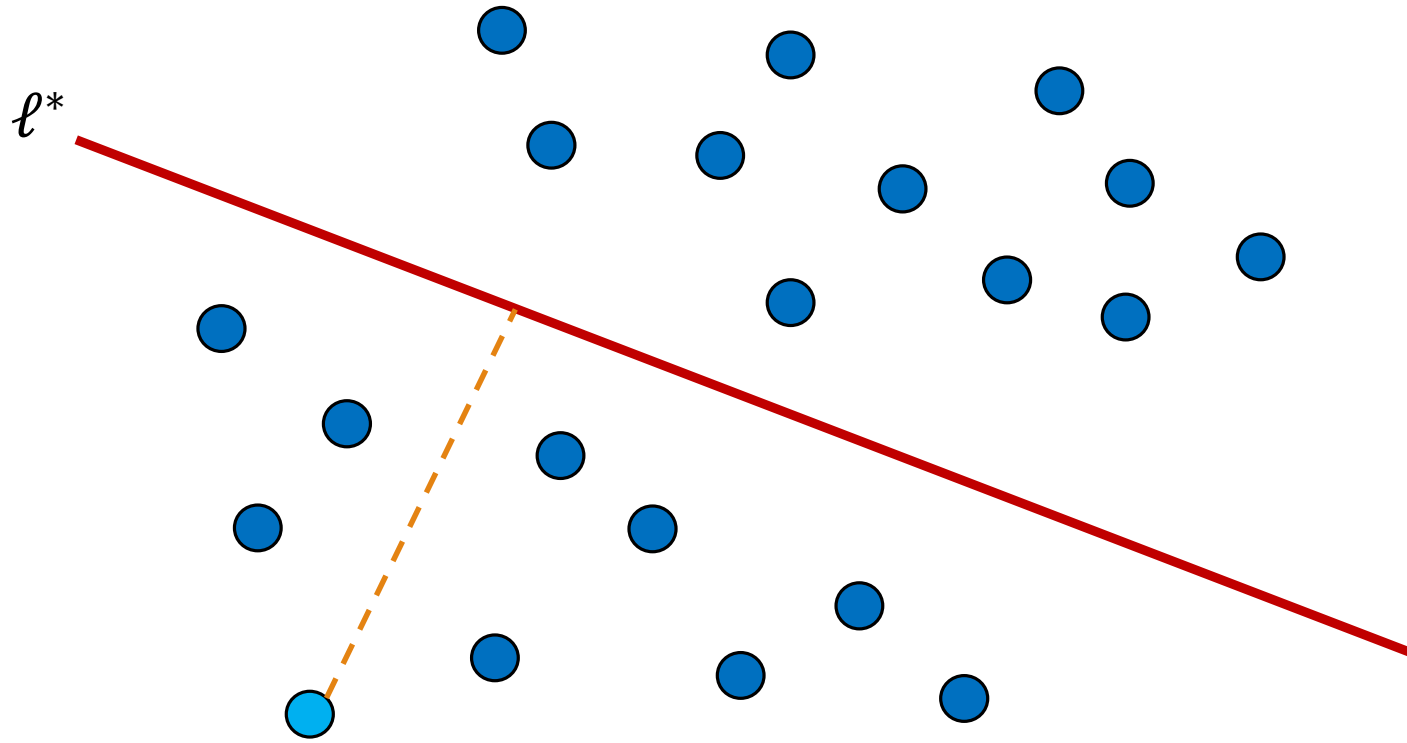


# Coreset for 1-Line in $R^2$



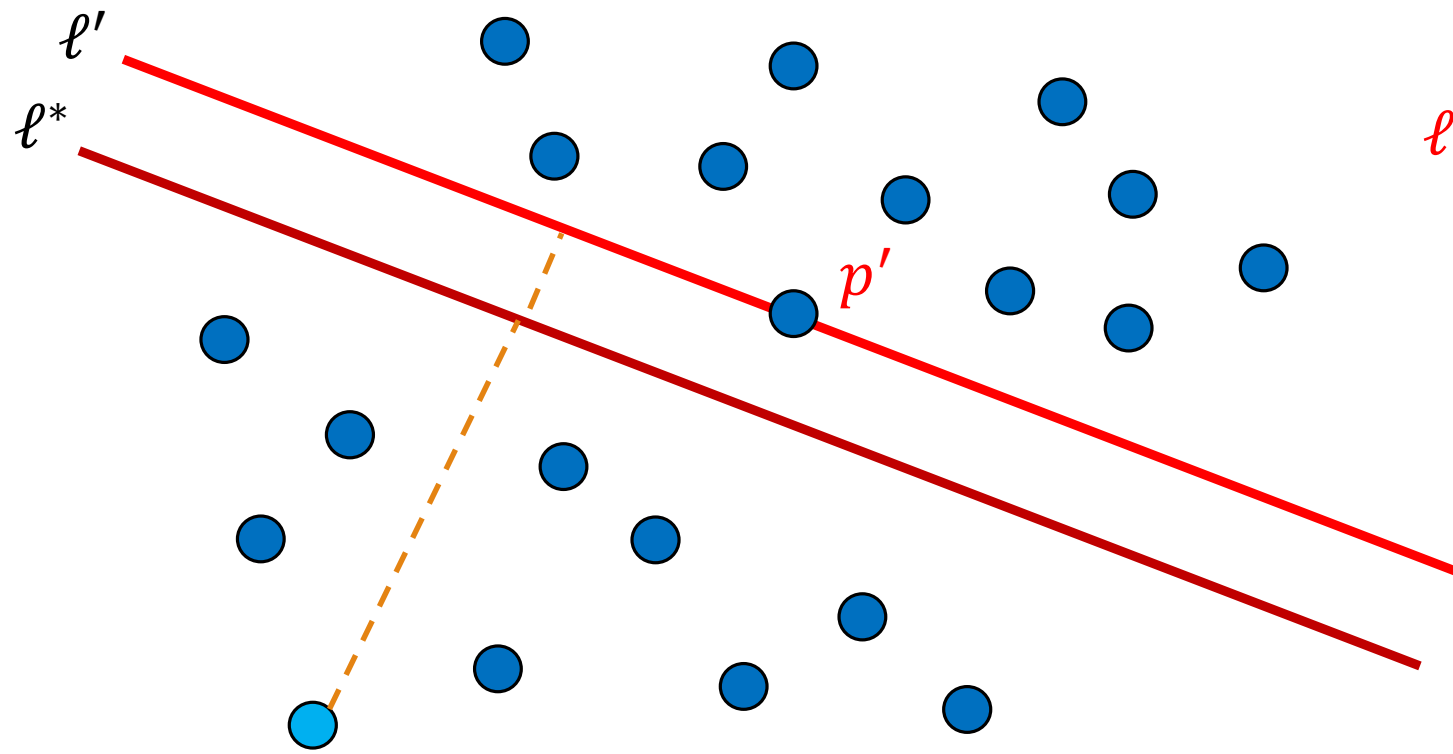
# Coreset for 1-Line in $R^2$

$\ell^*$  is the line that minimizes  
 $\max_{p \in P} \text{dist}(p, \ell)$



$$p^* = \arg \max_{p \in P} \text{dist}(p, \ell^*)$$

# Coreset for 1-Line in $R^2$



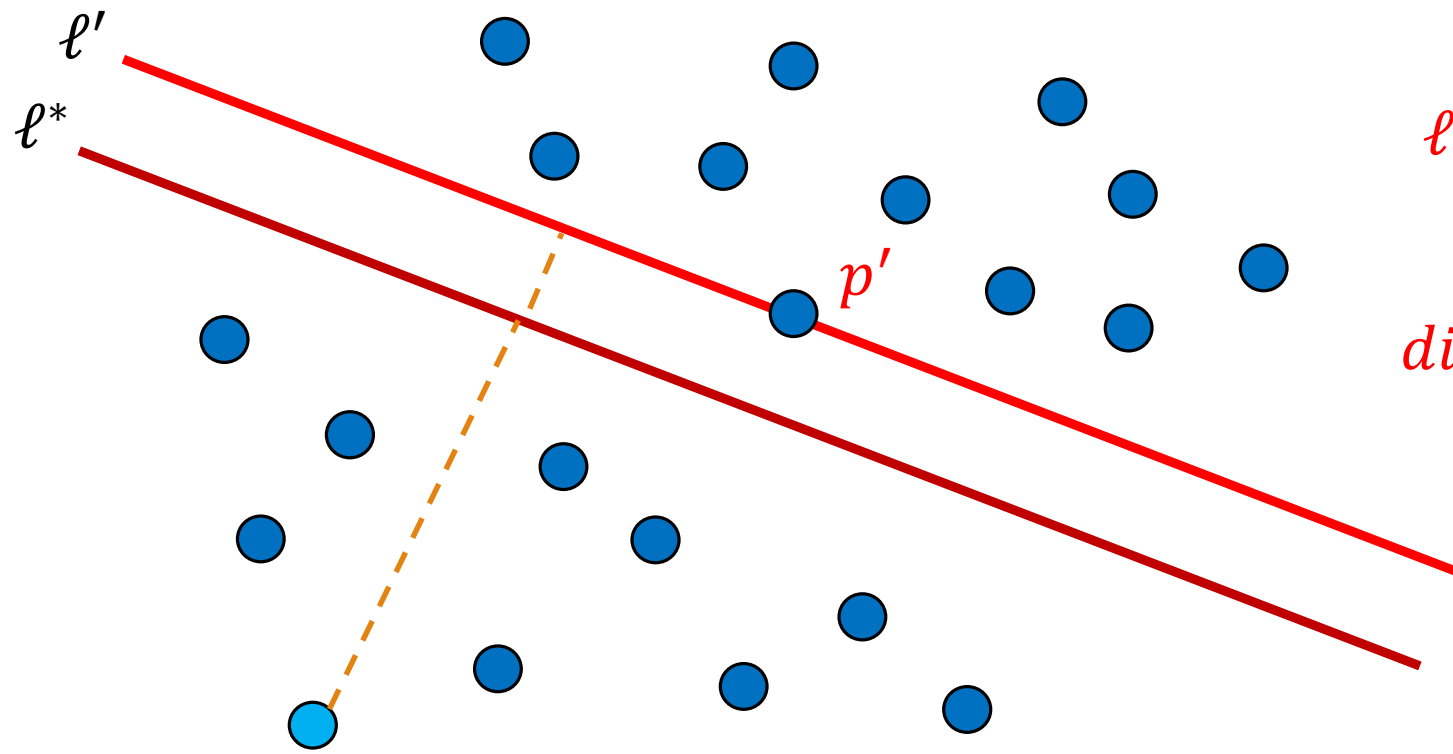
$\ell^*$  is the line that minimizes  
 $\max_{p \in P} \text{dist}(p, \ell)$

$\ell'$  is the translation of  $\ell^*$  to  
 $\ell^*$ 's closest point  $p'$

$$p^* = \arg \max_{p \in P} \text{dist}(p, \ell^*)$$



# Coreset for 1-Line in $R^2$



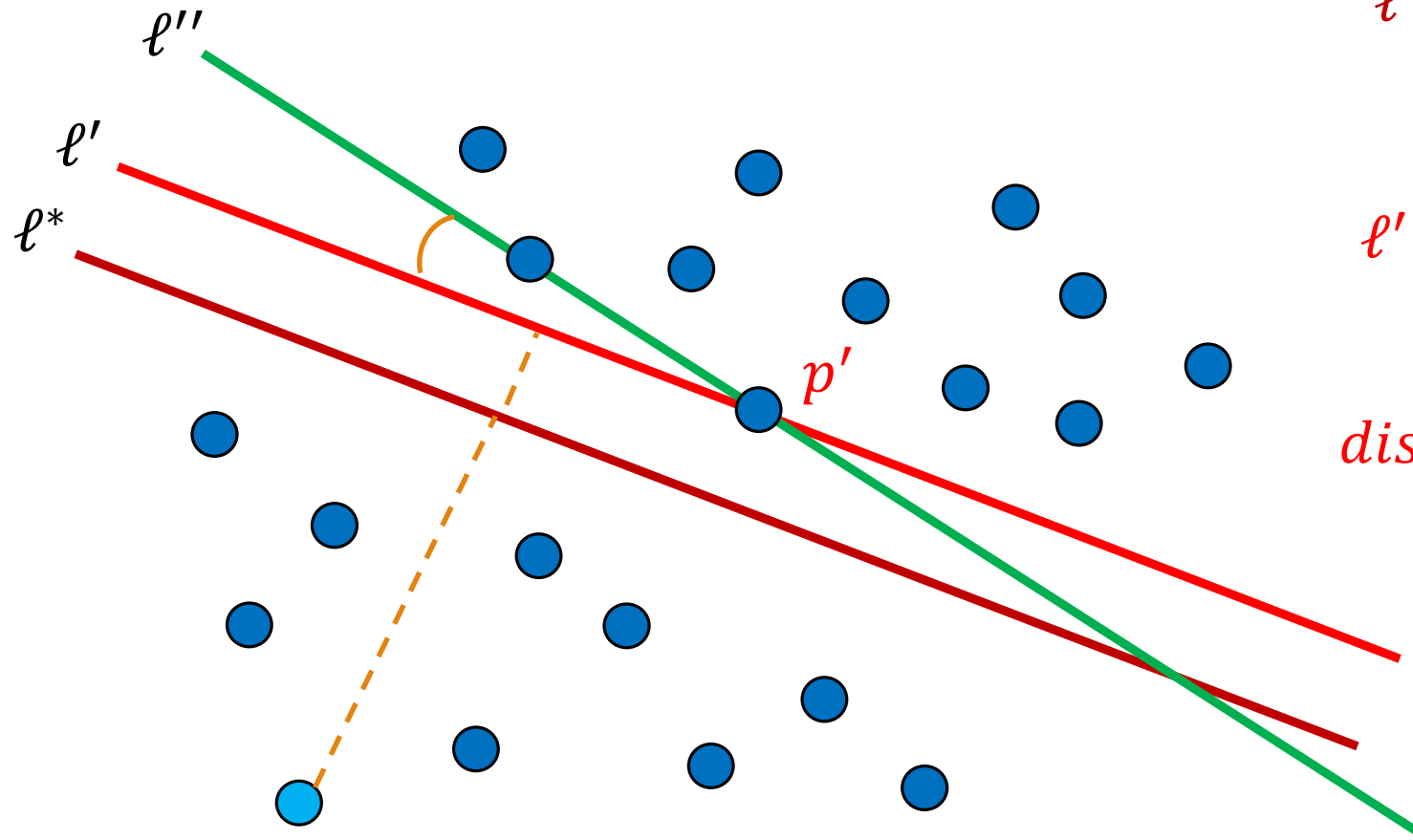
$l^*$  is the line that minimizes  
 $\max_{p \in P} \text{dist}(p, l)$

$l'$  is the translation of  $l^*$  to  
 $l^*$ 's closest point  $p'$

$$\text{dist}(p, l') \leq 2 \cdot \text{dist}(p, l^*)$$

$$p^* = \arg \max_{p \in P} \text{dist}(p, l^*)$$

# Coreset for 1-Line in $R^2$



$l^*$  is the line that minimizes  
 $\max_{p \in P} \text{dist}(p, l)$

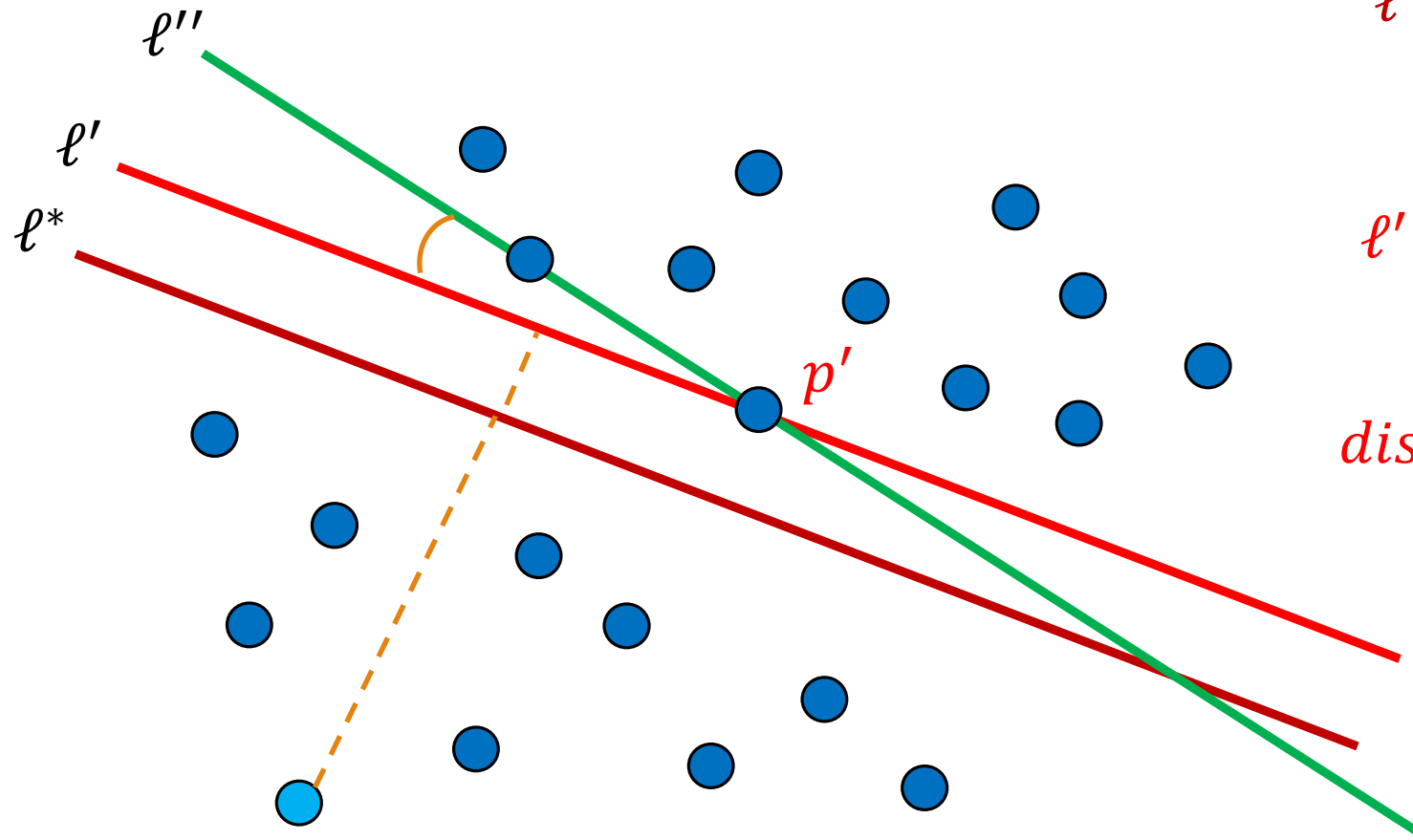
$l'$  is the translation of  $l^*$  to  
 $l^*$ 's closest point  $p'$

$$\text{dist}(p, l') \leq 2 \cdot \text{dist}(p, l^*)$$

$l''$  is the rotation of  $l'$   
around  $p'$  to  $l''$ 's closest  
point

$$p^* = \arg \max_{p \in P} \text{dist}(p, l^*)$$

# Coreset for 1-Line in $R^2$



$l^*$  is the line that minimizes  
 $\max_{p \in P} \text{dist}(p, l)$

$l'$  is the translation of  $l^*$  to  
 $l^*$ 's closest point  $p'$

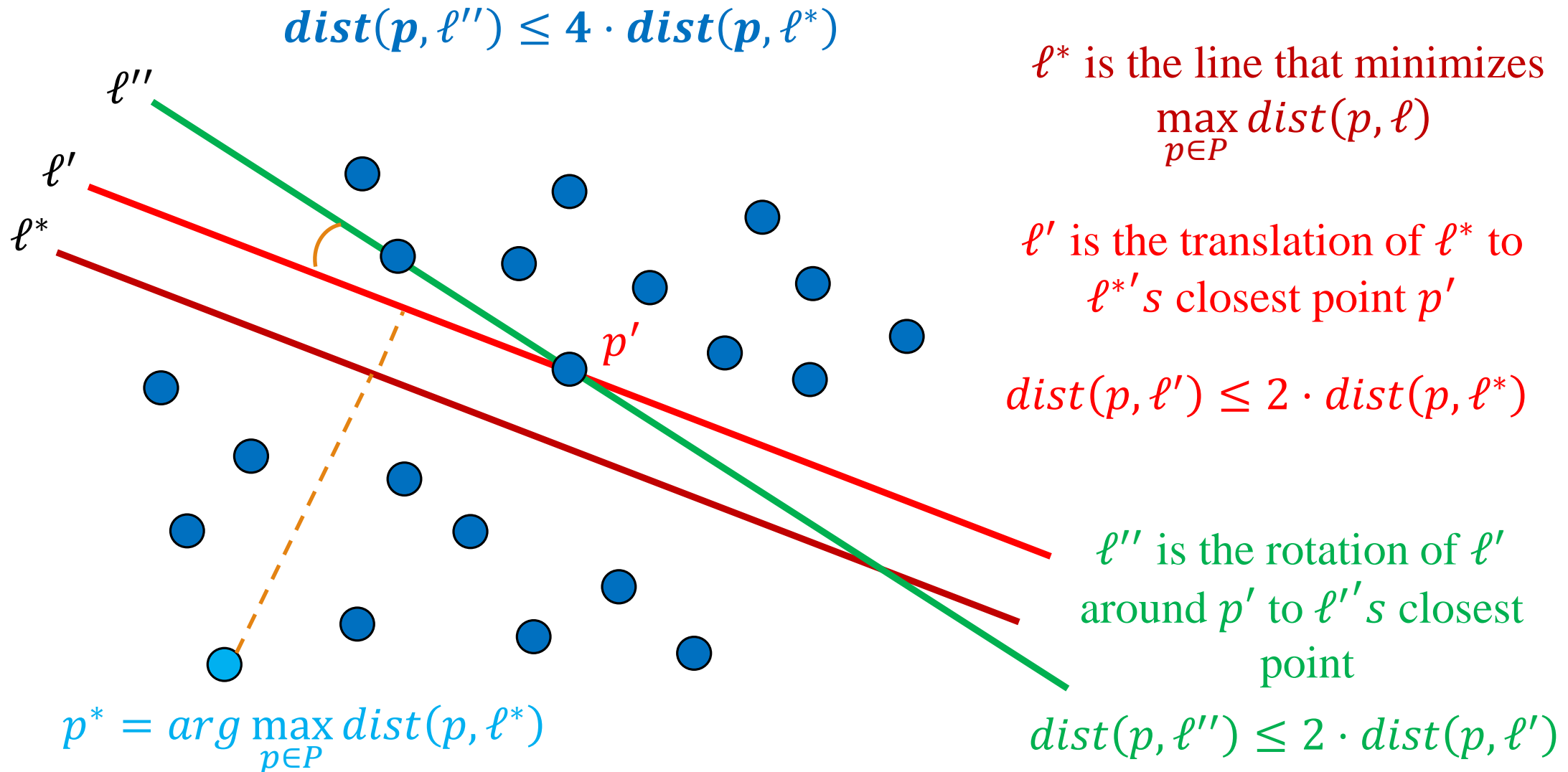
$$\text{dist}(p, l') \leq 2 \cdot \text{dist}(p, l^*)$$

$l''$  is the rotation of  $l'$   
around  $p'$  to  $l''$ 's closest  
point

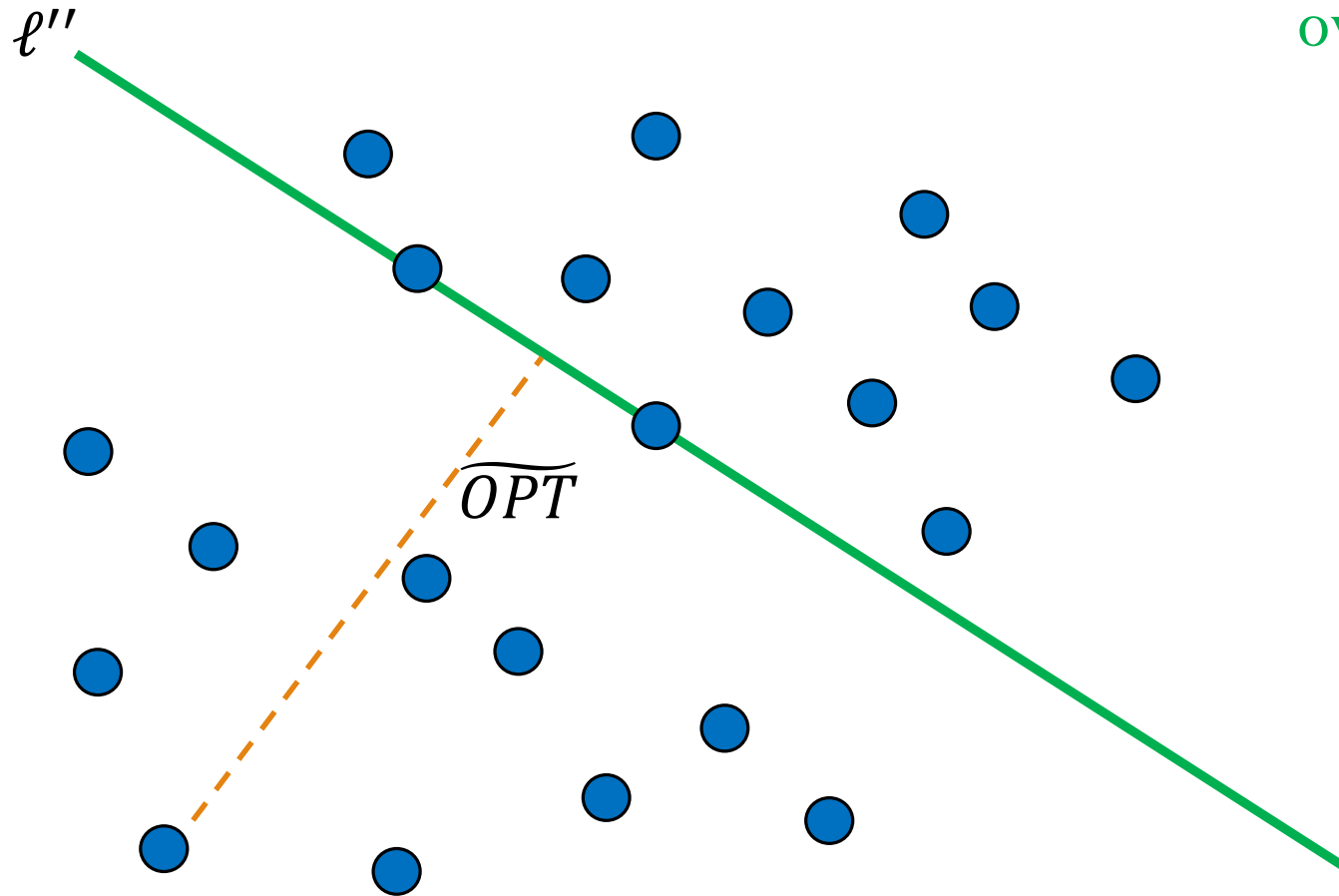
$$\text{dist}(p, l'') \leq 2 \cdot \text{dist}(p, l')$$

