# Robustness of OLS to sample removals: theoretical analysis and implications

#### Boaz Nadler

Weizmann Institute of Science, Israel

Joint work with
Michael Feldman and Eyar Azar

Vienna, Sep. 2025

- Trust in a learned model; specifically its robustness to removal of few samples

- Trust in a learned model; specifically its robustness to removal of few samples
- Most influential subset selection problem / robustness auditing

- Trust in a learned model; specifically its robustness to removal of few samples
- Most influential subset selection problem / robustness auditing
- A real data example

- Trust in a learned model; specifically its robustness to removal of few samples
- Most influential subset selection problem / robustness auditing
- A real data example
- Theoretical analysis of robustness auditing for ordinary least squares

- Trust in a learned model; specifically its robustness to removal of few samples
- Most influential subset selection problem / robustness auditing
- A real data example
- Theoretical analysis of robustness auditing for ordinary least squares
- Revisit the real data example

- Trust in a learned model; specifically its robustness to removal of few samples
- Most influential subset selection problem / robustness auditing
- A real data example
- Theoretical analysis of robustness auditing for ordinary least squares
- Revisit the real data example
- Insights and implications

Standard workflow of supervised learning:

Standard workflow of supervised learning:

#### Input:

training set of n samples  $(\mathbf{x}_i, y_i)$ , i = 1, ..., n

Standard workflow of supervised learning:

Input:

training set of n samples  $(\mathbf{x}_i, y_i)$ , i = 1, ..., n

Goal:

construct a predictor for the response y for new  $\mathbf{x}$ 's whose responses y are not observed

Standard workflow of supervised learning:

Input:

training set of n samples  $(\mathbf{x}_i, y_i)$ ,  $i = 1, \ldots, n$ 

Goal:

construct a predictor for the response y for new  $\mathbf{x}$ 's whose responses y are not observed

Focus in this talk: regression setting

Standard workflow of supervised learning:

Input:

training set of 
$$n$$
 samples  $(\mathbf{x}_i, y_i)$ ,  $i = 1, \ldots, n$ 

Goal:

construct a predictor for the response y for new  $\mathbf{x}$ 's whose responses y are not observed

Focus in this talk: regression setting

$$\mathbf{x} \in \mathbb{R}^p$$
  $y \in \mathbb{R}$ 

Standard workflow of supervised learning:

Input:

training set of 
$$n$$
 samples  $(\mathbf{x}_i, y_i)$ ,  $i = 1, ..., n$ 

Goal:

construct a predictor for the response y for new  $\mathbf{x}$ 's whose responses y are not observed

Focus in this talk: regression setting

$$\mathbf{x} \in \mathbb{R}^p$$
  $y \in \mathbb{R}$ 

Assume:

p = number of features < n = number of samples

X -  $n \times p$  matrix of all sample is of full rank p



Predictor f often found by choosing a loss function  $\ell$ , and minimizing empirical risk over some class of functions  $\mathcal{F}$ ,

Predictor f often found by choosing a loss function  $\ell$ , and minimizing empirical risk over some class of functions  $\mathcal{F}$ ,

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(\mathbf{x}_i))$$

Predictor f often found by choosing a loss function  $\ell$ , and minimizing empirical risk over some class of functions  $\mathcal{F}$ ,

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(\mathbf{x}_i))$$

**Example:** Squared loss combined with class  $\mathcal{F}$  of linear functions gives *ordinary least squares*,

$$\hat{f}(\mathbf{x}) = \widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}$$

Predictor f often found by choosing a loss function  $\ell$ , and minimizing empirical risk over some class of functions  $\mathcal{F}$ ,

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(\mathbf{x}_i))$$

**Example:** Squared loss combined with class  $\mathcal{F}$  of linear functions gives *ordinary least squares*,

$$\hat{f}(\mathbf{x}) = \widehat{\boldsymbol{\beta}}^{\top} \mathbf{x}$$

where

$$\widehat{\boldsymbol{\beta}} = (X^{\top}X)^{-1}X^{\top}\boldsymbol{y}$$



Can we trust learned model / its predictions?

Can we trust learned model / its predictions?

Accuracy of predictions: cross validation, conformal prediction

Can we trust learned model / its predictions?

Accuracy of predictions: cross validation, conformal prediction

**Accuracy of estimated model parameters:** asymptotic confidence intervals, bootstrap.

Can we trust learned model / its predictions?

Accuracy of predictions: cross validation, conformal prediction

**Accuracy of estimated model parameters:** asymptotic confidence intervals, bootstrap.

**Robustness:** detection of outliers, or very influential (individual) samples.

Can we trust learned model / its predictions?

Accuracy of predictions: cross validation, conformal prediction

**Accuracy of estimated model parameters:** asymptotic confidence intervals, bootstrap.

**Robustness:** detection of outliers, or very influential (individual) samples.

In the statistics literature: long history of works under umbrella of Robustness Diagnostics

In past few years, several researchers noted that in various datasets, removing a *small* number  $k \ll n$  of (specifically chosen) training samples leads to *large* changes in the learned model,

In past few years, several researchers noted that in various datasets, removing a *small* number  $k \ll n$  of (specifically chosen) training samples leads to *large* changes in the learned model, very different predictions and/or estimated coefficients

In past few years, several researchers noted that in various datasets, removing a *small* number  $k \ll n$  of (specifically chosen) training samples leads to *large* changes in the learned model,

very different predictions and/or estimated coefficients

**Example:** For OLS, let  $S \subset [n]$  be remaining set of samples after removal of k specific samples, |S| = n - k,

$$\widehat{\boldsymbol{\beta}}_{\mathcal{S}} = \left(\boldsymbol{X}_{\mathcal{S}}^{\top}\boldsymbol{X}_{\mathcal{S}}\right)^{-1}\boldsymbol{X}_{\mathcal{S}}^{\top}\boldsymbol{y}_{\mathcal{S}}.$$

