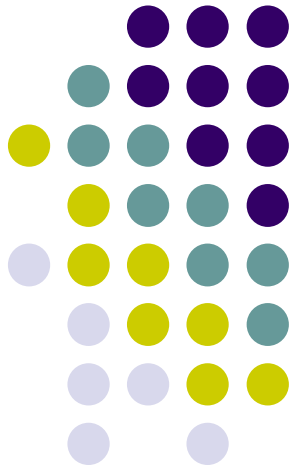


# Fast and Accurate Least-Mean-Squares Solvers

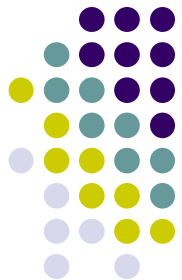


Alaa  
Maalouf

Ibrahim  
Jubran

Dan  
Feldman

The Robotics & Big Data Lab,  
Department of Computer Science,  
University of Haifa,  
Israel



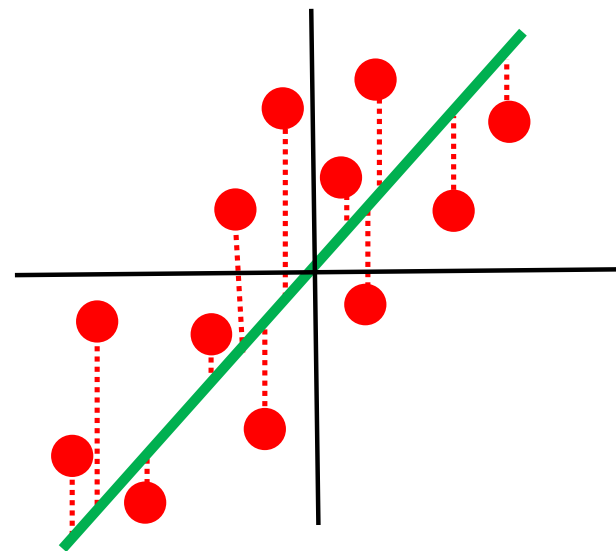
# Least-Mean-Squares Solvers

Input: A matrix

$$A = \begin{pmatrix} | & & | \\ a_1 & \dots & a_n \\ | & & | \end{pmatrix}^T \in \mathbb{R}^{n \times d}$$

and a vector

$$b = (b_1, \dots, b_n)^T \in \mathbb{R}^n$$



Output:

$$x^* \in \mathbb{R}^d$$

that minimizes

$$f(\|Ax - b\|_2^2) + g(x)$$

constraint

over every  $x \in \mathbb{R}^d$ , where  $f: \mathbb{R} \rightarrow \mathbb{R}$  and  $g: \mathbb{R}^d \rightarrow \mathbb{R}$ .



# Least-Mean-Squares Solvers

Input: A matrix

Linear regression:

$$\|Ax - b\|_2^2$$

$$\begin{pmatrix} | \\ a_n \\ | \end{pmatrix}^T$$

$$\in \mathbb{R}^{n \times d}$$

and a vector

$$b = (b_1, \dots, b_n)^T$$

Ridge regression:

$$\|Ax - b\|_2^2 + \alpha \|x\|_2^2$$

Lasso regression:

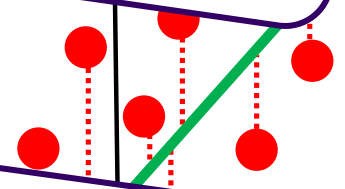
$$\frac{1}{2n} \|Ax - b\|_2^2 + \alpha \|x\|_1$$

Elastic-net regression:

$$\|Ax - b\|_2^2 + \rho \alpha \|x\|_2^2 + \frac{(1 - \rho)}{2} \alpha \|x\|_1$$

over every  $x \in \mathbb{R}^d$ , where

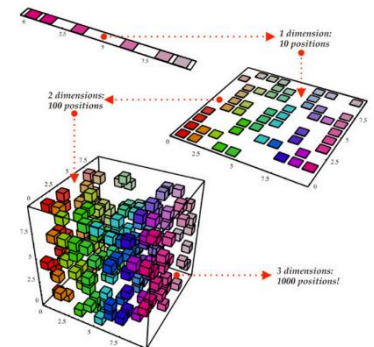
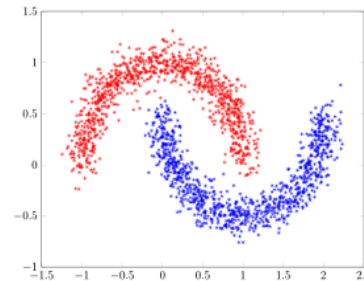
constraint  
 $\mathbb{R}^d \rightarrow \mathbb{R}$ .





# Applications

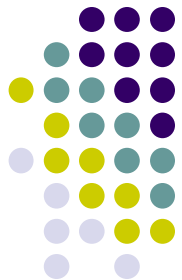
- Data analysis
- Prediction
- Feature selection
- Spectral clustering
- Dimensionality reduction





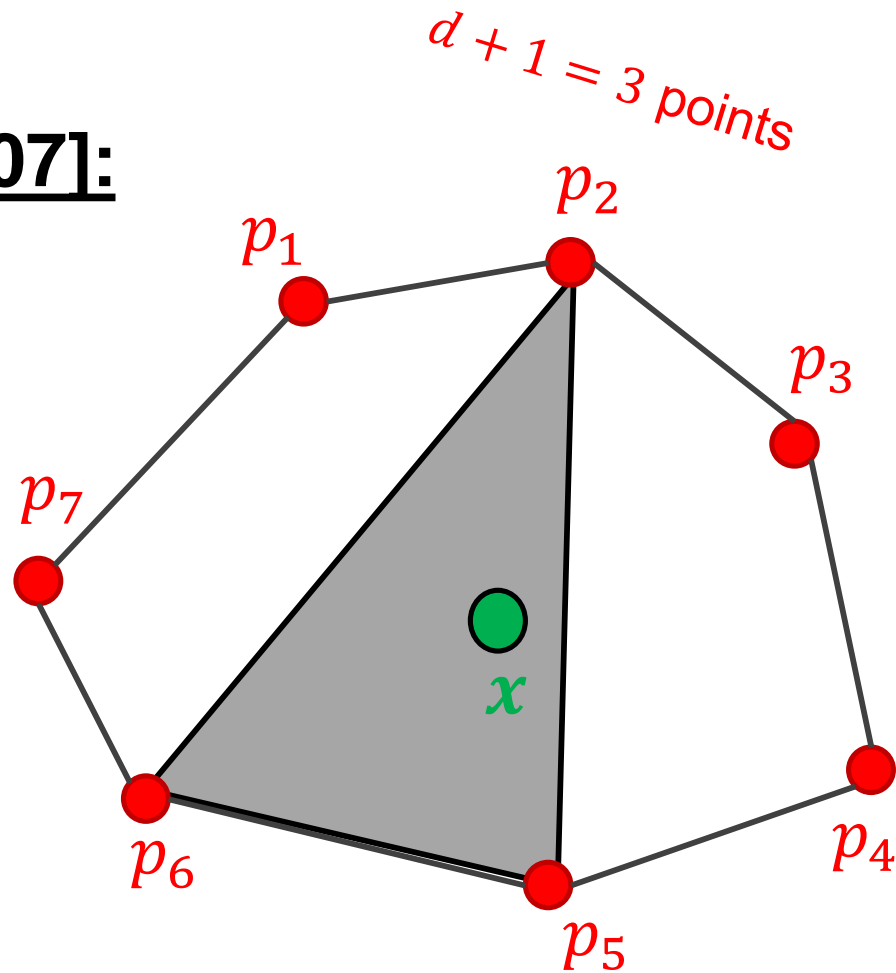
# Practical Considerations

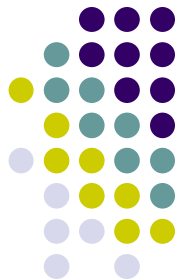
- **Computation time** - The use of cross validations (CV) with a huge number of hyperparameter is time consuming.
- **Space complexity** - The use of SVD or other factorizations on massive input leads to extensive memory usage.
- **Numerical stability** - There are faster yet numerically unstable solutions.



