# Two-Stage Robust and Sparse Distributed Statistical Inference for Large-Scale Data

Emadaldin Mozafari-Majd, *Student Member, IEEE,* Visa Koivunen, *Fellow, IEEE*

*Abstract*—In this paper, we address the problem of conducting statistical inference in settings involving large-scale data that may be high-dimensional and contaminated by outliers. The high volume and dimensionality of the data require distributed processing and storage solutions. We propose a two-stage distributed and robust statistical inference procedures coping with high-dimensional models by promoting sparsity. In the first stage, known as model selection, relevant predictors are locally selected by applying robust Lasso estimators to the distinct subsets of data. The variable selections from each computation node are then fused by a voting scheme to find the sparse basis for the complete data set. It identifies the relevant variables in a robust manner. In the second stage, the developed statistically robust and computationally efficient bootstrap methods are employed. The actual inference constructs confidence intervals, finds parameter estimates and quantifies standard deviation. Similar to stage 1, the results of local inference are communicated to the fusion center and combined there. By using analytical methods, we establish the favorable statistical properties of the robust and computationally efficient bootstrap methods including consistency for a fixed number of predictors, and robustness. The proposed two-stage robust and distributed inference procedures demonstrate reliable performance and robustness in variable selection, finding confidence intervals and bootstrap approximations of standard deviations even when data is high-dimensional and contaminated by outliers.

*Index Terms*—High-dimensional, large-scale data, big data, sparsity, robust estimator, Lasso, distributed computation and storage, fixed-point equations, information fusion.

## I. INTRODUCTION

**M**ASSIVE quantities of ubiquitous and heterogeneous data are generated by social media, smart phones, IoT, environmental monitoring, astronomical imaging devices and financial markets. Harnessing information from such large-scale data provides enterprises with meaningful insights into their performance and offers tremendous business opportunities. However, these benefits come with formidable challenges in handling storage, processing, acquisition and privacy concerns of high-speed and high volume data [1]. In order to remedy the storage and processing issues, distributed computation and storage solutions are preferred. In a variety of statistical inference applications, one is dealing with high-dimensional data where the number of explaining variables $p$ may be comparable or much larger than the number of observations $n$. Often, high-dimensional problems exhibit a lower dimensional structure such as sparsity or low-rank. Regularization may be necessary to address the ill-posed high-dimensional problems

The authors are with the Department of Signal Processing and Acoustics, Aalto University, FI-00076 Aalto, Finland (email: emadaldin.mozafarimajd@aalto.fi; visa.koivunen@aalto.fi).

and capture the parsimonious representation. In the context of linear regression, regularization by $\ell_1$-norm is known as the celebrated Lasso and performs simultaneous parameter estimation and model selection [2]. It improves the prediction accuracy by introducing some bias while reducing the variance. There are certain scenarios where regularization by $\ell_1$-norm may cause significant bias in estimated coefficients. Moreover, the exact characterization of the limiting distribution for the Lasso estimator is a challenging task. In recent years, several valid statistical inference procedures have been introduced to characterize the construction of confidence intervals and hypothesis testing for high-dimensional problems [3], [4], [5]. In general, the state-of-art procedures to address the uncertainty associated with parameter estimates belong to three main categories. First category concerns inference procedures based on bootstrapping. However, the conventional bootstrap methods fail to provide a reliable approximation to the distribution of Lasso estimator. In order to address this issue, Chatterjee and Lahiri [3], [6] proposed two alternative solutions, modified residual bootstrap Lasso and residual bootstrap adaptive Lasso. These solutions consistently estimate their limiting distributions and provide valid approximation of confidence intervals. The second category includes post-Lasso inference methods where the first stage involves model selection using Lasso and the actual inference is made in the second stage using the selected variables. This category includes sample splitting [7], bootstrap Lasso-OLS [5], bootstrap Lasso-partial Ridge [8] and post-selection inference [9], [10] methods. Herein, we extend the concept of Post-Lasso estimator to conduct statistical inference for large-scale data sets using distributed computation and storage. The third category is based on the debiased-Lasso method [4], [11] where the key idea is to remove the bias introduced by regularization from Lasso solution. These methods offer a concrete and general framework to quantify the uncertainty associated with parameters in high-dimensional settings, e.g. hypothesis testing and construction of confidence intervals. Other alternatives are covariance test [12], knockoff filter [13], the ridge projection and bias correction [14].