$$p^* = \arg \max_{p \in P} \text{dist}(p, l^*)$$

# Coreset for 1-Line in $R^2$

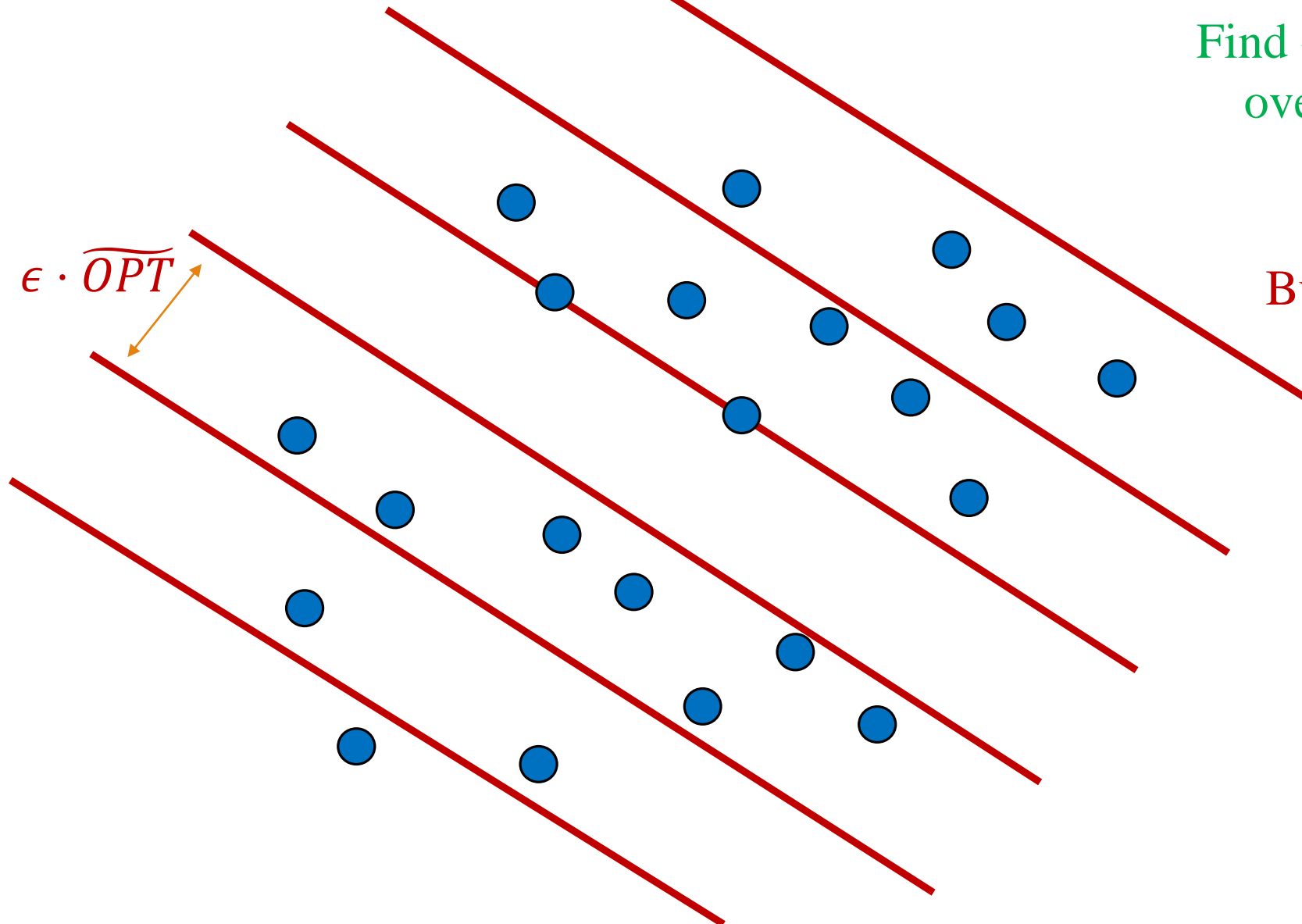


# Coreset for 1-Line in $R^2$



Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$

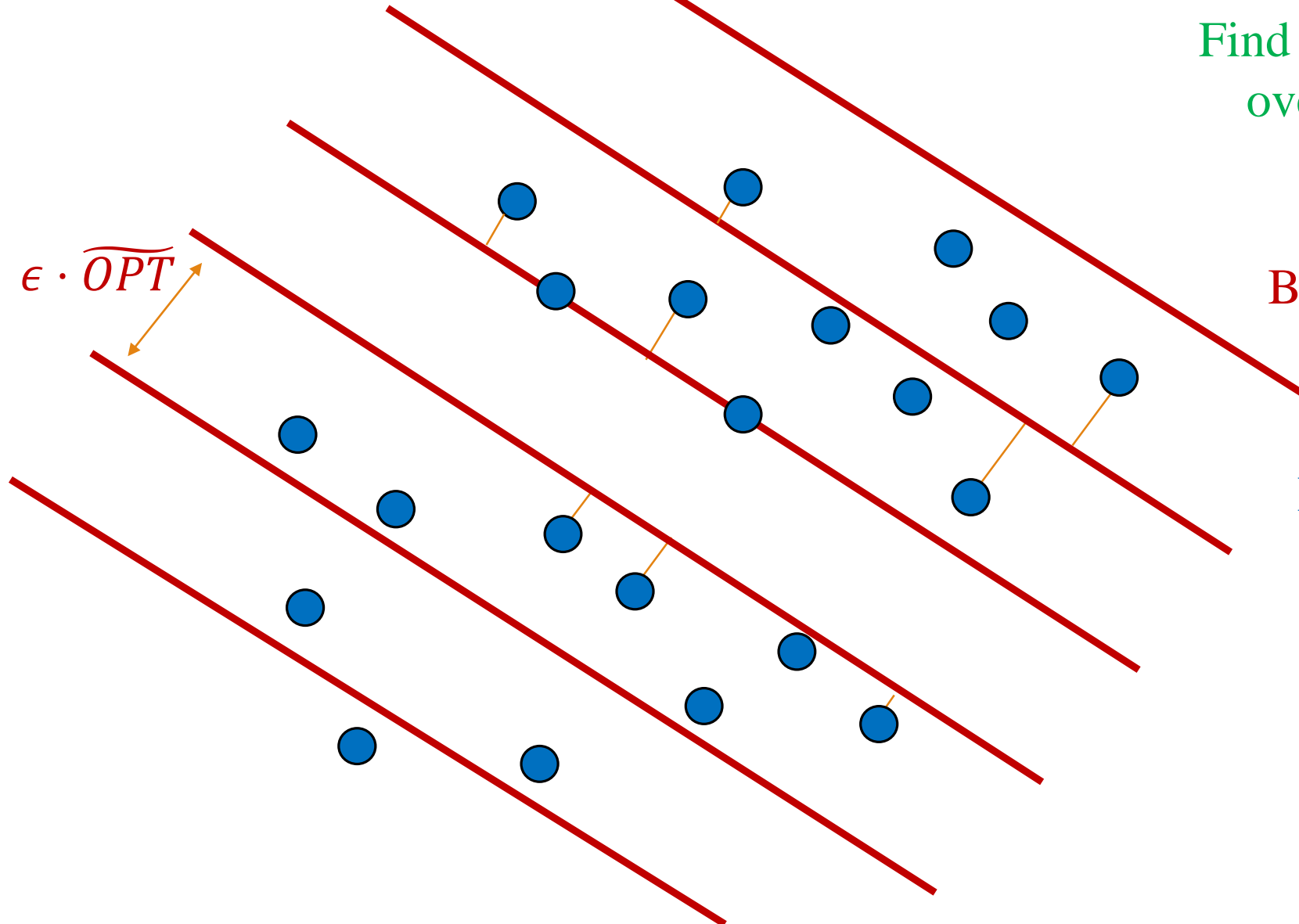
# Coreset for 1-Line in $R^2$



Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$

Build a grid of lines with  
 $\epsilon \cdot \overline{OPT}$  distance

# Coreset for 1-Line in $R^2$

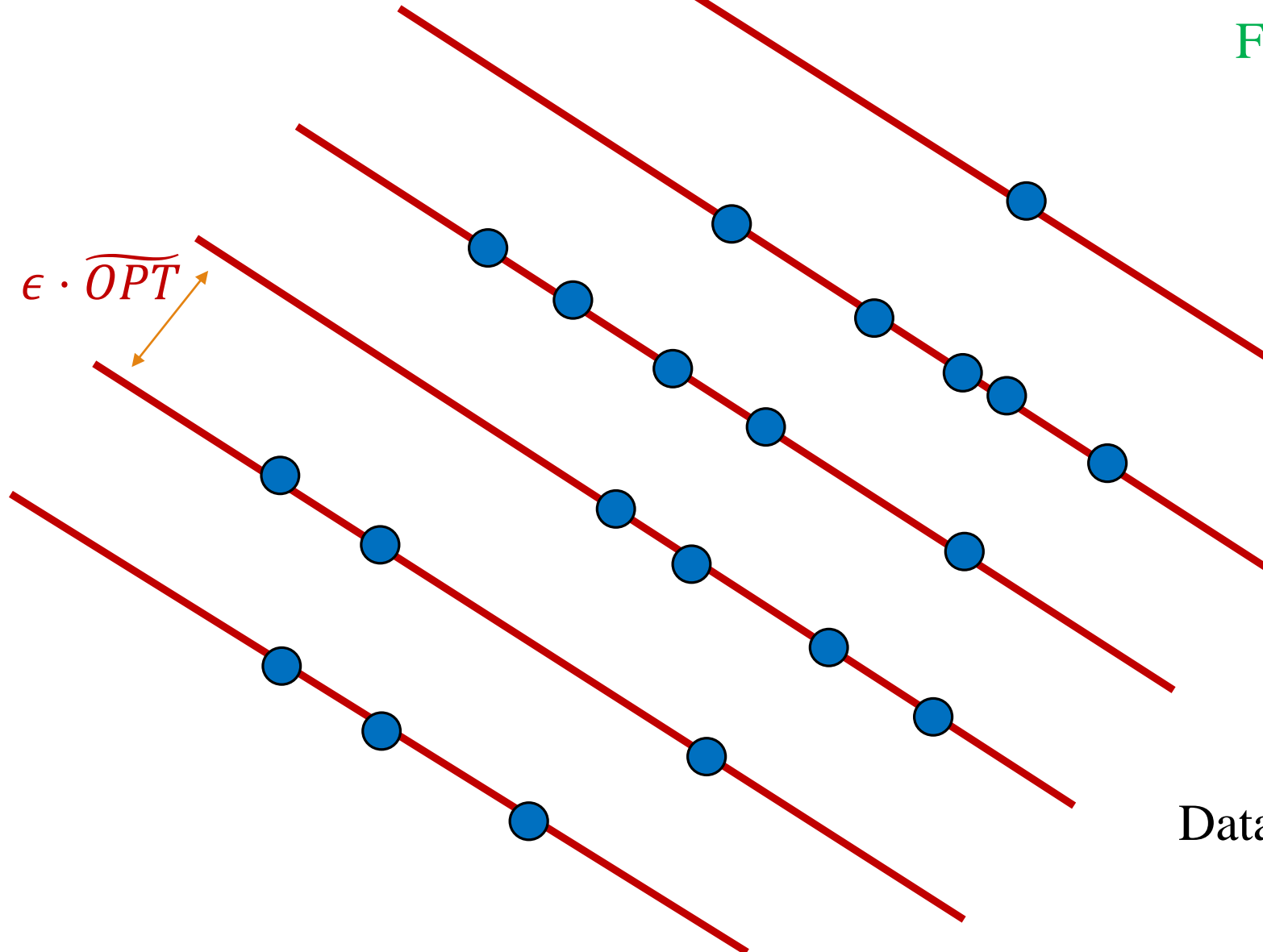


Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$

Build a grid of lines with  
 $\epsilon \cdot \overline{OPT}$  distance

Project each point onto  
it's closest line

# Coreset for 1-Line in $R^2$



Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$

Build a grid of lines with  
 $\epsilon \cdot \overline{OPT}$  distance

Project each point onto  
it's closest line

Data dimension is now reduced.



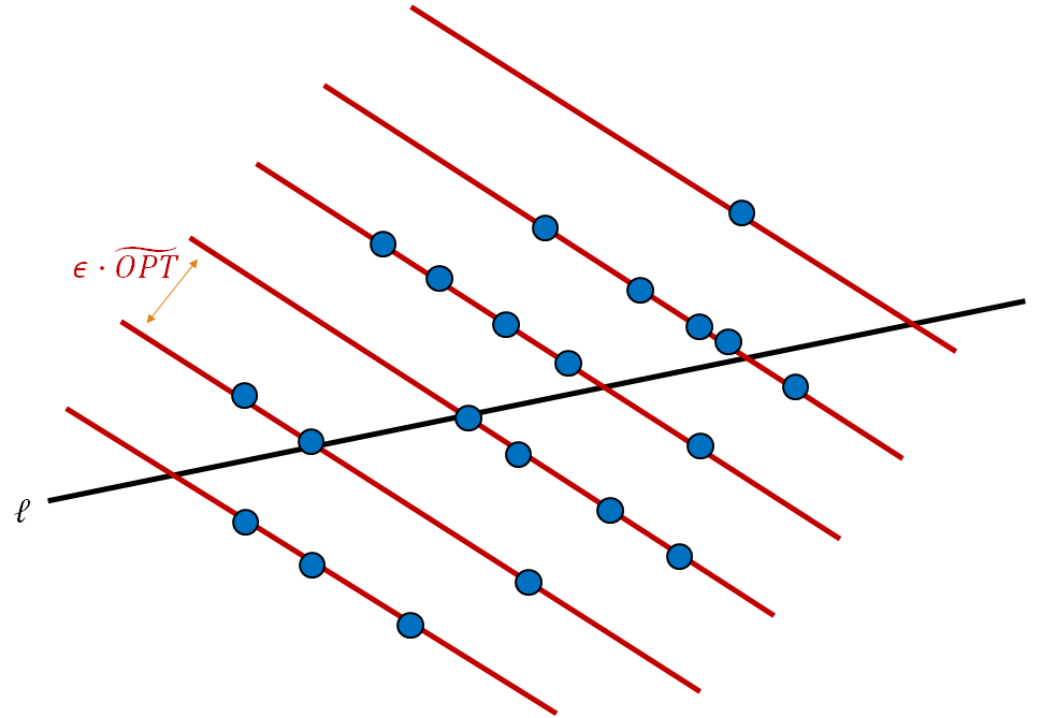
# Coreset for 1-Line in $R^2$

Claim: The projected  $n$  points  $P'$  are a “coreset” (not part of the input data) for any line query:

$$\max_{p \in P} \text{dist}(p, \ell) - \max_{p \in P'} \text{dist}(p, \ell) \leq \epsilon \cdot \widetilde{OPT}$$

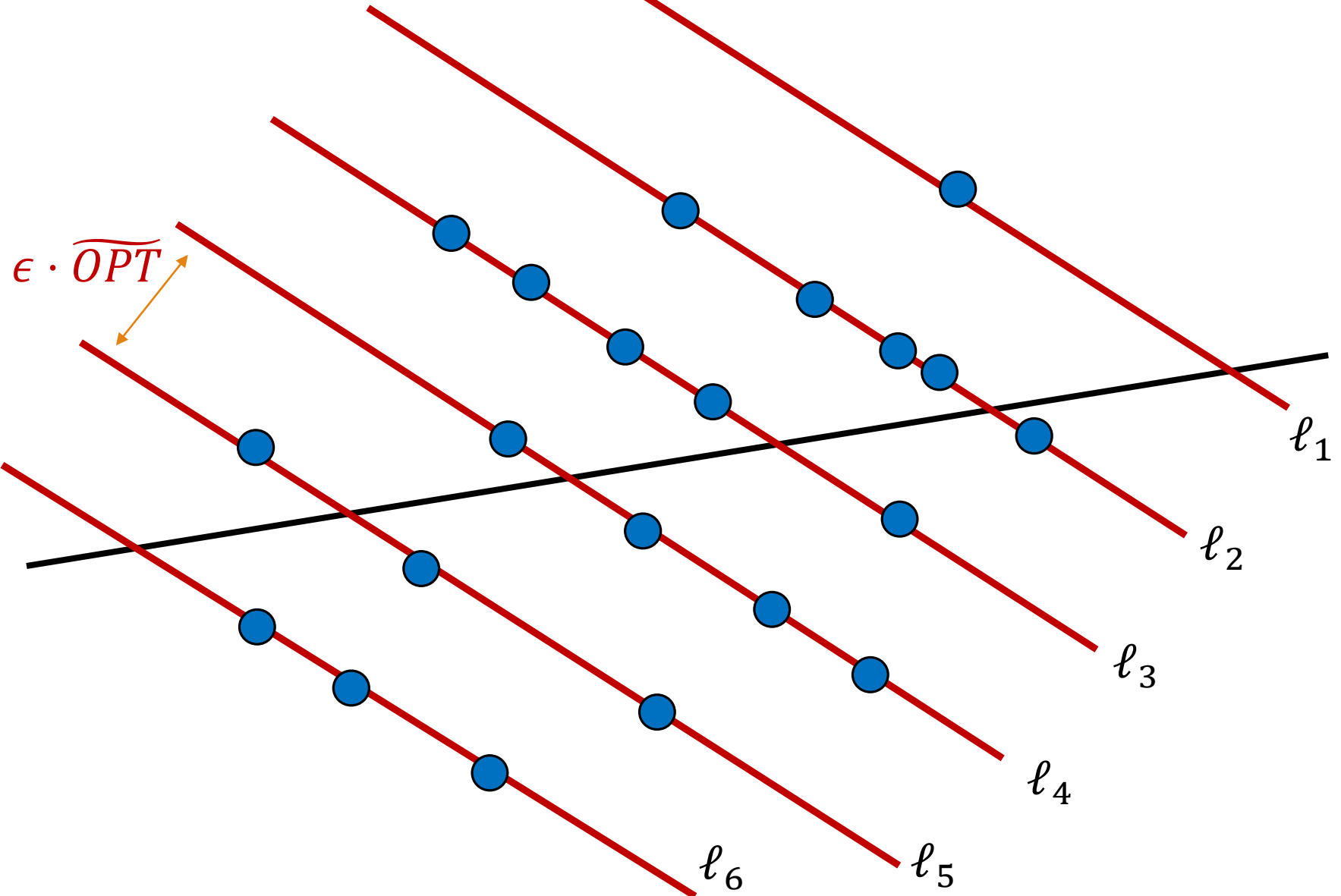
$$\leq 4\epsilon \cdot OPT$$

$$\leq 4\epsilon \cdot \max_{p \in P} \text{dist}(p, \ell)$$



→ Run with  $\epsilon' = \frac{\epsilon}{4}$

# Coreset for 1-Line in $R^2$



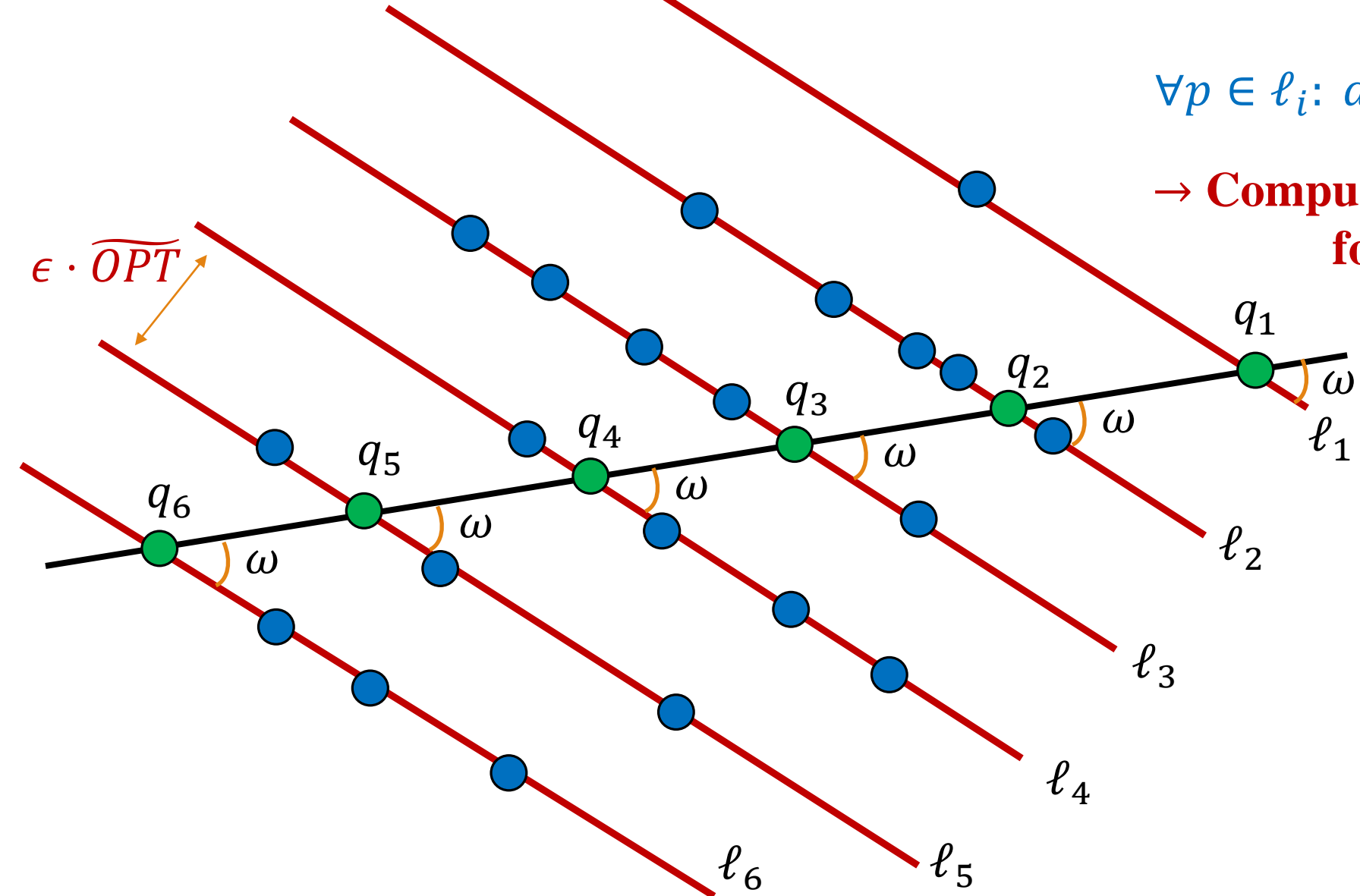
# Coreset for 1-Line in $R^2$

Has no effect since it is the same weight for all points

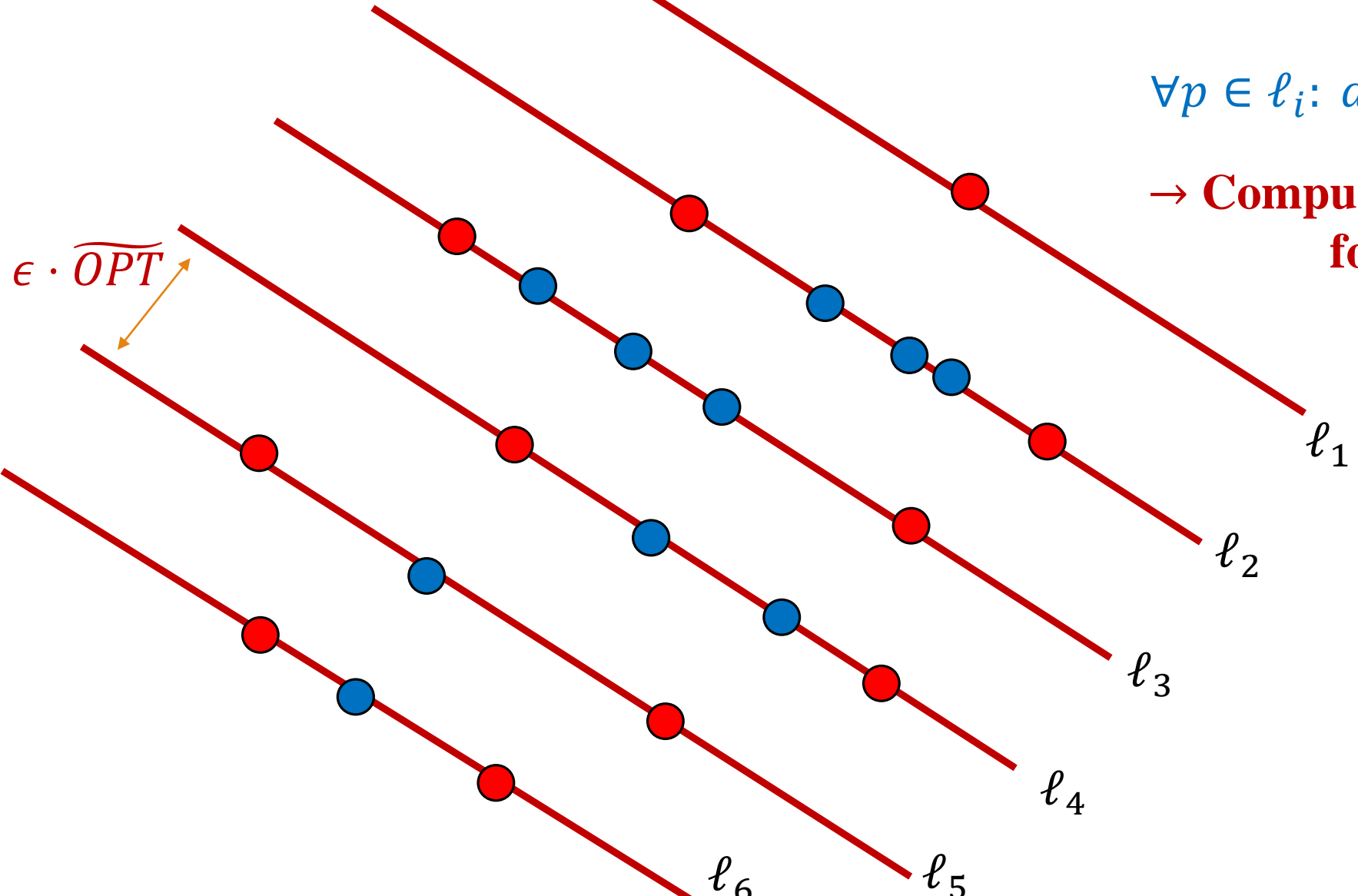
$$\forall p \in \ell_i: \text{dist}(p, \ell) = \omega \cdot \text{dist}(p, q_i)$$

→ Compute a 1-Center coreset  $C_i$  for each line  $\ell_i$ !

$\epsilon \cdot \overline{OPT}$



# Coreset for 1-Line in $R^2$



Has no effect since it is the same weight for all points

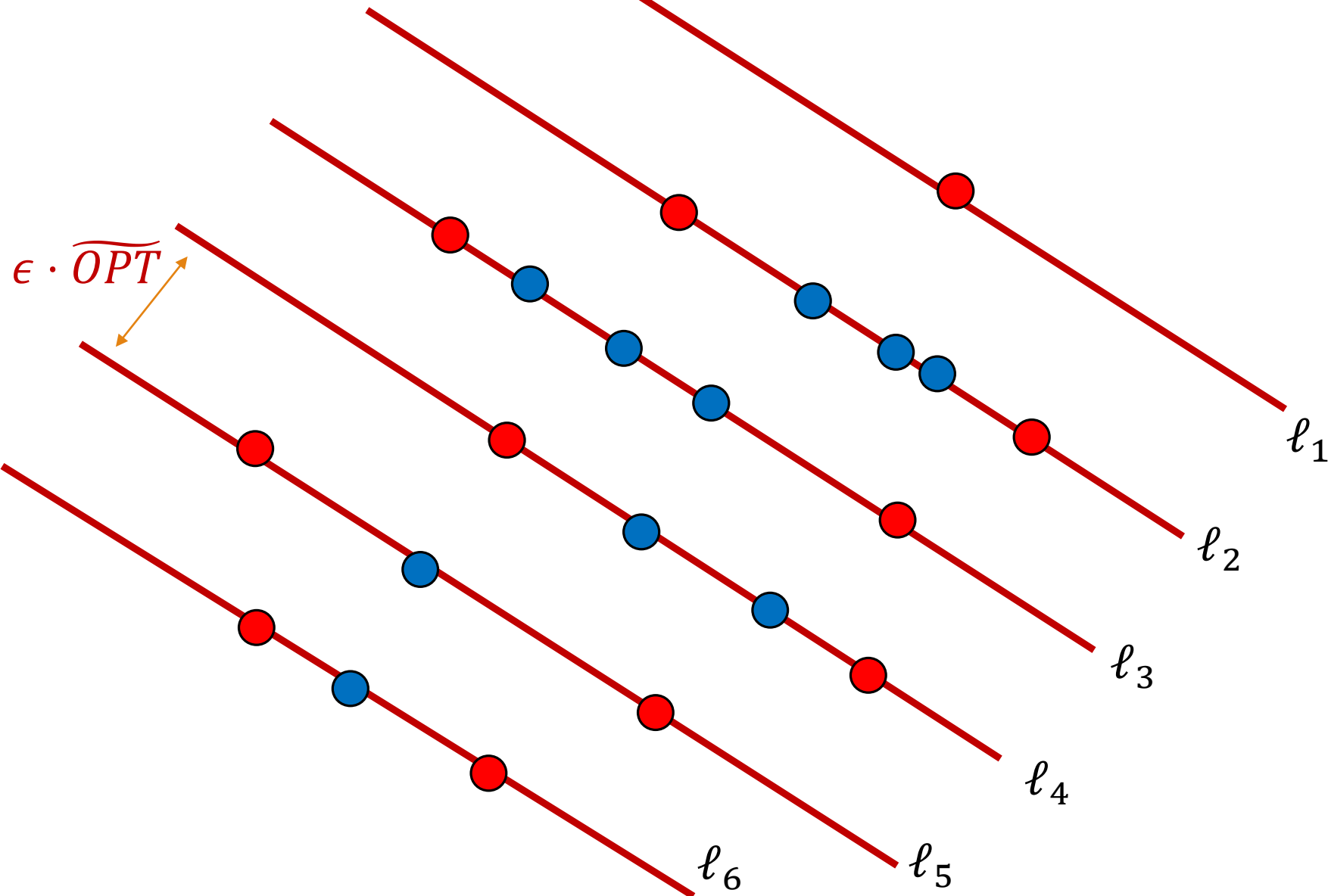
$$\forall p \in \ell_i: \text{dist}(p, \ell) = \omega \cdot \text{dist}(p, q_i)$$

→ Compute a 1-Center coreset  $C_i$  for each line  $\ell_i$ !

$$C = \bigcup C_i$$

since a union of two coresets is a coreset.

# Coreset for 1-Line in $R^2$



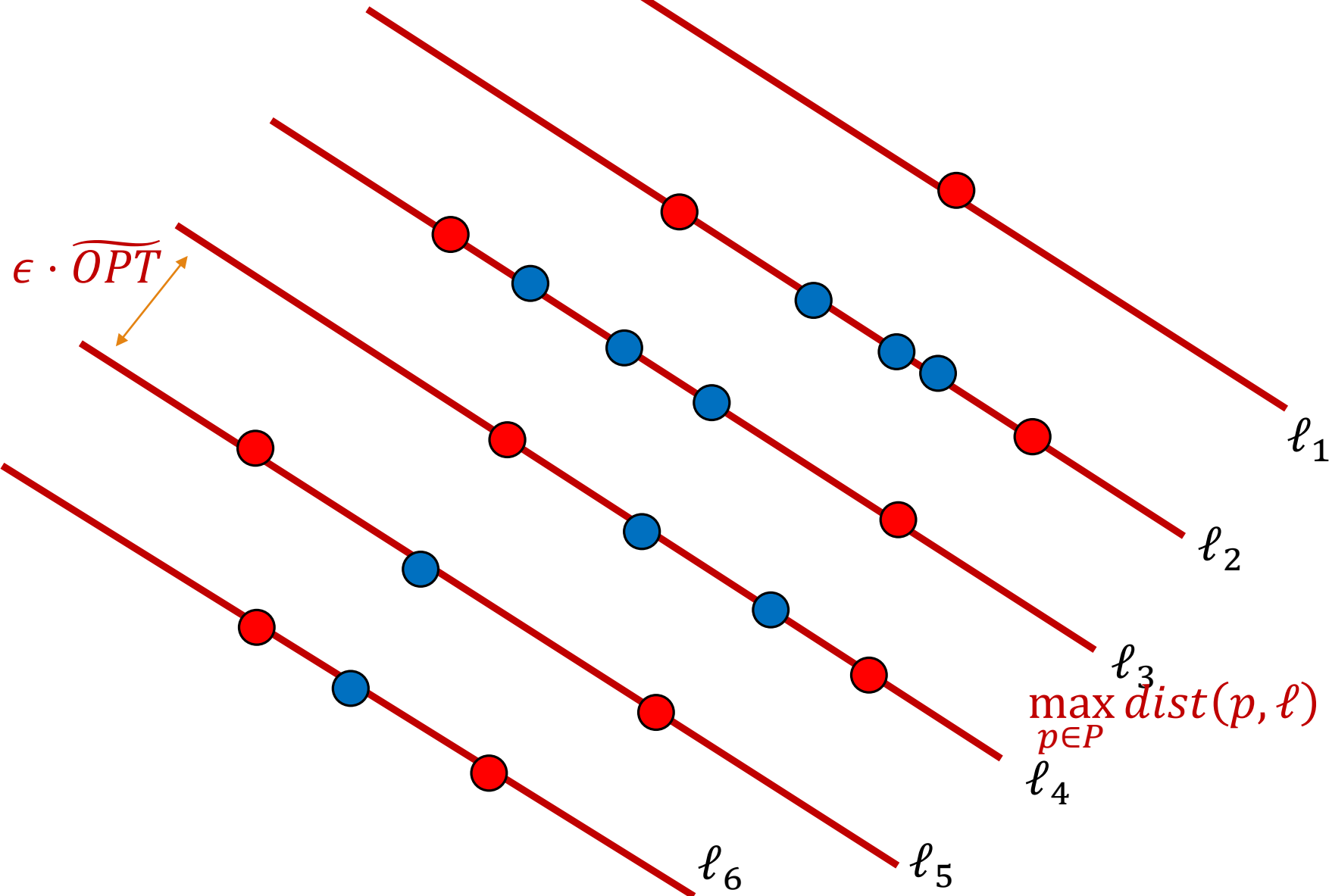
Problem:

The coreset is not part of the input data.