In past few years, several researchers noted that in various datasets, removing a *small* number  $k \ll n$  of (specifically chosen) training samples leads to *large* changes in the learned model,

very different predictions and/or estimated coefficients

**Example:** For OLS, let  $S \subset [n]$  be remaining set of samples after removal of k specific samples, |S| = n - k,

$$\widehat{\boldsymbol{\beta}}_{\mathcal{S}} = \left( X_{\mathcal{S}}^{\top} X_{\mathcal{S}} \right)^{-1} X_{\mathcal{S}}^{\top} \mathbf{y}_{\mathcal{S}}.$$

If for a specific coefficient  $j \in [p]$ , with  $k \ll n$  samples removed,

$$\widehat{oldsymbol{eta}}_j > 0$$
 but  $(\widehat{oldsymbol{eta}}_{\mathcal{S}})_j < 0$ 

our trust in the model may be questionable



In recent years several authors emphasized the need to assess how subsets of training samples collectively affect a learned model,

In recent years several authors emphasized the need to assess how subsets of training samples collectively affect a learned model,

- Koh et al, On the accuracy of influence functions for measuring group effects, NeurIPS 19'
- Basu et al, On second-order group influence functions for black-box predictions, ICML 20'
- Hu et al, Most influential subset selection: Challenges, promises, and beyond, NeurIPS 24'

In recent years several authors emphasized the need to assess how subsets of training samples collectively affect a learned model,

- Koh et al, On the accuracy of influence functions for measuring group effects, NeurIPS 19'
- Basu et al, On second-order group influence functions for black-box predictions, ICML 20'
- Hu et al, Most influential subset selection: Challenges, promises, and beyond, NeurIPS 24'

Broderick et al, 2020, considered a more stringent (worst-case) form of robustness, called robustness auditing:

In recent years several authors emphasized the need to assess how subsets of training samples collectively affect a learned model,

- Koh et al, On the accuracy of influence functions for measuring group effects, NeurIPS 19'
- Basu et al, On second-order group influence functions for black-box predictions, ICML 20'
- Hu et al, Most influential subset selection: Challenges, promises, and beyond, NeurIPS 24'

Broderick et al, 2020, considered a more stringent (worst-case) form of robustness, called robustness auditing:

estimated parameters / predictions need to be robust to removal of *any* subset of k samples

[Rubinstein and Hopkins 25']

**Definition (OLS)**: Robustness of  $\widehat{\beta}$  to k sample removals in a fixed direction  $\mathbf{v}$  is measured by

$$\Delta_k(\mathbf{v}) = \max_{\mathcal{S} \subseteq [n], |\mathcal{S}| = n - k} \langle \widehat{eta} - \widehat{eta}_{\mathcal{S}}, \mathbf{v} \rangle$$

[Rubinstein and Hopkins 25']

**Definition (OLS)**: Robustness of  $\widehat{\beta}$  to k sample removals in a fixed direction  $\mathbf{v}$  is measured by

$$\Delta_k(\mathbf{v}) = \max_{\mathcal{S} \subseteq [n], |\mathcal{S}| = n-k} \langle \widehat{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}}_{\mathcal{S}}, \mathbf{v} \rangle$$

In particular  $\Delta_k(e_j)$  measures sensitivity of j-th coefficient to removal of k samples.

[Rubinstein and Hopkins 25']

**Definition (OLS)**: Robustness of  $\widehat{\beta}$  to k sample removals in a fixed direction  $\mathbf{v}$  is measured by

$$\Delta_k(\mathbf{v}) = \max_{\mathcal{S} \subseteq [n], |\mathcal{S}| = n-k} \langle \widehat{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}}_{\mathcal{S}}, \mathbf{v} \rangle$$

In particular  $\Delta_k(e_j)$  measures sensitivity of j-th coefficient to removal of k samples.

#### **Key Questions:**

- Practical: For a given dataset and a given k, how large can  $\Delta_k(\mathbf{v})$  be ?



[Rubinstein and Hopkins 25']

**Definition (OLS)**: Robustness of  $\widehat{\beta}$  to k sample removals in a fixed direction  $\mathbf{v}$  is measured by

$$\Delta_k(\mathbf{v}) = \max_{\mathcal{S} \subseteq [n], |\mathcal{S}| = n - k} \langle \widehat{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}}_{\mathcal{S}}, \mathbf{v} \rangle$$

In particular  $\Delta_k(e_j)$  measures sensitivity of j-th coefficient to removal of k samples.

#### **Key Questions:**

- Practical: For a given dataset and a given k, how large can  $\Delta_k(\mathbf{v})$  be ?
- Theoretical: How large should we expect  $\Delta_k(\mathbf{v})$  to be under reasonable assumptions about the data?