Another challenge in dealing with large-scale data is that the probability of observing outliers may increase as the dimensionality ($p$) and sample size ($n$) grows larger. Outliers may become masked and hence difficult to detect. Outliers may severely deteriorate the performance of ordinary least square and regularized least square estimators. Robust multivariate statistical procedures are required to cope with outlying observations and ensure the veracity of estimation, classification and decision making. In the context of linear regression, robust

regularized estimators are employed to ensure robustness and find sparse solutions simultaneously [15], [16], [17], [18], [19], [11].

Modern statistical inference procedures need to accommodate distributed storage and parallel computations to deal with high volume and high-dimensional data. Bootstrap is a powerful tool for quantifying the uncertainty of estimates, i.e., confidence intervals and hypothesis tests. However, the applicability of conventional bootstrap techniques in large-scale settings may not be feasible because of computational constraints. The bootstrap resamples may have the same size as the original large-scale data, and repeating estimation for each bootstrap replicate becomes prohibitively expensive. The Bag of Little Bootstraps (BLB) [20] offers a scalable and computationally efficient inference procedure for quantification of the uncertainty associated with estimates that accommodates distributed and parallel computing architectures. It subdivides the complete large-scale data into smaller distinct subsets. It can be considered as sampling without replacement from the complete data set. Then, bootstrap is applied to each subset and combines the inference results from each subset. However, the BLB procedure is highly sensitive to outliers. In order to overcome this issue, a statistically robust BLFRB approach that extends the idea of fast fixed-point computations of FRB method in [21] to MM-estimators was introduced in [22]. The modified bootstrap replicates are calculated by applying a linear correction factor to the one-step approximation of bootstrap replicates. Despite its robustness and computational efficiency, BLFRB performs unreliably in high-dimensional settings.

This paper focuses on statistical inference in large-scale settings with distributed architecture where data may be high-dimensional and contaminated by outliers. This problem has not been thoroughly studied in open literature. In particular, we propose two-stage statistical inference procedures that are robust to outlying observations and allow for using distributed storage and processing architectures for scalability. In the first stage, the relevant predictors are selected in two steps. First robust Lasso estimator is applied to distinct subsets of data in order to perform local variable selection. The variables for the whole data set are then found by applying a fusion rule to the selections from individual nodes at the fusion center or cloud. In the second stage, we conduct inference by constructing confidence intervals, finding point estimates of the selected parameters and their standard deviations based on the large-scale data. The developed distributed and low-complexity inference procedures that use linearly corrected one-step robust estimators are employed. In special cases with very high dimensionality ($p/n \approx 1$ or $p \gg b$, $p$ number of predictors, $n$ sample size and $b$ subsample size), one may accelerate the model selection by a preprocessing stage excluding the majority of irrelevant variables via a robust variable screening procedure on the distinct subsets of data stored at each node. We address this issue in **Supplemental Materials**. The methods proposed in this paper extend our previous work on two-stage robust and distributed inference procedures [23], [24] by deriving computationally more efficient estimation methods and establishing statistical properties

of the inference methods using analytical tools. We emphasize that our asymptotic analysis is restricted to the classical fixed $p$ setting. In [5], the asymptotic properties were established for a special case of post-Lasso estimator Lasso+mLS and the valid bootstrap approximation while allowing $p$ to grow at an exponential rate in $n$. However, their analysis does not cover the distributed and robust settings.

The main new contributions of the paper are summarized as follows:

- **T**wo-**S**tage **R**obust and **D**istributed inference employing the class of $\tau$-estimators called TSRD-$\tau$ is introduced. A robust $\tau$-Lasso sparsity promoting estimator is employed in the first stage.
- The proposed TSRD-$\tau$ employs a novel **R**obust and **S**calable linearly corrected **O**ne-step **B**ootstrap procedure using $\boldsymbol{\tau}$-estimator (RSOB-$\tau$) for performing the actual inference. Its computational complexity is reduced by efficient estimation of bootstrap replicates.
- **T**wo-**S**tage **R**obust and **D**istributed inference employing the class of MM-estimators called TSRD-MM is developed. The sparsity promoting estimator in the first stage is robust MM-Lasso. It extends the BLFRB procedure [22] by a robust variable selection stage promoting sparsity.
- Analytical results proving the robustness and consistency of the TSRD-$\tau$ method are derived for fixed $p$ and diverging $n$. In order to formally show the quantitative robustness of $\tau$-Lasso, its finite-sample breakdown is characterized.
- Extensive simulations are conducted to assess the performance of robust and distributed two-stage procedures in variable selection, bootstrap estimation of standard deviation, robustness of confidence intervals and computational complexity of the RSOB-$\tau$ procedure.
- Analytical results on the consistency and asymptotic normality of the RSOB-$\tau$ procedure are verified through computer simulations. Furthermore, the favorable theoretical findings on robustness of bootstrap replications are confirmed by extensive simulation studies and comparing them to [22].