Solution:

Pick the closest points in the input data to the points of  $C$ .

# Coreset for 1-Line in $R^2$



Problem:

The coreset is not part of the input data.

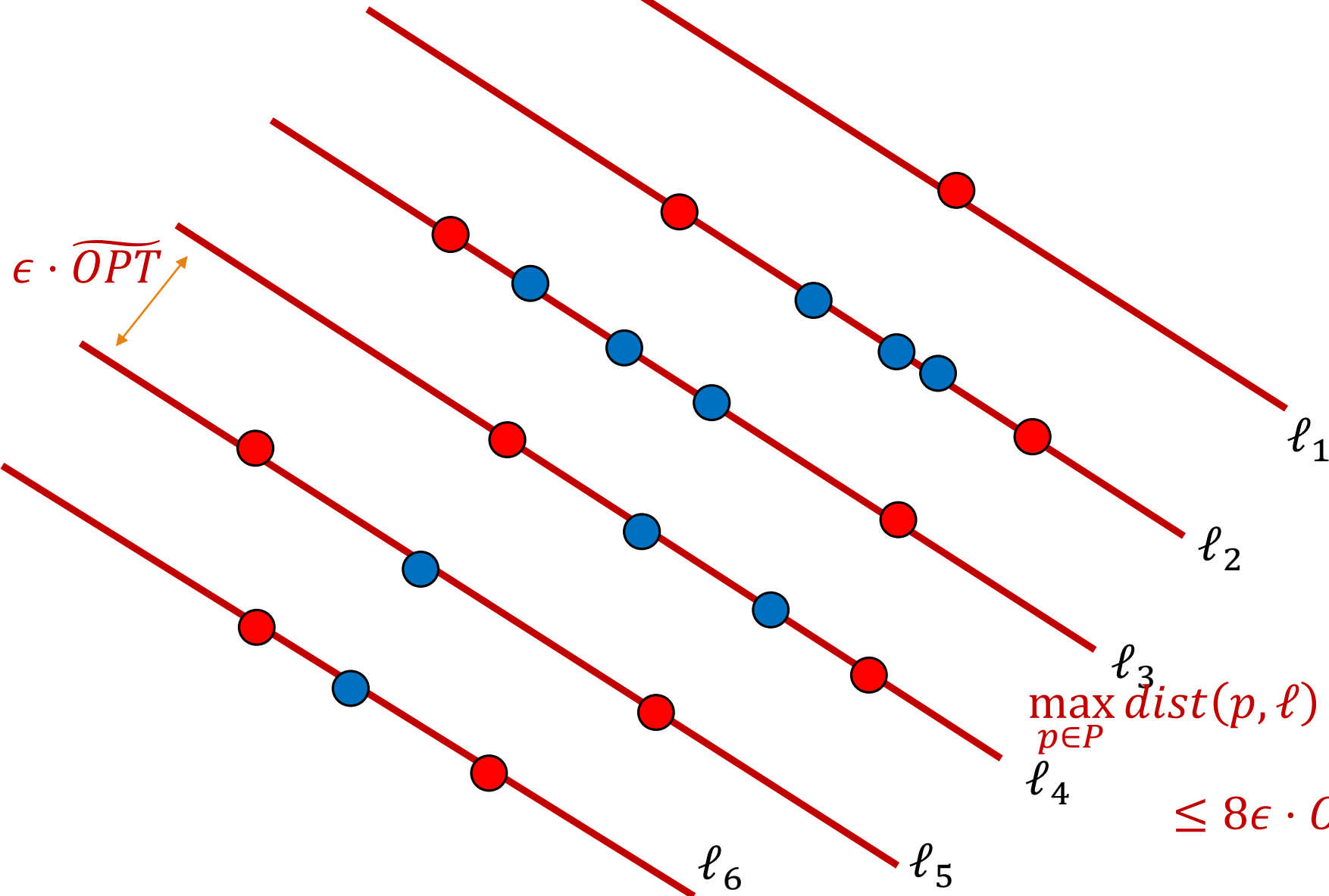
Solution:

Pick the closest points in the input data to the points of  $C$ .

→ **This adds another error of  $\epsilon \cdot \widetilde{OPT}$**

$$\begin{aligned} \max_{p \in P} \text{dist}(p, \ell) &\leq \max_{p \in P'} \text{dist}(p, \ell) + 2\epsilon \cdot \widetilde{OPT} \\ &\leq (1 + 8\epsilon) \cdot \max_{p \in P'} \text{dist}(p, \ell) \end{aligned}$$

# Coreset for 1-Line in $R^2$



Problem:

The coreset is not part of the input data.

Solution:

Pick the closest points in the input data to the points of  $C$ .

→ **This adds another error of  $\epsilon \cdot \widetilde{OPT}$**

$$\begin{aligned} \max_{p \in P} \text{dist}(p, \ell) - \max_{p \in P'} \text{dist}(p, \ell) &\leq 2\epsilon \cdot \widetilde{OPT} \\ &\leq 8\epsilon \cdot OPT \leq 8\epsilon \cdot \max_{p \in P} \text{dist}(p, \ell) \end{aligned}$$

# Coreset for 1-Line in $R^2$

Total time:

$O(n^2)$ .

Coreset size:

$$|C| \leq 2 \cdot \#lines = 2 \cdot \frac{2}{\epsilon} = \frac{4}{\epsilon}.$$



# Coreset for 1-Line in $R^2$

Total time:

$O(n^2)$ .

Coreset size:

$$|C| \leq 2 \cdot \#lines = 2 \cdot \frac{2}{\epsilon} = \frac{4}{\epsilon}.$$

Improvement:

Run the above algorithm using the streaming tree.

Run on batches of size  $2 \cdot |C| = \frac{8}{\epsilon}$ .

Total time:

$$O(n \cdot TimeForBatch) = O\left(n \cdot \left(\frac{8}{\epsilon}\right)^2\right).$$

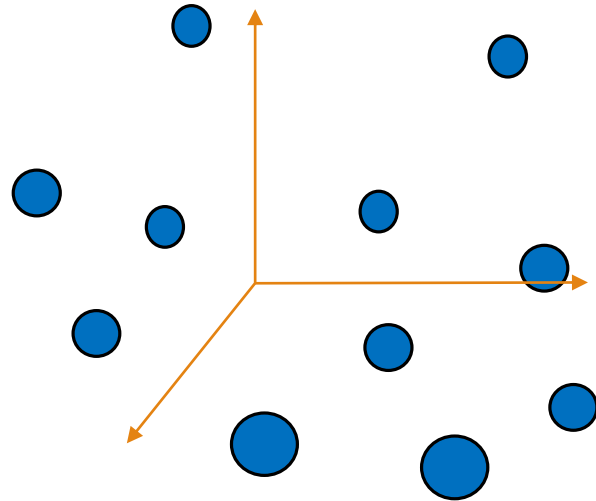
Error for streaming tree:

The error increases to  $(1 + \epsilon)^{\log n} \sim (1 + \epsilon \log n)$

→ Run with  $\epsilon' = \frac{\epsilon}{\log n}$ .

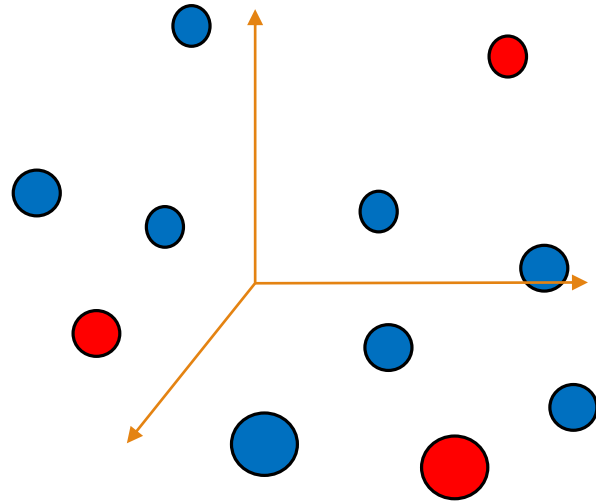
# Coreset for 1-Plane in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\pi \mid \pi \text{ is a plane in } R^3\}$
- Cost function:  $dist(p, \pi) = \min_{x \in \pi} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \pi \in Q: \max_{p \in P} dist(p, \pi) - \max_{c \in C} dist(c, \pi) \leq \epsilon \cdot \max_{p \in P} dist(p, \pi)$



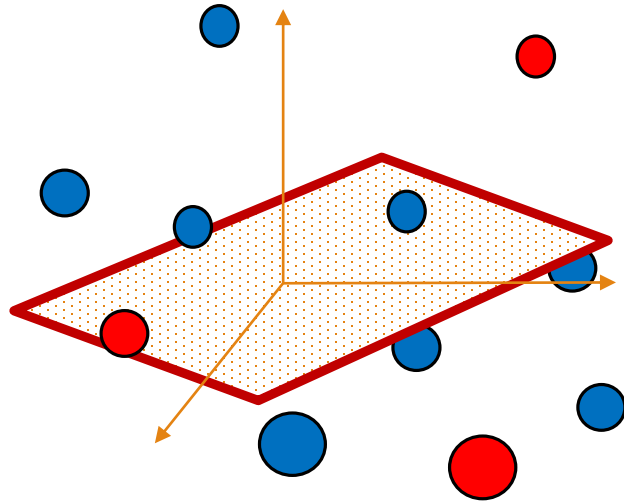
# Coreset for 1-Plane in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\pi \mid \pi \text{ is a plane in } R^3\}$
- Cost function:  $dist(p, \pi) = \min_{x \in \pi} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \pi \in Q: \max_{p \in P} dist(p, \pi) - \max_{c \in C} dist(c, \pi) \leq \epsilon \cdot \max_{p \in P} dist(p, \pi)$



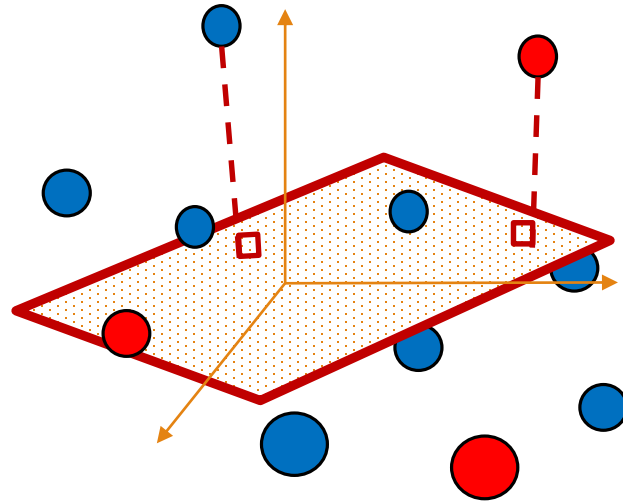
# Coreset for 1-Plane in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\pi \mid \pi \text{ is a plane in } R^3\}$
- Cost function:  $dist(p, \pi) = \min_{x \in \pi} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \pi \in Q: \max_{p \in P} dist(p, \pi) - \max_{c \in C} dist(c, \pi) \leq \epsilon \cdot \max_{p \in P} dist(p, \pi)$

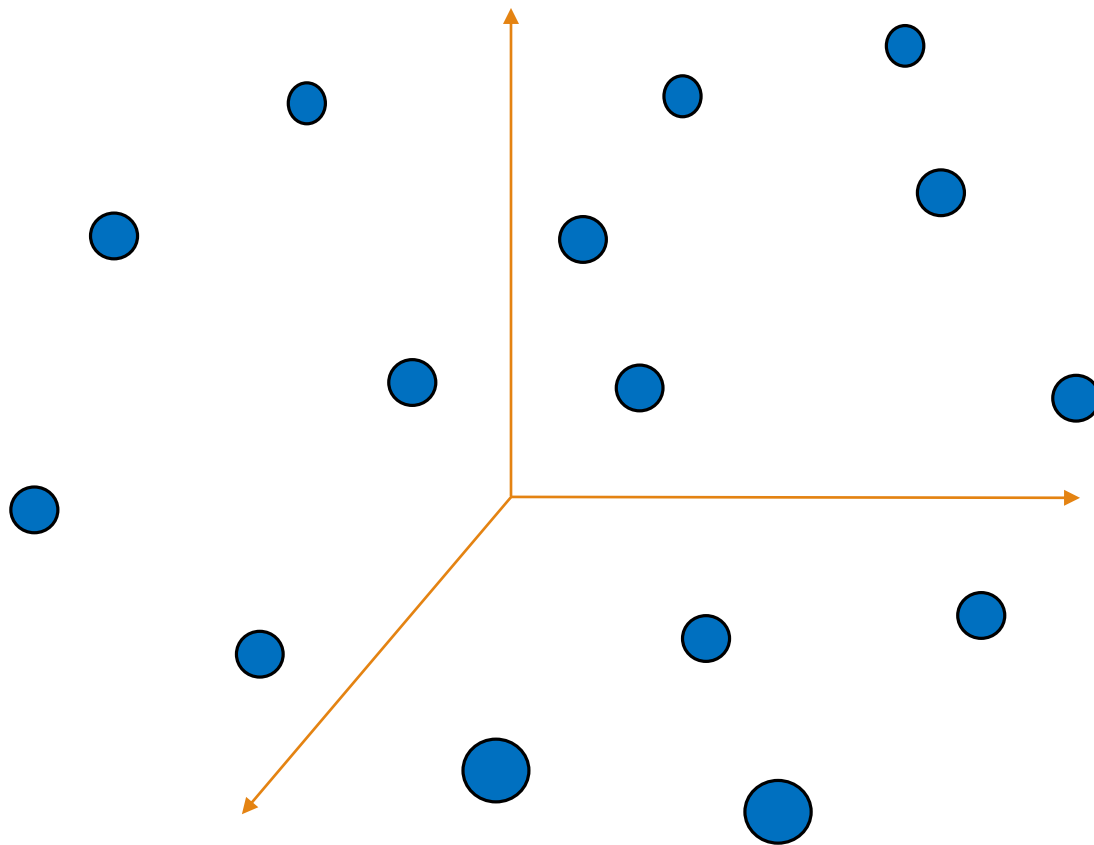


# Coreset for 1-Plane in $R^3$

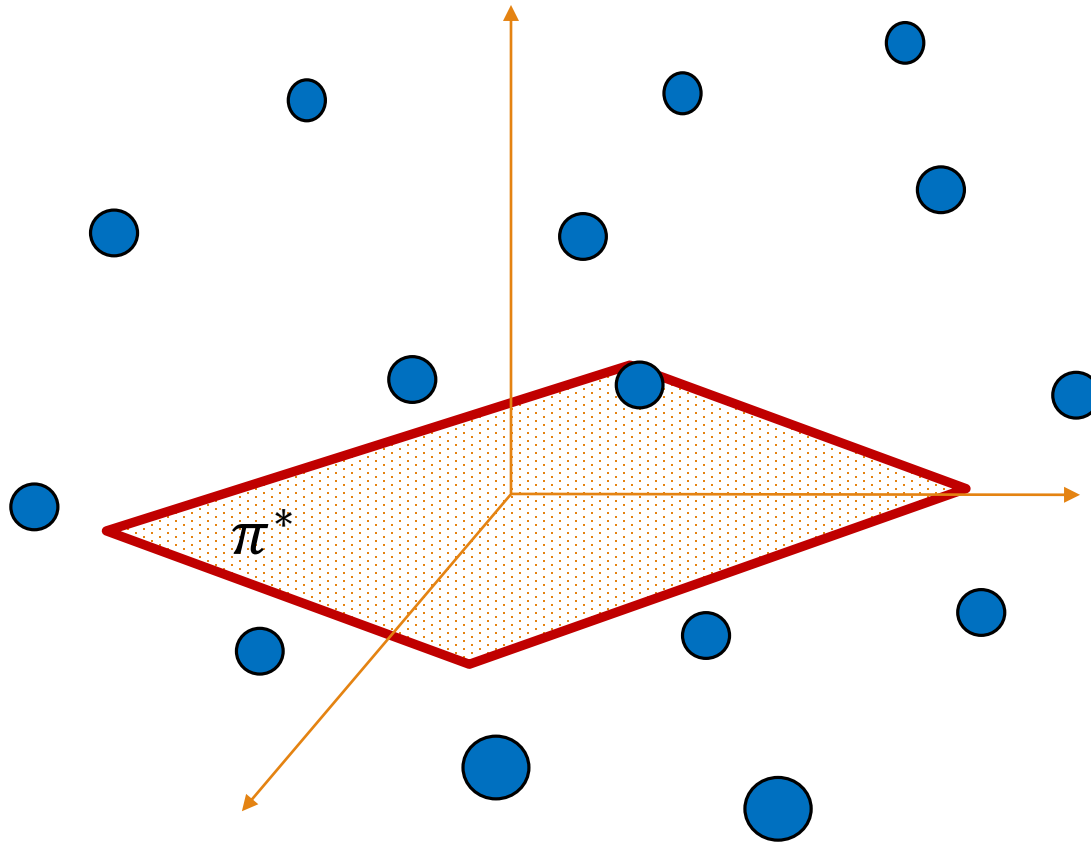
- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\pi \mid \pi \text{ is a plane in } R^3\}$
- Cost function:  $dist(p, \pi) = \min_{x \in \pi} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \pi \in Q: \max_{p \in P} dist(p, \pi) - \max_{c \in C} dist(c, \pi) \leq \epsilon \cdot \max_{p \in P} dist(p, \pi)$



# Coreset for 1-Plane in $R^3$

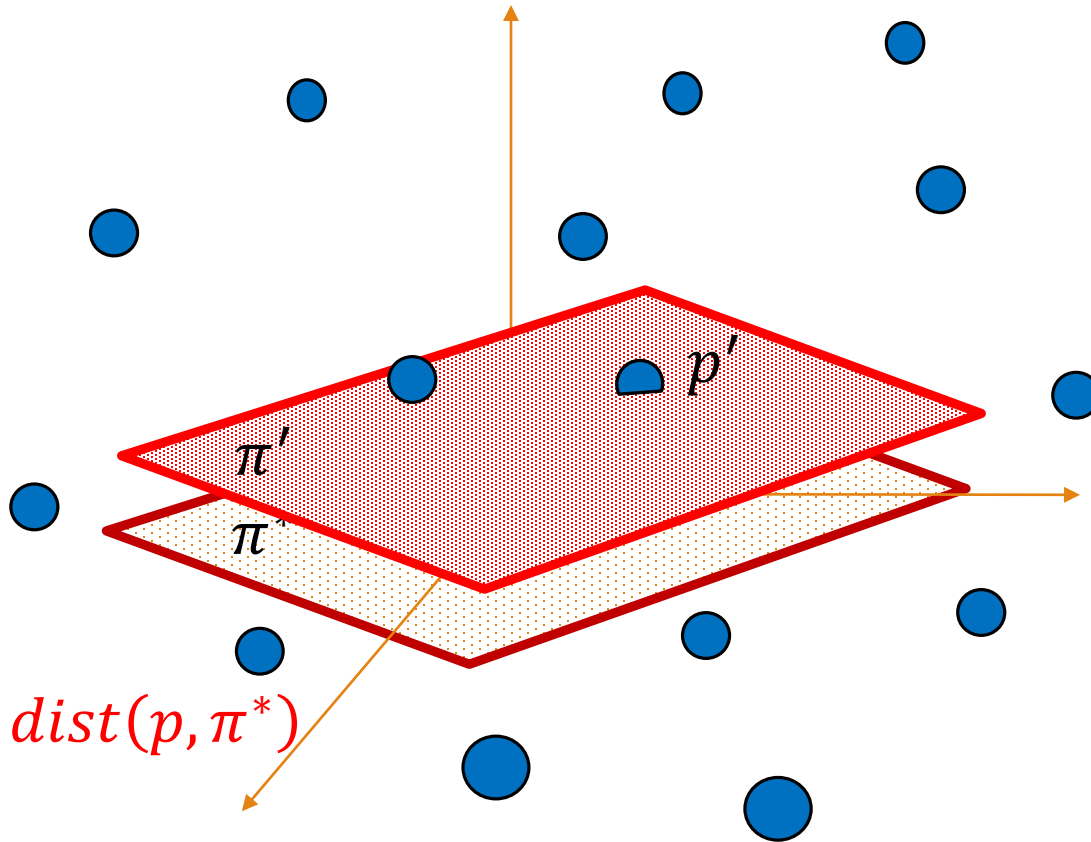


# Coreset for 1-Plane in $R^3$



$\pi^*$  is the plane that minimizes  
 $\max_{p \in P} \text{dist}(p, \pi)$

# Coreset for 1-Plane in $R^3$



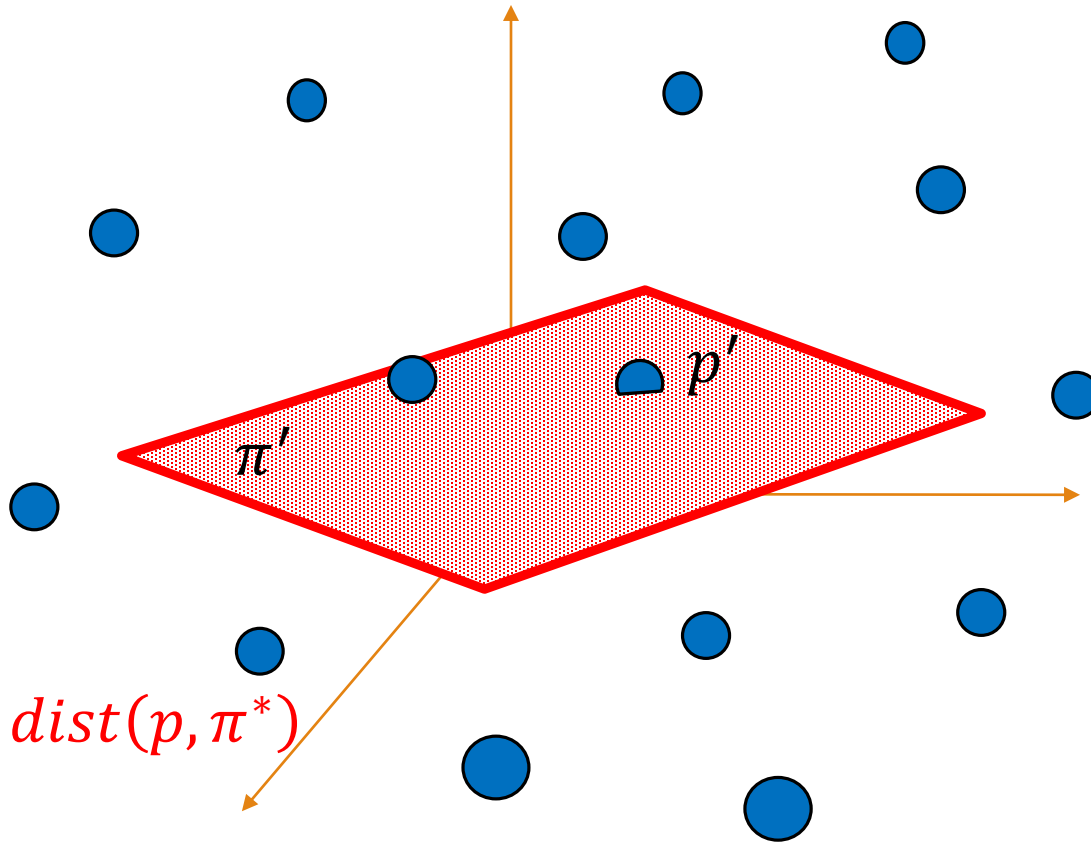
$\pi^*$  is the plane that minimizes  
 $\max_{p \in P} dist(p, \pi)$

$\pi'$  is the translation of  $\pi^*$  to  
 $\pi^*$ 's closest point  $p'$

$$dist(p, \pi') \leq 2 \cdot dist(p, \pi^*)$$



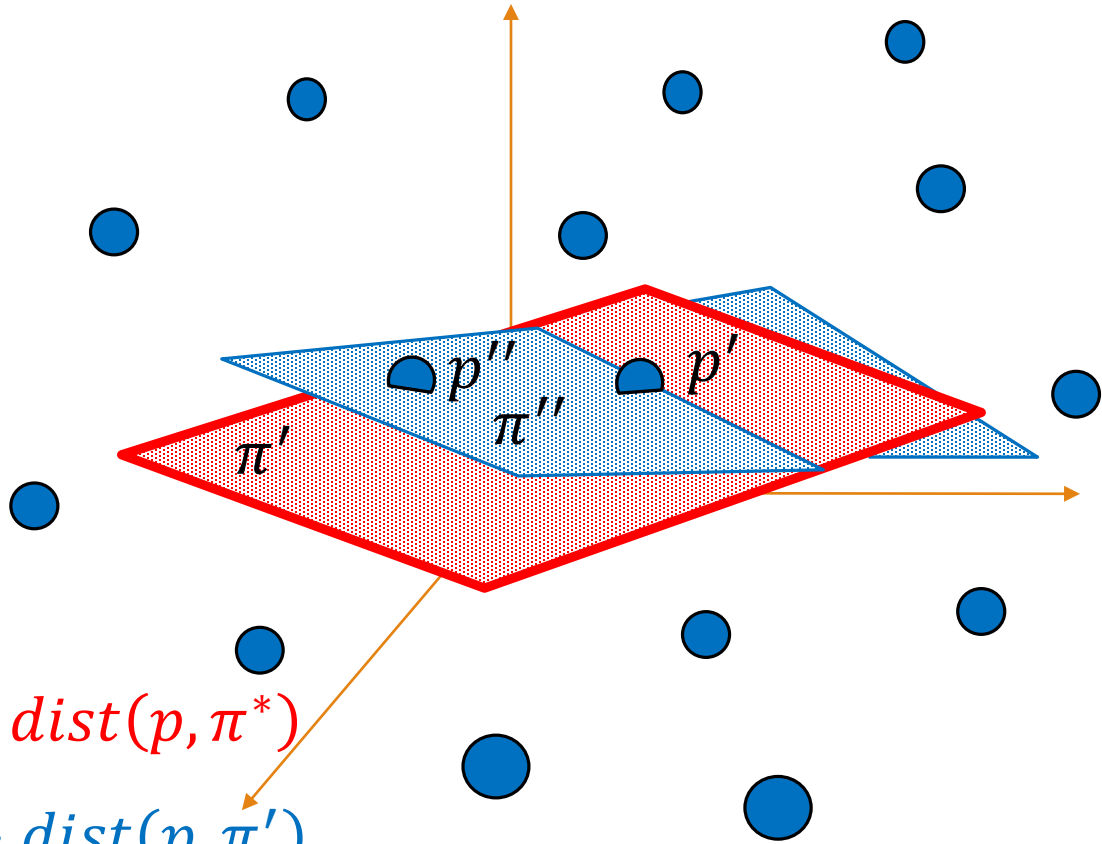
# Coreset for 1-Plane in $R^3$



$\pi'$  is the translation of  $\pi^*$  to  $\pi^*$ 's closest point  $p'$

$$\text{dist}(p, \pi') \leq 2 \cdot \text{dist}(p, \pi^*)$$

# Coreset for 1-Plane in $R^3$



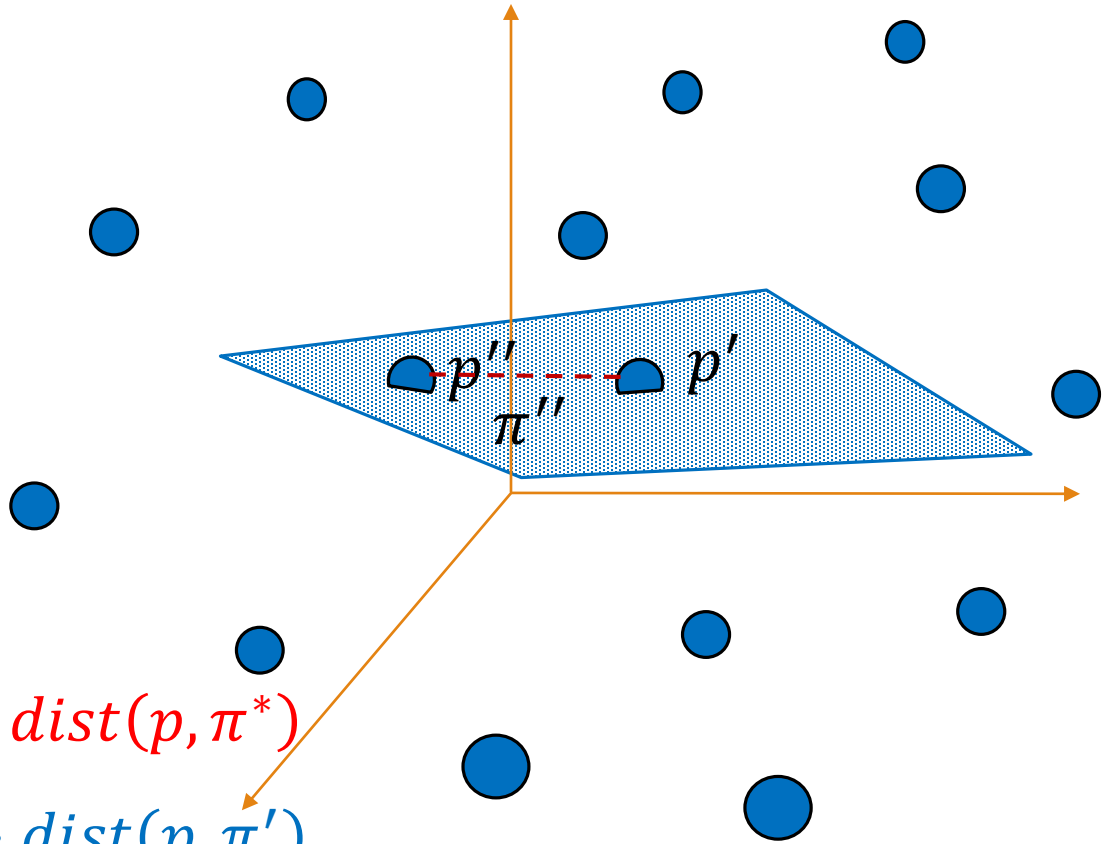
$\pi'$  is the translation of  $\pi^*$  to  $\pi^*$ 's closest point  $p'$