- Exact computation of  $\Delta_k(\mathbf{v})$ :  $\binom{n}{k}$  subsets is computationally intractable

- Exact computation of  $\Delta_k(\mathbf{v})$ :  $\binom{n}{k}$  subsets is computationally intractable

approximations as well as upper/lower bounds for  $\Delta_k(\mathbf{v})$ 

## Robustness Auditing

- Exact computation of  $\Delta_k(\mathbf{v})$ :  $\binom{n}{k}$  subsets is computationally intractable

approximations as well as upper/lower bounds for  $\Delta_k(\mathbf{v})$ 

- Broderick, T., Giordano, R., and Meager, AMIP, 20'
- Kuschnig, Zens and Cuaresma, 21'.
- Moitra and Rohatgi, 23'
- Rubinstein and Hopkins, ACRE, 25'

## Robustness Auditing

- Exact computation of  $\Delta_k(\mathbf{v})$ :  $\binom{n}{k}$  subsets is computationally intractable

approximations as well as upper/lower bounds for  $\Delta_k(\mathbf{v})$ 

- Broderick, T., Giordano, R., and Meager, AMIP, 20'
- Kuschnig, Zens and Cuaresma, 21'.
- Moitra and Rohatgi, 23'
- Rubinstein and Hopkins, ACRE, 25'

Several of these works show that on various datasets, OLS is *not* robust to removal of even just a handful of samples.

[Angelucci and De Giorgi 2009']

Econometric study on Mexico's progresa aid program

[Angelucci and De Giorgi 2009']

Econometric study on Mexico's progresa aid program 506 rural villages: 320 participated in program, 186 control group.

[Angelucci and De Giorgi 2009']

Econometric study on Mexico's progresa aid program

506 rural villages: 320 participated in program, 186 control group.

Poor households in participating villages received financial support

[Angelucci and De Giorgi 2009']

Econometric study on Mexico's progresa aid program 506 rural villages: 320 participated in program, 186 control group. Poor households in participating villages received financial support Effect of program estimated by a linear regression

$$y = \beta_0 + \beta_1 x_1 + \sum_{j=2}^{17} \beta_j x_j,$$

[Angelucci and De Giorgi 2009']

Econometric study on Mexico's progresa aid program 506 rural villages: 320 participated in program, 186 control group. Poor households in participating villages received financial support Effect of program estimated by a linear regression

$$y = \beta_0 + \beta_1 x_1 + \sum_{j=2}^{17} \beta_j x_j,$$

response y - total household consumption in pesos,

[Angelucci and De Giorgi 2009']

Econometric study on Mexico's progresa aid program 506 rural villages: 320 participated in program, 186 control group. Poor households in participating villages received financial support Effect of program estimated by a linear regression

$$y = \beta_0 + \beta_1 x_1 + \sum_{j=2}^{17} \beta_j x_j,$$

response y - total household consumption in pesos,  $x_1$  - binary treatment variable (1 if village participated in Progresa),  $x_2, \ldots, x_{17}$  additional covariates of each household.

[Angelucci and De Giorgi 2009']

Econometric study on Mexico's progresa aid program 506 rural villages: 320 participated in program, 186 control group. Poor households in participating villages received financial support Effect of program estimated by a linear regression

$$y = \beta_0 + \beta_1 x_1 + \sum_{j=2}^{17} \beta_j x_j,$$

response y - total household consumption in pesos,  $x_1$  - binary treatment variable (1 if village participated in Progresa),  $x_2, \ldots, x_{17}$  additional covariates of each household. In non-poor households,  $\beta_1$  captures (indirect) effect of the

Progresa aid program.

### Cash Transfers Dataset

Period	Poor	n	$\widehat{eta}_1$	AMIP
8	Υ	10781	16.53	225
8	Ν	4543	-5.53	5
9	Υ	9489	28.65	321
9	Ν	3769	23.19	21
10	Υ	10368	32.52	570
10	N	4191	21.12	26

6 rows: three time periods  $\times$  two groups (poor/non poor)

### Cash Transfers Dataset

Period	Poor	n	$\widehat{eta}_{1}$	AMIP
8	Υ	10781	16.53	225
8	Ν	4543	-5.53	5
9	Υ	9489	28.65	321
9	Ν	3769	23.19	21
10	Υ	10368	32.52	570
10	N	4191	21.12	26

6 rows: three time periods  $\times$  two groups (poor/non poor) Column 5: size of smallest subset found by AMIP whose removal changes the sign of  $\widehat{\beta}_1$ 

How come only k = 5 - 30 samples out of n > 4000 can change the sign of an OLS coefficient? (even when it is considered as statistically significant by standard tests)

How come only k = 5 - 30 samples out of n > 4000 can change the sign of an OLS coefficient? (even when it is considered as statistically significant by standard tests)

This talk:

How come only k=5-30 samples out of n>4000 can change the sign of an OLS coefficient? (even when it is considered as statistically significant by standard tests)

#### This talk:

How come only k=5-30 samples out of n>4000 can change the sign of an OLS coefficient? (even when it is considered as statistically significant by standard tests)

#### This talk:

Theoretical analysis: robustness auditing of OLS

- Derive conditions (and understanding) when would OLS be provably robust to k-sample removals

How come only k=5-30 samples out of n>4000 can change the sign of an OLS coefficient? (even when it is considered as statistically significant by standard tests)

#### This talk:

- Derive conditions (and understanding) when would OLS be provably robust to k-sample removals
- conversely, when would OLS be provably non-robust to k sample removals

How come only k=5-30 samples out of n>4000 can change the sign of an OLS coefficient? (even when it is considered as statistically significant by standard tests)

#### This talk:

- Derive conditions (and understanding) when would OLS be provably robust to k-sample removals
- conversely, when would OLS be provably non-robust to k sample removals
- Revisit the cash transfer dataset and identify potential causes

How come only k=5-30 samples out of n>4000 can change the sign of an OLS coefficient? (even when it is considered as statistically significant by standard tests)

#### This talk:

- Derive conditions (and understanding) when would OLS be provably robust to k-sample removals
- conversely, when would OLS be provably non-robust to k sample removals
- Revisit the cash transfer dataset and identify potential causes
- Implications for practitioners



### Question:

When would OLS be robust to removal of k samples?

