

# Robust Least Squares / Adversarial training in linear models

Andersen Ang  
U.Southampton UK

[angms.science](http://angms.science)

March 24, 2025

1st ver. March 21, 2025

## Content

Ordinary Least Squares

Motivation of Robust Least Squares

Robust Least Squares

Min-max optimization

Alternating optimization

$$\begin{aligned}(\text{RLS}) : \quad & \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \frac{1}{2} \|(\mathbf{A} + \mathbf{E})\mathbf{x} - \mathbf{b}\|_2^2 \\ & = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \sum_{i=1}^m \frac{1}{2} \left( \langle \mathbf{a}_i + \mathbf{e}_i, \mathbf{x} \rangle - b_i \right)^2\end{aligned}$$

# Contents

Ordinary Least Squares

Motivation of Robust Least Squares

Robust Least Squares

Min-max optimization

Alternating optimization

# Ordinary Least Squares

$$\text{(OLS)} : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

- Consider

$$\text{(OLS)} : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

- $\mathbf{A} : \mathbb{R}^{m \times n}$  is a given matrix, possibly not full rank
  - $\mathbf{b} : \mathbb{R}^m$  is a given vector
  - $\mathbf{x} \in \mathbb{R}^n$  is the optimization variable
  - The objective function  $f$  is a  $\mathbb{R}^n \rightarrow \mathbb{R}$  mapping
- 
- OLS is
    - the most basic regression model in statistics
    - the most basic linear model in machine learning
    - the most basic convex model in optimization

# Ordinary Least Squares

$$\text{(OLS)} : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

- Summation Form in statistics/machine learning

$$\text{(OLS)} \iff \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m \frac{1}{2} (\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i)^2 = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m \frac{1}{2} \ell(\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i)$$

- $\mathbf{a}_i$  is the  $i$ th row of  $\mathbf{A}$
  - $b_i$  is the  $i$ th element (a scalar) of  $\mathbf{b}$
  - $\langle \mathbf{a}_i, \mathbf{x} \rangle$  is the inner product between the  $i$ th row of  $\mathbf{A}$  and  $\mathbf{x}$
  - $\ell = \frac{1}{2}(\cdot)^2$  is called a lost function, and usually  $\ell_i(\cdot) := \frac{1}{2}(\langle \mathbf{a}_i, \cdot \rangle - b_i)^2$
- Quadratic Programming Form

$$\text{(OLS)} \iff \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \langle \mathbf{Q}\mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{p}, \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{b}\|_2^2$$

- $\mathbf{Q} = \mathbf{A}^\top \mathbf{A}$
- $\mathbf{p} = \mathbf{A}^\top \mathbf{b}$

here are 90 pages more details

# Contents

Ordinary Least Squares

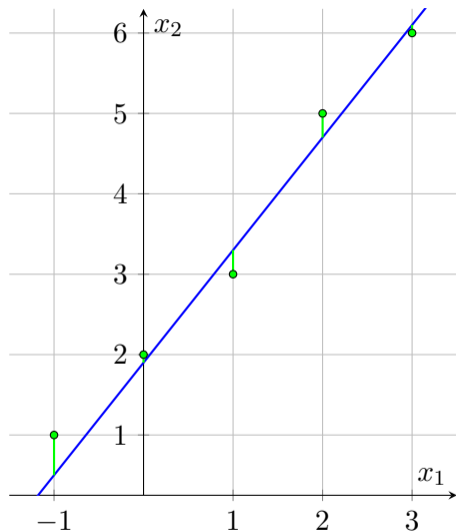
Motivation of Robust Least Squares

Robust Least Squares

Min-max optimization

Alternating optimization

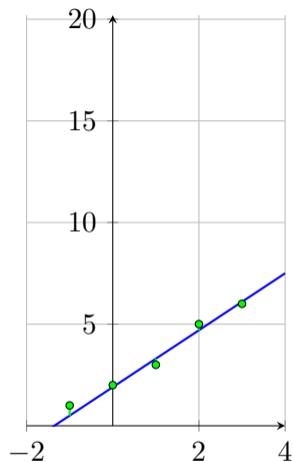
## Picture of OLS



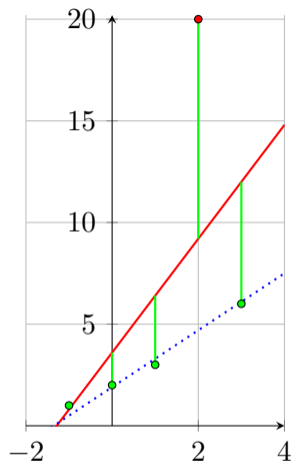
$$(\text{OLS}) \iff \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^m \frac{1}{2} \left( \langle \mathbf{a}_i, \mathbf{x} \rangle - b_i \right)^2$$