$\pi''$  is the rotation of  $\pi'$  around  $p'$  to  $\pi'$ 's closest point  $p''$

$$\text{dist}(p, \pi') \leq 2 \cdot \text{dist}(p, \pi^*)$$

$$\text{dist}(p, \pi'') \leq 2 \cdot \text{dist}(p, \pi')$$

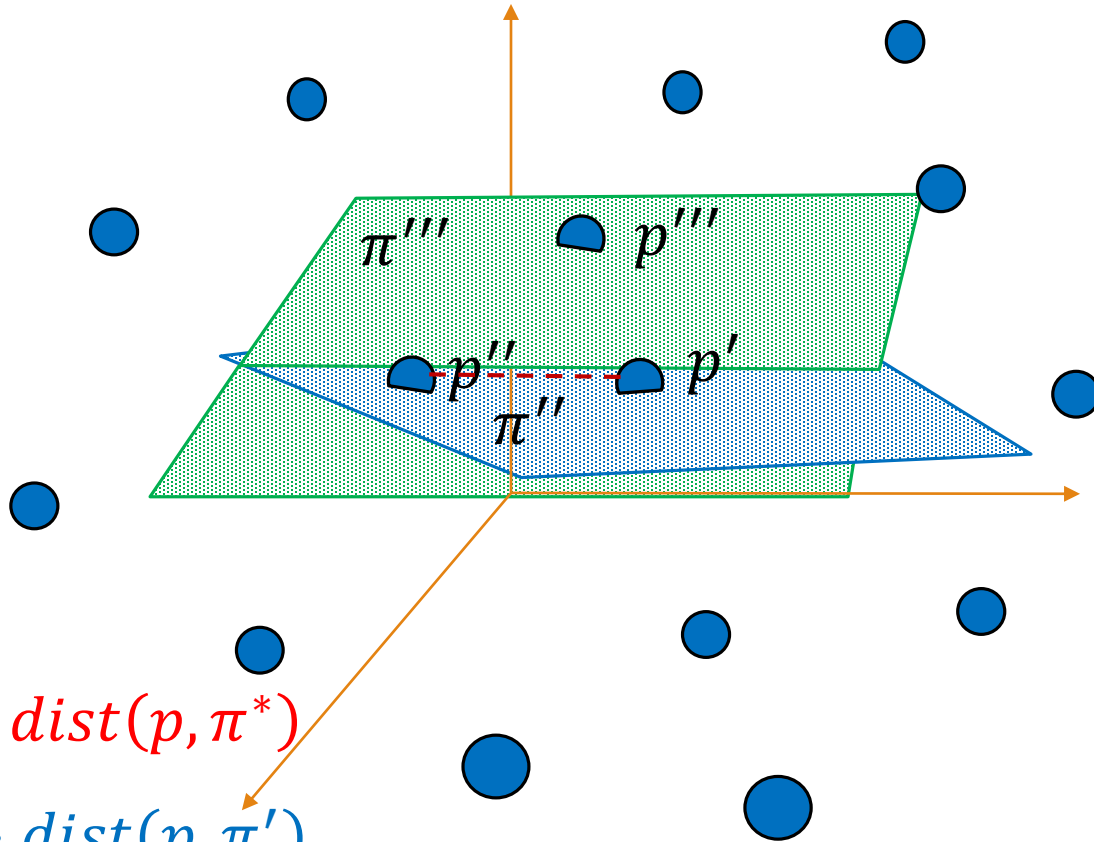
# Coreset for 1-Plane in $R^3$



$\pi''$  is the rotation of  $\pi'$  around  $p'$  to  $\pi'$ 's closest point  $p''$

$\text{dist}(p, \pi') \leq 2 \cdot \text{dist}(p, \pi^*)$   
 $\text{dist}(p, \pi'') \leq 2 \cdot \text{dist}(p, \pi')$

# Coreset for 1-Plane in $R^3$



$\pi''$  is the rotation of  $\pi'$  around  $p'$  to  $\pi'$ 's closest point  $p''$

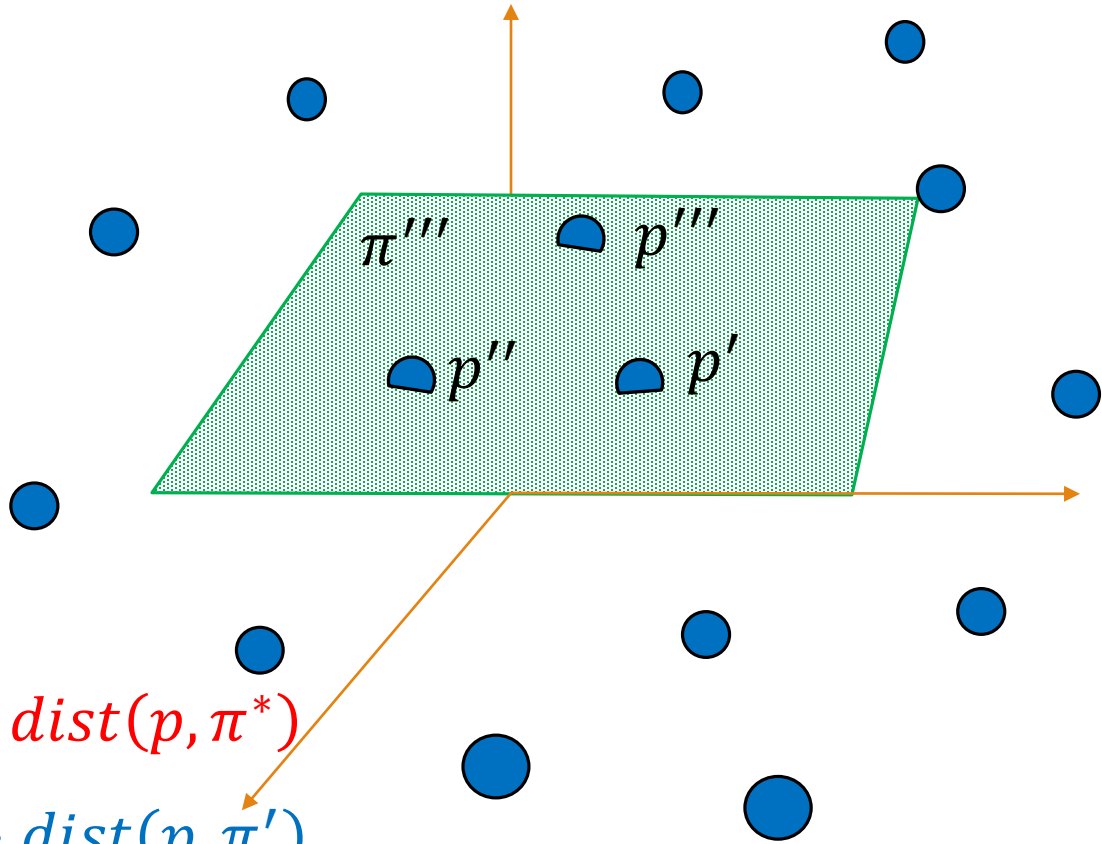
$\pi'''$  is the rotation of  $\pi''$  around  $p' - p''$  to  $\pi''$ 's closest point  $p'''$

$$\text{dist}(p, \pi') \leq 2 \cdot \text{dist}(p, \pi^*)$$

$$\text{dist}(p, \pi'') \leq 2 \cdot \text{dist}(p, \pi')$$

$$\text{dist}(p, \pi''') \leq 2 \cdot \text{dist}(p, \pi'')$$

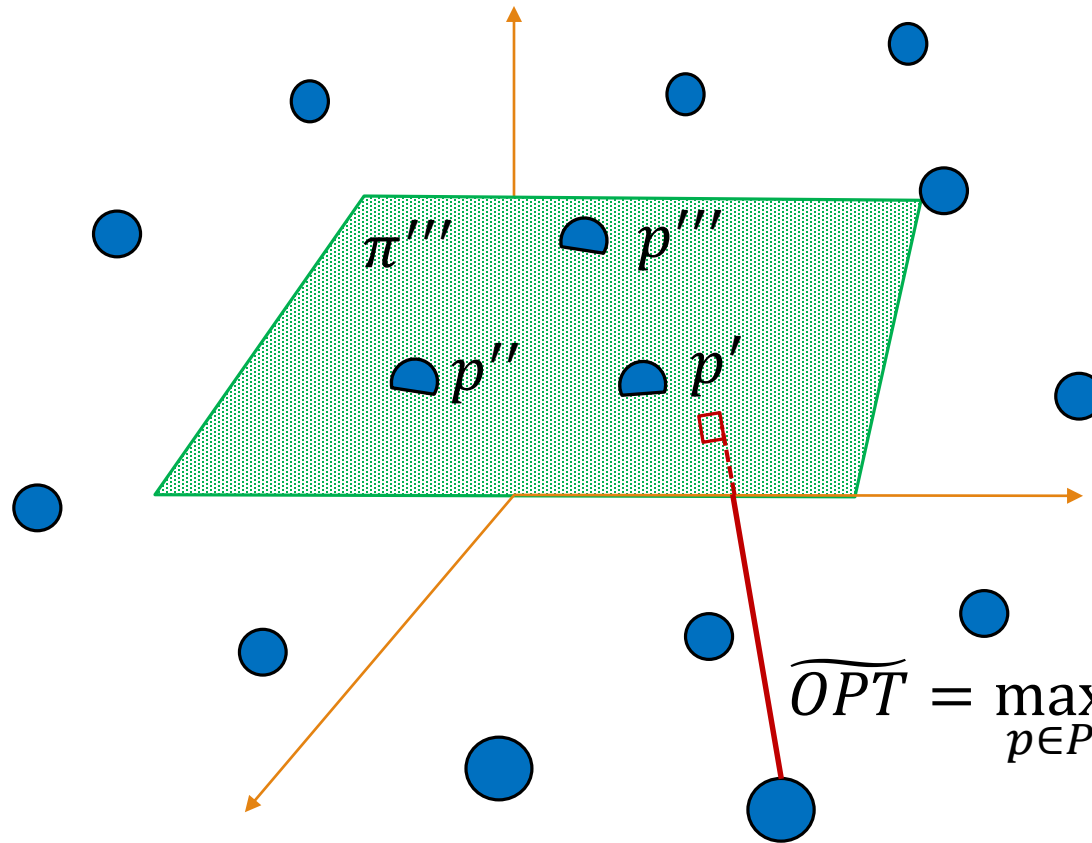
# Coreset for 1-Plane in $R^3$



$$\begin{aligned} \text{dist}(p, \pi') &\leq 2 \cdot \text{dist}(p, \pi^*) \\ \text{dist}(p, \pi'') &\leq 2 \cdot \text{dist}(p, \pi') \\ \text{dist}(p, \pi''') &\leq 2 \cdot \text{dist}(p, \pi'') \\ \text{dist}(p, \pi''') &\leq 8 \cdot \text{dist}(p, \pi^*) \end{aligned}$$

$\pi'''$  is the rotation of  $\pi''$  around  $p' - p''$  to  $\pi''$ 's closest point  $p'''$

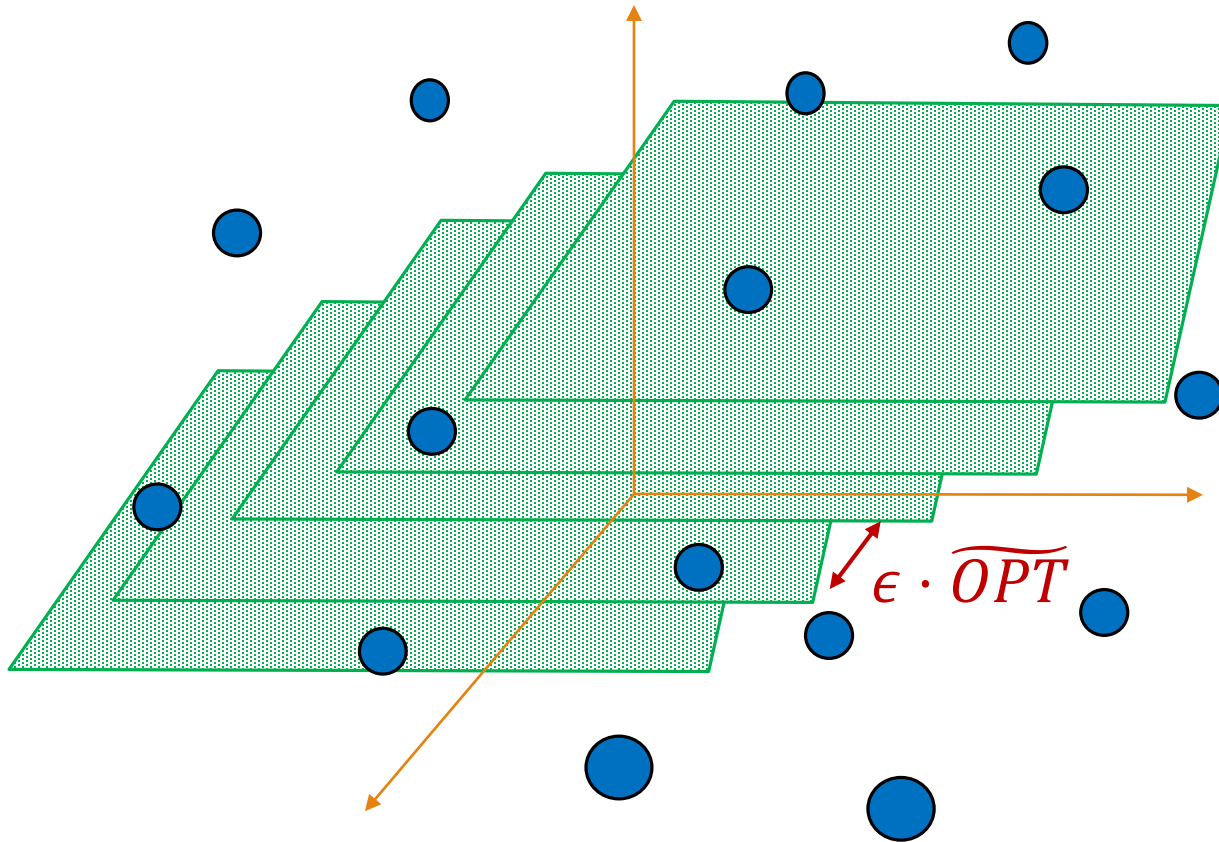
# Coreset for 1-Plane in $R^3$



Find  $\pi'''$  by exhaustive search  
over every triplet of points.  
 $O(n^3)$

$$\widetilde{OPT} = \max_{p \in P} \text{dist}(p, \pi''')$$

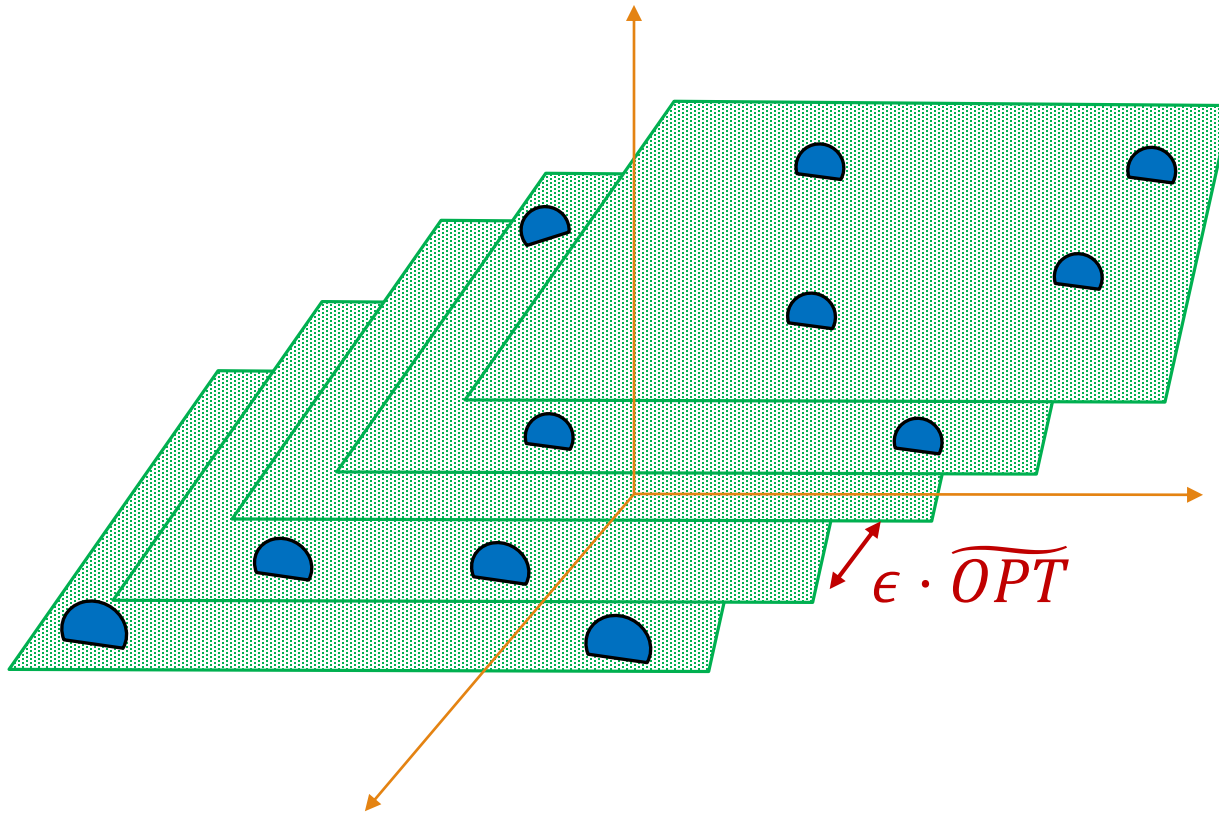
# Coreset for 1-Plane in $R^3$



Find  $\pi'''$  by exhaustive search  
over every triplet of points.  
 $O(n^3)$

Build a grid of planes with  
 $\epsilon \cdot \overline{OPT}$  distance

# Coreset for 1-Plane in $R^3$



Find  $\pi'''$  by exhaustive search  
over every triplet of points.  
 $O(n^3)$

Build a grid of planes with  
 $\epsilon \cdot \overline{OPT}$  distance

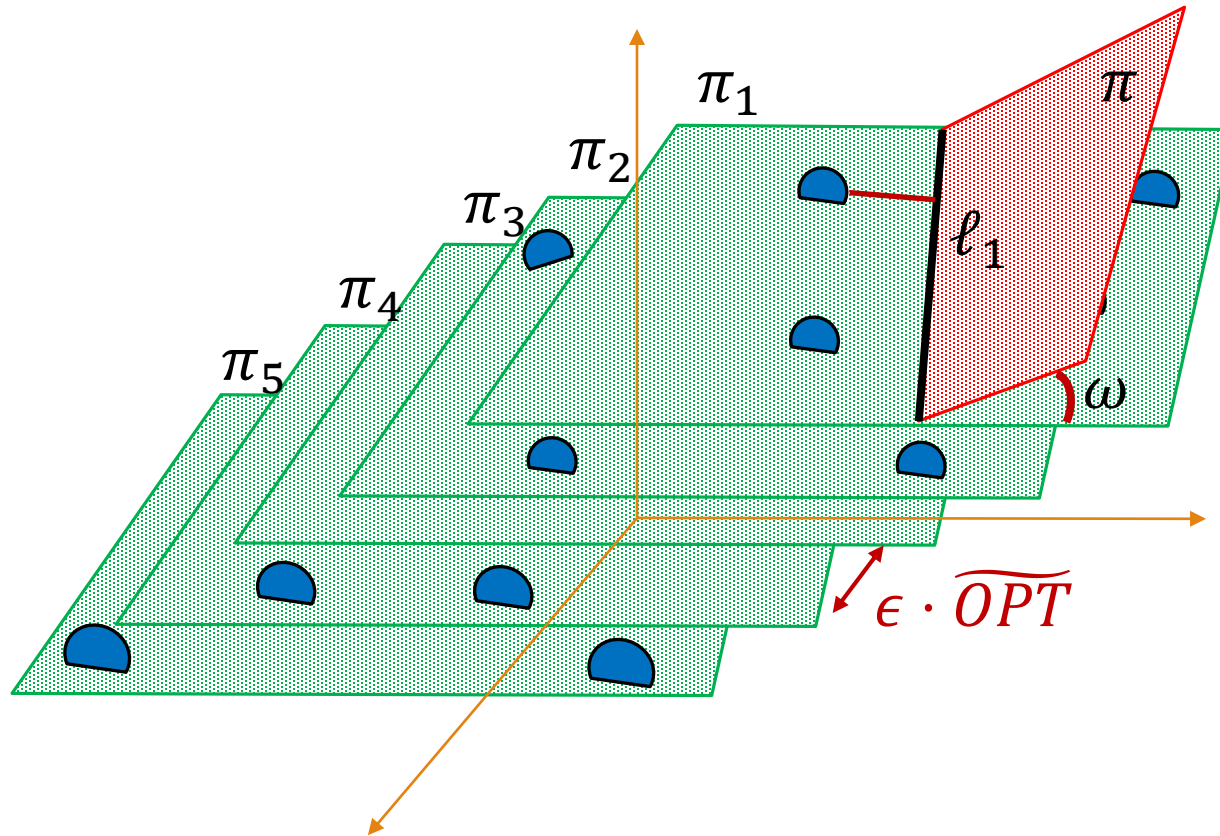
Project each point onto  
it's closest plane



# Coreset for 1-Plane in $R^3$

$$\forall p \in \pi_i: \text{dist}(p, \pi) = \omega \cdot \text{dist}(p, \ell_i)$$

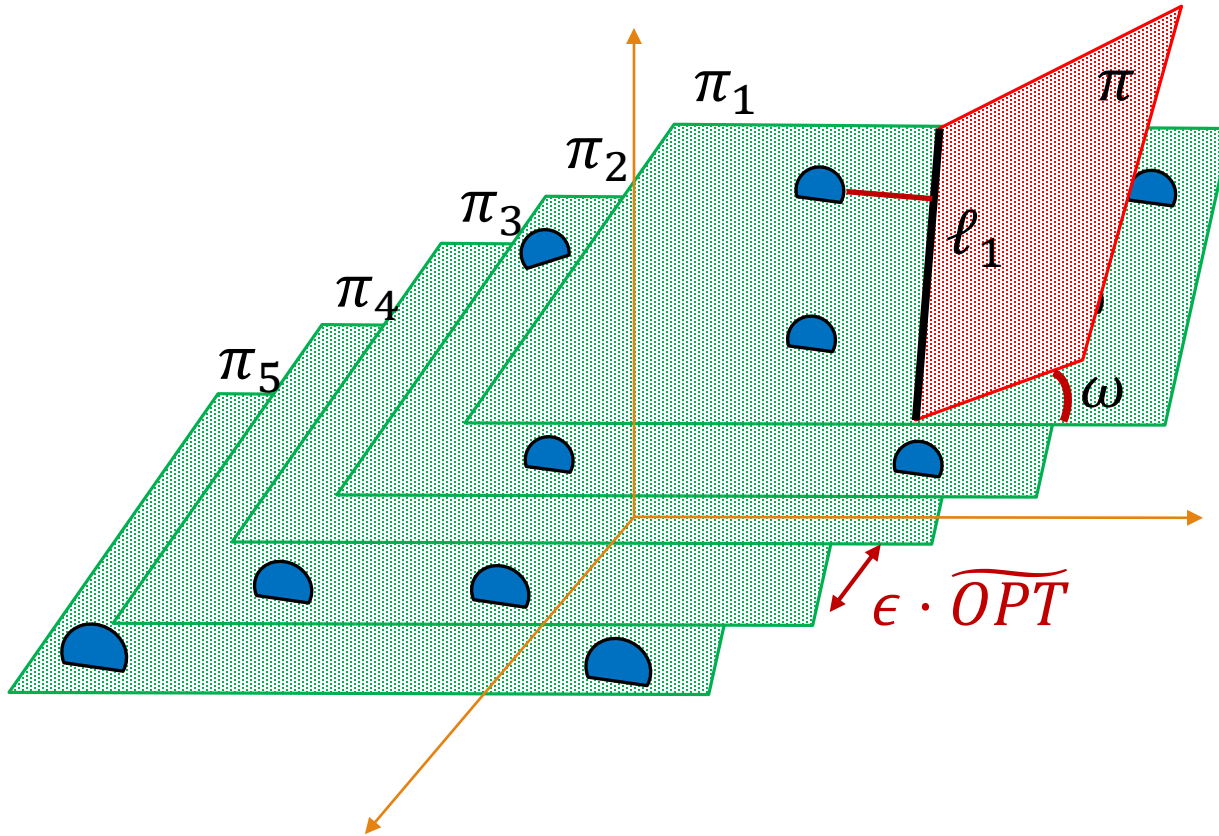
$$\ell_i = \pi_i \cap \pi$$



# Coreset for 1-Plane in $R^3$

$$\forall p \in \pi_i: \text{dist}(p, \pi) = \omega \cdot \text{dist}(p, \ell_i)$$

$$\ell_i = \pi_i \cap \pi$$



→ Compute a **1-Line** coreset  $C_i$   
for each **plane**  $\pi_i$ !

$$C = \bigcup C_i$$

since a union of two  
coresets is a coreset.

# Coreset for *Hyperplane* in $R^d$

*HyperplaneCoreset*( **$P$** ,  $d$ ):

# Coreset for *Hyperplane* in $R^d$

*HyperplaneCoreset*( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .

# Coreset for *Hyperplane* in $R^d$

*HyperplaneCoreset*( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .

# Coreset for *Hyperplane* in $R^d$

## *HyperplaneCoreset*( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the vector that is orthogonal to  $h'$ .

# Coreset for *Hyperplane* in $R^d$

## HyperplaneCoreset( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the vector that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .

# Coreset for *Hyperplane* in $R^d$

## HyperplaneCoreset( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the vector that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct a hyperplane parallel to  $h'$ . ( $\#Hyperplanes = \frac{2}{\epsilon}$ )



# Coreset for *Hyperplane* in $R^d$

## HyperplaneCoreset( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the vector that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct a hyperplane parallel to  $h'$ . ( $\#Hyperplanes = \frac{2}{\epsilon}$ )
- Compute the projection  $p'$  of each point  $p \in P$  onto its closest hyperplane  $h_p$ .

# Coreset for *Hyperplane* in $R^d$

## HyperplaneCoreset( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the vector that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct a hyperplane parallel to  $h'$ . ( $\#Hyperplanes = \frac{2}{\epsilon}$ )
- Compute the projection  $p'$  of each point  $p \in P$  onto its closest hyperplane  $h_p$ .
- $H_p \leftarrow$  an  $R^{d \times d-1}$  matrix whose columns span  $h_p$  and  $H_p^T H_p = I$ .

# Coreset for *Hyperplane* in $R^d$

## HyperplaneCoreset( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the vector that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct a hyperplane parallel to  $h'$ . ( $\#Hyperplanes = \frac{2}{\epsilon}$ )
- Compute the projection  $p'$  of each point  $p \in P$  onto its closest hyperplane  $h_p$ .
- $H_p \leftarrow$  an  $R^{d \times d-1}$  matrix whose columns span  $h_p$  and  $H_p^T H_p = I$ .
- $P' = \{H_p p' \mid p \in P\}$ .

# Coreset for *Hyperplane* in $R^d$

## *HyperplaneCoreset*( $P, d$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the optimal hyperplane of  $P$ .
- $\overline{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the vector that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \overline{OPT}$ .
- Through each grid point construct a hyperplane parallel to  $h'$ . ( $\#Hyperplanes = \frac{2}{\epsilon}$ )
- Compute the projection  $p'$  of each point  $p \in P$  onto its closest hyperplane  $h_p$ .
- $H_p \leftarrow$  an  $R^{d \times d-1}$  matrix whose columns span  $h_p$  and  $H_p^T H_p = I$ .
- $P' = \{H_p p' \mid p \in P\}$
- Call *HyperplaneCoreset*( $P', d - 1$ ).

# Coreset for *Hyperplane* in $R^d$

Total time:

$$O(n^d)$$

Coreset size:

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

# Coreset for *Hyperplane* in $R^d$

Total time:

$$O(n^d)$$

Coreset size:

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

Improvement:

Run the above algorithm using the streaming tree.

Run on batches of size  $2 \cdot |C| \leq 2 \left(\frac{1}{\epsilon}\right)^{O(d)}$ .

Total time:

$$O(n \cdot \text{TimeForBatch}) = O\left(n \cdot \left(\frac{1}{\epsilon}\right)^{O(d^2)}\right).$$

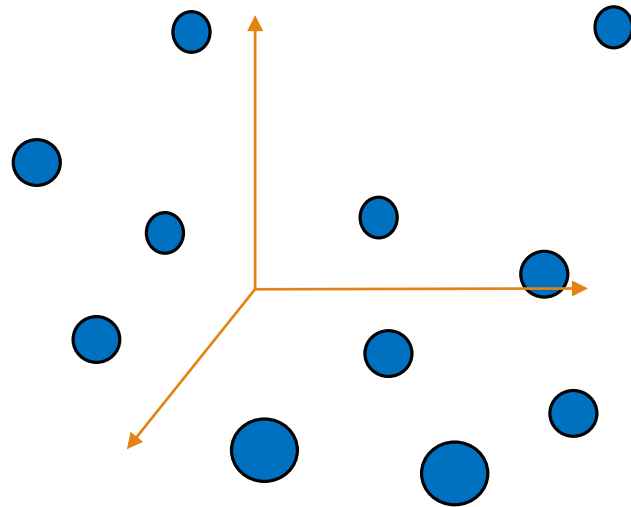
Error for streaming tree:

The error increases to  $(1 + \epsilon)^{\log n} \sim (1 + \epsilon \log n)$

→ Run with  $\epsilon' = \frac{\epsilon}{\log n}$ .