### Question:

When would OLS be robust to removal of k samples?

Assume samples  $(\mathbf{x}_i, y_i)$  are i.i.d. from joint distribution  $P(\mathbf{x}, y)$ ,

### Question:

When would OLS be robust to removal of k samples?

Assume samples  $(\mathbf{x}_i, y_i)$  are i.i.d. from joint distribution  $P(\mathbf{x}, y)$ ,

Analyze two models for  $P(\mathbf{x}, y)$ :

### Question:

When would OLS be robust to removal of k samples?

Assume samples  $(\mathbf{x}_i, y_i)$  are i.i.d. from joint distribution  $P(\mathbf{x}, y)$ ,

Analyze two models for  $P(\mathbf{x}, y)$ :

- Model 1: General P with mild regularity conditions

### Question:

When would OLS be robust to removal of k samples?

Assume samples  $(\mathbf{x}_i, y_i)$  are i.i.d. from joint distribution  $P(\mathbf{x}, y)$ ,

Analyze two models for  $P(\mathbf{x}, y)$ :

- Model 1: General P with mild regularity conditions
- Model 2: Gaussian Linear Model.

### OLS Robustness

### Model 1 - General

 $P(\mathbf{x}, y)$  - probability distribution over  $\mathbb{R}^{p+1}$  such that

- 1.  ${\bf x}$  is sub-Gaussian, zero mean, and has positive-definite covariance  ${\bf \Sigma} = \mathbb{E}[{\bf x}{\bf x}^{\top}]$
- 2. y is sub-Gaussian

### OLS Robustness

### Model 2 - Gaussian Linear

 $y = \boldsymbol{\beta}^{\top} \cdot \mathbf{x} + \varepsilon$ , where  $\boldsymbol{\beta} \in \mathbb{R}^p$  is deterministic and

- 1.  $\mathbf{x} \sim \mathcal{N}(0, \Sigma)$ ,  $\Sigma \succ 0$  is positive definite
- 2.  $\varepsilon$  is sub-Gaussian, mean zero, and independent of x

Optimal prediction under squared loss is conditional mean

Optimal prediction under squared loss is *conditional mean*  $\mathbb{E}[y \,|\, \mathbf{x}]$ 

Optimal prediction under squared loss is conditional mean

$$\mathbb{E}[y \mid \mathbf{x}]$$

Optimal OLS solution is

$$oldsymbol{eta}^{ extsf{OLS}} = \Sigma^{-1} \cdot \mathbb{E}[\mathbf{x}y]$$

Optimal prediction under squared loss is conditional mean

$$\mathbb{E}[y \mid \mathbf{x}]$$

Optimal OLS solution is

$$oldsymbol{eta}^{ exttt{OLS}} = \Sigma^{-1} \cdot \mathbb{E}[\mathbf{x} y]$$

The general model 1 is *mis-specified*, since  $\mathbb{E}[y \mid \mathbf{x}]$  may be a non-linear function of  $\mathbf{x}$ 



Optimal prediction under squared loss is conditional mean

$$\mathbb{E}[y \mid \mathbf{x}]$$

Optimal OLS solution is

$$oldsymbol{eta}^{ ext{OLS}} = \Sigma^{-1} \cdot \mathbb{E}[\mathbf{x} y]$$

The general model 1 is *mis-specified*, since  $\mathbb{E}[y \mid \mathbf{x}]$  may be a non-linear function of  $\mathbf{x}$ 

Under Gaussian-Linear Model 2,  $oldsymbol{eta}^{ extsf{oLS}}=oldsymbol{eta}$ 



### Theoretical Results under Model 1

### **Sub-Gaussian norm of** *y*:

$$||y||_{\psi_2} = \inf \{t > 0 : \mathbb{E}[\exp(|y|^2/t^2)] \le e\}.$$

### Theoretical Results under Model 1

### Sub-Gaussian norm of y:

$$||y||_{\psi_2} = \inf \{ t > 0 : \mathbb{E}[\exp(|y|^2/t^2)] \le e \}.$$

#### Sub-Gaussian norm of vector x:

$$\|\mathbf{x}\|_{\psi_2} = \sup_{\mathbf{v} \in \mathbb{S}^{p-1}} \|\mathbf{v}^{\top}\mathbf{x}\|_{\psi_2},$$

# Non-asymptotic bound

 $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  i.i.d. from Model 1.

## Non-asymptotic bound

 $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  i.i.d. from Model 1.

#### $\mathsf{Theorem}$

Assume  $k \le n/2$ ,  $p \le n - k$ , and define

$$\eta = \|\boldsymbol{\Sigma}^{-1/2}\| \Big(1 + \|\boldsymbol{\Sigma}^{-1/2}\mathbf{x}\|_{\psi_2}^4 \Big) \|\boldsymbol{y}\|_{\psi_2}^2.$$

There exist constants C, c > 0 such that with probability at least  $1 - 7(n/k)^{-ck}$ .

$$\max_{\mathcal{S} \subseteq [n]: |\mathcal{S}| \ge n-k} \|\widehat{\boldsymbol{\beta}}_{\mathcal{S}} - \widehat{\boldsymbol{\beta}}\| \le C\eta \sqrt{\frac{k \log(en/k)}{n-k}}$$

for all p, n, k satisfying

$$\|C\|\Sigma^{-1/2}\mathbf{x}\|_{\psi_2}^2\sqrt{rac{p}{n-k}}\leq 1.$$

18/36

### OLS Robustness Under Model 1

Asymptotics as  $n \to \infty$ ,  $k = k_n$  and  $p = p_n$  may tend to infinity.