The two-stage inference solutions presented in this work achieve scalability by performing inference over smaller distinct subsets of data in parallel, using multinomial weighting and discarding irrelevant predictors. Considerable computational gains are achieved by using the low-complexity procedures to calculate bootstrap replications while retaining the consistency. The model selection stage facilitates preventing the impact of undesirable bias introduced by regularization and allows for the inference free of regularization parameter tuning.

This paper is organized as follows: In section 2, we describe the proposed two-stage inference methods and the employed data model. In section 3, the robust model selection methods and the distributed inference procedures are explained in more detail. In section 4, theoretical characterization of finite-sample breakdown point of robust $\tau$-Lasso is provided. Moreover, the details of RSOB-$\tau$ are explained and its consistency and robustness properties are established. Section 5 studies the

performance of the two-stage robust inference procedures in simulations. Their consistency, robustness of confidence intervals, computational complexity, variable selection and standard deviation are considered. Section 6 concludes the paper. Detailed derivations and explanations of the theorems and their proofs can be found in the **Appendix** and **Supplemental Material**. Moreover, we present a robust variable screening procedure [25] prior to model selection, aiming to reduce the computational complexity in very high-dimensional data sets. We refer the interested readers to **Supplemental Materials**.

## II. OVERVIEW

In this section, we will define the employed data model and briefly describe the proposed two-stage inference procedures.

### A. Data Model

Consider a large-scale data with $n$ independent and identically distributed (i.i.d.) observations following a linear regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{v} \tag{1}$$

where $\mathbf{X} = (\mathbf{x}_{[1]}, \mathbf{x}_{[2]}, \ldots, \mathbf{x}_{[n]})^T \in \mathbb{R}^{n \times p}$ denotes a regression matrix, $\mathbf{y} = (y_1, y_2, \ldots, y_n)^T \in \mathbb{R}^n$ is a response vector, $\mathbf{v} = (v_1, v_2, \ldots, v_n)^T \in \mathbb{R}^n$ is a measurement noise and the errors $v_l$ are assumed to be independent of predictors $\mathbf{x}_{[l]}$. $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ denotes a sparse parameter vector with $k_s = |\mathcal{S}|$ non-zero entries and $\mathcal{S} = \{j : \mathbf{1}([\boldsymbol{\beta}_0]_{(j)} \neq 0)\}$. In order to deal with sheer volume of data, $\mathbf{Y} = (\mathbf{y}, \mathbf{X})$ is split into $s$ smaller distinct subsets of data $\check{\mathbf{Y}}^{(i)} = (\check{\mathbf{y}}^{(i)}, \check{\mathbf{X}}^{(i)}) \in \mathbb{R}^{b \times (p+1)}, i = 1, \cdots, s$ that can be stored and processed separately. The subsets may be formed by resampling without replacement from rows of the complete data set where $b = \{\lfloor n^\gamma \rfloor | \gamma \in [0.6, 1)\}$. The same situation would occur if subsets of data are stored on $s$ storage and computing nodes and each node contains $b$ observations.

### B. Assumptions

Suppose the measurement noise or errors $v_l$ follow some distribution $F_0$ and the distribution of the observed predictors $\mathbf{x}_{[l]}$ is $G_0$. Then, the joint distribution $H_0$ of $(y_l, \mathbf{x}_{[l]})$ is

$$H_0(\mathbf{x}_{[l]}, y_l) = G_0(\mathbf{x}_{[l]})F_0(y_l - \mathbf{x}_{[l]}^T\boldsymbol{\beta}_0). \tag{2}$$

We make the following assumptions on the distribution of errors and predictors.