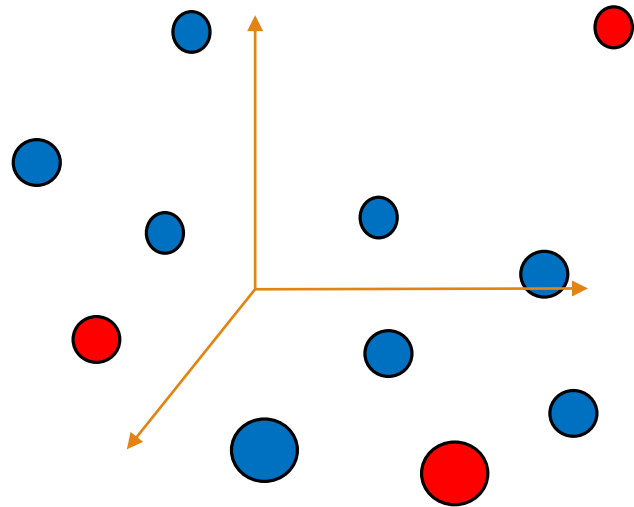
# Coreset for 1-Line parallel to z-axis in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3 \text{ parallel to the z-axis}\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$



# Coreset for 1-Line parallel to z-axis in $R^3$

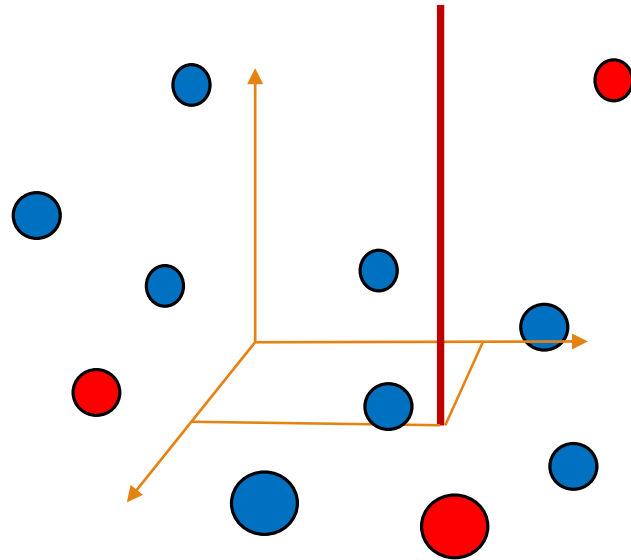
- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3 \text{ parallel to the z-axis}\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$





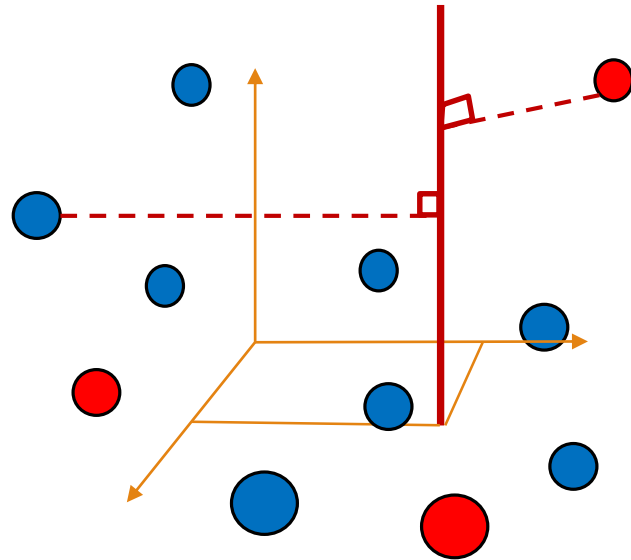
# Coreset for 1-Line parallel to z-axis in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3 \text{ parallel to the z-axis}\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$

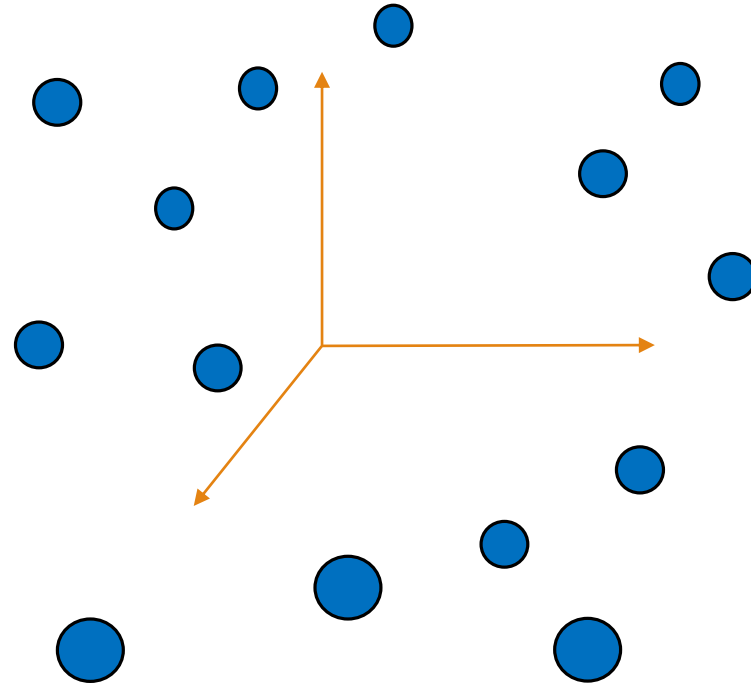


# Coreset for 1-Line parallel to z-axis in $R^3$

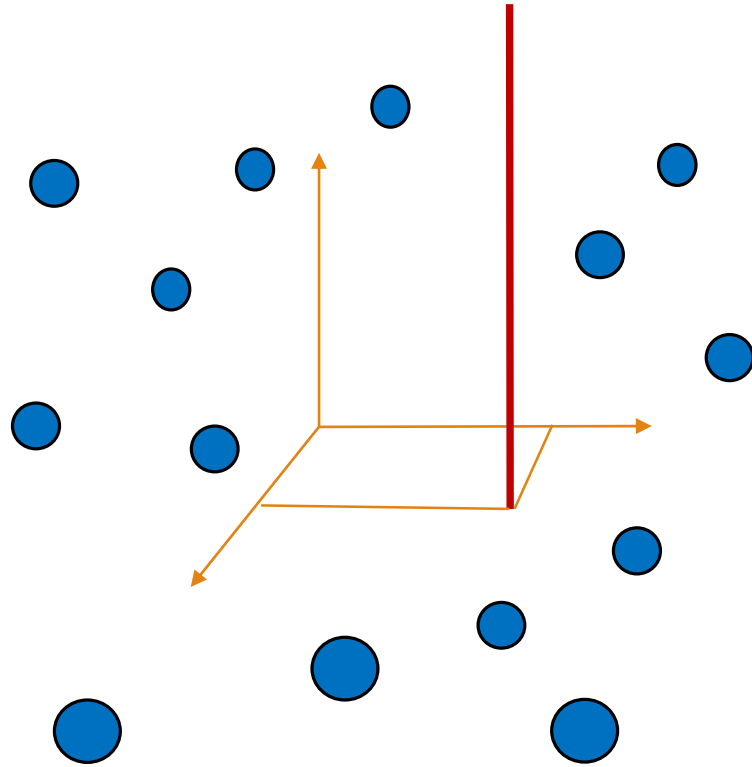
- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3 \text{ parallel to the z-axis}\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$



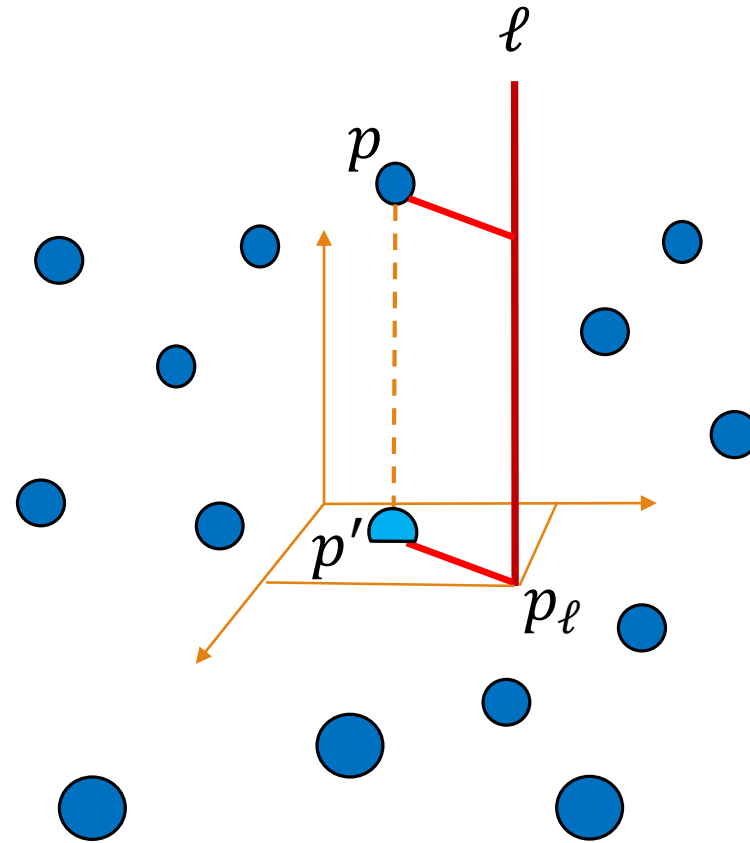
# Coreset for 1-Line parallel to $z$ -axis in $R^3$



# Coreset for 1-Line parallel to $z$ -axis in $R^3$



# Coreset for 1-Line parallel to z-axis in $R^3$



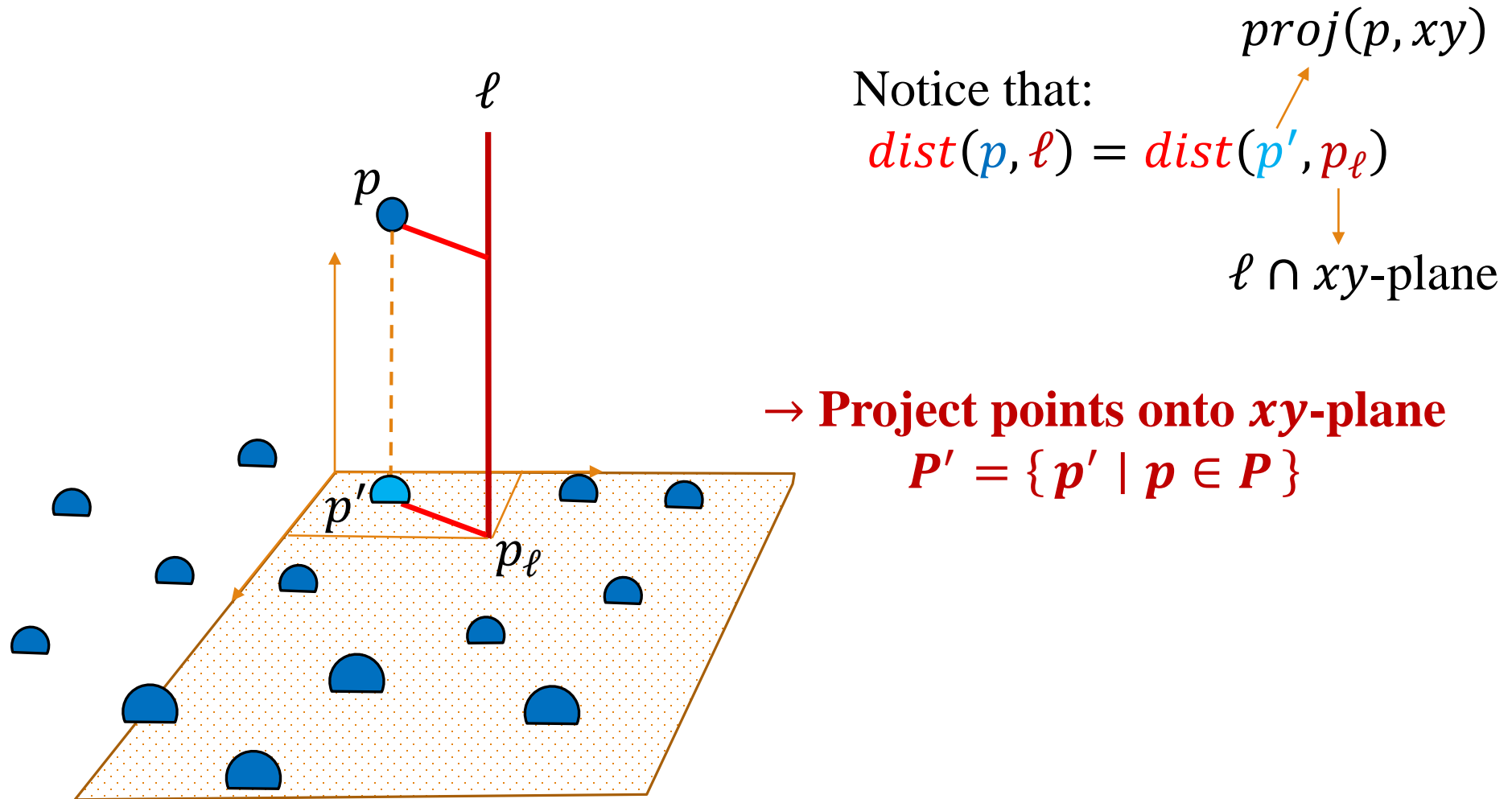
Notice that:

$$\text{dist}(p, \ell) = \text{dist}(p', p_\ell)$$

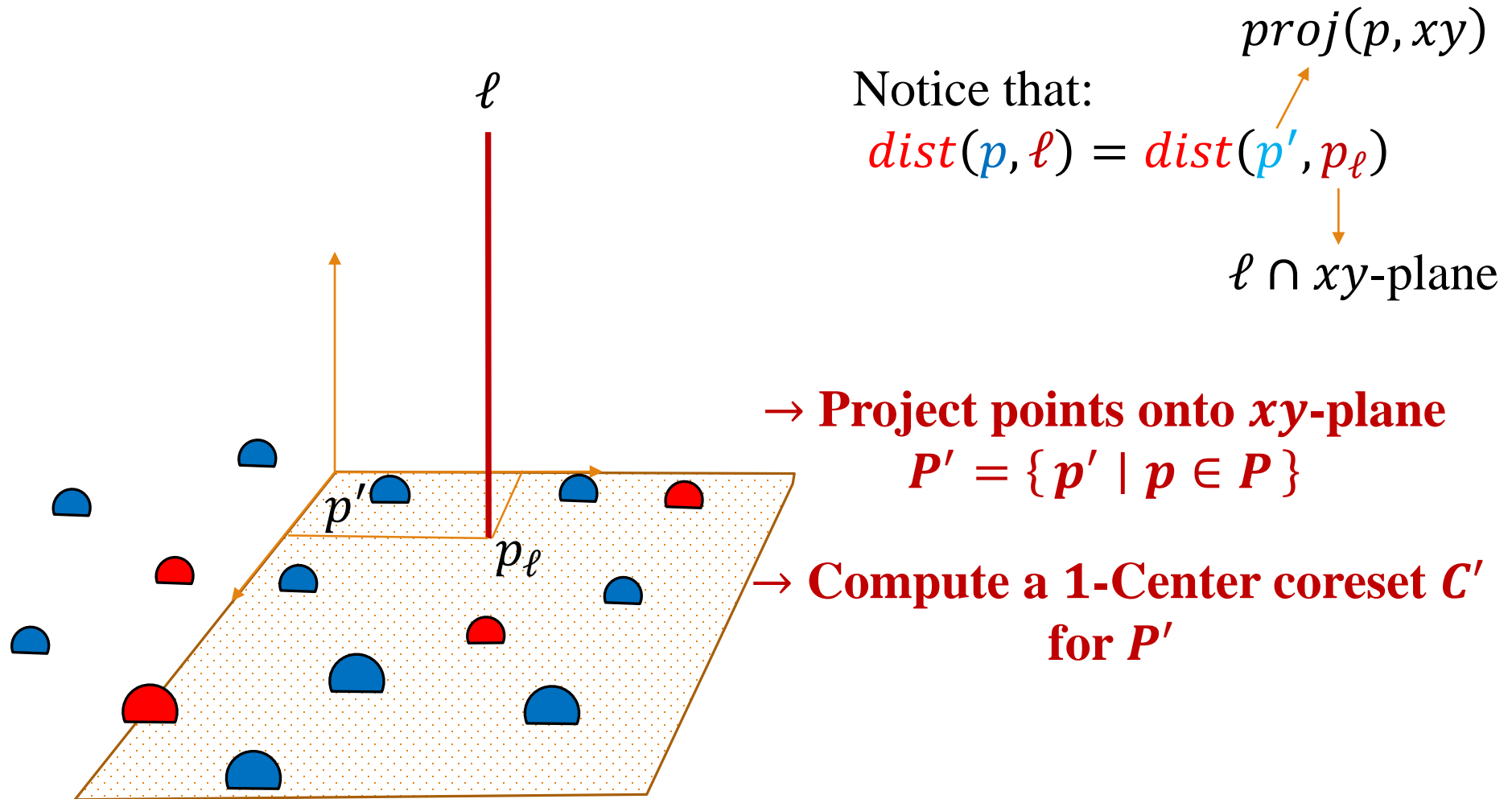
$\text{proj}(p, xy)$

$\ell \cap xy\text{-plane}$

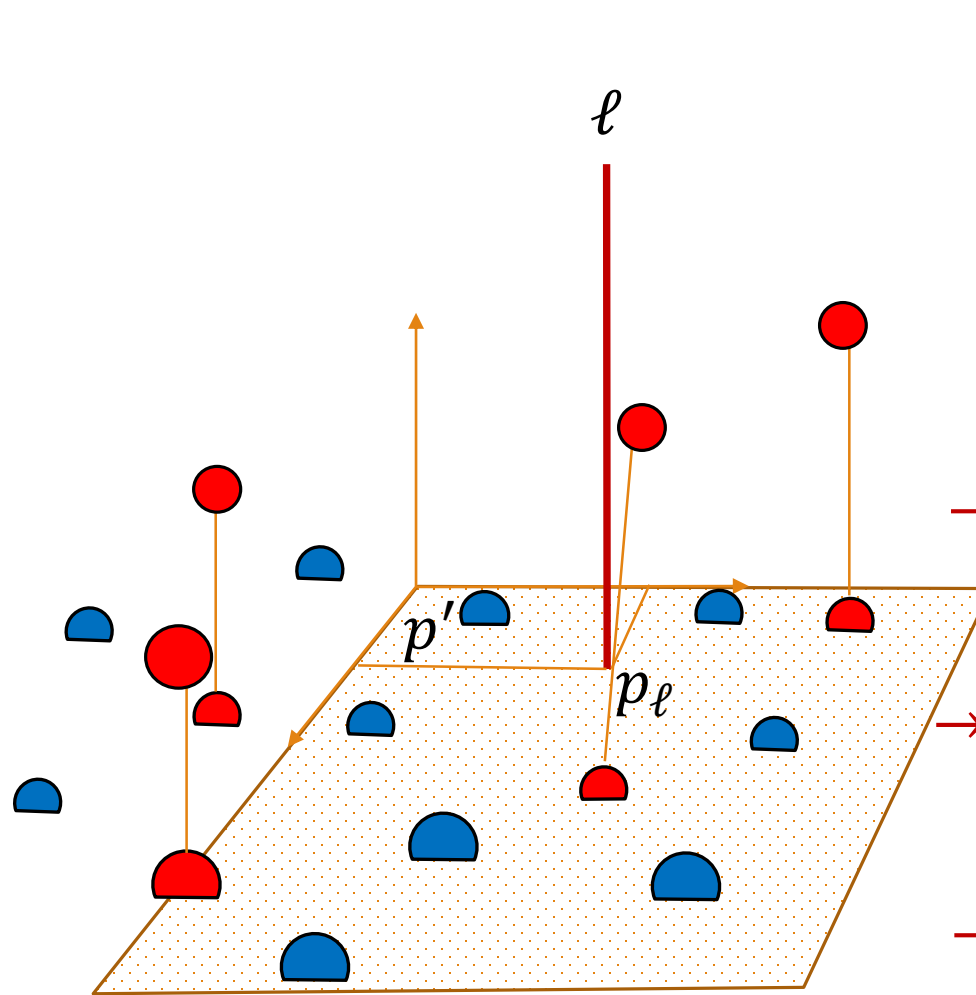
# Coreset for 1-Line parallel to z-axis in $R^3$



# Coreset for 1-Line parallel to z-axis in $R^3$



# Coreset for 1-Line parallel to z-axis in $R^3$



Notice that:

$$\text{dist}(p, \ell) = \text{dist}(p', p_\ell)$$

$\text{proj}(p, xy)$

$\ell \cap xy\text{-plane}$

→ Project points onto  $xy$ -plane

$$P' = \{p' \mid p \in P\}$$

→ Compute a 1-Center coreset  $C'$  for  $P'$

→ Return to original data

$$C' \rightarrow C$$



# Coreset for 1-Line parallel to z-axis in $R^3$

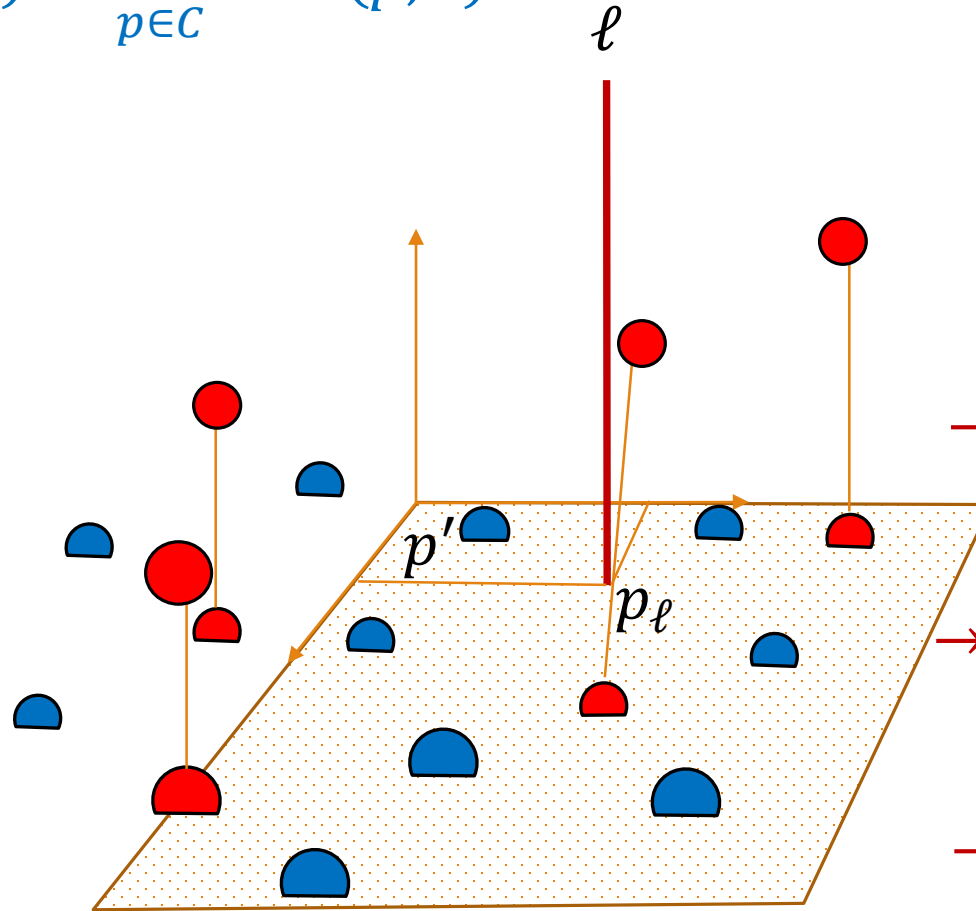
$$\forall \ell \in Q: \max_{p \in P} \text{dist}(p, \ell) - \max_{p \in C} \text{dist}(p, \ell)$$

Notice that:

$$\text{dist}(p, \ell) = \text{dist}(p', p_\ell)$$

$\text{proj}(p, xy)$

$\ell \cap xy\text{-plane}$



→ Project points onto  $xy\text{-plane}$

$$P' = \{p' \mid p \in P\}$$

→ Compute a 1-Center coreset  $C'$  for  $P'$

→ Return to original data

$$C' \rightarrow C$$

# Coreset for 1-Line parallel to z-axis in $R^3$

$$\forall \ell \in Q: \max_{p \in P} \text{dist}(p, \ell) - \max_{p \in C} \text{dist}(p, \ell)$$

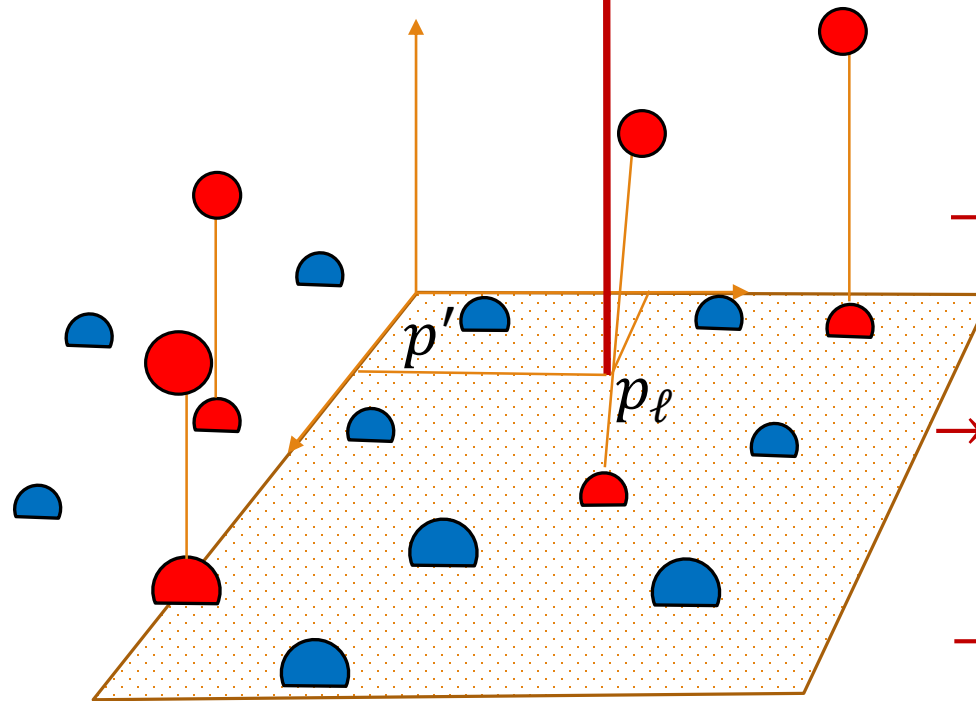
$$= \max_{p' \in P'} \text{dist}(p', p_\ell) - \max_{p' \in C'} \text{dist}(p', p_\ell)$$

Notice that:

$$\text{dist}(p, \ell) = \text{dist}(p', p_\ell)$$

$\text{proj}(p, xy)$

$\ell \cap xy\text{-plane}$



→ Project points onto  $xy\text{-plane}$

$$P' = \{p' \mid p \in P\}$$

→ Compute a 1-Center coreset  $C'$  for  $P'$

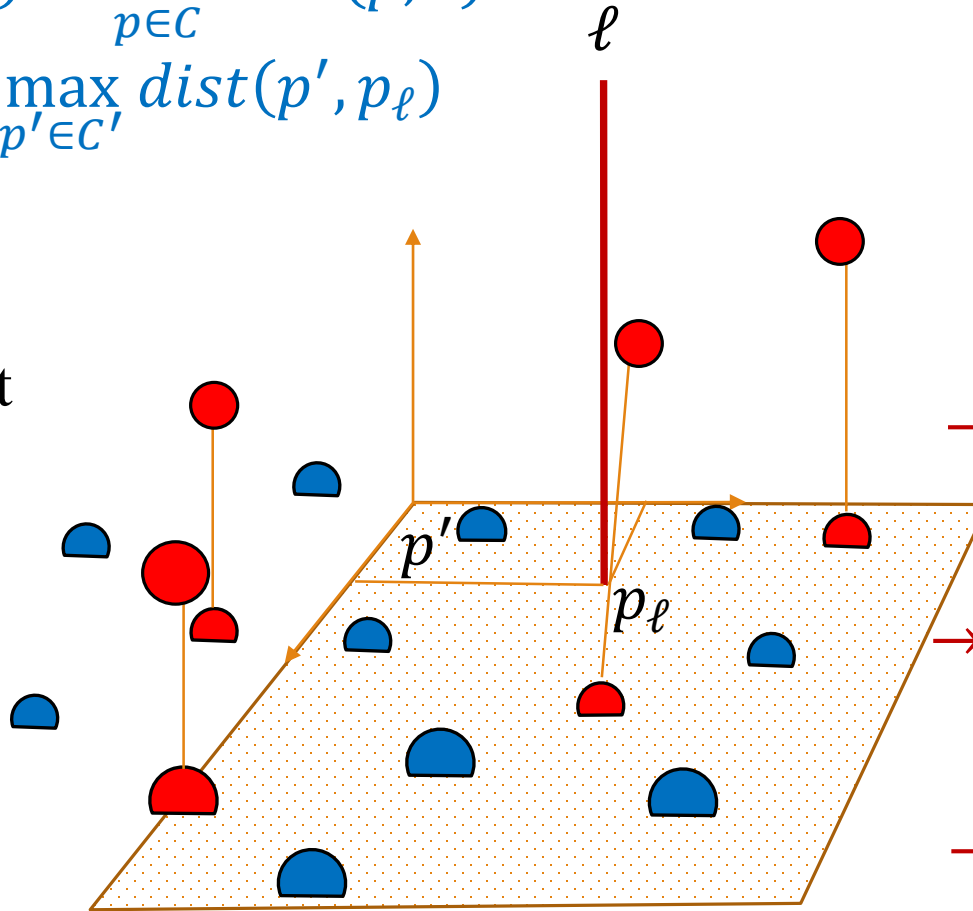
→ Return to original data

$$C' \rightarrow C$$

# Coreset for 1-Line parallel to z-axis in $R^3$

$$\begin{aligned} & \forall \ell \in Q: \max_{p \in P} \text{dist}(p, \ell) - \max_{p \in C} \text{dist}(p, \ell) \\ &= \max_{p' \in P'} \text{dist}(p', p_\ell) - \max_{p' \in C'} \text{dist}(p', p_\ell) \\ &\leq \epsilon \cdot \max_{p' \in P'} \text{dist}(p', p_\ell) \end{aligned}$$

$C'$  is a 1-Center coreset



Notice that:

$$\text{dist}(p, \ell) = \text{dist}(proj(p, xy), \ell \cap xy\text{-plane})$$

→ Project points onto  $xy$ -plane

$$P' = \{p' \mid p \in P\}$$

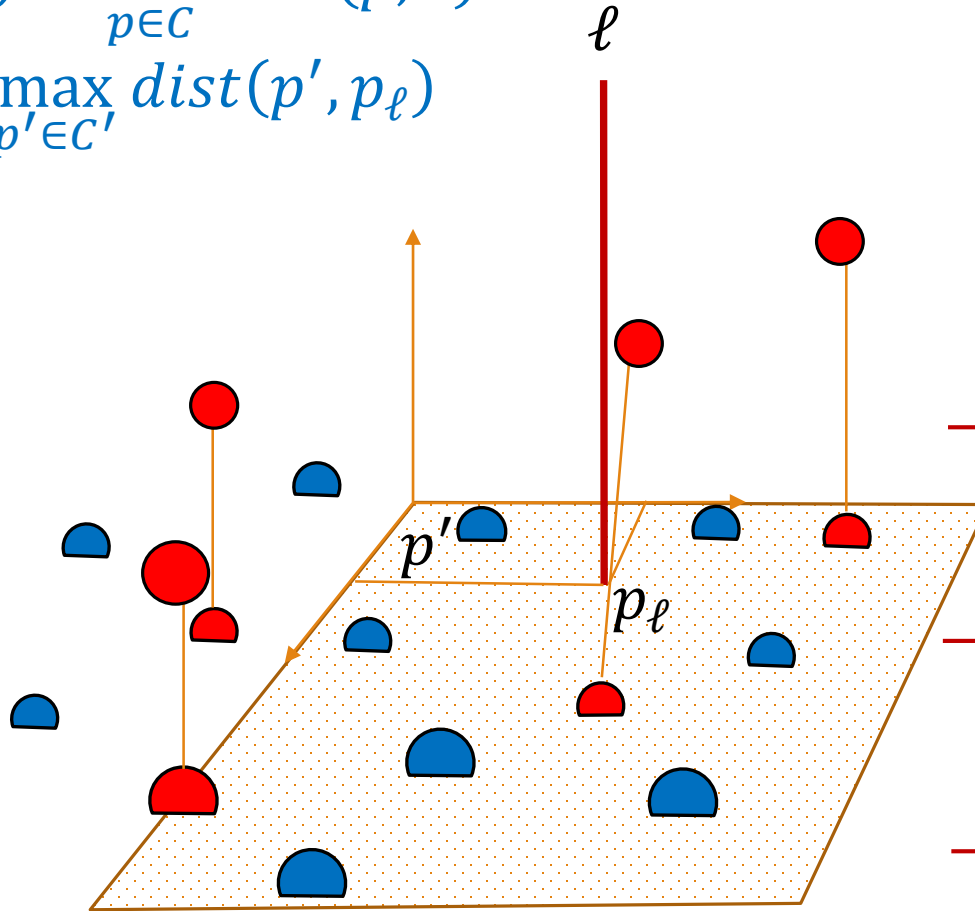
→ Compute a 1-Center coreset  $C'$  for  $P'$