# Main Technique

Theorem [Caratheodory 1907]:





# Main Technique

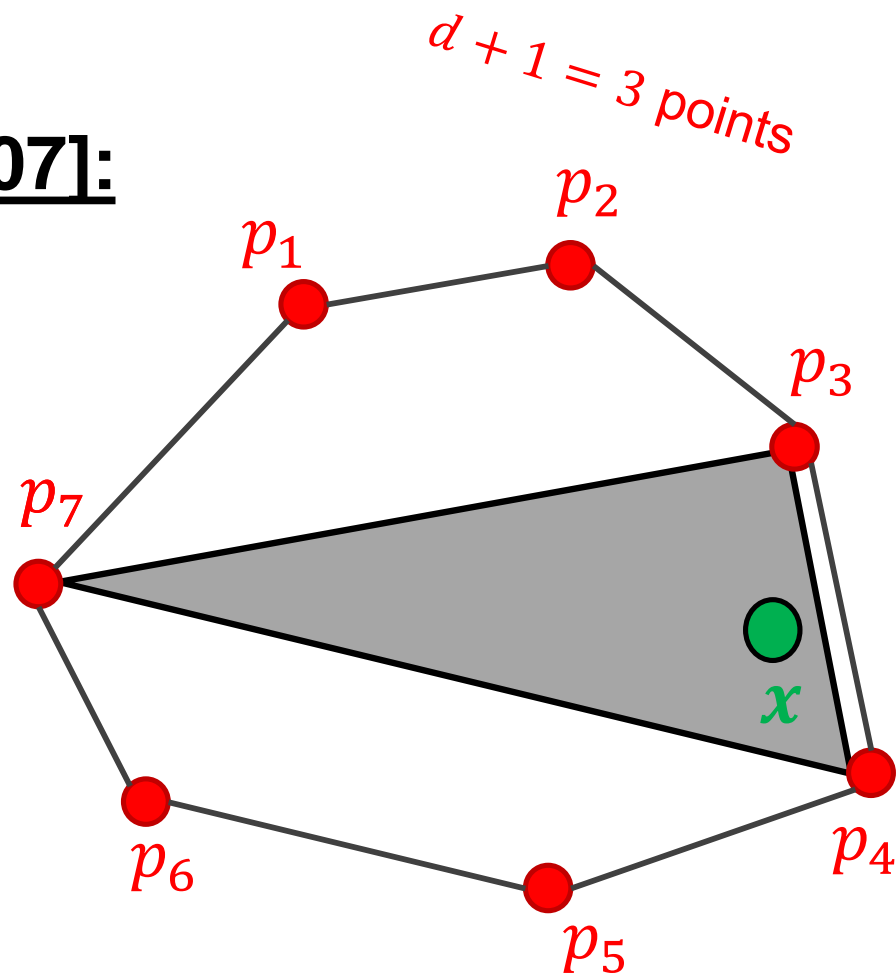
## Theorem [Caratheodory 1907]:

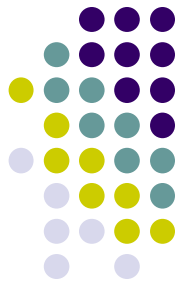
This *Caratheodory set* can be computed in  $O(n^2 d^2)$  time.

**Bottleneck**

↓ **Goal**

$O(nd)$

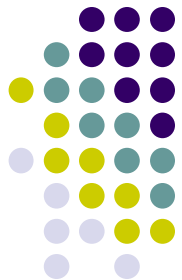




# Our Contribution

- Computing the *Caratheodory set* in  $O(nd + \log(n) d^4)$  time.
- Practically and provably boosting:
  - **time complexity** of LMS solvers.
  - **space complexity** of LMS solvers.
  - **numerical accuracy** of LMS solvers.
- Evaluation on **real-world** data.
- Full **open source code**.





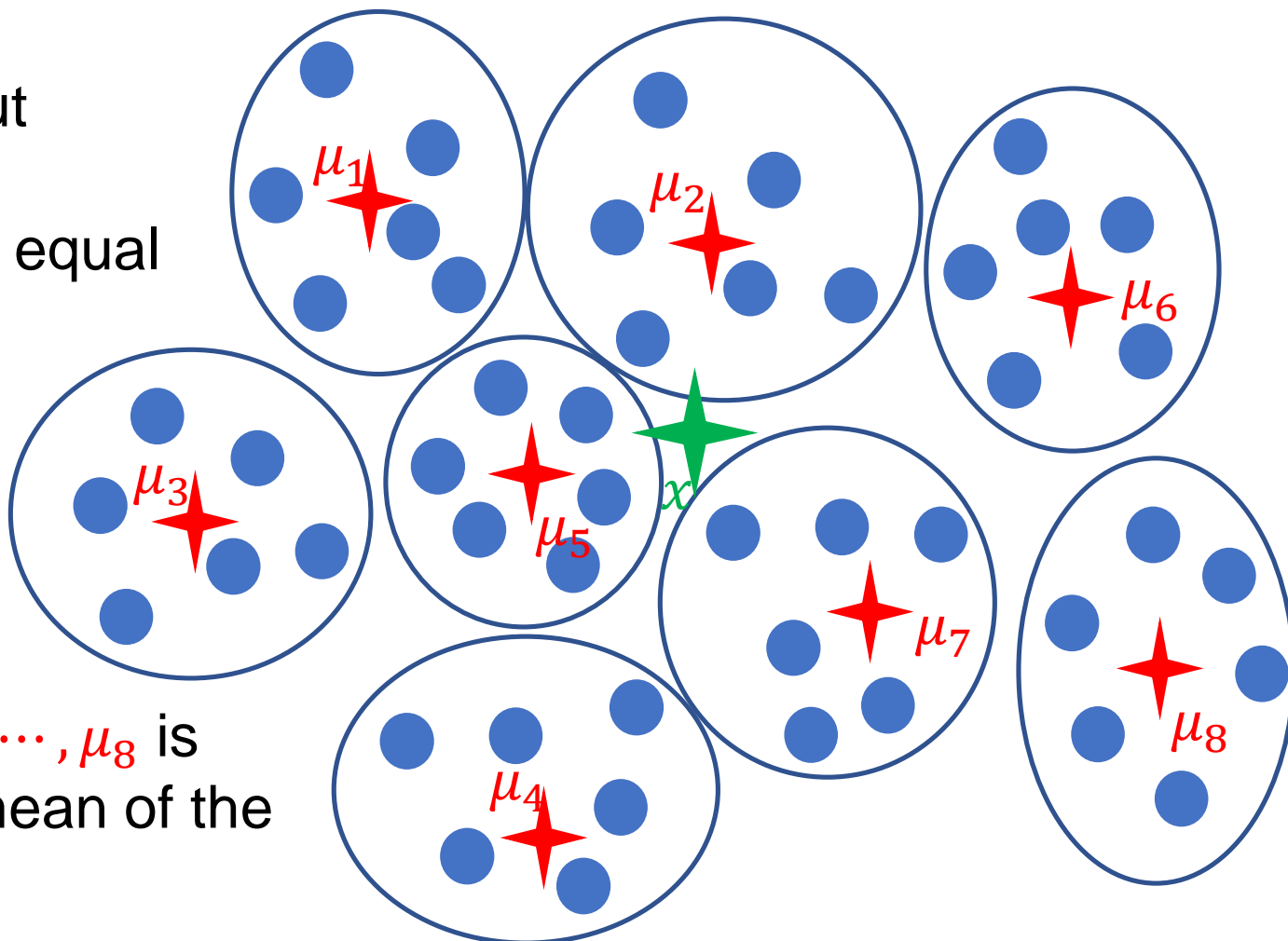
# Caratheodory Booster Illustration

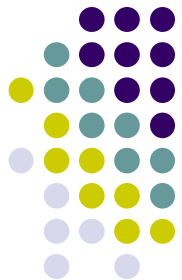
## Step 1:

Partition the input  $d$ -dimensional points into  $k = 8$  equal sized subsets.

$\mu_1, \dots, \mu_8$  are the subset means.

The mean of  $\mu_1, \dots, \mu_8$  is equal to  $x$ , the mean of the input.