- The probability density $f_0(u)$ associated with probability distribution $F_0$ of the errors $v_l$ has the following properties: even, monotone decreasing in $|v|$ and strictly decreasing in a neighborhood of 0.
- $\mathbb{P}(\mathbf{x}_{[l]}^T\boldsymbol{\beta} = \mathbf{0}) < 1 - \delta$ for all non-zero $\boldsymbol{\beta}$ and $\delta$ as defined by equation (5).
- $G_0$ has a finite second moment and $\mathbb{E}_{G_0}[\mathbf{x}_{[l]}\mathbf{x}_{[l]}^T]$ is non-singular.

The condition 1 generalizes the result established in this work to extremely heavy-tailed errors by imposing no moment conditions on the residual distribution $F_0$. The condition 2 guarantees the probability that observed values of explanatory variable are concentrated on a hyperplane does not get too large. The condition 3 concerns the second moment of explanatory variables and very common in the asymptotic analysis of regression estimators.

### C. Proposed Two-Stage Inference Methods

In order to perform inference on potentially high-dimensional models in the presence of sparsity and outlying observations, we propose two-stage robust and distributed statistical inference methods where robust variable selection is performed in the first stage and then selected variables from each distinct subset of data are combined by using a fusion rule in the fusion center or cloud. The actual inference is done in the second stage by using robust and low-complexity bootstrapping procedures on the variables selected in the first stage.

The robust variable selection incorporates features of split-and-conquer approach [26], Bolasso [27] and robust lasso estimators. The Bolasso discusses that Lasso fails to produce consistent variable selection results under certain decays of regularization parameter. More specifically, the probability of selecting irrelevant variables is strictly positive. In order to overcome this issue, Bolasso recommends to generate sufficient number of bootstrap samples of the original data set and intersect the support of the parameter vector estimated based on each bootstrap. Thus, the irrelevant variables would randomly be selected by Lasso and could be eliminated from the support during intersection. The majority voting scheme is often preferred over intersection because if few of the relevant variables are erroneously not selected by one bootstrap replication, the AND-rule excludes those variables from the model. In the split-and-conquer approach, the data set is split into $s$ subsets and distinct subsets are processed separately. In large-scale settings with distributed storage and processing architecture, multiple distinct subsets of i.i.d. data are stored on nodes. One can estimate the support of the parameter vector based on distinct subsets of data stored at each node and combine them in the fusion center by using a voting scheme. Due to page limitations, we chose to discuss only TSRD-$\tau$ and its related equations. The TSRD-MM method basically replaces the class of $\tau$-estimators in TSRD-$\tau$ with MM-estimators.

### D. TSRD-$\tau$

We describe the core components of two-stage robust and distributed inference procedure employing class of $\tau$-estimators. In the first stage, the robust variable selection exploits the $\tau$-Lasso method [28], [15] to select the sparse basis for each of the $s$ distinct data sets of $b$ observations. The selection results from each node are fused in a cloud or fusion center by using a percentage-based voting rule to choose the variables for the complete large data set. The chosen basis is communicated to each distributed computing and storage node and used in the second stage of the inference. A new robust and scalable technique using one-step bootstrap approximation of robust $\tau$-estimators is proposed to find parameter estimates for the selected basis and their confidence intervals in the presence

of outliers. Bootstrap replicates are computed by using a robust low-dimensional $\tau$-estimator of regression [29]. The bootstrap percentile method is used to estimate confidence intervals associated with selected variables. The estimated confidence intervals from each node are communicated to the cloud or fusion center for the inference on complete large scale data. The confidence intervals estimated for each subset of data are combined at the fusion center by coordinate-wise averaging over the lower and upper bound of confidence intervals.

In addition, we show that the distribution of the proposed RSOB-$\tau$ asymptotically converges to the same limiting distribution as the sampling distribution of $\tau$-estimator for distinct subsets of data. Combining this with the assumption of $\mathbb{P}(\hat{\mathcal{S}} = \mathcal{S}) = 1$ as $n \to \infty$ implies that the asymptotic distribution of bootstrap $\tau$-estimates for non-zero variables would converge to the same limiting normal distribution of $\tau$-estimator as if the true non-zero variables were known a priori. The robustness properties of $\tau$-Lasso and RSOB-$\tau$ are established by using analytical methods and computer simulations. A simple schematic of the proposed two-stage inference method is presented in **Fig. 1**. The details of the proposed TSDR-$\tau$ method are presented in the following sections.

## III. The Proposed Robust Inference Procedure

### A. Stage I: Model Selection

The proposed distributed and robust variable selection algorithm employed in the first stage of the inference procedure is described in detail. The main steps of the algorithm are summarized in **Fig. 2**.