→ Return to original data

$$C' \rightarrow C$$

# Coreset for 1-Line parallel to z-axis in $R^3$

$$\begin{aligned}
 & \forall \ell \in Q: \max_{p \in P} \text{dist}(p, \ell) - \max_{p \in C} \text{dist}(p, \ell) \\
 &= \max_{p' \in P'} \text{dist}(p', p_\ell) - \max_{p' \in C'} \text{dist}(p', p_\ell) \\
 &\leq \epsilon \cdot \max_{p' \in P'} \text{dist}(p', p_\ell) \\
 &= \epsilon \cdot \max_{p \in P} \text{dist}(p, \ell)
 \end{aligned}$$



Notice that:

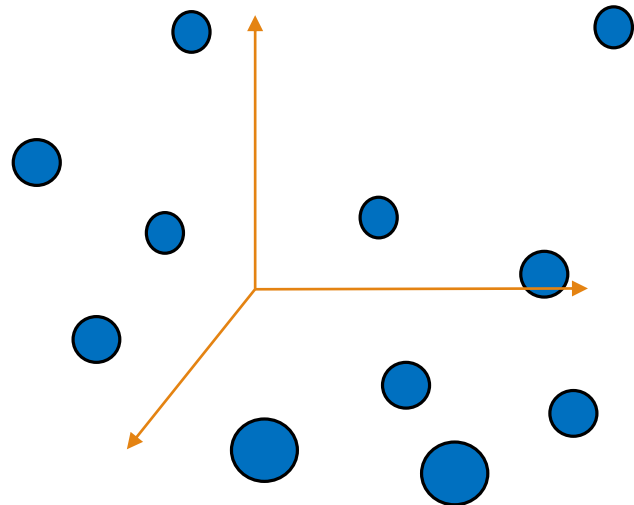
$$\text{dist}(p, \ell) = \text{dist}(p', p_\ell)$$

$\text{proj}(p, xy)$   
 $\ell \cap xy\text{-plane}$

- Project points onto **xy-plane**  
 $P' = \{p' \mid p \in P\}$
- Compute a **1-Center coreset C**  
for  $P'$
- Return to original data  
 $C' \rightarrow C$

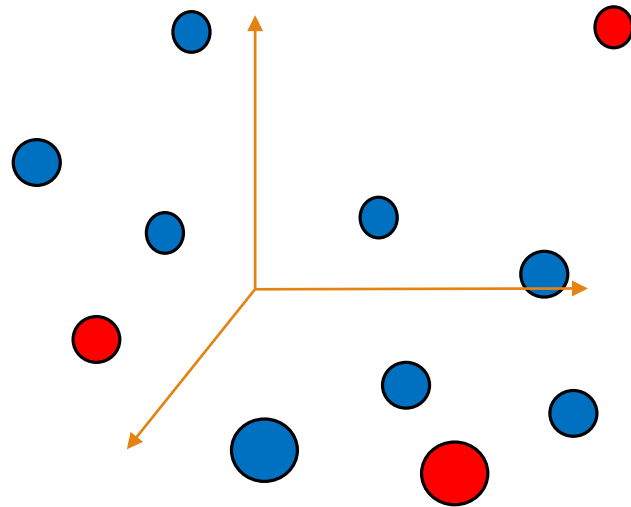
# Coreset for 1-Line in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$



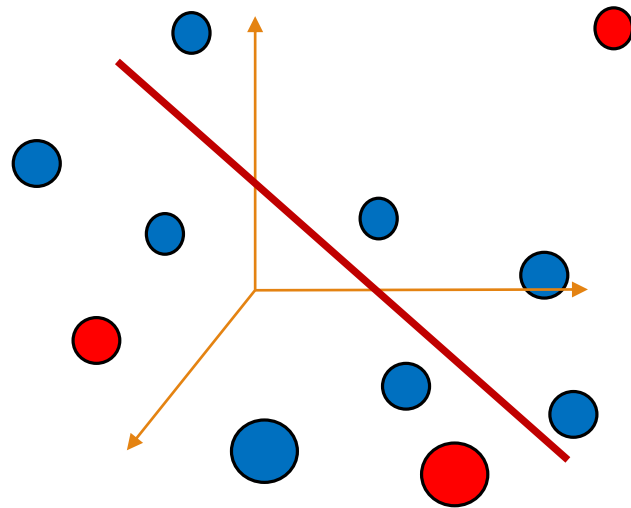
# Coreset for 1-Line in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$



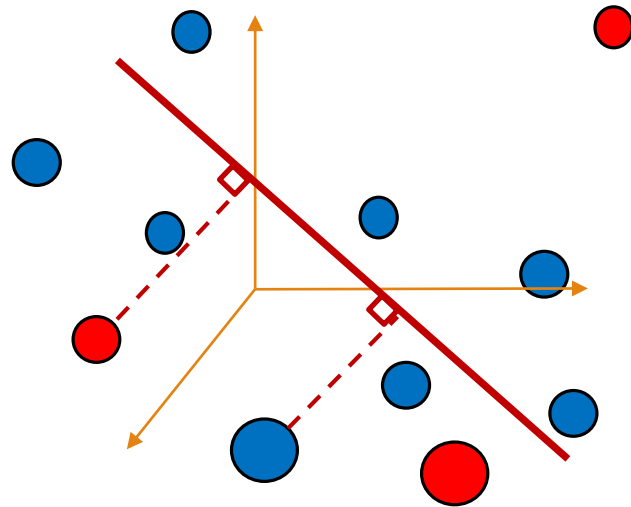
# Coreset for 1-Line in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$



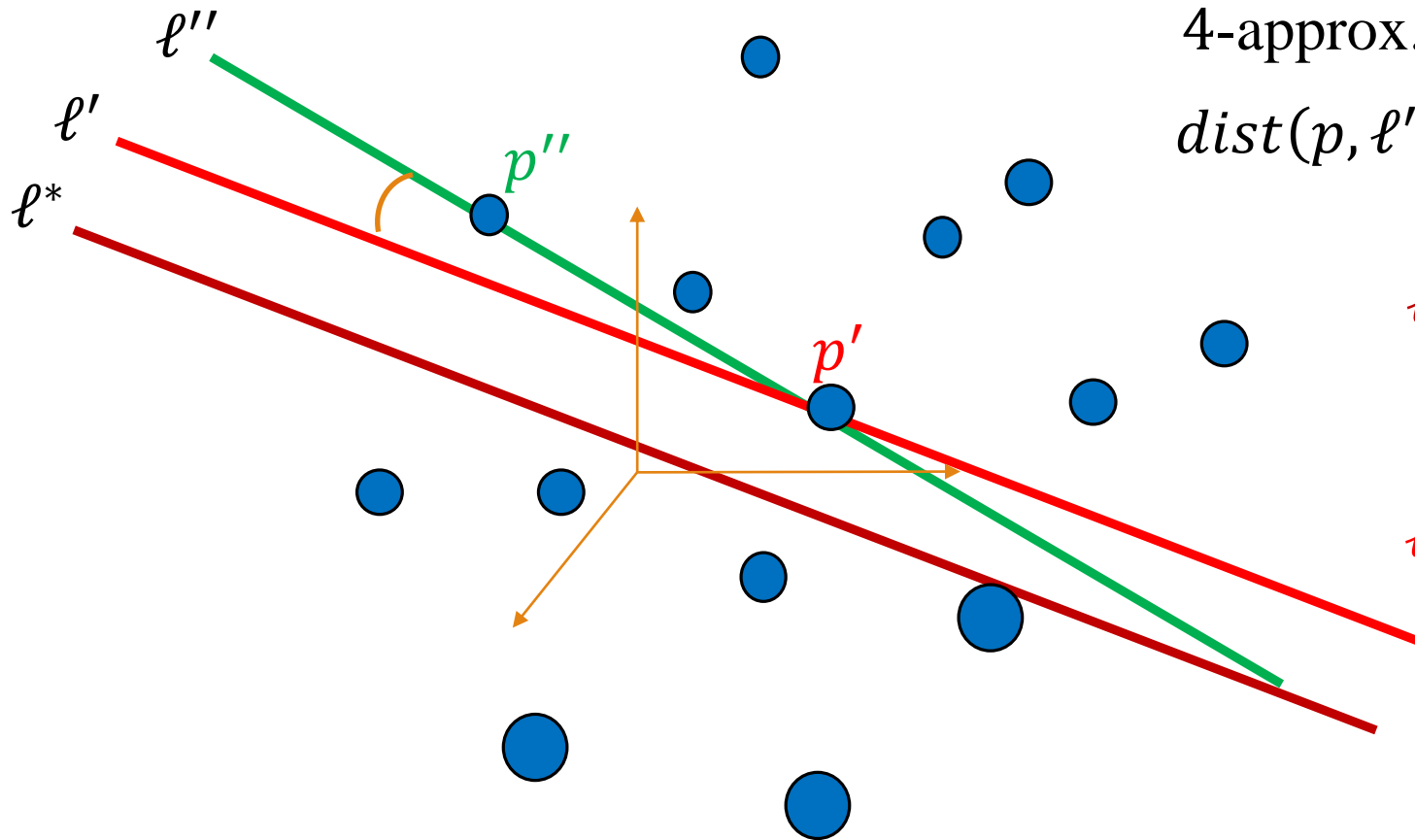
# Coreset for 1-Line in $R^3$

- Input:  $P \subseteq R^3$
- Query space:  $Q = \{\ell \mid \ell \text{ is a line in } R^3\}$
- Cost function:  $dist(p, \ell) = \min_{x \in \ell} \|p - x\|_2$
- Output:  $C \subseteq P$  s. t.  $\forall \ell \in Q: \max_{p \in P} dist(p, \ell) - \max_{c \in C} dist(c, \ell) \leq \epsilon \cdot \max_{p \in P} dist(p, \ell)$





# Coreset for 1-Line in $R^3$



Similar to the problem in  $R^2$ , there is a line  $\ell''$  that passes through 2 points of the data and is a 4-approx. to the optimal line  $\ell^*$ .

$$\text{dist}(p, \ell'') \leq 4 \cdot \text{dist}(p, \ell^*)$$

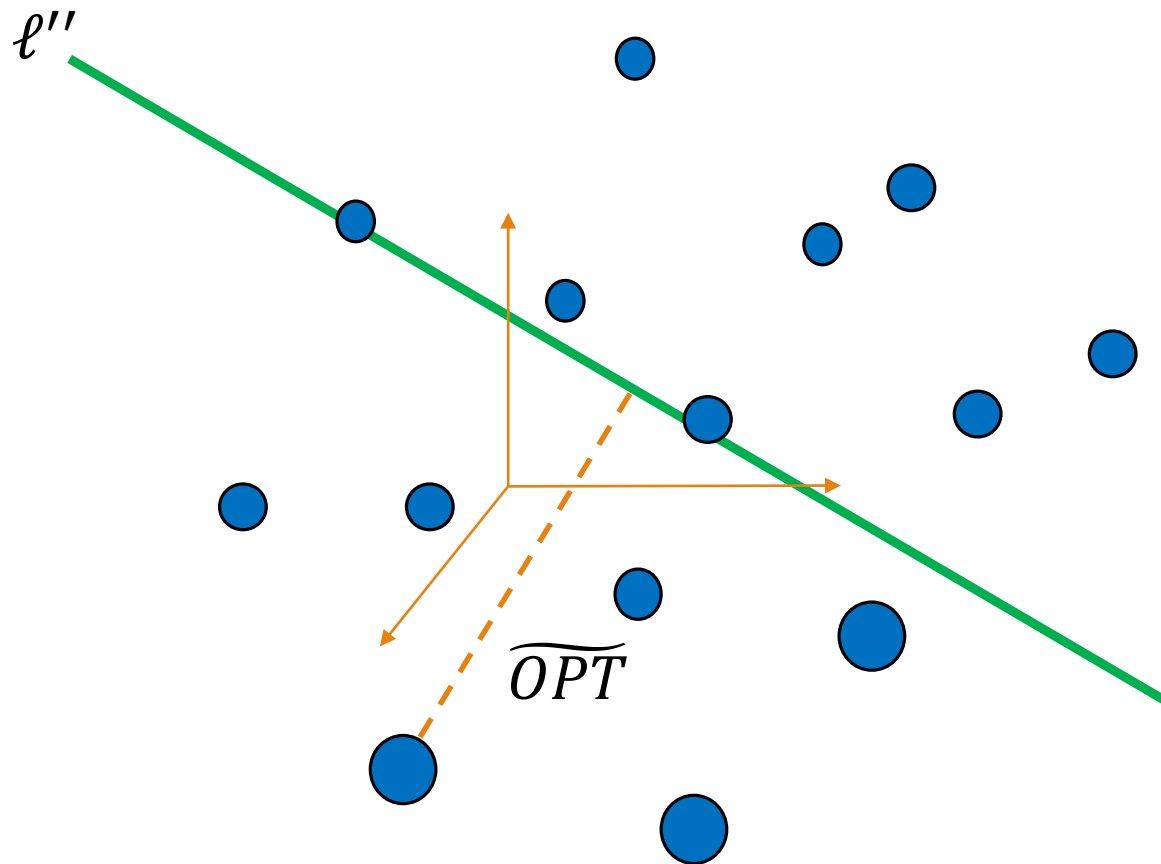
$\ell^*$  is the line that minimizes  $\max_{p \in P} \text{dist}(p, \ell)$

$\ell'$  is the translation of  $\ell^*$  to  $\ell^*$ 's closest point  $p'$

$\ell''$  is the rotation of  $\ell'$  around  $p'$  to  $\ell''$ 's closest point  $p''$

# Coreset for 1-Line in $R^3$

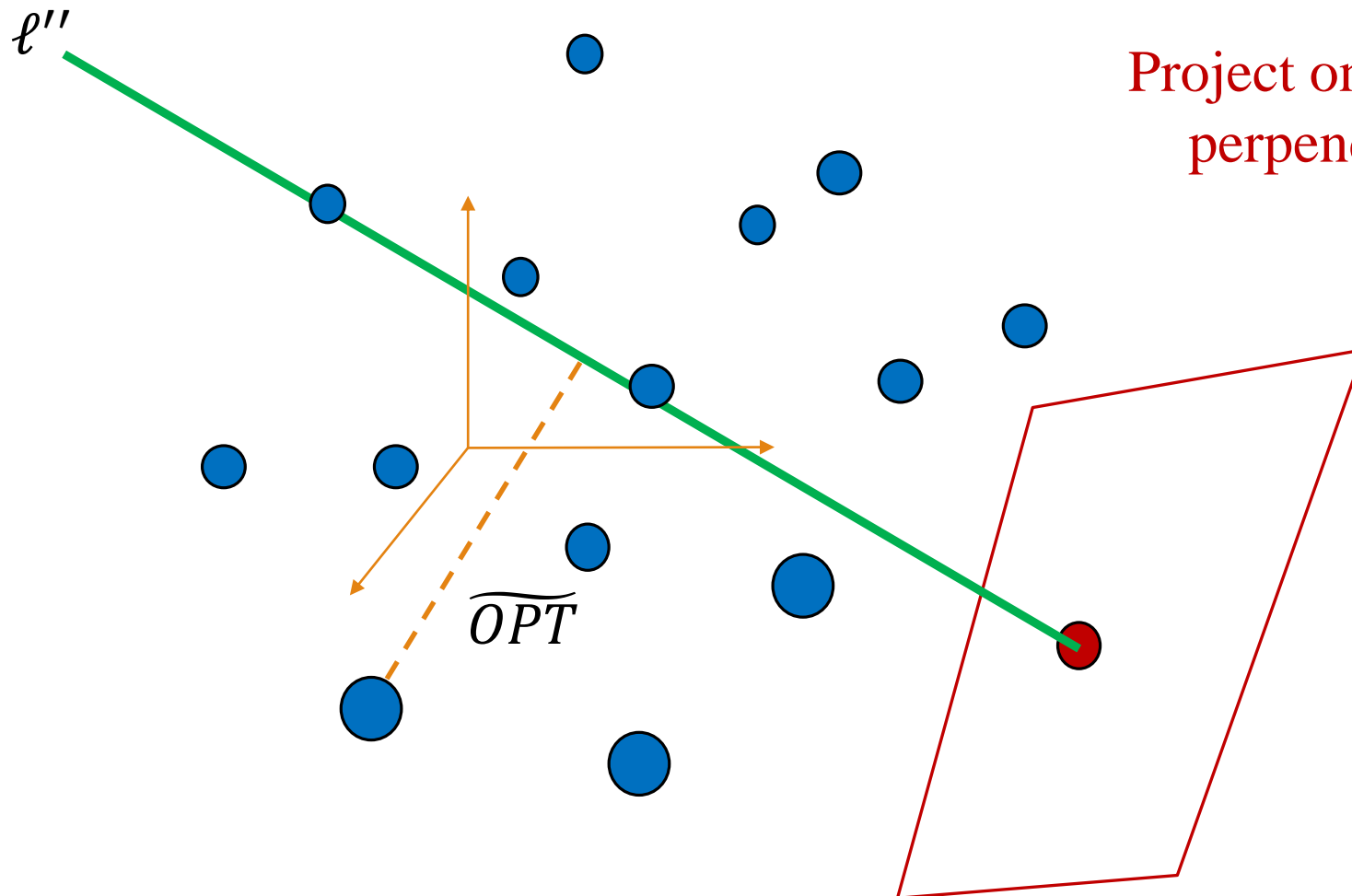
Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$



# Coreset for 1-Line in $R^3$

Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$

Project onto the plane  $\pi$   
perpendicular to  $\ell''$

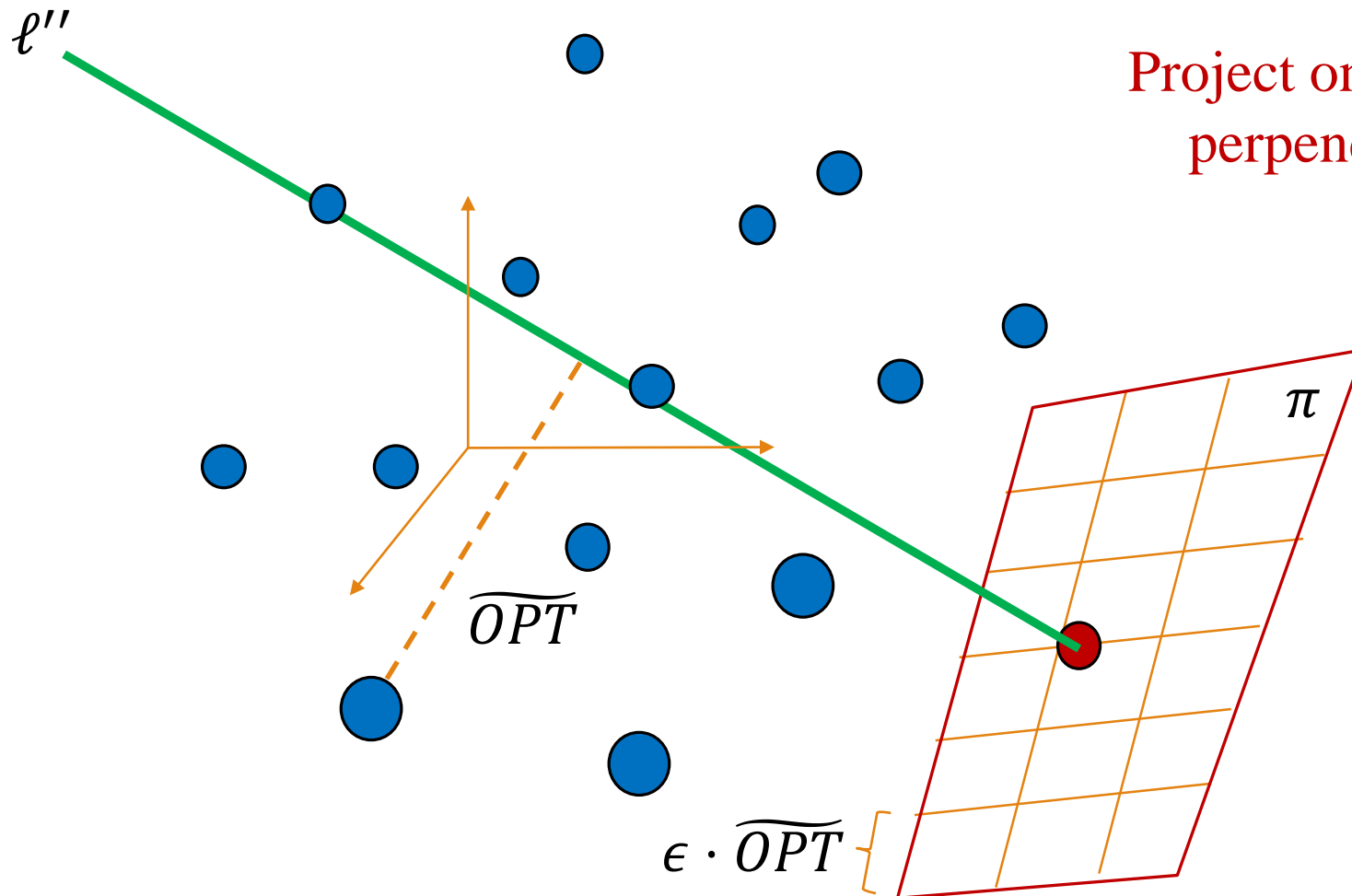


# Coreset for 1-Line in $R^3$

Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$

Project onto the plane  $\pi$   
perpendicular to  $\ell''$

Build a grid with  
distances  $\epsilon \cdot \overline{OPT}$



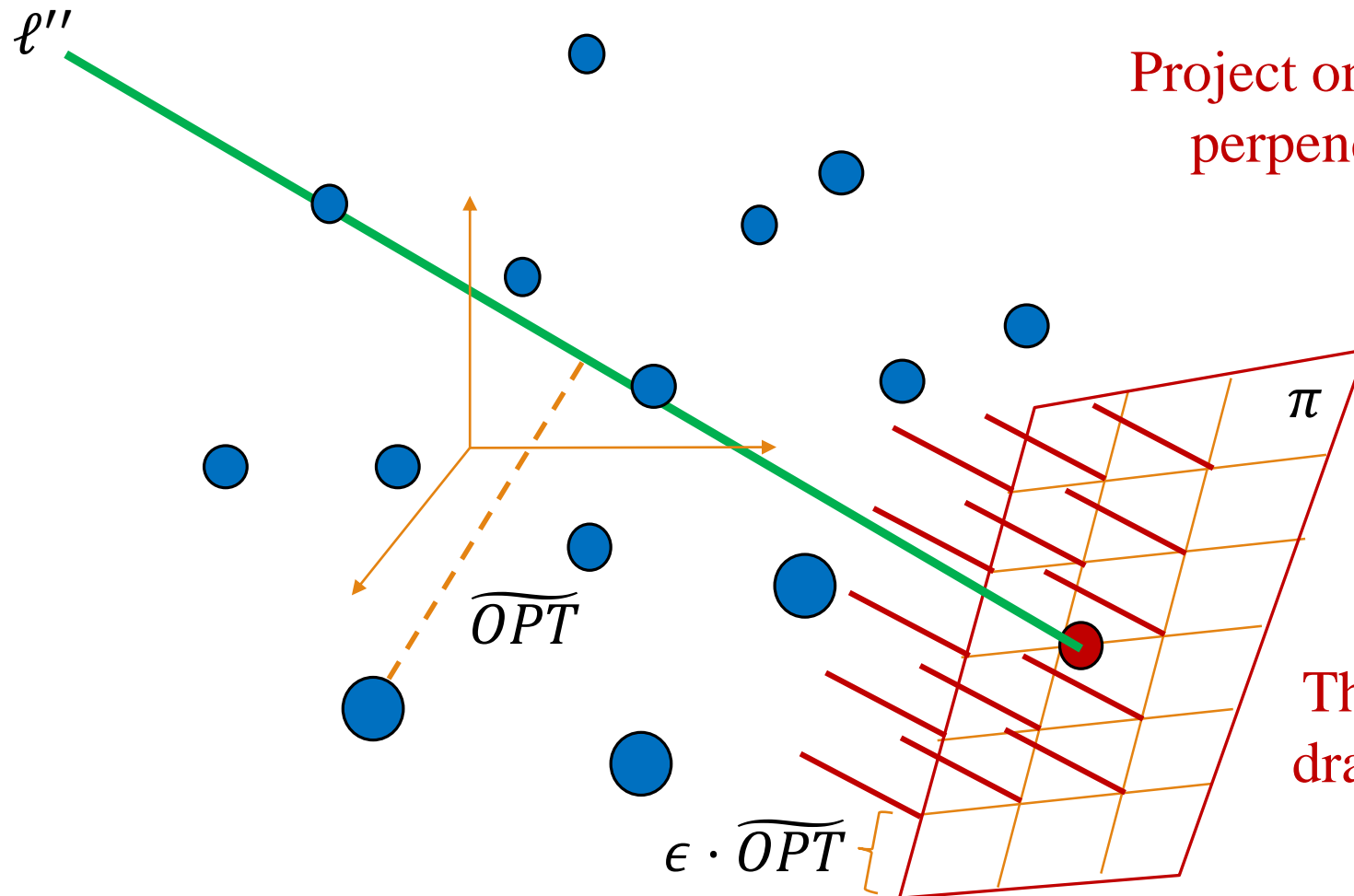
# Coreset for 1-Line in $R^3$

Find  $\ell''$  by exhaustive search  
over every pair of points.  
 $O(n^2)$

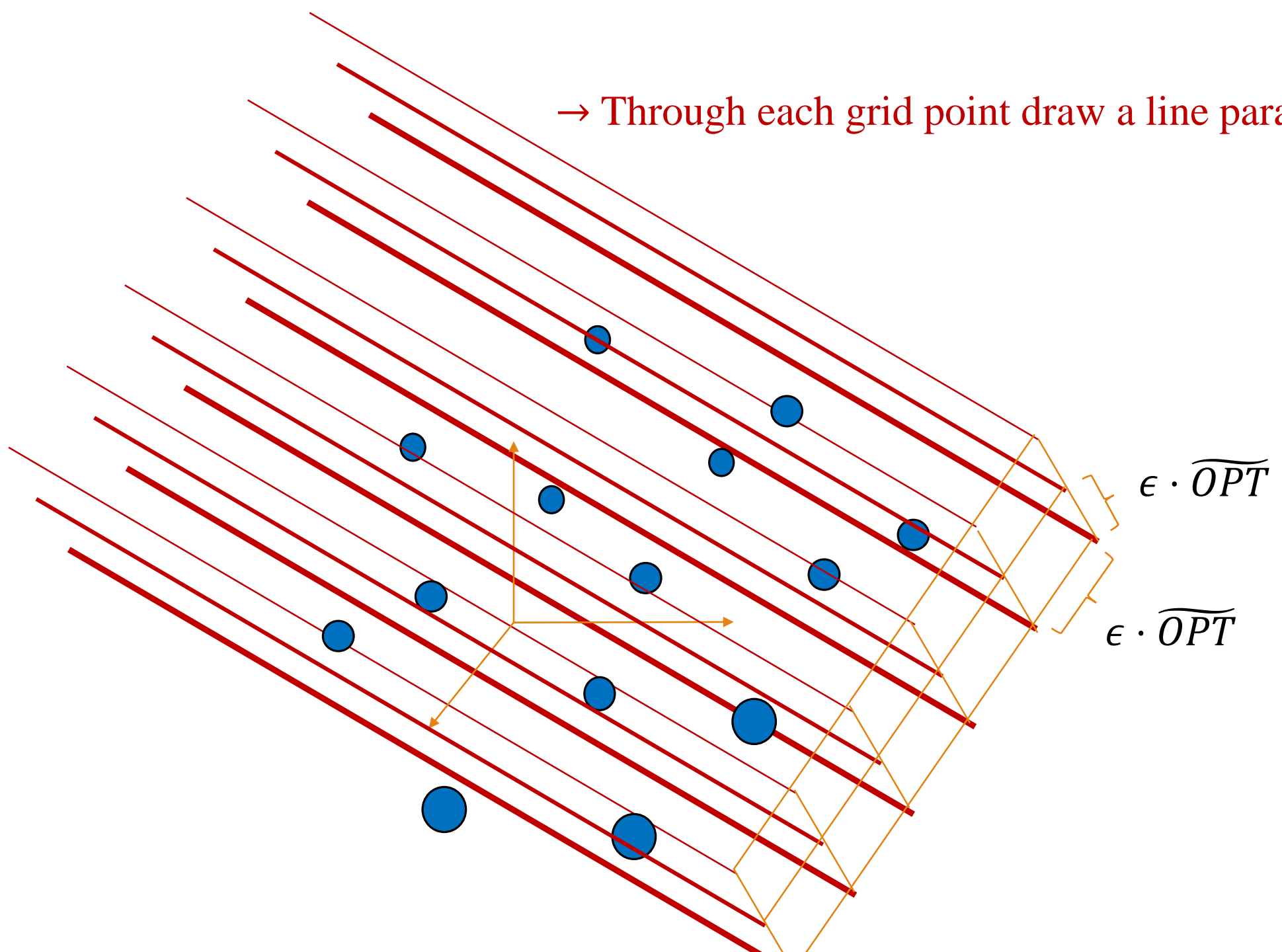
Project onto the plane  $\pi$   
perpendicular to  $\ell''$