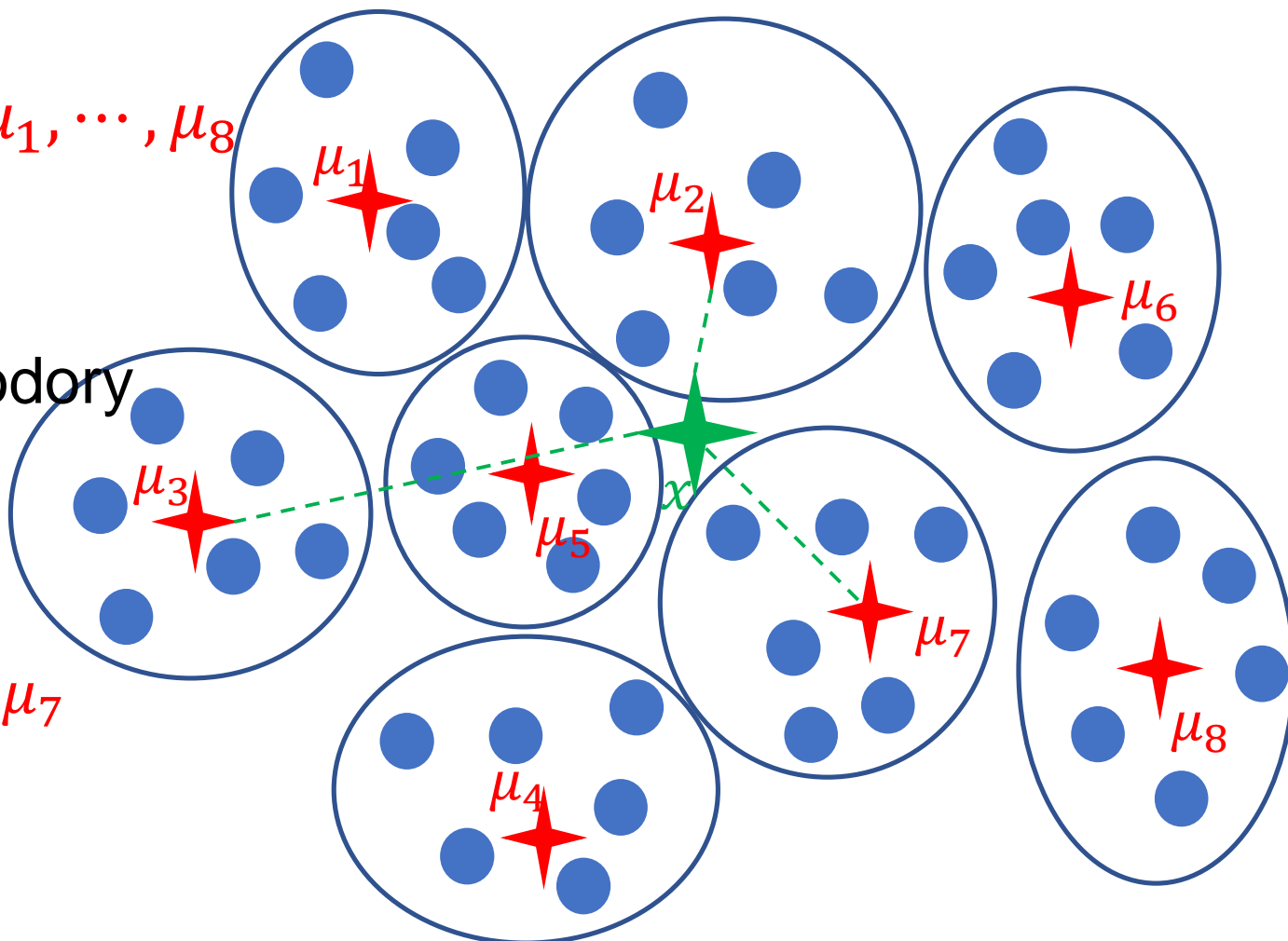


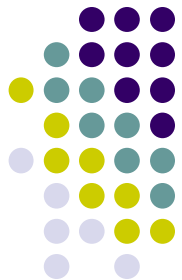
# Caratheodory Booster Illustration

## Step 2:

Consider only  $\mu_1, \dots, \mu_8$  and  $x$ .

Apply Caratheodory to represent  $x$  as a convex combination of a subset  $\mu_2, \mu_3, \mu_7$  of the means.



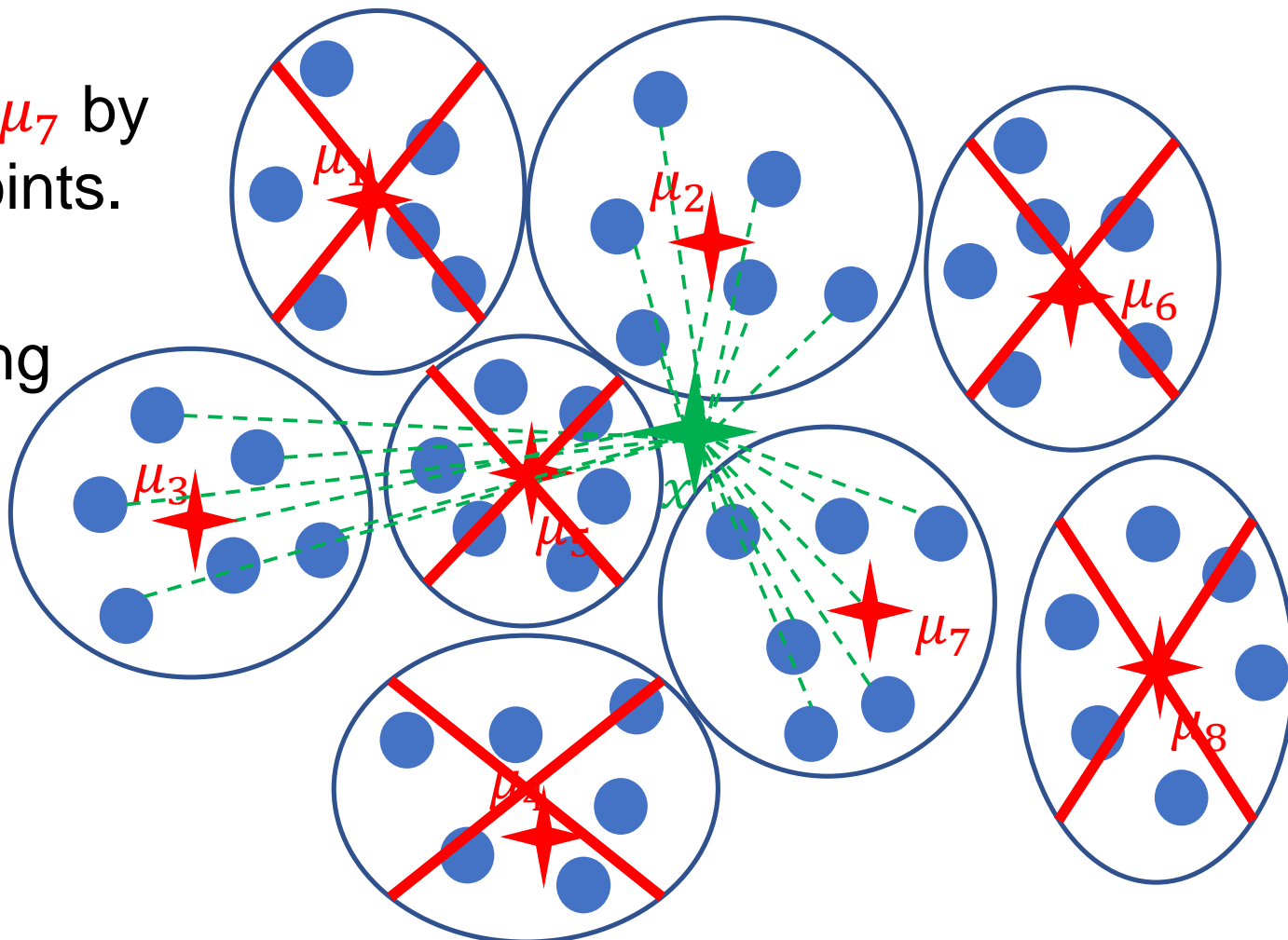


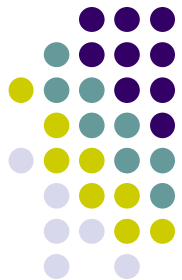
# Caratheodory Booster Illustration

## Step 3:

Replace  $\mu_2, \mu_3, \mu_7$  by their original points.

Delete remaining points.