- We are giving a bunch of points (2D points here)
- Solving OLS is to find the blue line
- The blue line is the best fit to these points
- The word “best” refers to the sum of the length of the green segments is the shortest
- Geometrically,  
solve OLS  $\iff$  move the **line** for the shortest **sum**

## Motivation of Robust LS: OLS is sensitive to outlier



- Data is “ok”
- The line fits all points well “on average”



- A point is an **outlier**
- The **line** is dragged by that outlier
- To fit the line we should “ignore” the outlier

# Contents

Ordinary Least Squares

Motivation of Robust Least Squares

**Robust Least Squares**

Min-max optimization

Alternating optimization



## Robust Least Squares

$$\text{(OLS)} : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

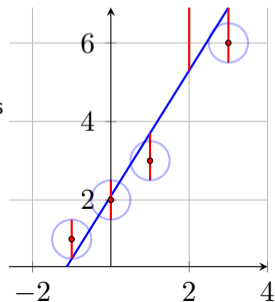
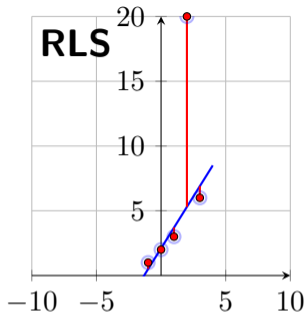
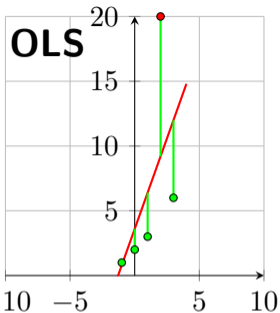
$$\text{(RLS)} : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|e_i\| \leq \delta} \frac{1}{2} \|(\mathbf{A} + \mathbf{E})\mathbf{x} - \mathbf{b}\|_2^2$$

- What's new:  $\max$  and  $\mathbf{E}$ 
  - $\mathbf{E} : \mathbb{R}^{m \times n}$  is an unknown matrix, also an optimization variable here
  - $e_i : \mathbb{R}^n$  is the  $i$ th row of  $\mathbf{E}$
  - $\delta \in \mathbb{R}$  is the max perturbation radius, it tells the max magnitude of the perturbation
- May be easier to understand RLS in summation form

$$\text{(RLS)} \iff \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|e_i\| \leq \delta} \sum_{i=1}^m \frac{1}{2} \left( \langle \mathbf{a}_i + e_i, \mathbf{x} \rangle - b_i \right)^2$$

- think of  $\mathbf{a}_i$  is "true value" and  $e_i$  is "additive noise"
- think of  $\hat{\mathbf{a}}_i = \mathbf{a}_i + e_i$  is  $\mathbf{a}_i$  corrupted by noise
- $\|e_i\| \leq \delta$  means how large is the noise's size
- we don't know the exact value of noise, so we consider all the possible  $e_i$
- $\max_{\|e_i\| \leq \delta}$  means we take the worst case among all the possible noise

## OLS vs RLS

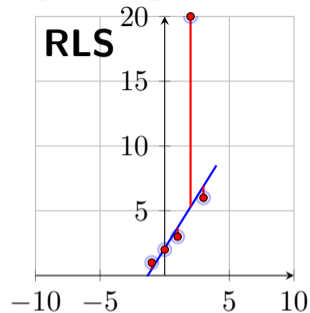


- Data points have “zero volume” (I draw them thicker for visibility)
- The vertical line touches the points

In RLS

- Data points are “points with volume”
- Circle radius = perturbation radius  $\delta$
- We are now considering the largest sum of red length touching circles
- “largest sum” contains the meaning of the “worst case”
- using worst case formulation, we magically ignore the outlier !