Build a grid with  
distances  $\epsilon \cdot \overline{OPT}$

Through each grid point  
draw a line parallel to  $\ell''$

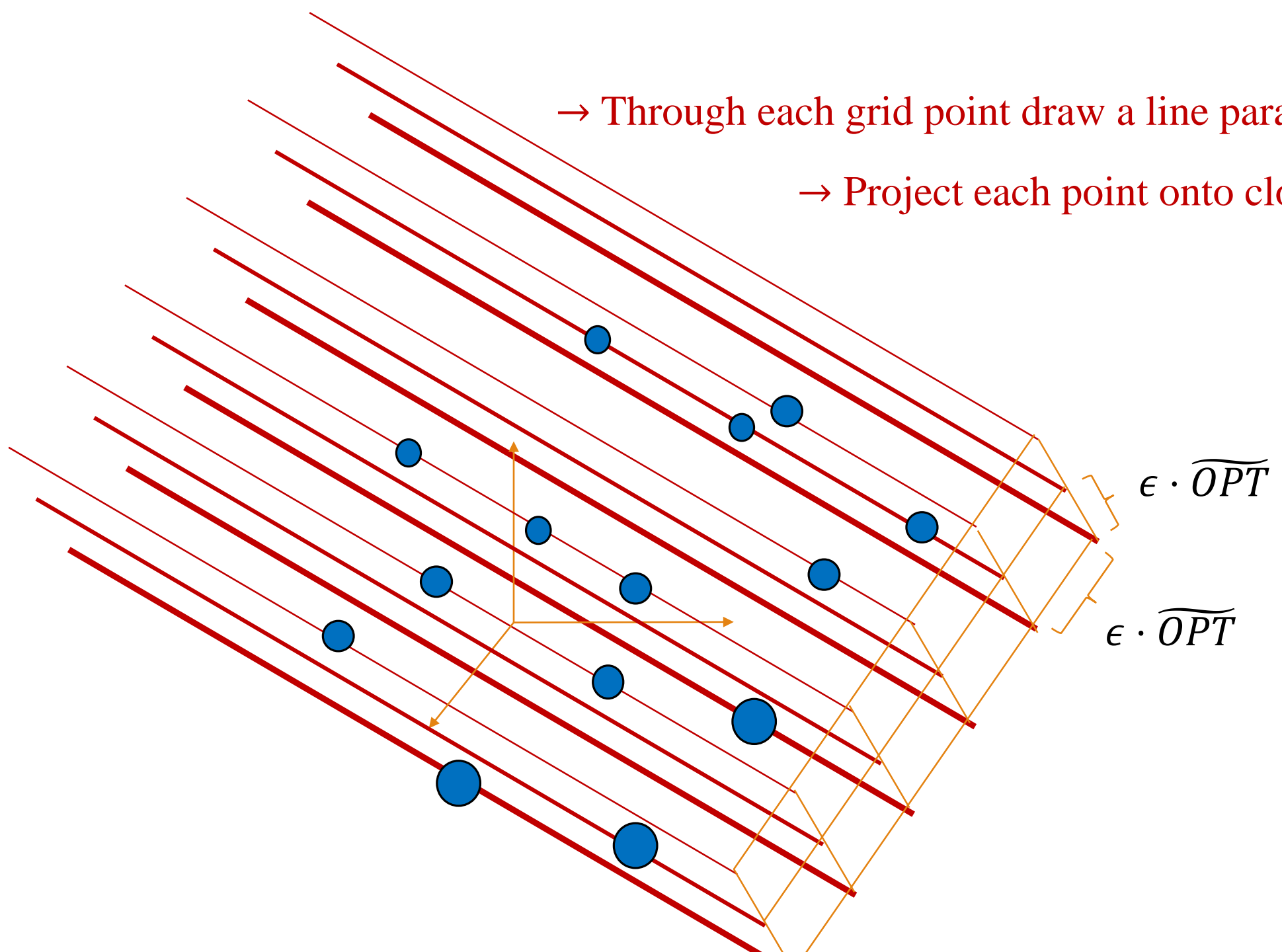


→ Through each grid point draw a line parallel to  $\ell''$



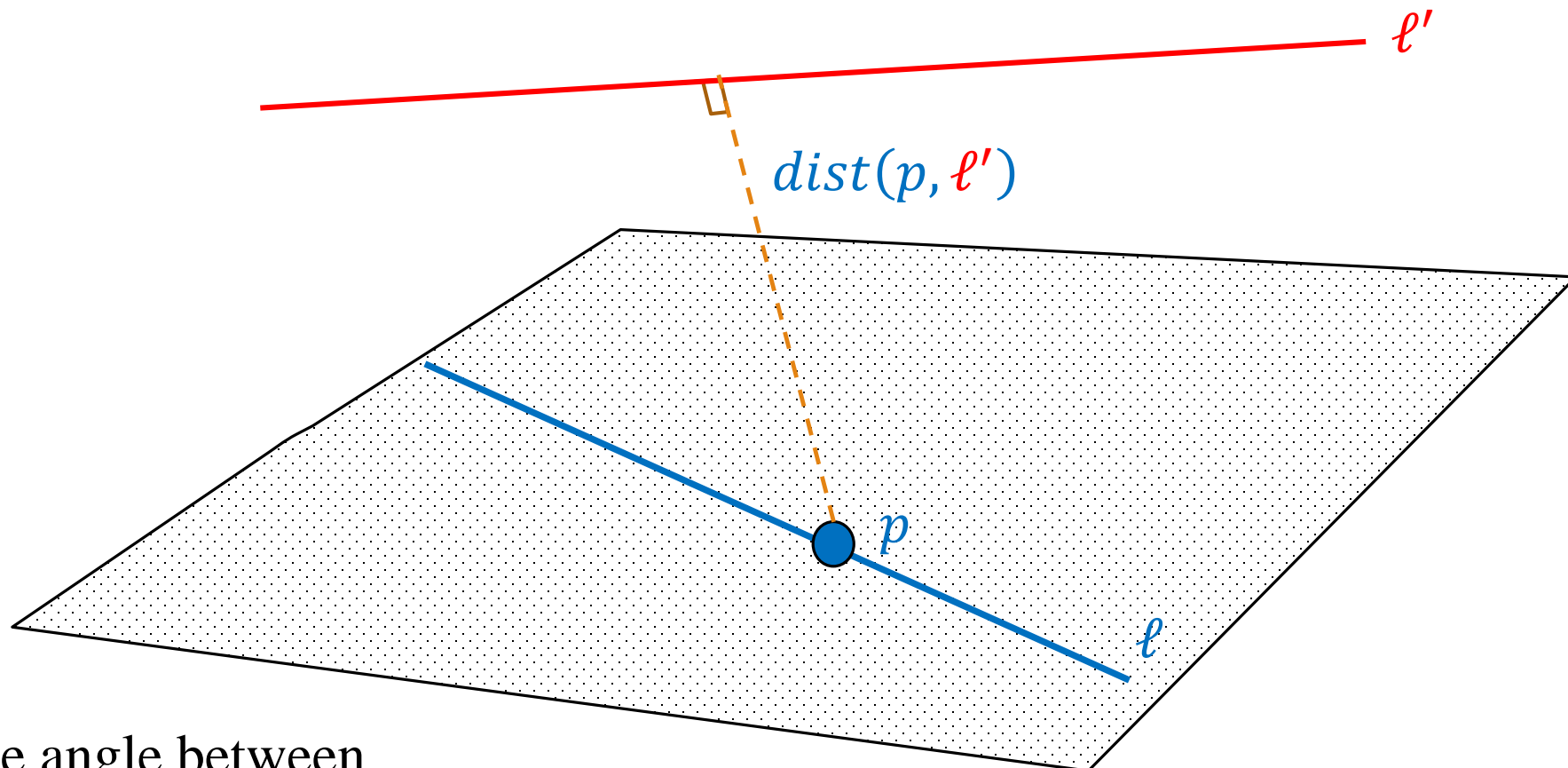
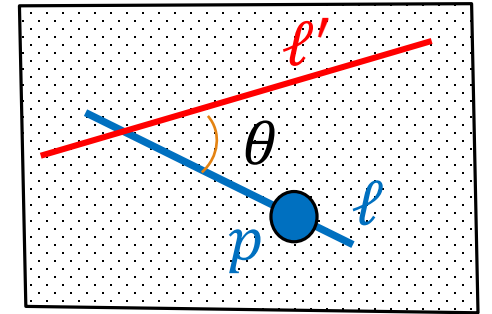
→ Through each grid point draw a line parallel to  $\ell''$

→ Project each point onto closest line



# Distance between $p \in \ell$ and $\ell'$

Top view

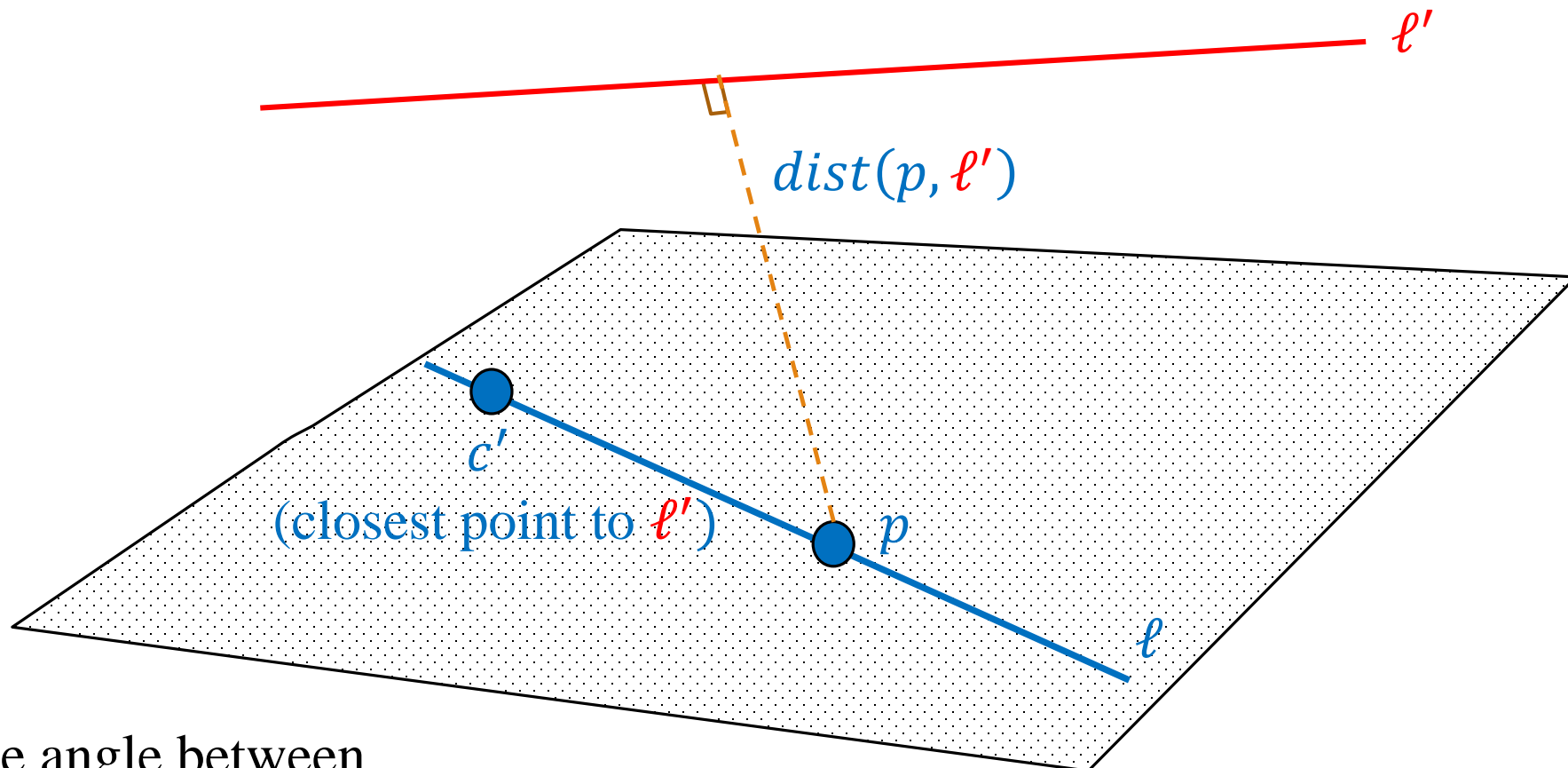
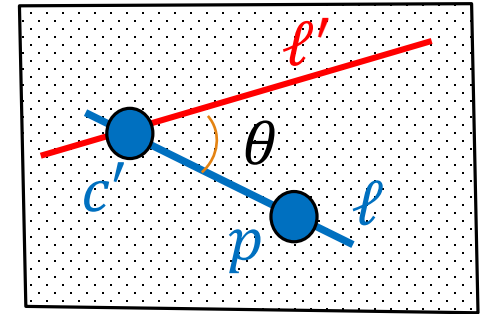


$\theta$  is the angle between  
the line directions



# Distance between $p \in \ell$ and $\ell'$

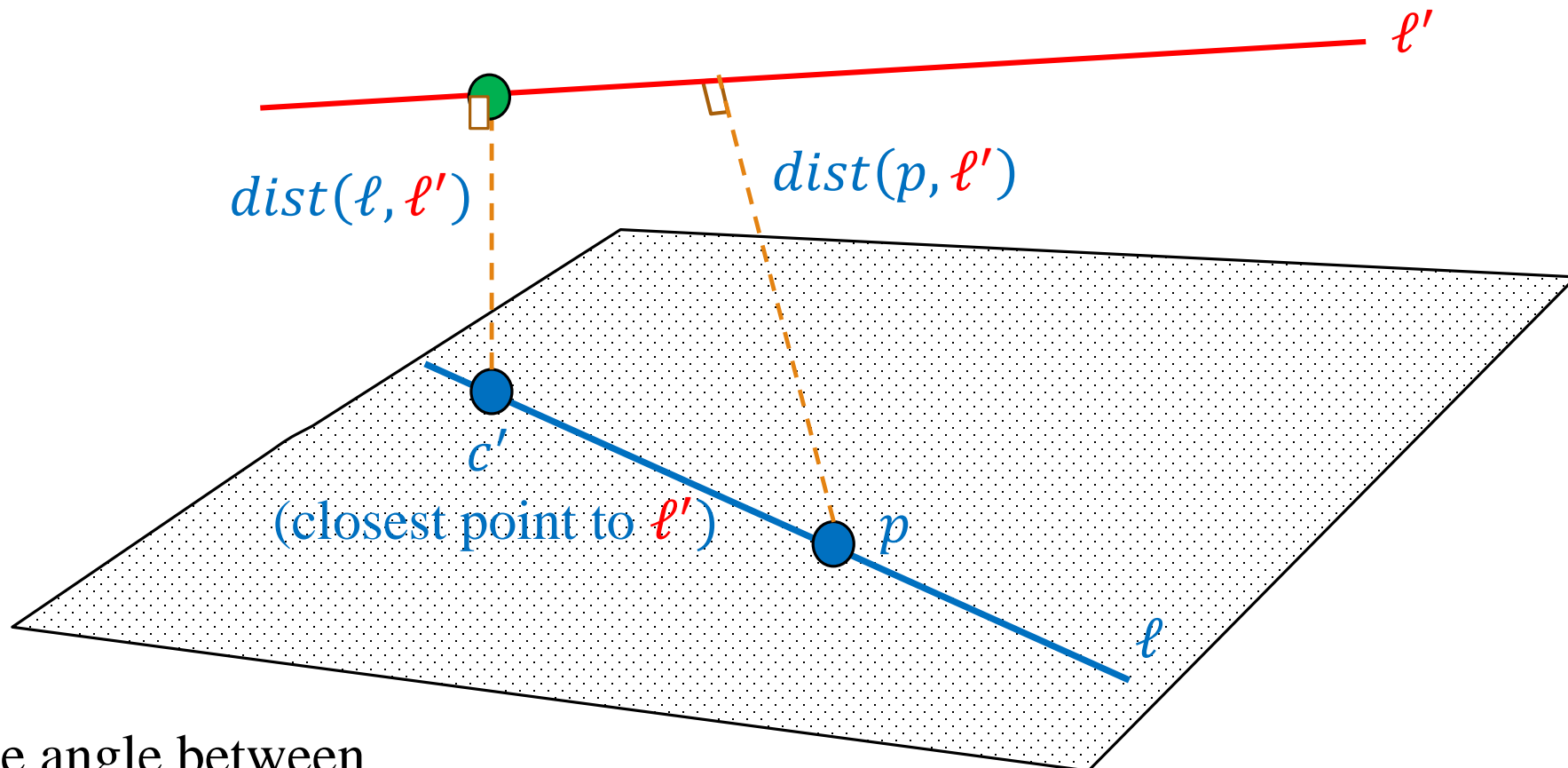
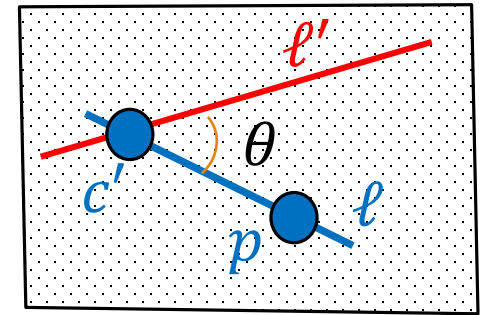
Top view



$\theta$  is the angle between  
the line directions

# Distance between $p \in \ell$ and $\ell'$

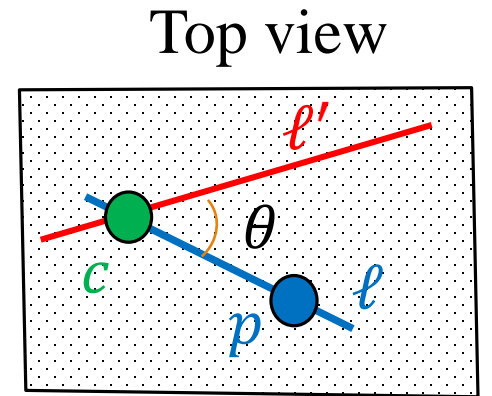
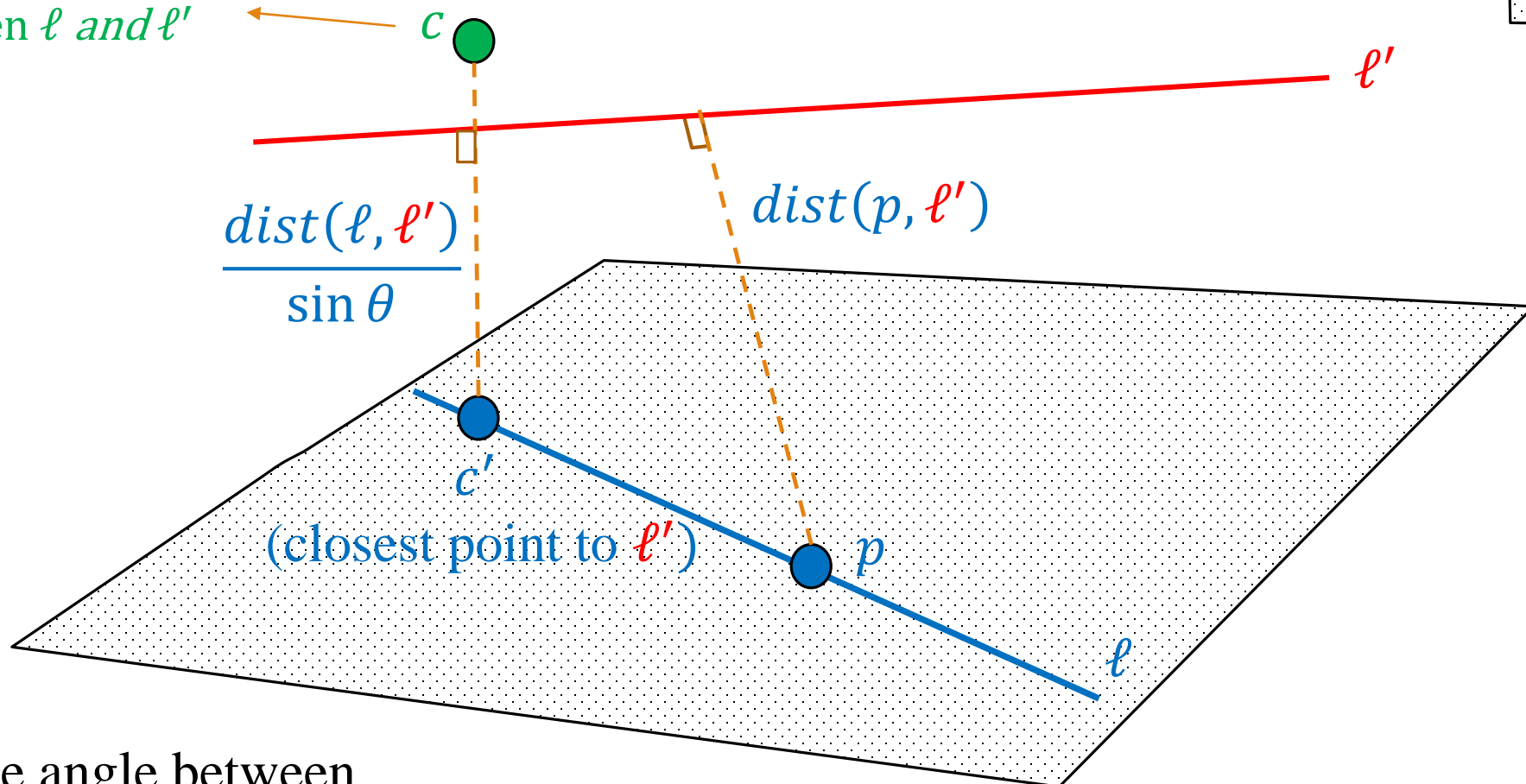
Top view



$\theta$  is the angle between  
the line directions

# Distance between $p \in \ell$ and $\ell'$

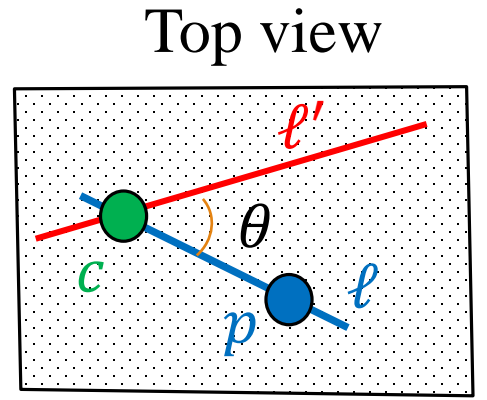
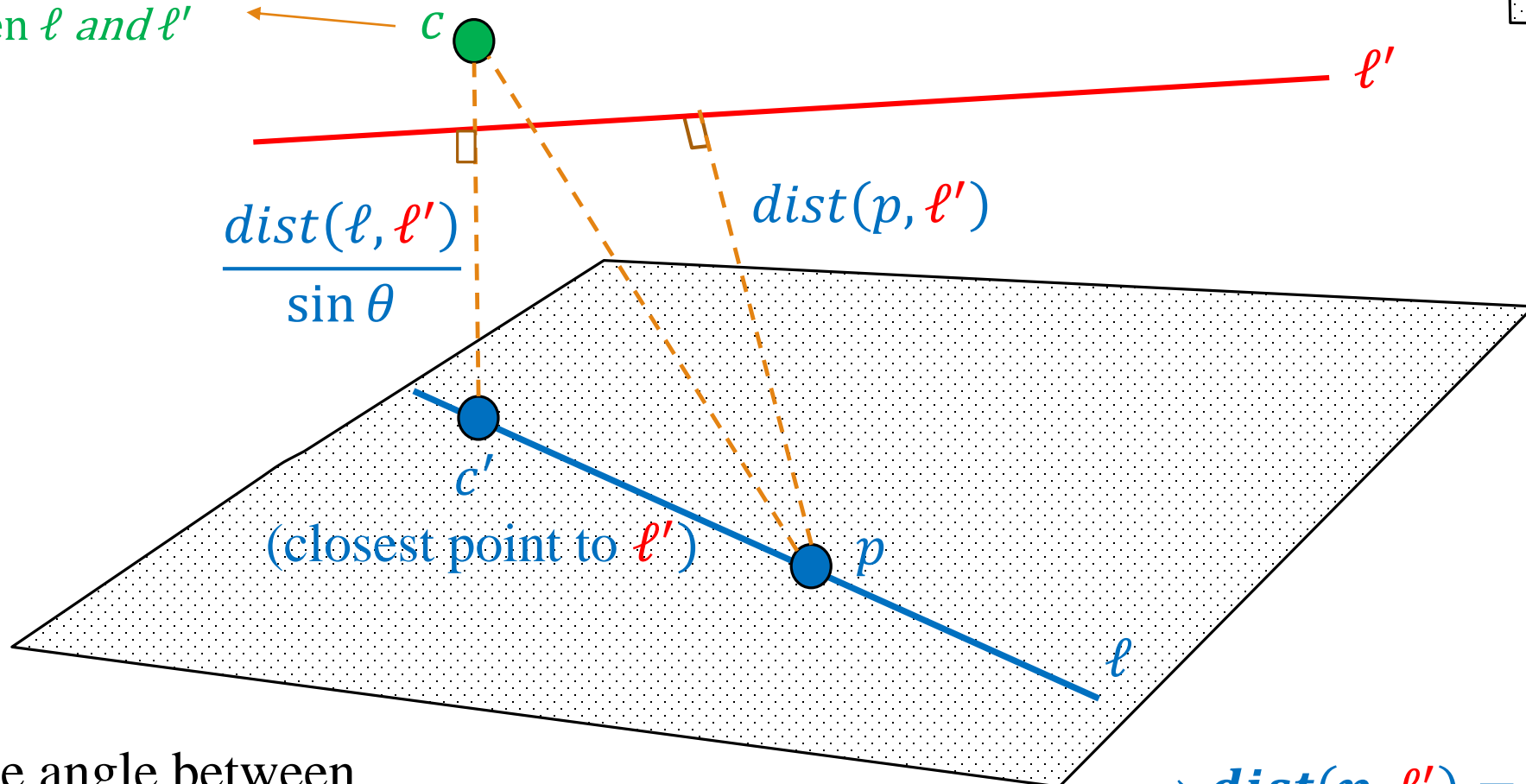
point on the line that spans the shortest distance between  $\ell$  and  $\ell'$