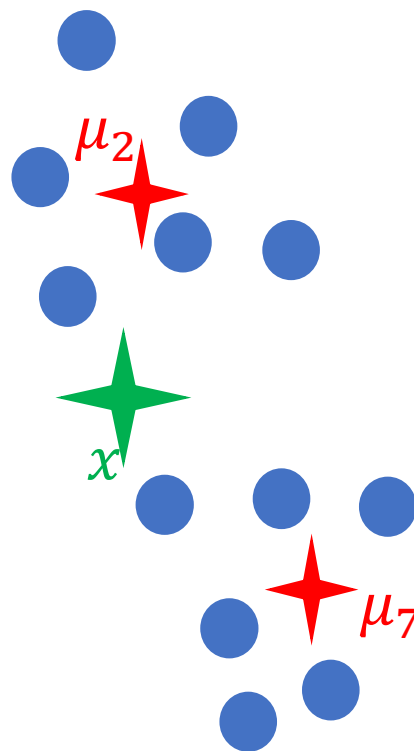
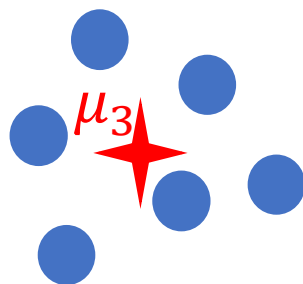


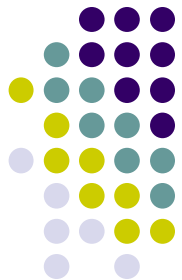


# Caratheodory Booster Illustration

## Step 4:

Repeat steps above until only few points remain.



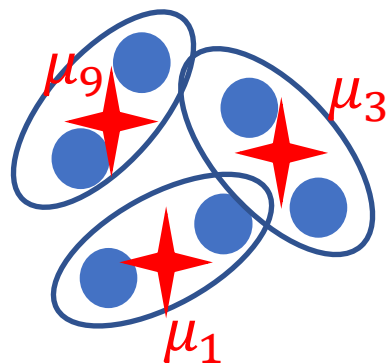


# Caratheodory Booster Illustration

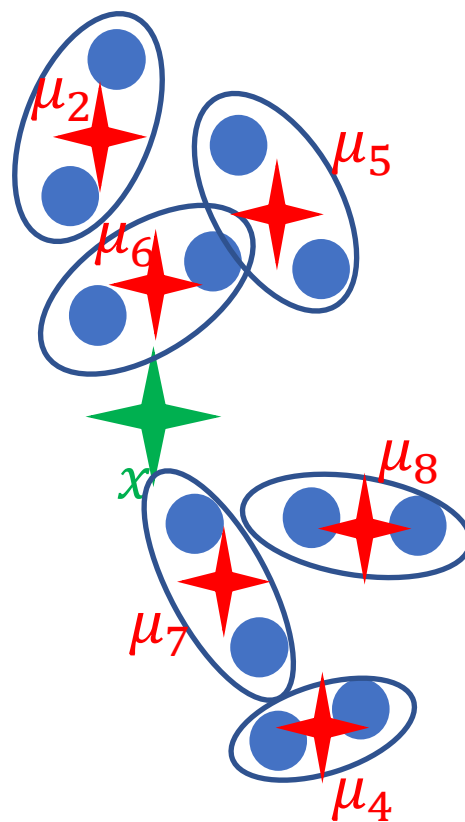
## Step 1:

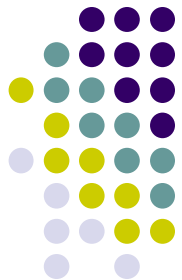
Partition the input  $d$ -dimensional points into  $k = 8$  equal sized subsets.

$\mu_1, \dots, \mu_8$  are the subset means.



The mean of  $\mu_1, \dots, \mu_8$  is equal to  $x$ , the mean of the input.



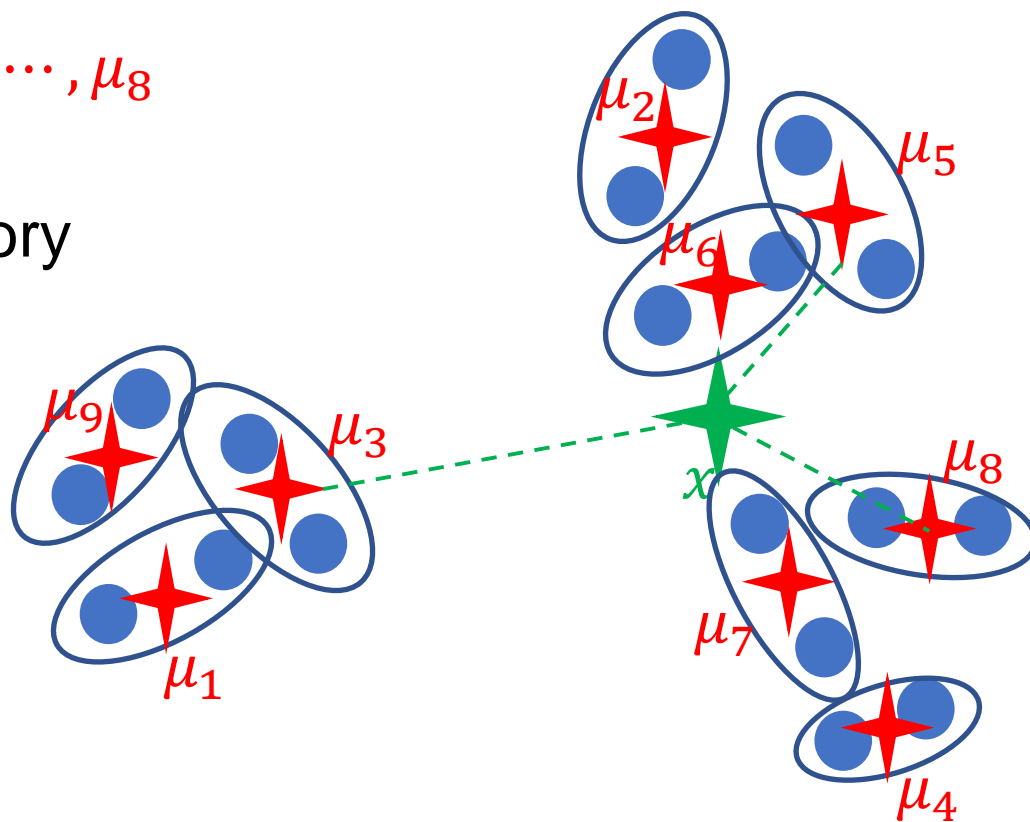


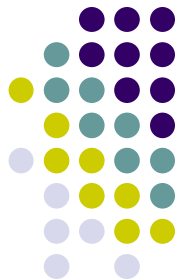
# Caratheodory Booster Illustration

## Step 2:

Consider only  $\mu_1, \dots, \mu_8$  and  $x$ .

Apply Caratheodory to represent  $x$  as a convex combination of a subset  $\mu_2, \mu_3, \mu_7$  of the means.