### OLS Robustness Under Model 1

Asymptotics as  $n \to \infty$ ,  $k = k_n$  and  $p = p_n$  may tend to infinity.

Assume  $\|\Sigma^{-1/2}\|$ ,  $\|\mathbf{x}\|_{\psi_2}$ , and  $\|y\|_{\psi_2}$  remain bounded as  $n \to \infty$ .

Asymptotics as  $n \to \infty$ ,  $k = k_n$  and  $p = p_n$  may tend to infinity.

Assume  $\|\Sigma^{-1/2}\|$ ,  $\|\mathbf{x}\|_{\psi_2}$ , and  $\|y\|_{\psi_2}$  remain bounded as  $n \to \infty$ .

#### Theorem

 $(\mathbf{x}_1,y_1),\ldots,(\mathbf{x}_n,y_n)$  i.i.d. from Model 1. Assume that as  $n\to\infty$ ,  $\limsup \|\Sigma^{-1/2}\mathbf{x}\|_{\psi_2}^2 \sqrt{p/n} < c$ , where c>0 is an absolute constant. Then, if  $\mathbf{k}/n\to0$ ,

$$\max_{\mathcal{S}\subseteq[n],|\mathcal{S}|\geq n-k}\|\widehat{\boldsymbol{\beta}}_{\mathcal{S}}-\widehat{\boldsymbol{\beta}}\|\xrightarrow{p}0$$



In simple words: for "well-behaved" data, if  $k/n \to 0$ , OLS is robust to k sample removals, regardless of model mis-specification, and regardless of data dimension.

In simple words: for "well-behaved" data, if  $k/n \to 0$ , OLS is robust to k sample removals, regardless of model mis-specification, and regardless of data dimension.

In fact, theorem allows  $p \to \infty$  provided that  $p/n \to 0$ , as in this case the condition of the theorem is satisfied.

In simple words: for "well-behaved" data, if  $k/n \to 0$ , OLS is robust to k sample removals, regardless of model mis-specification, and regardless of data dimension.

In fact, theorem allows  $p \to \infty$  provided that  $p/n \to 0$ , as in this case the condition of the theorem is satisfied.

Robustness measure  $\Delta_k(\mathbf{v})$  converges to zero, uniformly in  $\mathbf{v}$ :

$$\Delta_k(\mathbf{v}) = \max_{\mathcal{S} \subseteq [n], |\mathcal{S}| = n-k} \langle \widehat{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}}_{\mathcal{S}}, \mathbf{v} \rangle \leq \max_{\mathcal{S} \subseteq [n], |\mathcal{S}| \geq n-k} \|\widehat{\boldsymbol{\beta}}_{\mathcal{S}} - \widehat{\boldsymbol{\beta}}\|.$$

Hence, if  $k/n \rightarrow 0$ ,

$$\sup_{\boldsymbol{v}\in\mathbb{S}^{p-1}}\Delta_k(\boldsymbol{v})\stackrel{p}{\longrightarrow} 0.$$



## Consistency under sample removals / Model 1

#### **Theorem**

Under same conditions above, and additional assumption that  $\kappa(\Sigma)$  and  $\|\beta\|$  remain bounded, as  $n \to \infty$  and  $(p+k)/n \to 0$ ,

$$\max_{\mathcal{S}\subseteq [n], |\mathcal{S}|\geq n-k} \|\widehat{\boldsymbol{\beta}}_{\mathcal{S}} - \boldsymbol{\beta}\| \xrightarrow{p} 0.$$

## Consistency under sample removals / Model 1

#### Theorem

Under same conditions above, and additional assumption that  $\kappa(\Sigma)$  and  $\|\beta\|$  remain bounded, as  $n \to \infty$  and  $(p+k)/n \to 0$ ,

$$\max_{\mathcal{S}\subseteq[n],|\mathcal{S}|\geq n-k}\|\widehat{\boldsymbol{\beta}}_{\mathcal{S}}-\boldsymbol{\beta}\|\xrightarrow{p} 0.$$

Maximum number of samples that can be removed while keeping OLS error rate:

$$k \log \left(\frac{n}{k}\right) \ll \sqrt{p}$$



# Robustness / Model 2

### Robustness / Model 2

 $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  i.i.d. from Gaussian Linear Model 2.

### Robustness / Model 2

 $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  i.i.d. from Gaussian Linear Model 2.

#### Theorem

There exist constants C, c > 0 such that if

$$Ck \le n, \quad t \ge 0, \quad \text{and} \quad \sqrt{\frac{p}{n-k}} + t < c,$$

then with probability at least  $1 - 4e^{-c(n-k)t^2}$ ,

$$\max_{\substack{\mathcal{S}\subseteq[n],\\|\mathcal{S}|\geq n-k}}\|\widehat{\boldsymbol{\beta}}_{\mathcal{S}}-\boldsymbol{\beta}\|\leq C\|\Sigma^{-1/2}\|\|\boldsymbol{\varepsilon}\|_{\psi_2}\left(\frac{k\log n}{n-k}+\sqrt{\frac{p}{n-k}}+t\right)$$

Boaz Nadler

### Robustness under Linear Model

We conjecture theorem is optimal up to logarithmic factors.

#### Robustness under Linear Model

We conjecture theorem is optimal up to logarithmic factors.

The error rate of  $\max_{S\subseteq[n],|S|\geq n-k}\|\widehat{\beta}_S-\beta\|$  matches that of OLS on the full dataset if

$$k \ll \frac{\sqrt{np}}{\log n}$$

#### Robustness under Linear Model

We conjecture theorem is optimal up to logarithmic factors.