The large-scale data is first split into $s$ distinct subsets via resampling without replacement. Then, the robust model selection is carried out in two steps in a distributed manner as follows:

#### 1) Variable Selection Using $\tau$-Lasso

At each node, the relevant variables are selected by applying robust $\tau$-Lasso estimator [28], [15] to each distinct subset of observations $\check{\mathbf{Y}}^{(i)} = (\check{\mathbf{y}}^{(i)}, \check{\mathbf{X}}^{(i)}), i = 1, ..., s$. Hence, the robust variable selection is performed by solving a set of $s$ estimation sub-problems defined as follows:

$$\hat{\boldsymbol{\beta}}^{(i)} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}}\, \mathcal{L}(\boldsymbol{\beta}) = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left( [\hat{\sigma}_\tau^{(i)}]^2 + \lambda \|\boldsymbol{\beta}\|_{\ell_1} \right) \quad (3)$$

where $\lambda$ controls the level of sparsity imposed by the $\ell_1$-norm penalty term. $\hat{\sigma}_\tau^{(i)}$ is a shorthand for $\hat{\sigma}_\tau(\check{\mathbf{r}}^{(i)}(\boldsymbol{\beta}))$ an efficient estimate of scale defined as follows:

$$[\hat{\sigma}_\tau^{(i)}]^2 = \frac{[\hat{\sigma}_b^{(i)}]^2}{b} \sum_{l=1}^{b} \rho_1\left( \frac{\check{r}_l^{(i)}(\boldsymbol{\beta})}{\hat{\sigma}_b^{(i)}} \right) \quad (4)$$

where $\check{\mathbf{r}}^{(i)}(\boldsymbol{\beta}) = \check{\mathbf{y}}^{(i)} - \check{\mathbf{X}}^{(i)}\boldsymbol{\beta}$ and $\hat{\sigma}_b^{(i)}$ denotes a shorthand for $\hat{\sigma}_M(\check{\mathbf{r}}^{(i)}(\boldsymbol{\beta}))$ an M-scale estimate of residuals $\check{\mathbf{r}}^{(i)}(\boldsymbol{\beta})$ defined as the solution to

$$\frac{1}{b} \sum_{l=1}^{b} \rho_0\left( \frac{\check{r}_l^{(i)}(\boldsymbol{\beta})}{\hat{\sigma}_b^{(i)}} \right) = \delta_1 \quad (5)$$

where $\delta_1$ is a tuning constant controlling the asymptotic breakdown point of the estimator. $\rho_0(\cdot)$ and $\rho_1(\cdot)$ are even and bounded functions satisfying the properties of $\rho$-function defined by Maronna et al. [30]. Tukey's bisquare $\rho$-function is considered as a popular choice in robust regression and defined as $\rho_i(t) = 1 - \left( 1 - (t/c_i)^2 \right)^3 \mathbf{1}(|t| \leqslant c_i), i = 0,1$ where $c_0$ and $c_1$ are chosen so that the desired normal efficiency $\zeta^*$ and breakdown point $\delta^*$ are attained for $\lambda = 0$, respectively. This can be achieved by finding $c_0$ and $c_1$ that satisfy $\mathbb{E}[\rho_0(t)] = \delta^*$ and $\left( \mathbb{E}[\psi'(t)] \right)^2 / \mathbb{E}[\psi^2(t)] = \zeta^*$ under the normality assumption of errors $t \sim \mathcal{N}(0,1)$ when $\lambda = 0$, simultaneously. $\psi(t)$, $\psi_0(t)$, $\psi_1(t)$ and $W$ are defined as follows:

$$\psi(t) = W\psi_0(t) + \psi_1(t),$$
$$\psi_0(t) = \partial\rho_0(t)/\partial t, \quad \psi_1(t) = \partial\rho_1(t)/\partial t, \quad (6)$$
$$W = \left( 2\mathbb{E}[\rho_1(t)] - \mathbb{E}[\psi_1(t)t] \right)/\mathbb{E}[\psi_0(t)t].$$