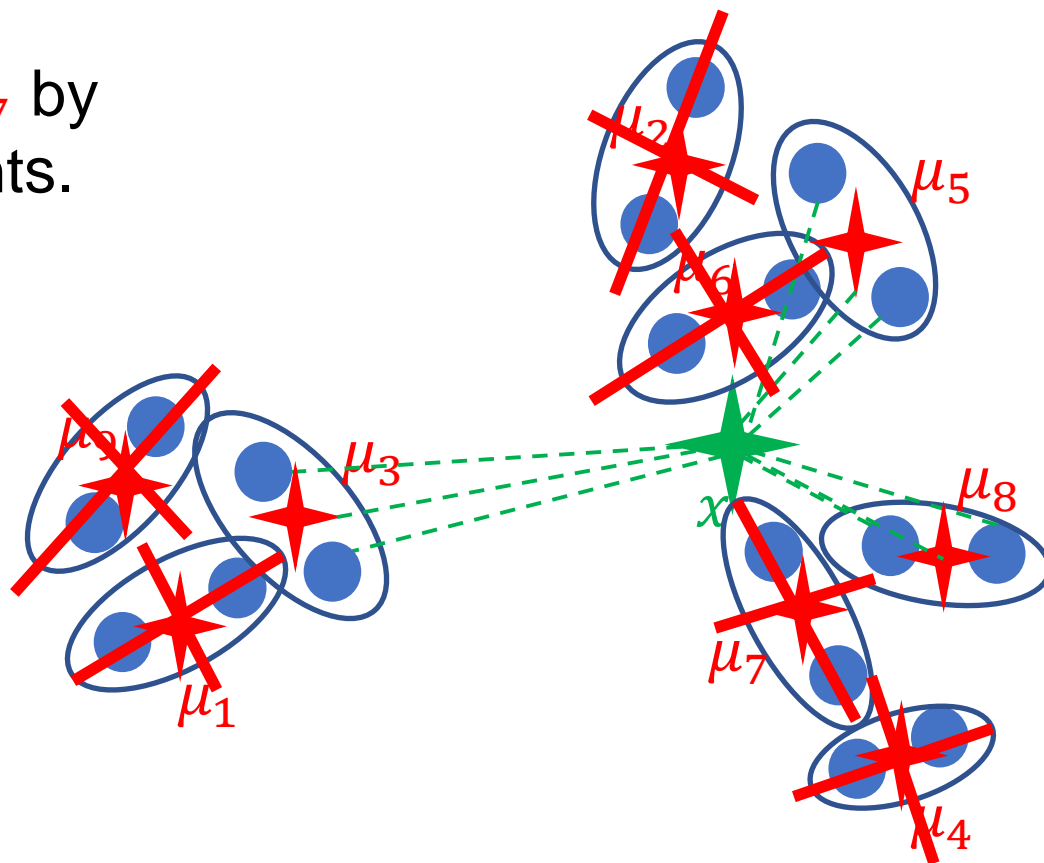


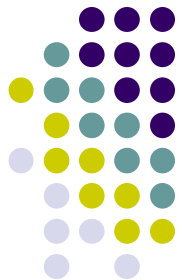
# Caratheodory Booster Illustration

## Step 3:

Replace  $\mu_2, \mu_3, \mu_7$  by their original points.

Delete remaining points.

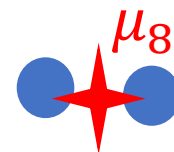
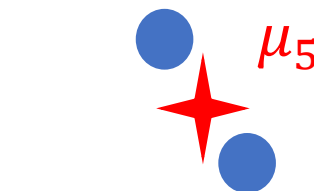




# Caratheodory Booster Illustration

## Step 4:

Repeat steps above until only few points remain.



$O(nd + d^4 \log n)$ ,  
i.e.,  $O(nd)$  time for  
sufficiently large  $n$ .





# Key Observation

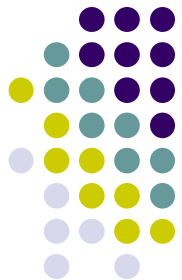
$$\begin{aligned}\|Ax - b\|_2^2 &= \left\| [A \mid b] \begin{matrix} x' \\ x \\ -1 \end{matrix} \right\|_2^2 \\ &= \left\| [A \mid b] x' \right\|_2^2 \\ &= x'^T \underbrace{[A \mid b]^T [A \mid b]}_{\text{Covariance matrix}} x'\end{aligned}$$

Therefore, for any  $m \geq 1$ ,  $C \in \mathbb{R}^{m \times d}$  and  $y \in \mathbb{R}^m$  such that

$$\underbrace{[A \mid b]^T [A \mid b]}_{P^T P} = \underbrace{[C \mid y]^T [C \mid y]}_{Z^T Z}$$

we have that:

$$\begin{aligned}f(\|Ax - b\|_2^2) + g(x) &= f(x'^T P^T P x') + g(x) \\ &= f(x'^T Z^T Z x') + g(x) \\ &= f(\|Cx - y\|_2^2) + g(x)\end{aligned}$$



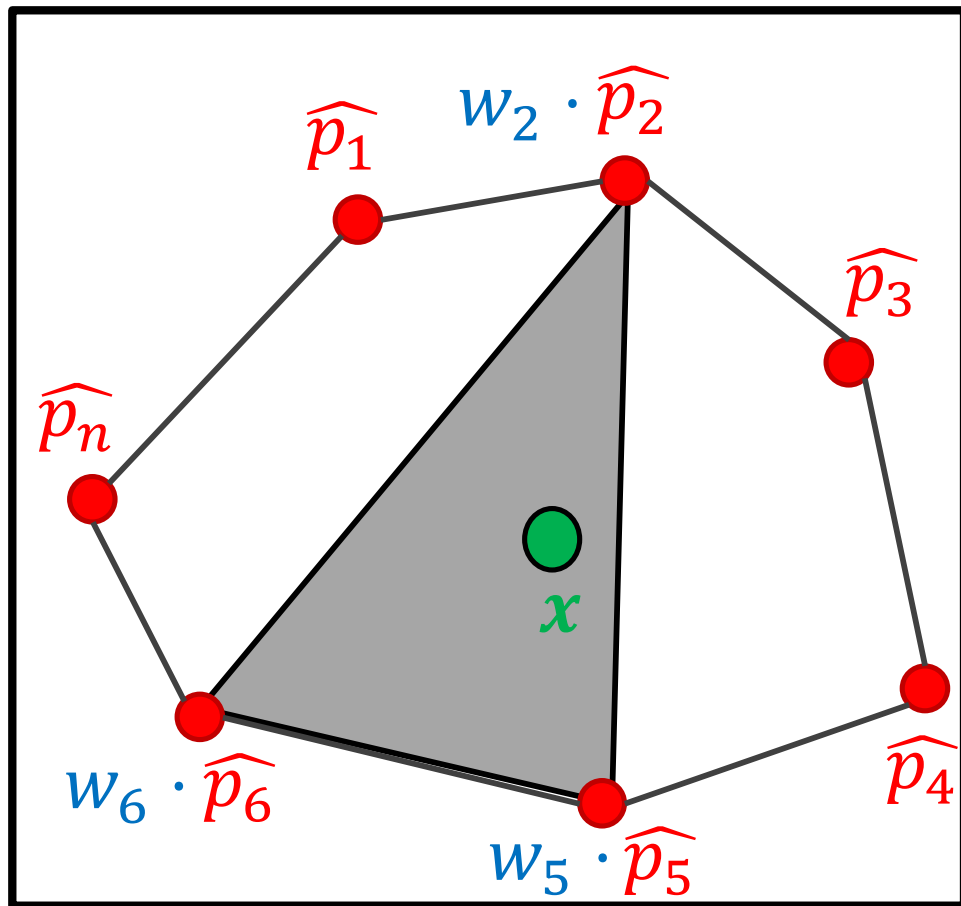
# Maintaining the Covariance Matrix

Input:

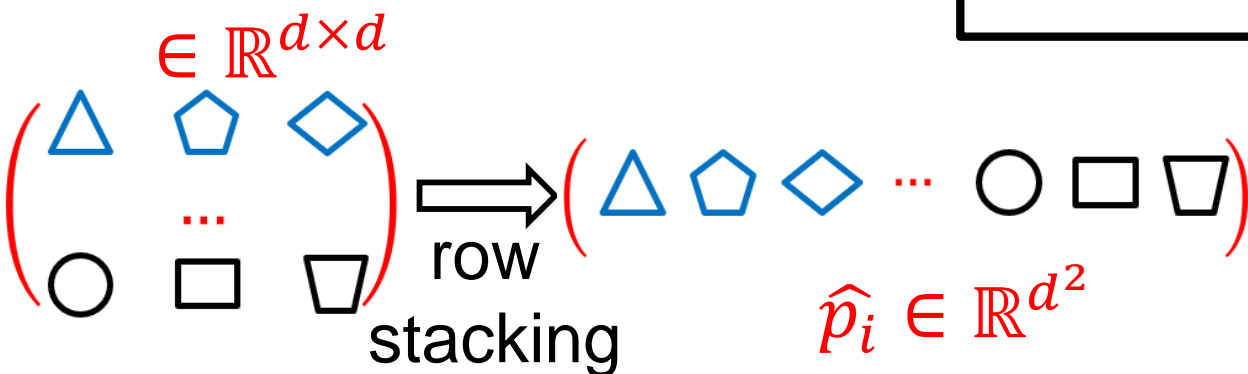
$$P = \begin{pmatrix} | & & | \\ p_1 & \dots & p_n \\ | & & | \end{pmatrix}^T \in \mathbb{R}^{n \times d}$$