The error rate of  $\max_{S\subseteq[n],|S|\geq n-k}\|\widehat{\beta}_S-\beta\|$  matches that of OLS on the full dataset if

$$k \ll \frac{\sqrt{np}}{\log n}$$

Under linear model, OLS can tolerate the removal of significantly more samples than under the general model

[Rubinstein and Hopkins, ICLR, 25'] ACRE= Algorithm for Certifying Robustness Efficiently

[Rubinstein and Hopkins, ICLR, 25']

ACRE= Algorithm for Certifying Robustness Efficiently

For a fixed  $\mathbf{v}$ , ACRE computes upper and lower bounds  $U_k(\mathbf{v})$  and  $L_k(\mathbf{v})$  such that without any modeling assumptions,

$$L_k(\mathbf{v}) \leq \Delta_k(\mathbf{v}) \leq U_k(\mathbf{v}).$$

Rubinstein and Hopkins derived following theoretical guarantee:

Rubinstein and Hopkins derived following theoretical guarantee: Under a variant of Model 2, with less restrictive conditions, there exists a threshold

$$K = \widetilde{\Theta}\left(\min\left(\frac{n}{\sqrt{p}}, \frac{n^2}{p^2}\right)\right)$$

such that for all  $k \leq K$ , with high probability,

$$\frac{U_k(\mathbf{v})}{L_k(\mathbf{v})} = 1 + \widetilde{O}\left(\frac{p + k\sqrt{p}}{n}\right).$$

Rubinstein and Hopkins derived following theoretical guarantee: Under a variant of Model 2, with less restrictive conditions, there exists a threshold

$$K = \widetilde{\Theta}\left(\min\left(\frac{n}{\sqrt{p}}, \frac{n^2}{p^2}\right)\right)$$

such that for all  $k \leq K$ , with high probability,

$$\frac{U_k(\mathbf{v})}{L_k(\mathbf{v})} = 1 + \widetilde{O}\left(\frac{p + k\sqrt{p}}{n}\right).$$

When  $p + k\sqrt{p} \ll n$ , the upper and lower bounds  $U_k(\mathbf{v})$  and  $L_k(\mathbf{v})$  are tight, so ACRE accurately measures robustness to removals.



Comparison to our theoretical results:

Comparison to our theoretical results:

Under model 2, OLS is *provably* robust to removals in a *broader* parameter regime  $p + k \ll n$ .

Comparison to our theoretical results:

Under model 2, OLS is *provably* robust to removals in a *broader* parameter regime  $p + k \ll n$ .

**Open Question:** whether the upper and lower bounds of ACRE remain tight in more general regimes, in particular where OLS is non-robust.

For the general model 1:

For the general model 1:

proof is "standard": based on concentration inequalities and union bounds.

For the general model 1:

proof is "standard": based on concentration inequalities and union bounds.

For Gaussian linear model 2:

For the general model 1:

proof is "standard": based on concentration inequalities and union bounds.

For Gaussian linear model 2:

sharper results require more involved proof.

For the general model 1:

proof is "standard": based on concentration inequalities and union bounds.

For Gaussian linear model 2:

sharper results require more involved proof.

Careful use of Gaussian comparison inequalities

For the general model 1:

proof is "standard": based on concentration inequalities and union bounds.

For Gaussian linear model 2:

sharper results require more involved proof.

Careful use of Gaussian comparison inequalities

#### **Open Question:**

For the general model 1:

proof is "standard": based on concentration inequalities and union bounds.

For Gaussian linear model 2:

sharper results require more involved proof.

Careful use of Gaussian comparison inequalities

#### **Open Question:**

Derive robustness guarantees for other models

### Non-robustness for $k \propto n$

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$$
 i.i.d. from Model 2. Set  $\alpha = k/n$ .

### Non-robustness for $k \propto n$

 $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  i.i.d. from Model 2. Set  $\alpha = k/n$ .

#### Theorem

Fix  $\mathbf{v} \in \mathbb{S}^{p-1}$ . Assume that p < k and  $\gamma = p/(n-k) < 1/4$ . There exist absolute constants C, c > 0 such that if  $\alpha = k/n \le c$ , then with probability at least  $1 - 17e^{-c(n-k)t^2}$ 

$$\Delta_k(\mathbf{v}) \geq \left\| \Sigma^{-1/2} \mathbf{v} \right\| \left( \mathbb{E}[\varepsilon z \, \mathbb{1}(\varepsilon z > q_{1-lpha+t)})] - \frac{C(t+\sqrt{\gamma})}{(C-\sqrt{\gamma}-t)^2} \right),$$

for any  $t \in (0, \min\{\alpha, 1/2 - \alpha\})$ . Here,  $z \sim \mathcal{N}(0, 1)$  is independent of  $\varepsilon$ , and  $q_{1-\alpha+t}$  is the  $(1-\alpha+t)$ -quantile of  $\varepsilon z$ .