## Why RLS ignores the outlier?



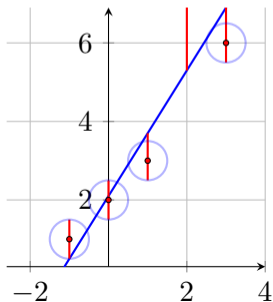
- By definition: outliers are far away  $\implies$  the **length** is large
- If  $\delta$  is small, the perturbation  $\|e_i\| \leq \delta$  has limited impact on the **length**
- For inliers (non-outliers), the **length** is smaller
- the perturbation  $\|e_i\| \leq \delta$  has larger impact on the **length**
- If the **line** is moving away from the inliers, the inliers error

$$\text{inliers error} = \max_{\|e_i\| \leq \delta} \sum_{i \in \text{inliers}} \frac{1}{2} \left( \langle \mathbf{a}_i + \mathbf{e}_i, \mathbf{x} \rangle - b_i \right)^2$$

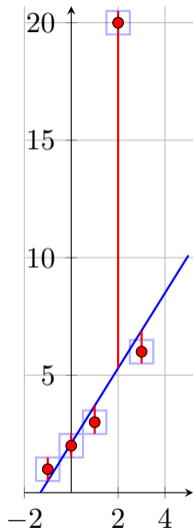
will go up much faster than the outliers error

$$\text{outliers error} = \max_{\|e_i\| \leq \delta} \sum_{i \notin \text{inliers}} \frac{1}{2} \left( \langle \mathbf{a}_i + \mathbf{e}_i, \mathbf{x} \rangle - b_i \right)^2 \approx \sum_{i \notin \text{inliers}} \frac{1}{2} \left( \langle \mathbf{a}_i, \mathbf{x} \rangle - b_i \right)^2$$

therefore the argmin will pull the **line** towards inliers



## $l_\infty$ -norm version of RLS



$$(\text{RLS}) : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \frac{1}{2} \|(\mathbf{A} + \mathbf{E})\mathbf{x} - \mathbf{b}\|_2^2$$

- Previously we used  $\|\mathbf{e}_i\|_2 \leq \delta$

$$\|\mathbf{v}\|_2 = \sqrt{v_1^2 + v_2^2 + \dots}$$

- Here we used  $\|\mathbf{e}_i\|_\infty \leq \delta$

$$\|\mathbf{v}\|_\infty = \max_i \{|v_1|, |v_2|, \dots\}$$

- There is also  $\|\mathbf{e}_i\|_1 \leq \delta$

$$\|\mathbf{v}\|_1 = |v_1| + |v_2| + \dots$$

# Contents

Ordinary Least Squares

Motivation of Robust Least Squares

Robust Least Squares

**Min-max optimization**

Alternating optimization

# Min-max Optimization

$$\text{(RLS)} : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \frac{1}{2} \|(\mathbf{A} + \mathbf{E})\mathbf{x} - \mathbf{b}\|_2^2 = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \sum_{i=1}^m \frac{1}{2} \left( \langle \mathbf{a}_i + \mathbf{e}_i, \mathbf{x} \rangle - b_i \right)^2$$

- RLS is an example of min-max optimization

$$\text{(Min-max)} : \min_{\mathbf{x}} \max_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$$

- Goal of solving min-max problem: find saddle point  $(\mathbf{x}^*, \mathbf{y}^*)$

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq f(\mathbf{x}, \mathbf{y}^*)$$

- Where will you see Min-max optimization
  - Robust optimization
  - Dual-based optimization
  - Game Theory
  - Machine Learning using robustness and adversarial training

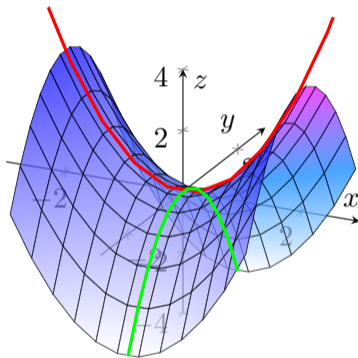
# What is a saddle point

$$(\text{Min-max}) : \min_x \max_y f(\mathbf{x}, \mathbf{y})$$

- Goal of solving min-max problem: find saddle point  $(\mathbf{x}^*, \mathbf{y}^*)$

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq f(\mathbf{x}, \mathbf{y}^*)$$

- $\mathbf{x}^*$  is minimizer if you view along x-axis
- $\mathbf{y}^*$  is maximizer if you view along y-axis
- $(\mathbf{x}^*, \mathbf{y}^*)$  is x-min y-max if we view x-axis y-axis together



# General theory of Min-max optimization

$$(\text{Min-max}) : \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}) \quad \text{Saddle point: } f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq f(\mathbf{x}, \mathbf{y}^*)$$