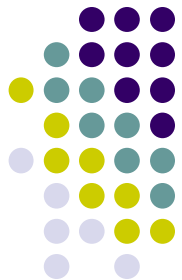
Immediate:

$$x = \frac{P^T P}{n} = \frac{1}{n} \sum_i p_i p_i^T$$



**Caratheodory**



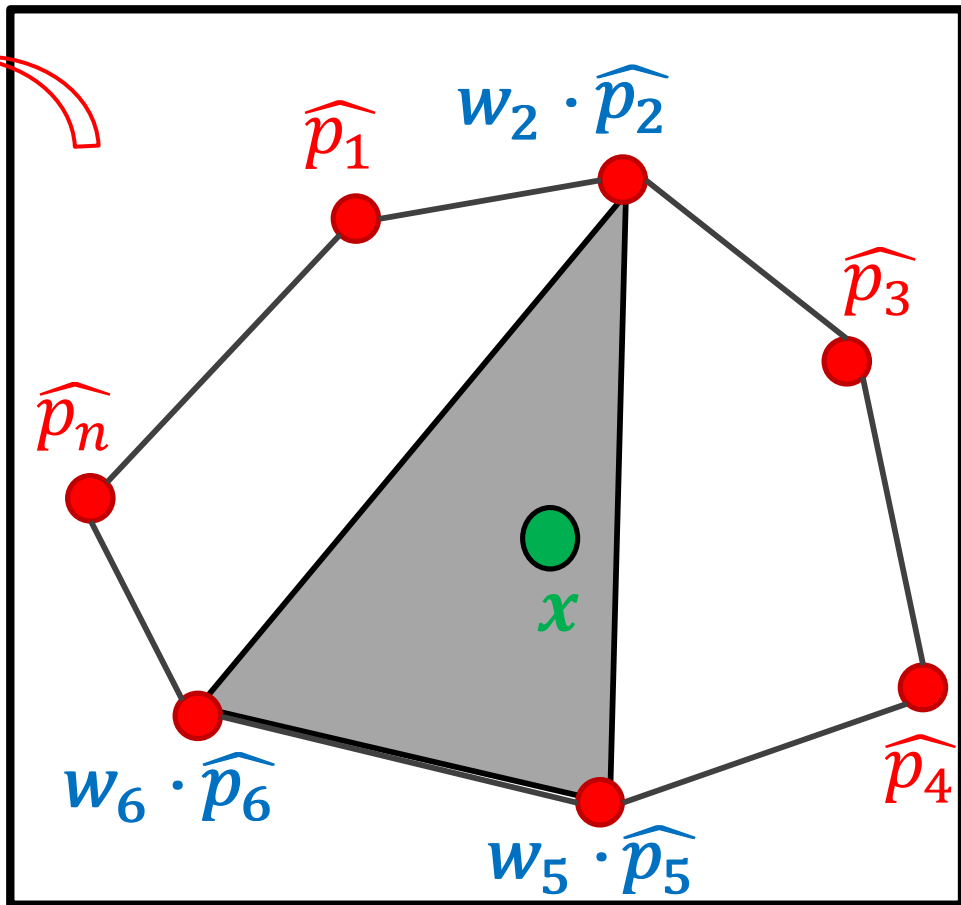


# Maintaining the Covariance Matrix

$$\begin{bmatrix} \sqrt{n \cdot w_2} p_2^T \\ \sqrt{n \cdot w_5} p_5^T \\ \sqrt{n \cdot w_6} p_6^T \end{bmatrix}^T \begin{bmatrix} \sqrt{n \cdot w_2} p_2^T \\ \sqrt{n \cdot w_5} p_5^T \\ \sqrt{n \cdot w_6} p_6^T \end{bmatrix}$$

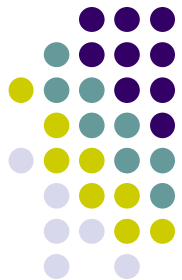
$$= Z^T Z$$

$$P^T P = Z^T Z$$

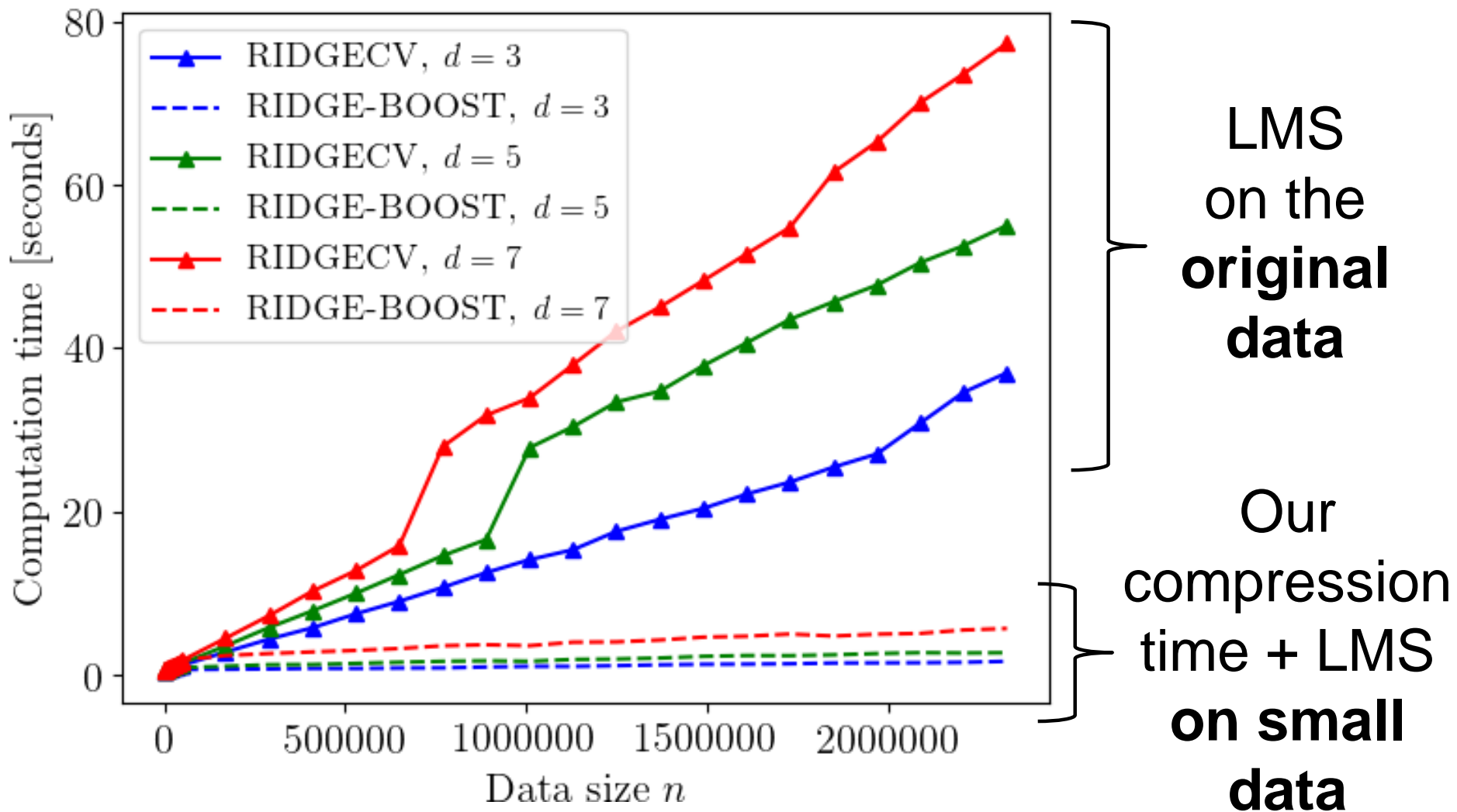


Caratheodory

Preserve covariance accurately!