Boaz Nadler Robustness Auditing

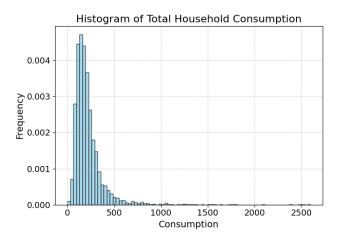
## Robustness and Consistency Regimes for OLS / Model 2

Region	$k, p \ll n$	$k \ll n, p \asymp n$	$k \asymp n, p \ll n$	$k \asymp p \asymp n$
Robust	$\checkmark$	$\checkmark$	×	×
Consistent	✓	×	✓	×

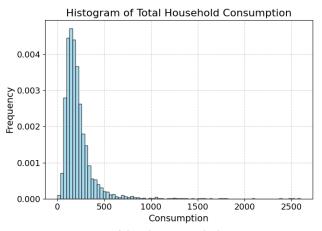
Response Y = total household consumption in pesos

Response Y= total household consumption in pesos In all 6 datasets,  $\mathbb{E}[Y]\approx 200$ , and  $\mathrm{std}(Y)$  comparable to  $\mathbb{E}[Y]$ 

Response Y= total household consumption in pesos In all 6 datasets,  $\mathbb{E}[Y]\approx 200$ , and  $\mathrm{std}(Y)$  comparable to  $\mathbb{E}[Y]$ 



Response Y = total household consumption in pesosIn all 6 datasets,  $\mathbb{E}[Y] \approx 200$ , and  $\operatorname{std}(Y)$  comparable to  $\mathbb{E}[Y]$ 

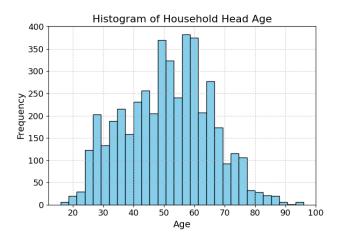


#### Cash Transfers Dataset

Explanatory variables  $x_1, \ldots, x_{17}$  are all "well-behaved", some are categorical,

#### Cash Transfers Dataset

Explanatory variables  $x_1,\ldots,x_{17}$  are all "well-behaved", some are categorical,



### Cash Transfers Dataset

Period	Poor	n	$\widehat{eta}_{1}$	AMIP	$\mu_{y}$	$\sigma_{y}$	$>$ 5 $\sigma_y$	$>$ $10\sigma_y$
8	Υ	10781	16.53	225	170	126	48	12
8	Ν	4543	-5.53	5	219	172	29	5
9	Υ	9489	28.65	321	176	182	48	15
9	Ν	3769	23.19	21	226	273	20	9
10	Υ	10368	32.52	570	172	156	56	13
10	N	4191	21.12	26	217	267	19	7

 $\mu_y$  and  $\sigma_y$ : empirical mean and standard deviation of response y. Last 2 columns: # samples larger than  $\mu_y$  by  $> 5\sigma_y$  or  $> 10\sigma_y$ 

### Cash Transfers Dataset

Period	Poor	n	$\widehat{eta}_{1}$	AMIP	$\mu_{y}$	$\sigma_{y}$	$>$ 5 $\sigma_y$	$>$ $10\sigma_y$
8	Υ	10781	16.53	225	170	126	48	12
8	N	4543	-5.53	5	219	172	29	5
9	Υ	9489	28.65	321	176	182	48	15
9	Ν	3769	23.19	21	226	273	20	9
10	Υ	10368	32.52	570	172	156	56	13
10	N	4191	21.12	26	217	267	19	7

 $\mu_y$  and  $\sigma_y$ : empirical mean and standard deviation of response y. Last 2 columns: # samples larger than  $\mu_y$  by  $> 5\sigma_y$  or  $> 10\sigma_y$  Y is extremely heavy tailed



Period	Poor	n	$\widehat{eta}_{1}$	AMIP	$\mu_{y}$	$\mu_{\mathbf{y}}$ amip	$y_{\sf max}^{\sf AMIP}$
8	Υ	10781	16.53	225	170	572	4380
8	Ν	4543	-5.53	5	219	2018	2483
9	Υ	9489	28.65	321	176	580	5117
9	Ν	3769	23.19	21	226	2670	5801
10	Υ	10368	32.52	570	172	412	5080
10	N	4191	21.12	26	217	2154	7470

Period	Poor	n	$\widehat{eta}_1$	AMIP	$\mu_{y}$	$\mu_{\mathbf{y}}$ amip	$y_{max}^{AMIP}$
8	Υ	10781	16.53	225	170	572	4380
8	Ν	4543	-5.53	5	219	2018	2483
9	Υ	9489	28.65	321	176	580	5117
9	Ν	3769	23.19	21	226	2670	5801
10	Υ	10368	32.52	570	172	412	5080
10	N	4191	21.12	26	217	2154	7470

AMIP removes samples with extreme Y values

Period	Poor	n	$\widehat{eta}_1$	AMIP	$\mu_{y}$	$\mu_{y}$ amip	$y_{max}^{AMIP}$
8	Υ	10781	16.53	225	170	572	4380
8	N	4543	-5.53	5	219	2018	2483
9	Υ	9489	28.65	321	176	580	5117
9	Ν	3769	23.19	21	226	2670	5801
10	Υ	10368	32.52	570	172	412	5080
10	N	4191	21.12	26	217	2154	7470

AMIP removes samples with extreme Y values

Perhaps not surprising few samples suffice to reverse sign  $\hat{eta}_1$ 

Suppose instead of OLS, we fit linear model under Huber loss

Suppose instead of OLS, we fit linear model under Huber loss

$$\widehat{oldsymbol{eta}}^{\mathsf{Huber}} = \operatorname*{\mathsf{argmin}}_{oldsymbol{eta} \in \mathbb{R}^p} \sum_{i=1}^n h_{ au}(oldsymbol{eta}^{ op} \mathbf{x}_i - y_i)$$

Suppose instead of OLS, we fit linear model under Huber loss

$$\widehat{oldsymbol{eta}}^{\mathsf{Huber}} = \operatorname*{\mathsf{argmin}}_{oldsymbol{eta} \in \mathbb{R}^p} \sum_{i=1}^n h_ au(oldsymbol{eta}^ op \mathbf{x}_i - y_i)$$

where  $h_{\tau}$  is the Huber loss function, given by

$$h_{ au}(z) = egin{cases} rac{z^2}{2} & |z| \leq au, \ au\Big(|z| - rac{ au}{2}\Big) & |z| > au. \end{cases}$$