$f$  is  $\mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$

- **Von Neumann's Minimax Theorem** If  $f$  is convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ , then

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{y} \in \mathcal{Y}} \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \mathbf{y})$$

- **Subgradient 1st-order optimality**

$$\mathbf{0}_n \in \partial_{\mathbf{x}} \left( f(\mathbf{x}^*, \mathbf{y}^*) + \iota_{\mathcal{X}}(\mathbf{x}^*, \mathbf{y}^*) \right), \quad \mathbf{0}_m \in \partial_{\mathbf{y}} \left( f(\mathbf{x}^*, \mathbf{y}^*) + \iota_{\mathcal{Y}}(\mathbf{x}^*, \mathbf{y}^*) \right)$$

- **Alternating Optimization: Gradient Descent-Ascent (GDA)**

$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)$  followed by proj, prox ect

$\mathbf{y}_{k+1} = \mathbf{y}_k + \alpha \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k)$  followed by proj, prox ect



# Contents

Ordinary Least Squares

Motivation of Robust Least Squares

Robust Least Squares

Min-max optimization

Alternating optimization

## How to solve RLS

$$\text{(RLS)} : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \frac{1}{2} \|(\mathbf{A} + \mathbf{E})\mathbf{x} - \mathbf{b}\|_2^2 = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \sum_{i=1}^m \frac{1}{2} \left( \langle \mathbf{a}_i + \mathbf{e}_i, \mathbf{x} \rangle - b_i \right)^2$$

- Alternating Optimization: repeat

1. Update  $\mathbf{x}$  with a fix  $\mathbf{E}$
2. Update  $\mathbf{E}$  with a fix  $\mathbf{x}$

- x-update: gradient descent ([here are 90 pages more details](#))

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \left( (\mathbf{A} + \mathbf{E})^\top (\mathbf{A} + \mathbf{E}) \mathbf{x}_k - (\mathbf{A} + \mathbf{E})^\top \mathbf{b} \right)$$

- E-update: projected gradient ascent

$$\mathbf{e}_{i,k+1} = \operatorname{proj}_{\|\mathbf{e}\| \leq \delta} \left( \mathbf{e}_{i,k} + \alpha \nabla_{\mathbf{e}} \frac{1}{2} \left( \langle \mathbf{a}_i + \mathbf{e}, \mathbf{x} \rangle - b_i \right)^2 \right)$$

- [projected gradient DEscent](#) (change DEscent to AScent)
- [proj onto  \$\ell\_2\$ -ball](#)
- [proj onto  \$\ell\_1\$ -ball](#)

**Find**  $\nabla_{\mathbf{x}} \left( (\mathbf{a} + \mathbf{x})^\top \mathbf{c} - b \right)^2$

$$f(\mathbf{x}) = \left( (\mathbf{a} + \mathbf{x})^\top \mathbf{c} - b \right)^2 = \left( \mathbf{a}^\top \mathbf{c} + \mathbf{x}^\top \mathbf{c} - b \right)^2$$

• Let  $u = \mathbf{a}^\top \mathbf{c} + \mathbf{x}^\top \mathbf{c} - b$ , then  $f(\mathbf{x}) = u^2$

• Chain rule

$$\frac{\partial f}{\partial \mathbf{x}} = 2u \frac{\partial u}{\partial \mathbf{x}}$$

• With  $\frac{\partial u}{\partial \mathbf{x}} = \mathbf{c}$

$$\frac{\partial f}{\partial \mathbf{x}} = 2u\mathbf{c} = 2 \left( \mathbf{a}^\top \mathbf{c} + \mathbf{x}^\top \mathbf{c} - b \right) \mathbf{c} = 2 \left( (\mathbf{a} + \mathbf{x})^\top \mathbf{c} - b \right) \mathbf{c}$$

• Final gradient:

$$\nabla_{\mathbf{x}} \frac{1}{2} \left( (\mathbf{a} + \mathbf{x})^\top \mathbf{c} - b \right)^2 = \left( (\mathbf{a} + \mathbf{x})^\top \mathbf{c} - b \right) \mathbf{c} = 2f(\mathbf{x})\mathbf{c}$$

## Alternating optimization solve RLS

$$(\text{RLS}) : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \frac{1}{2} \|(\mathbf{A} + \mathbf{E})\mathbf{x} - \mathbf{b}\|_2^2 = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|\mathbf{e}_i\| \leq \delta} \sum_{i=1}^m \frac{1}{2} \left( \langle \mathbf{a}_i + \mathbf{e}_i, \mathbf{x} \rangle - b_i \right)^2$$