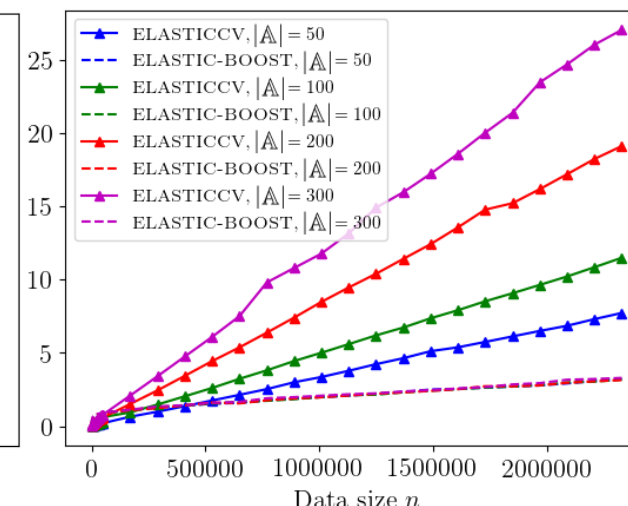
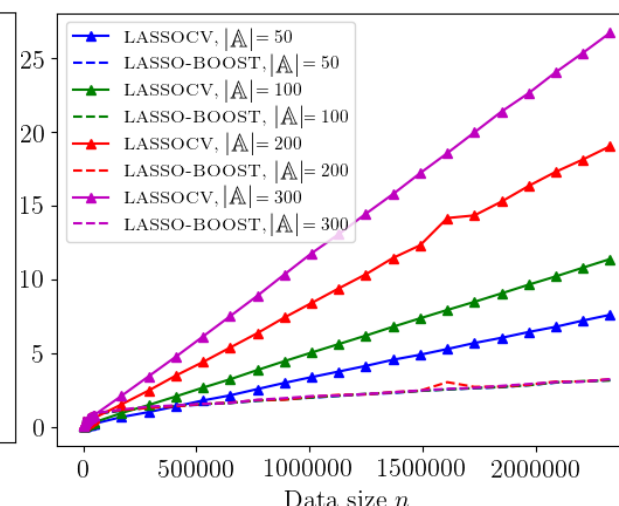
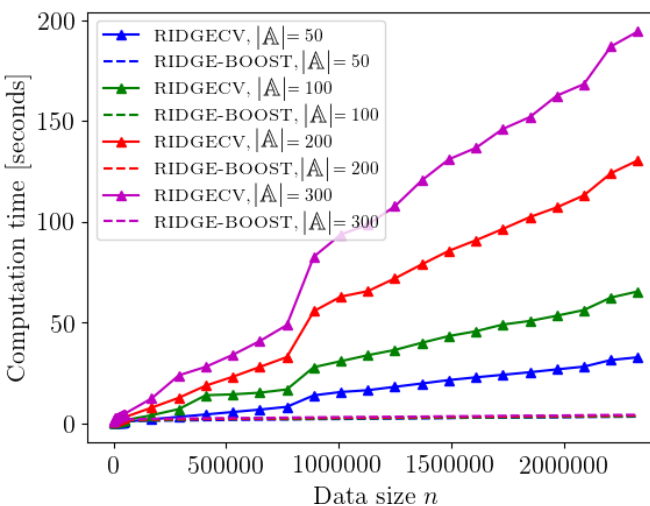
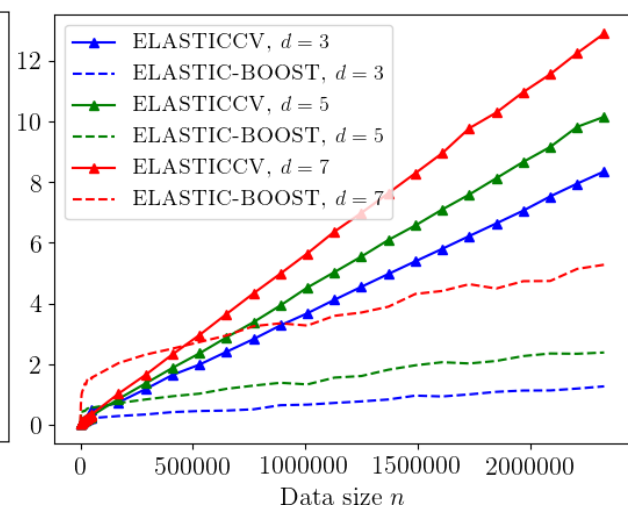
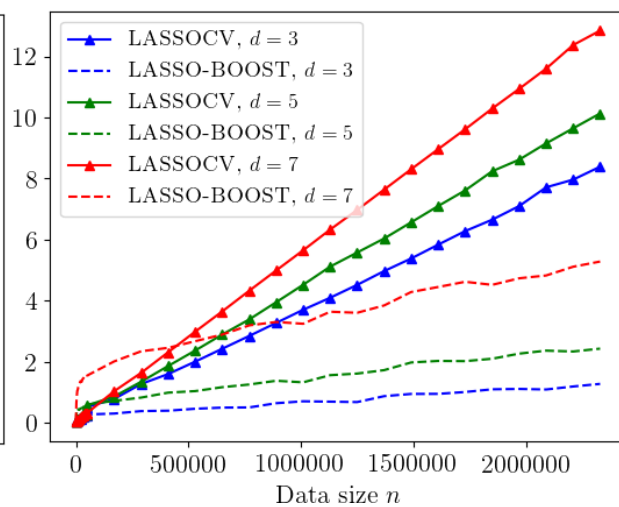
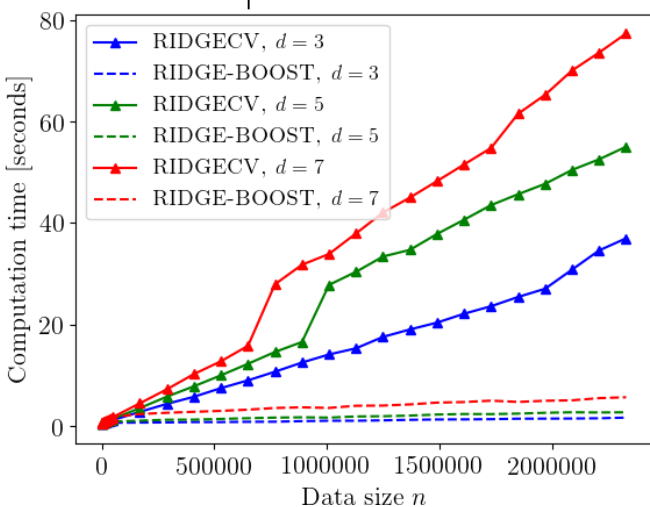


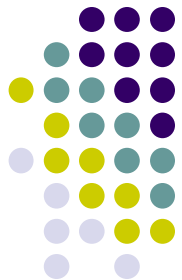
# Boosting Computation Time On Real-World Data (roughly x50)