Suppose instead of OLS, we fit linear model under Huber loss

$$\widehat{oldsymbol{eta}}^{\mathsf{Huber}} = \operatorname*{\mathsf{argmin}}_{oldsymbol{eta} \in \mathbb{R}^p} \sum_{i=1}^n h_ au(oldsymbol{eta}^ op \mathbf{x}_i - y_i)$$

where  $h_{\tau}$  is the Huber loss function, given by

$$h_{ au}(z) = egin{cases} rac{z^2}{2} & |z| \leq au, \ au\Big(|z| - rac{ au}{2}\Big) & |z| > au. \end{cases}$$

au>0 controls the transition from squared loss to absolute loss. In our experiments we took au=1.

Period	Poor	n	$\widehat{eta}_1$	AMIP	$\widehat{eta}_1^{Huber}$	AMIP Huber
8	Υ	10781	16.53	225	16.55	725
8	Ν	4543	-5.53	5	-5.53	30
9	Υ	9489	28.65	321	27.92	915
9	Ν	3769	23.19	21	22.15	228
10	Υ	10368	32.52	570	31.31	1242
10	N	4191	21.12	26	19.06	217

Period	Poor	n	$\widehat{eta}_1$	AMIP	$\widehat{eta}_1^{Huber}$	AMIP Huber
8	Υ	10781	16.53	225	16.55	725
8	Ν	4543	-5.53	5	-5.53	30
9	Υ	9489	28.65	321	27.92	915
9	Ν	3769	23.19	21	22.15	228
10	Υ	10368	32.52	570	31.31	1242
10	N	4191	21.12	26	19.06	217

 $\widehat{\beta}_1$  and  $\widehat{\beta}_1^{\mathsf{Huber}}$  -treatment effect under OLS and Huber regression

Period	Poor	n	$\widehat{eta}_{1}$	AMIP	$\widehat{eta}_1^{Huber}$	AMIP Huber
8	Υ	10781	16.53	225	16.55	725
8	Ν	4543	-5.53	5	-5.53	30
9	Υ	9489	28.65	321	27.92	915
9	Ν	3769	23.19	21	22.15	228
10	Υ	10368	32.52	570	31.31	1242
10	Ν	4191	21.12	26	19.06	217

 $\widehat{\beta}_1$  and  $\widehat{\beta}_1^{\text{Huber}}$  -treatment effect under OLS and Huber regression Columns "AMIP" and "AMIP Huber": size of the smallest subset identified by AMIP whose removal reverses the sign of  $\widehat{\beta}_1$  under each method.

Period	Poor	n	$\widehat{eta}_{1}$	AMIP	$\widehat{eta}_1^{Huber}$	AMIP Huber
8	Υ	10781	16.53	225	16.55	725
8	Ν	4543	-5.53	5	-5.53	30
9	Υ	9489	28.65	321	27.92	915
9	Ν	3769	23.19	21	22.15	228
10	Υ	10368	32.52	570	31.31	1242
10	N	4191	21.12	26	19.06	217

 $\widehat{eta}_1$  and  $\widehat{eta}_1^{\text{Huber}}$  -treatment effect under OLS and Huber regression Columns "AMIP" and "AMIP Huber": size of the smallest subset identified by AMIP whose removal reverses the sign of  $\widehat{eta}_1$  under each method.

Assuming AMIP approximation is accurate Huber regression substantially more robust

- Robustness Auditing: important to enhance trust in a learned model. Framework goes beyond influence function of individual samples.

- Robustness Auditing: important to enhance trust in a learned model. Framework goes beyond influence function of individual samples.
- Presented theoretical analysis of robustness auditing for OLS.

- Robustness Auditing: important to enhance trust in a learned model. Framework goes beyond influence function of individual samples.
- Presented theoretical analysis of robustness auditing for OLS.
- Well behaved data and  $k \ll n$ , OLS is provably robust to sample removals.

- Robustness Auditing: important to enhance trust in a learned model. Framework goes beyond influence function of individual samples.
- Presented theoretical analysis of robustness auditing for OLS.
- Well behaved data and  $k \ll n$ , OLS is provably robust to sample removals.
- Implications: If removal of  $k \ll n$  samples significantly changes linear model: need to carefully inspect potential reasons: heavy tails, outliers, non-i.i.d. data, etc.

- Robustness Auditing: important to enhance trust in a learned model. Framework goes beyond influence function of individual samples.
- Presented theoretical analysis of robustness auditing for OLS.
- Well behaved data and  $k \ll n$ , OLS is provably robust to sample removals.
- Implications: If removal of  $k \ll n$  samples significantly changes linear model: need to carefully inspect potential reasons: heavy tails, outliers, non-i.i.d. data, etc.
- Multiple future directions: other models, less restrictive assumptions, heavy tailed distributions, outliers, etc.

- Robustness Auditing: important to enhance trust in a learned model. Framework goes beyond influence function of individual samples.
- Presented theoretical analysis of robustness auditing for OLS.
- Well behaved data and  $k \ll n$ , OLS is provably robust to sample removals.
- Implications: If removal of  $k \ll n$  samples significantly changes linear model: need to carefully inspect potential reasons: heavy tails, outliers, non-i.i.d. data, etc.
- Multiple future directions: other models, less restrictive assumptions, heavy tailed distributions, outliers, etc.

#### Thank You!