- x-update:

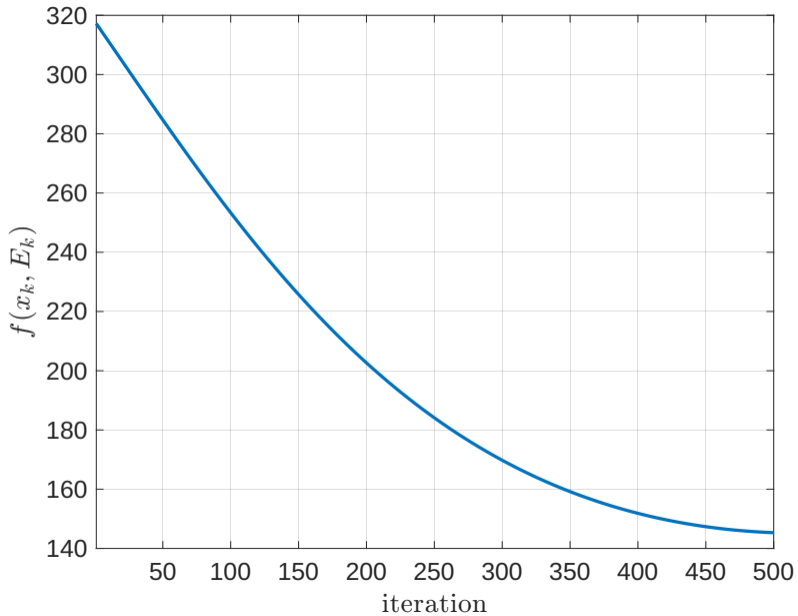
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \left( (\mathbf{A} + \mathbf{E})^\top (\mathbf{A} + \mathbf{E}) \mathbf{x}_k - (\mathbf{A} + \mathbf{E})^\top \mathbf{b} \right)$$

- E-update:

$$\hat{\mathbf{e}}_{i,k+1} = \mathbf{e}_{i,k} + \alpha \left( \langle \mathbf{a}_i + \mathbf{e}_{i,k}, \mathbf{x}_k \rangle - b_i \right) \mathbf{x}_k$$

$$\mathbf{e}_{i,k+1} = \operatorname{proj}_{\|\mathbf{e}\| \leq \delta} \left( \hat{\mathbf{e}}_{i,k+1} \right) = \frac{\hat{\mathbf{e}}_{i,k+1}}{\max \left\{ 1, \frac{\|\hat{\mathbf{e}}_{i,k+1}\|}{\delta} \right\}}$$

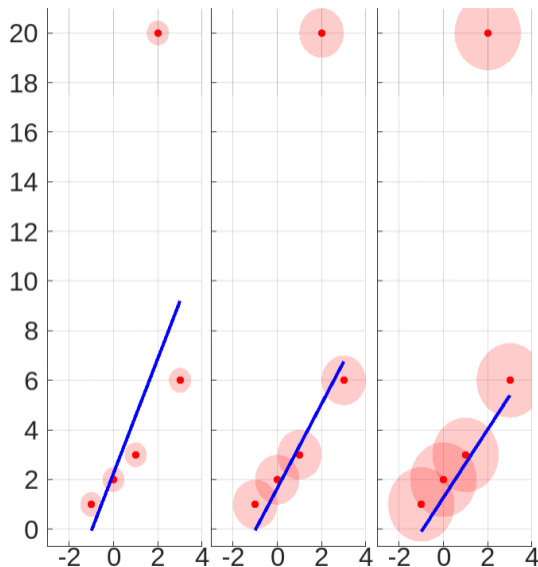
## Quick-dirty-done MATLAB code



## Effect of the radius $\delta$

$$(\text{RLS}) : \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \max_{\|e_i\| \leq \delta} \frac{1}{2} \|(\mathbf{A} + \mathbf{E})\mathbf{x} - \mathbf{b}\|_2^2$$

- If  $\delta \rightarrow 0$ , then  $\mathbf{E} \rightarrow \mathbf{0}$ , and RLS  $\rightarrow$  OLS
- Decreasing  $\delta$  makes RLS back to OLS
- Increasing  $\delta$  makes the inlier “more important”



## Last page - summary

Ordinary Least Squares

Motivation of Robust Least Squares

Robust Least Squares

Min-max optimization

Alternating optimization

End of document