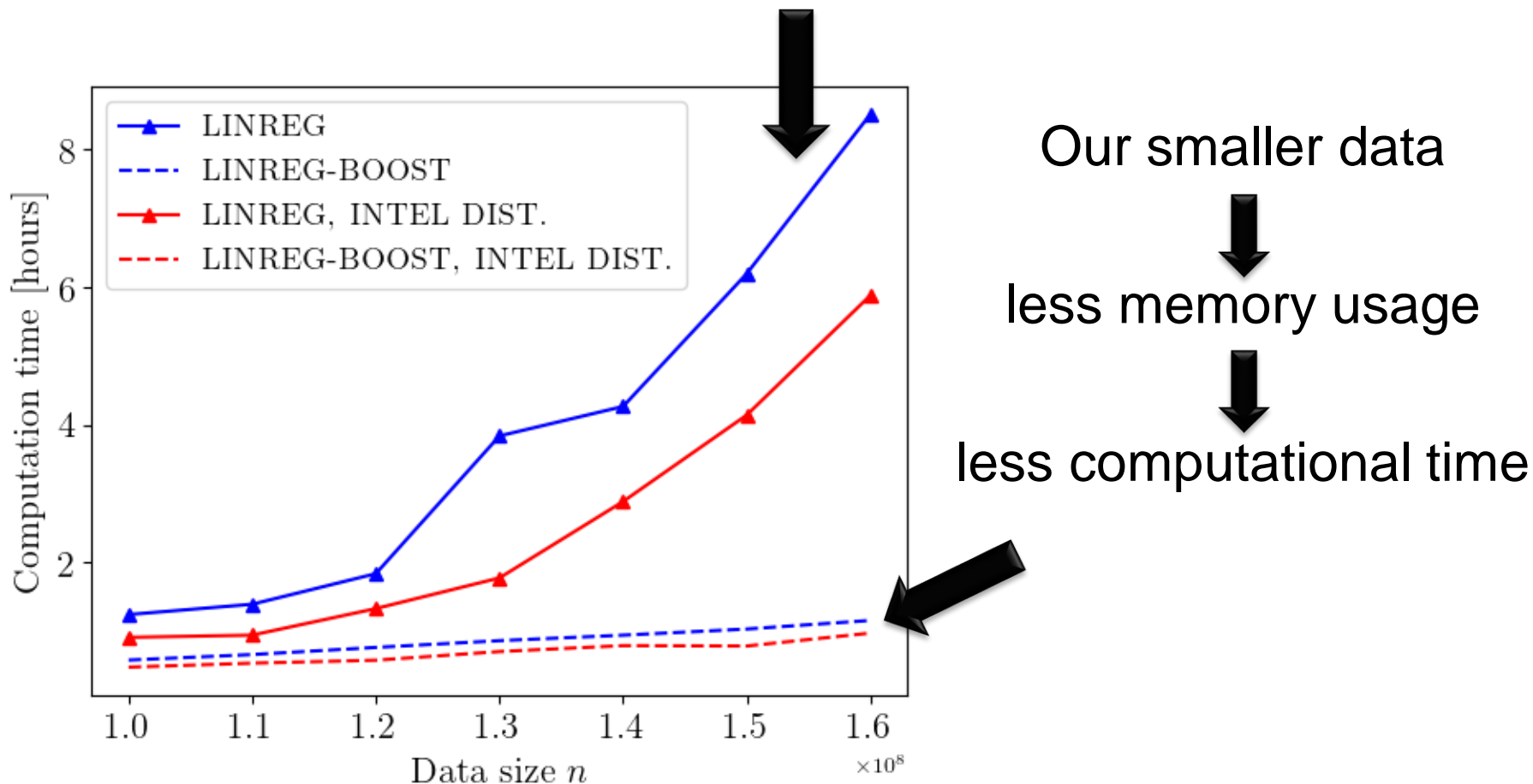
# Boosting Computation Time



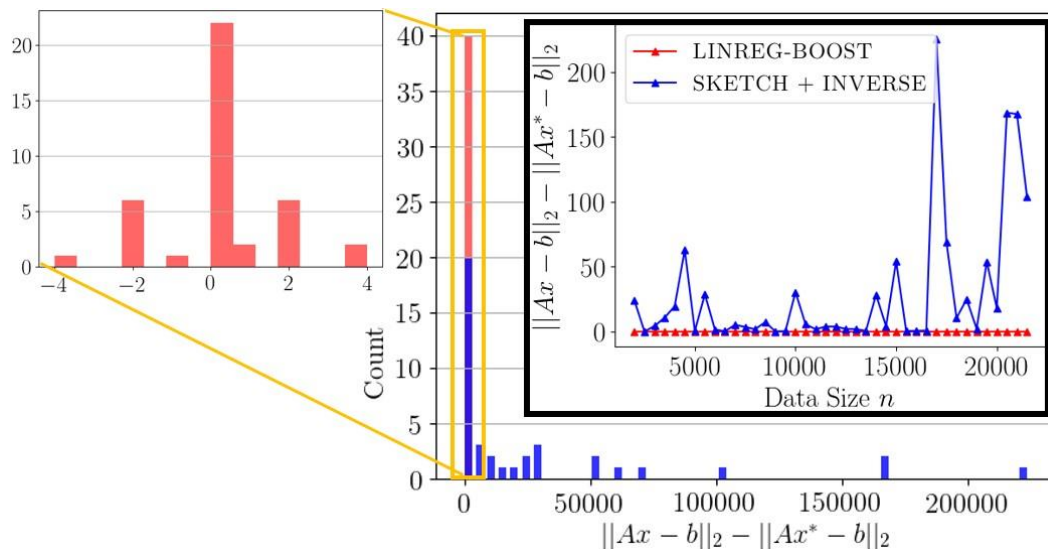
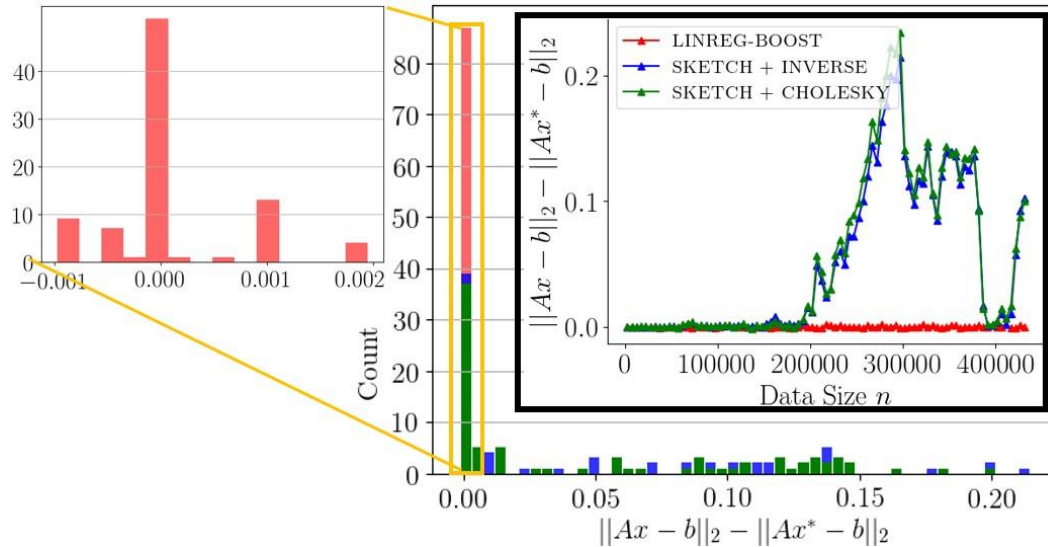


# Improving Space Complexity on Huge Data

Linear regression time is huge due to to memory overload.



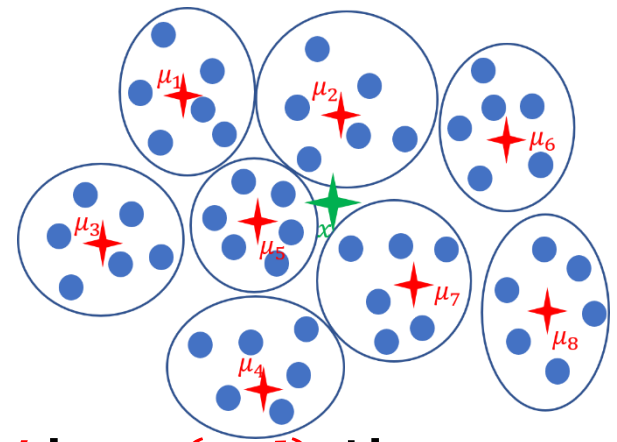
# Improving numerical stability vs. other compression schemes



Similar numerical improvement for Ridge / Lasso / Elastic-net

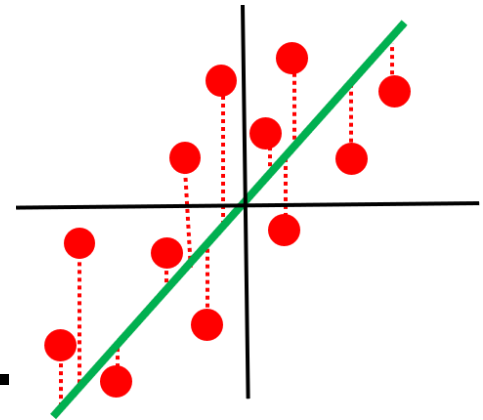


# Summary



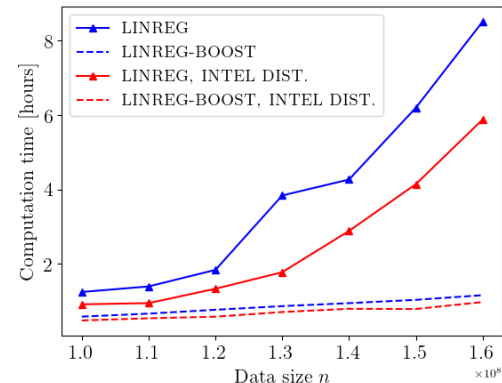
• Computing the *Caratheodory set* in  $O(nd)$  time.

- Practically and provably boosting:
- **time complexity** of LMS solvers.
  - **space complexity** of LMS solvers.
  - **numerical accuracy** of LMS solvers.



• Evaluation on **real-world** data.

• Full **open source code**.







Thank you 😊



Open source code