$\theta$  is the angle between the line directions

# Distance between $p \in \ell$ and $\ell'$

point on the line that spans the shortest distance between  $\ell$  and  $\ell'$

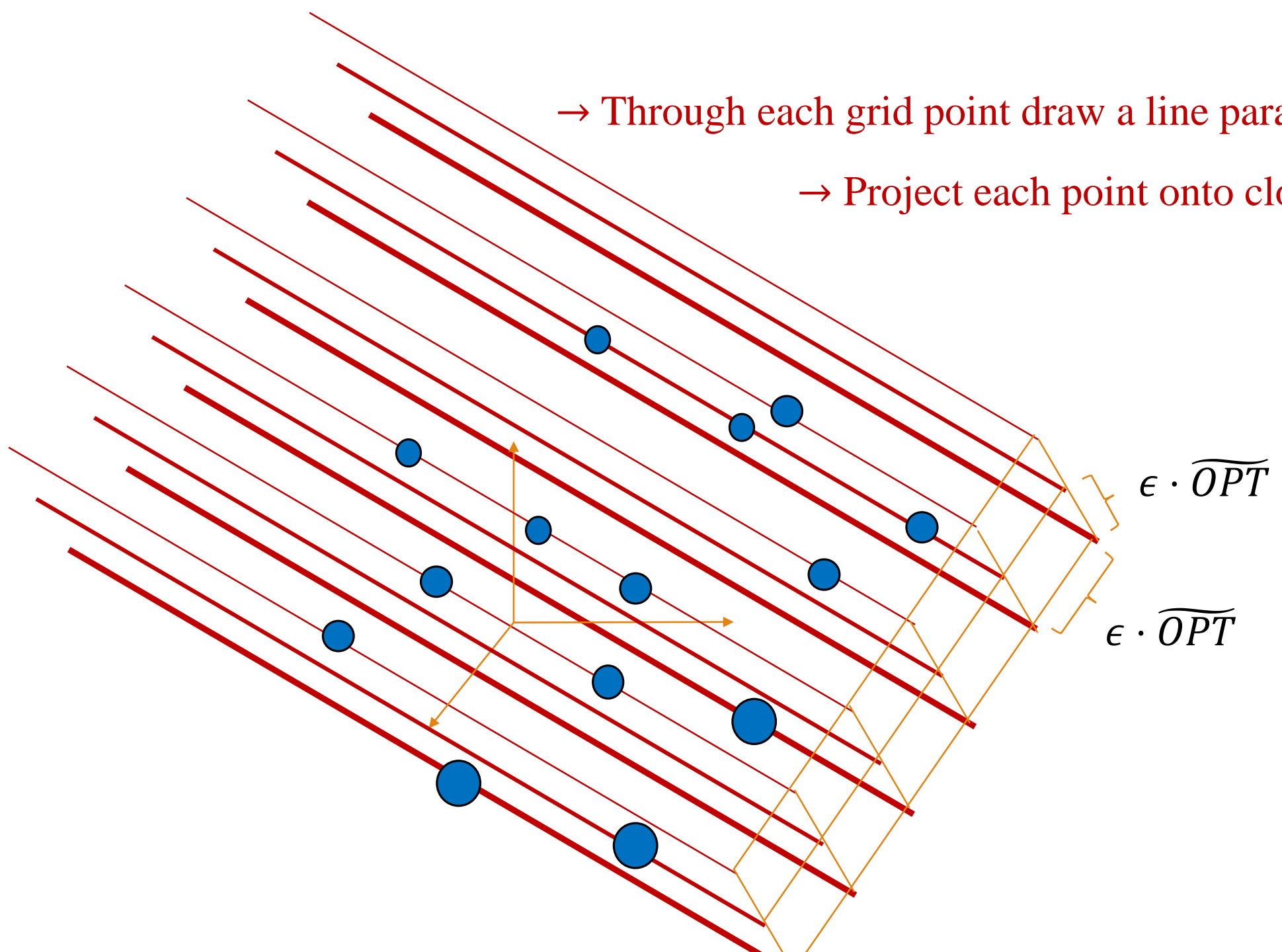


$\theta$  is the angle between the line directions

$$\rightarrow \text{dist}(p, \ell') = \sin \theta \cdot \text{dist}(p, c)$$

→ Through each grid point draw a line parallel to  $\ell''$

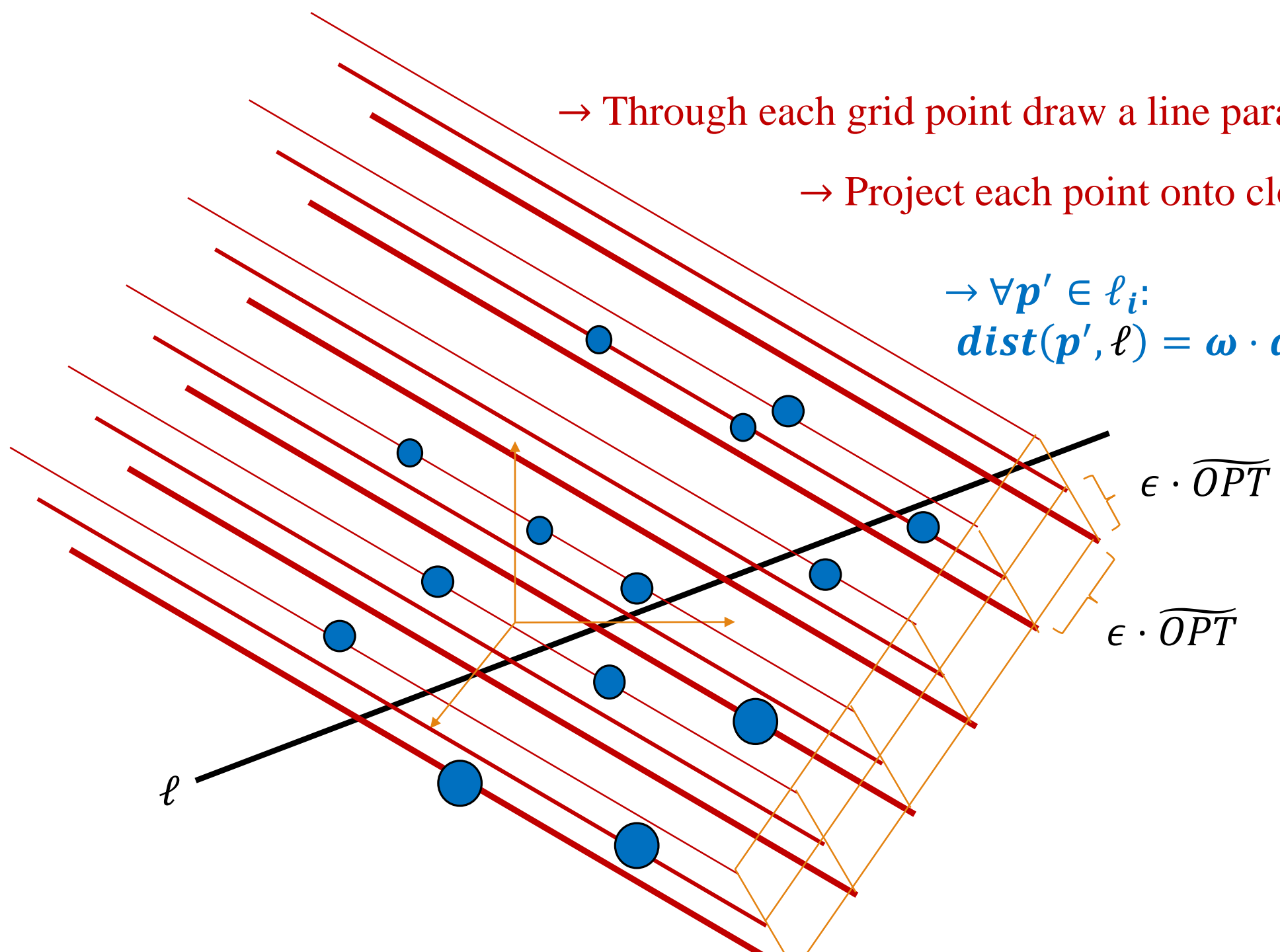
→ Project each point onto closest line



→ Through each grid point draw a line parallel to  $\ell''$

→ Project each point onto closest line

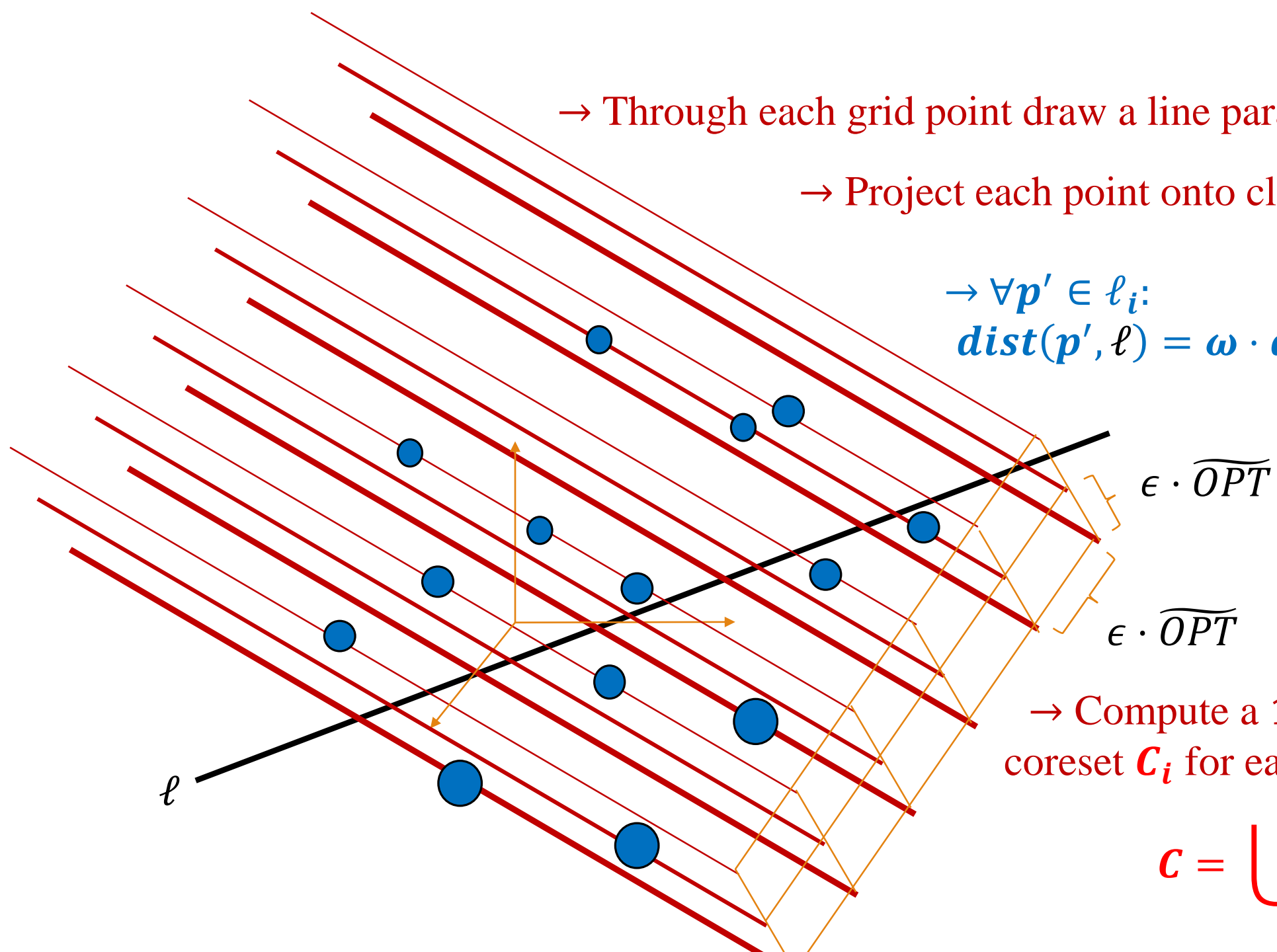
→  $\forall p' \in \ell_i$ :  
 $dist(p', \ell) = \omega \cdot dist(p', c_i)$



→ Through each grid point draw a line parallel to  $\ell''$

→ Project each point onto closest line

→  $\forall p' \in \ell_i$ :  
 $dist(p', \ell) = \omega \cdot dist(p', c_i)$



→ Compute a **1-Center** coresets  $C_i$  for each line  $\ell_i$ !

$$C = \bigcup C_i$$

# Coreset for $j$ -subspace in $R^d$

## Claim 1:

Let  $S$  be an  $r$ -dimensional subspace of  $R^d$  and let  $L$  be an  $(r + j)$ -dimensional subspace of  $R^d$  that contains  $S$ . Let  $V$  be a  $j$ -dimensional subspace of  $R^d$ . Then there is an orthogonal matrix  $U$  such that  $Ux = x$  for every  $x \in S$ , and  $Uc \in L$  for every  $c \in V$ .

## Claim 2:

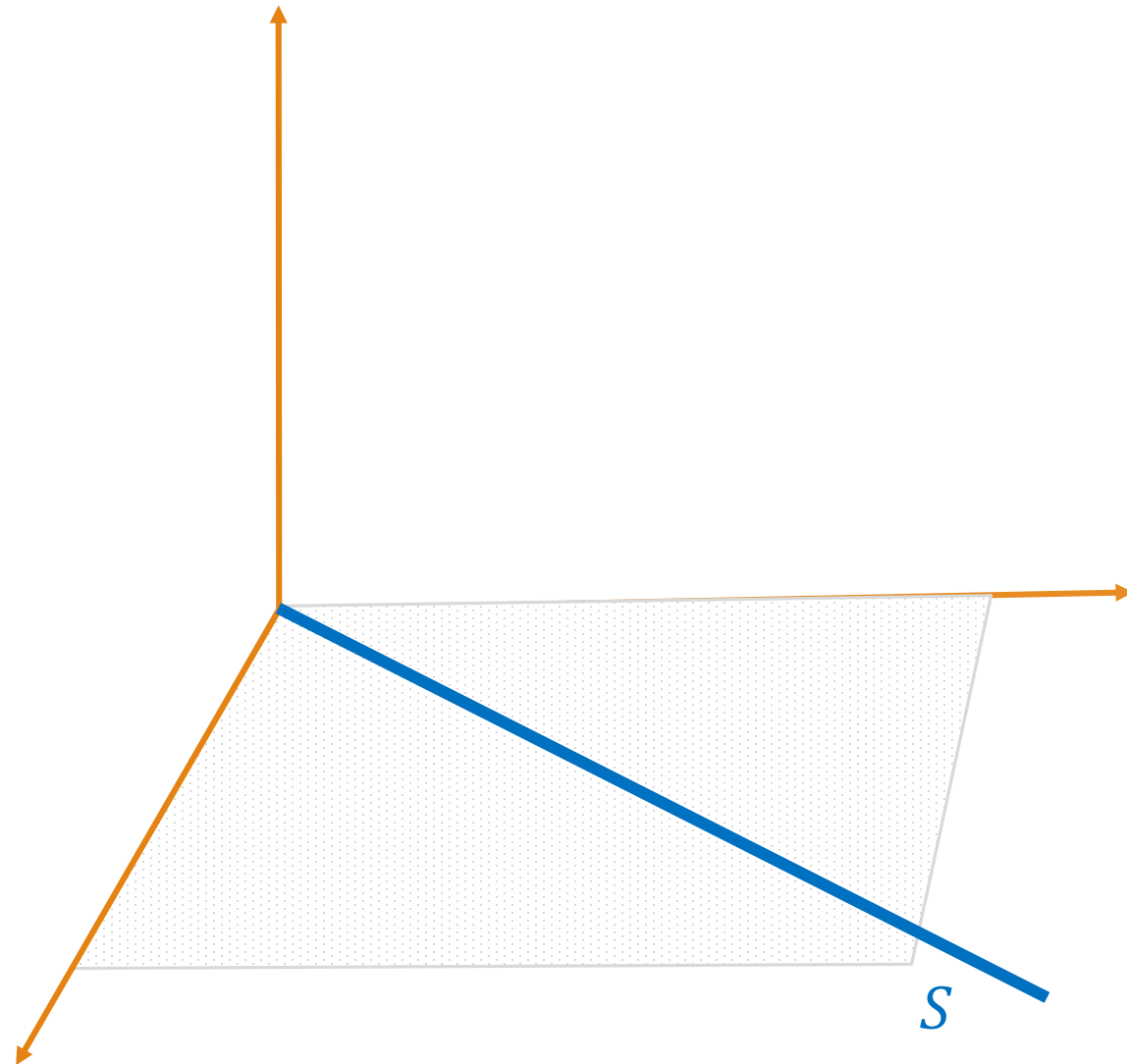
Let  $A \in R^{n \times d}$  be a matrix of rank  $r$  and let  $L$  be an  $(r + j + 1)$ -dimensional subspace of  $R^d$  that contains the row vectors  $(A_{i*})$  for every  $1 \leq i \leq n$ . Then for every affine  $j$ -dimensional subspace  $V$  of  $R^d$  there is a corresponding affine  $j$ -dimensional subspace  $V' \subseteq L$  such that for every  $i \in [n]$  we have

$$\text{dist}(A_{i*}, V) = \text{dist}(A_{i*}, V').$$



# Coreset for $j$ -subspace in $R^d$

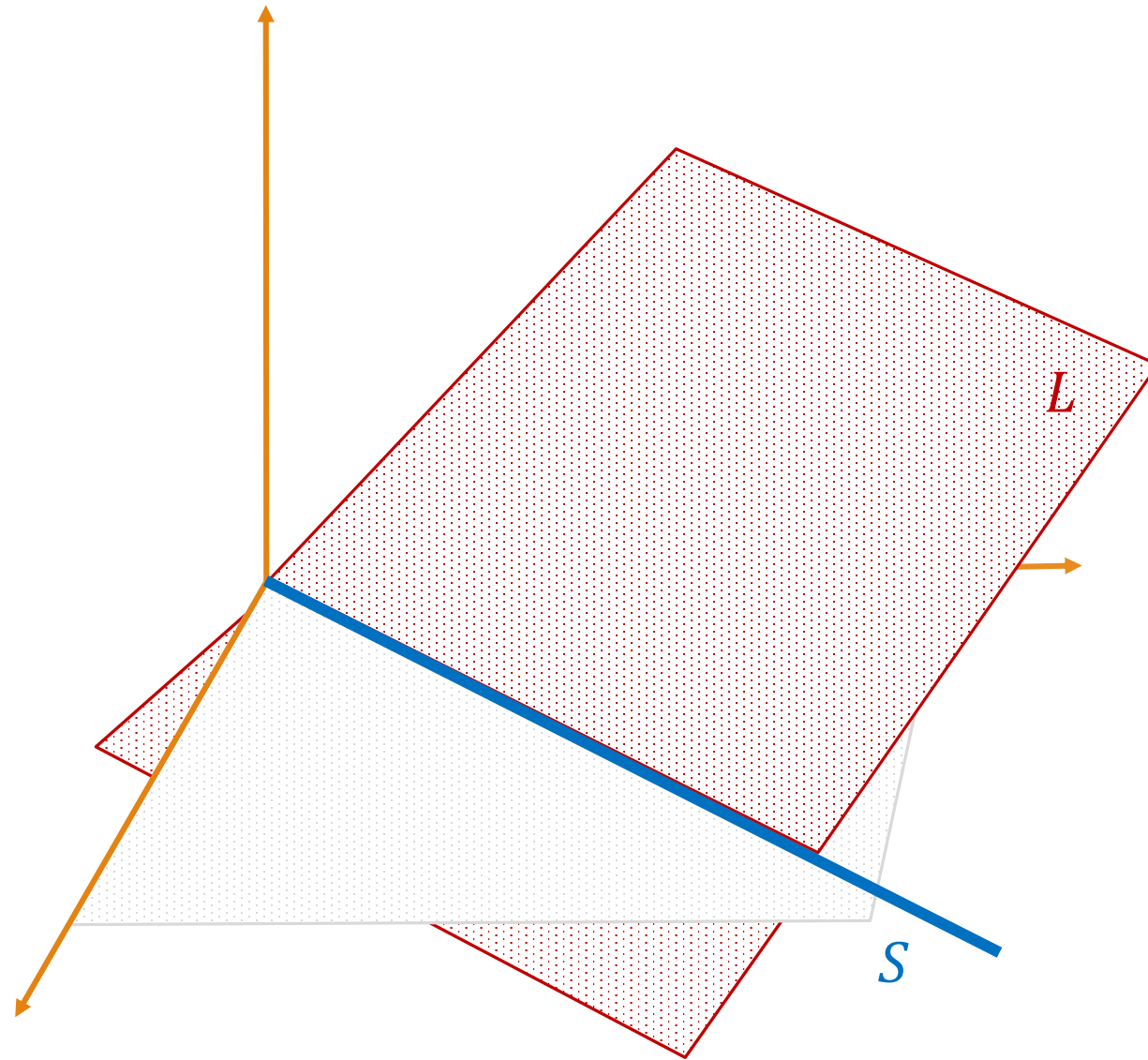
$$r = j = 1$$



# Coreset for $j$ -subspace in $R^d$

$$r = j = 1$$

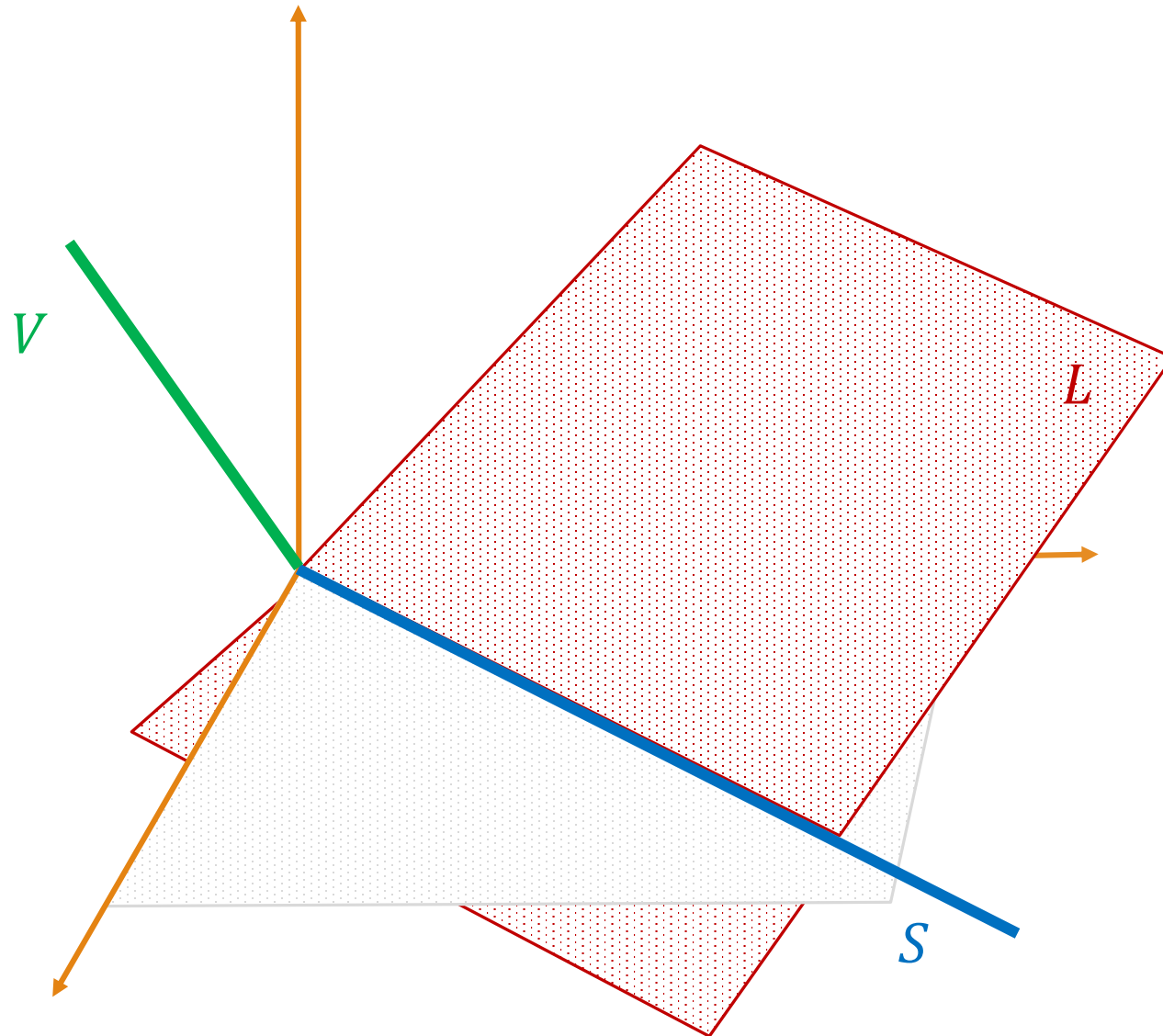
$L$  contains  $S$



# Coreset for $j$ -subspace in $R^d$

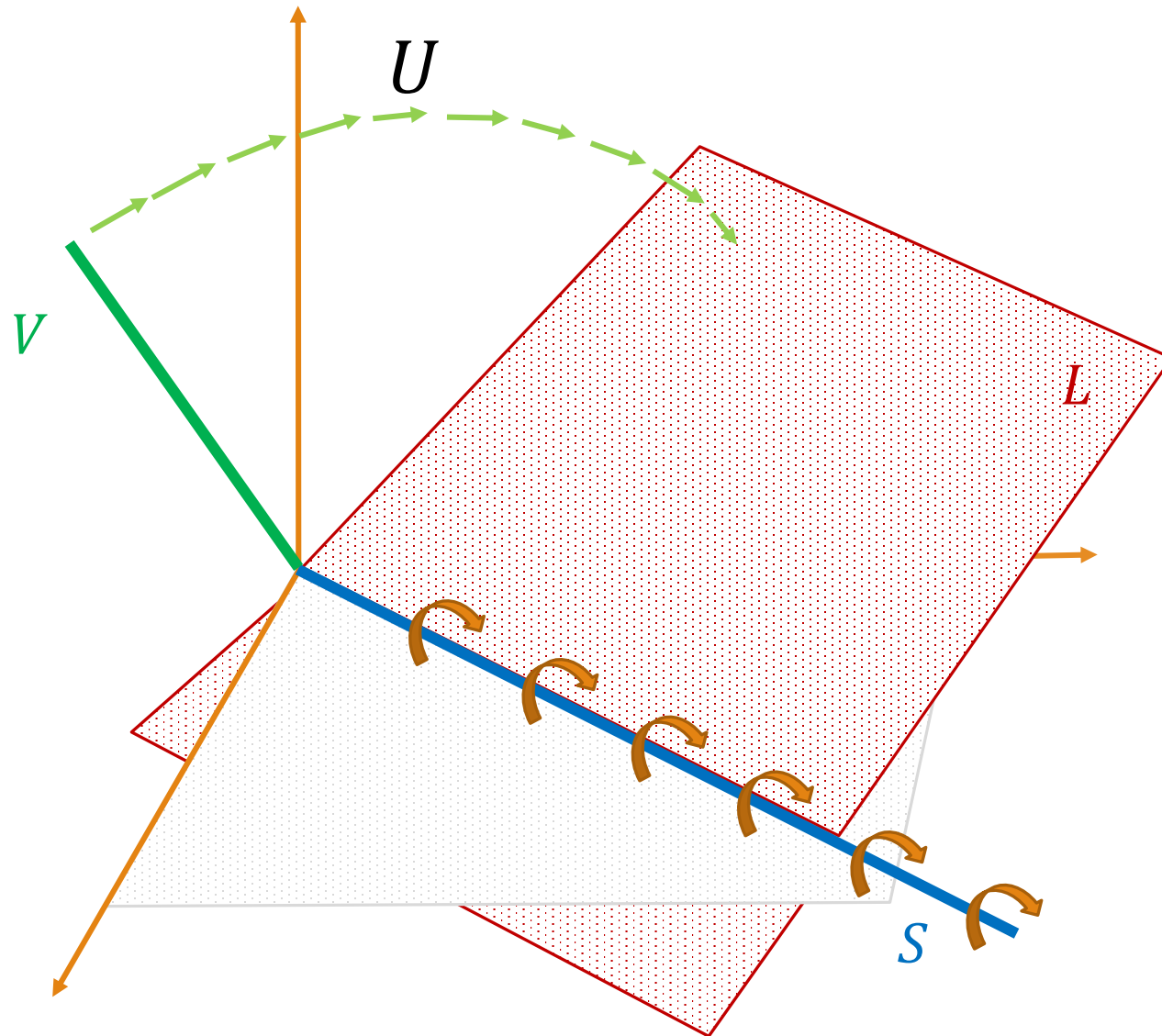
$$r = j = 1$$

$L$  contains  $S$



# Coreset for $j$ -subspace in $R^d$

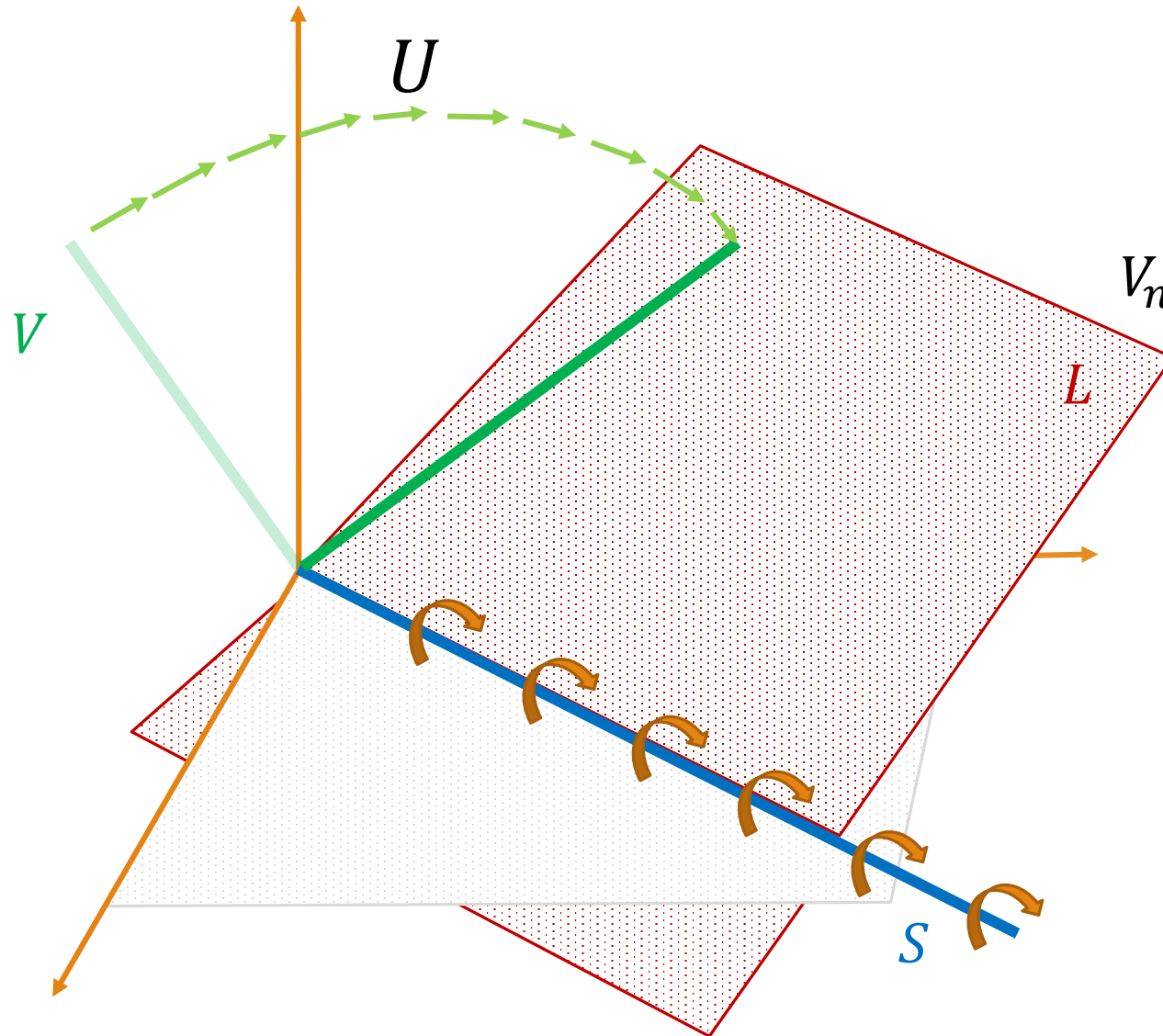
$r = j = 1$



$L$  contains  $S$

# Coreset for $j$ -subspace in $R^d$

$r = j = 1$



$L$  contains  $S$

$L$  contains

$$V_{new} = \{Uc \mid c \in V\}$$

# Coreset for $j$ -subspace in $R^d$

**JSubspaceCoreset( $P, j$ ):**

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .



# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the affine  $d - j$ -subspace that is orthogonal to  $h'$ .

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the affine  $d - j$ -subspace that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the affine  $d - j$ -subspace that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct an affine  $j$ -subspace parallel to  $h'$ .

$$\#JSubspaces = O\left(\left(\frac{2}{\epsilon}\right)^{d-j}\right).$$

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\overline{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the affine  $d - j$ -subspace that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \overline{OPT}$ .
- Through each grid point construct an affine  $j$ -subspace parallel to  $h'$ .

$$\#\text{JSubspaces} = O\left(\left(\frac{2}{\epsilon}\right)^{d-j}\right).$$

- Compute the projection  $p'$  of each point  $p \in P$  onto its closest affine  $j$ -subspace  $h_p$ .

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the affine  $d - j$ -subspace that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct an affine  $j$ -subspace parallel to  $h'$ .

$$\#\text{JSubspaces} = O\left(\left(\frac{2}{\epsilon}\right)^{d-j}\right).$$

- Compute the projection  $p'$  of each point  $p \in P$  onto its closest affine  $j$ -subspace  $h_p$ .
- $H_p \leftarrow$  an  $R^{d \times j}$  matrix whose columns span  $h_p$  and  $H_p^T H_p = I$ .

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the affine  $d - j$ -subspace that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct an affine  $j$ -subspace parallel to  $h'$ .

$$\#\text{JSubspaces} = O\left(\left(\frac{2}{\epsilon}\right)^{d-j}\right).$$

- Compute the projection  $p'$  of each point  $p \in P$  onto its closest affine  $j$ -subspace  $h_p$ .
- $H_p \leftarrow$  an  $R^{d \times j}$  matrix whose columns span  $h_p$  and  $H_p^T H_p = I$ .
- $P' = \{H_p p' \mid p \in P\}$ .

# Coreset for $j$ -subspace in $R^d$

## JSubspaceCoreset( $P, j$ ):

- $h' \leftarrow$  an  $\alpha$ -approximation for the affine  $j$ -subspace center of  $P$ .
- $\widetilde{OPT} = \max_{p \in P} \text{dist}(p, h')$ .
- $h^\perp \leftarrow$  the affine  $d - j$ -subspace that is orthogonal to  $h'$ .
- Construct a grid on  $h^\perp$  whose cell length is  $\epsilon \cdot \widetilde{OPT}$ .
- Through each grid point construct an affine  $j$ -subspace parallel to  $h'$ .

$$\#\text{JSubspaces} = O\left(\left(\frac{2}{\epsilon}\right)^{d-j}\right).$$

- Compute the projection  $p'$  of each point  $p \in P$  onto its closest affine  $j$ -subspace  $h_p$ .
- $H_p \leftarrow$  an  $R^{d \times j}$  matrix whose columns span  $h_p$  and  $H_p^T H_p = I$ .
- $P' = \{H_p p' \mid p \in P\}$ .
- Call **HyperplaneCoreset**( $P', j$ ).