*Computation*: In order to solve the optimization problem given in equation (38), we employ the generalized gradient to minimize the composite objective function consisting of a non-convex term and a non-smooth $\ell_1$-norm penalty term. The generalized gradient of the objective function is defined as $\partial_{\boldsymbol{\beta}}([\hat{\sigma}_\tau^{(i)}]^2 + \lambda\|\boldsymbol{\beta}\|_{\ell_1})$ where $\partial_{\boldsymbol{\beta}}[\hat{\sigma}_\tau^{(i)}]^2$ associated with the smooth, non-convex, continuously differentiable term is identical to its gradient $\nabla_{\boldsymbol{\beta}}[\hat{\sigma}_\tau^{(i)}]^2$ and $\partial_{\boldsymbol{\beta}}(\lambda\|\boldsymbol{\beta}\|_{\ell_1})$ associated with non-smooth, convex term coincides with its subdifferential [17], [31]. It follows from the local lipschitzity of the composite objective function, any point $\bar{\boldsymbol{\beta}} \in \mathbb{R}^p$ at which $\mathbf{0} \in \partial_{\boldsymbol{\beta}}([\hat{\sigma}_\tau^{(i)}]^2 + \lambda\|\boldsymbol{\beta}\|_{\ell_1})$ is a local minimum of the $\tau$-Lasso estimation problem. Therefore, the generalized gradient of the objective function wrt $\boldsymbol{\beta}$ may be leveraged to find the local minima of the given estimation problem. It can be shown the generalized gradient of the objective function is equivalent to the sub-gradient of the weighted least square penalized by $\ell_1$-norm except that the weights $w_l^{(i)}(\boldsymbol{\beta})$ here depend on the unknown $\boldsymbol{\beta}$. Hence, the original optimization problem may be reformulated as follows:

$$\hat{\boldsymbol{\beta}}^{(i)} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \|\boldsymbol{\Omega}^{(i)}(\check{\mathbf{y}}^{(i)} - \check{\mathbf{X}}^{(i)}\boldsymbol{\beta})\|_{\ell_2}^2 + \lambda'\|\boldsymbol{\beta}\|_{\ell_1}, \quad (7)$$

where $\lambda' = 2b\lambda/\hat{\sigma}_b^{(i)}$, $\boldsymbol{\Omega}^{(i)}$ is a diagonal matrix whose entries on diagonal are $\sqrt{w_l^{(i)}}$ and $w_l^{(i)}$ is given by,

$$w_l^{(i)} = \frac{\left[ w_\tau^{(i)}\psi_0(\tilde{r}_l^{(i)}) + \psi_1(\tilde{r}_l^{(i)}) \right]}{\check{r}_l^{(i)}},$$
$$w_\tau^{(i)} = \frac{\sum_{l=1}^{b} \left[ 2\rho_1(\tilde{r}_l^{(i)}) - \psi_1(\tilde{r}_l^{(i)})\tilde{r}_l^{(i)} \right]}{\sum_{l=1}^{b} \psi_0(\tilde{r}_l^{(i)})\tilde{r}_l^{(i)}}. \quad (8)$$

where the notation $\check{r}_l^{(i)}$ is a shorthand for $\check{r}_l^{(i)}(\boldsymbol{\beta})$ and $\tilde{r}_l^{(i)} = \check{r}_l^{(i)}/\hat{\sigma}_b^{(i)}$. In the spirit of iteratively reweighted least-squares (IRLS) [32]-[33], we use iteratively reweighted Lasso (IR-LASSO) alternating between finding the weight matrices $\boldsymbol{\Omega}^{(i)}$, refining $\hat{\boldsymbol{\beta}}^{(i)}$ and updating $\hat{\sigma}_b^{(i)}$. The M-scale estimates are calculated via fixed-point iterations at each step of IR-LASSO.
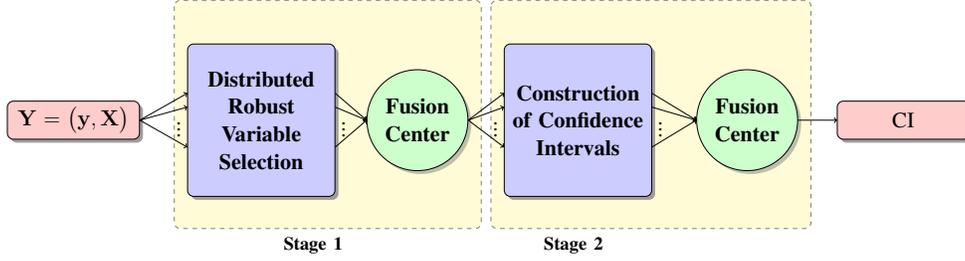
Fig. 1. A simple schematic of the proposed two-stage robust and distributed inference.
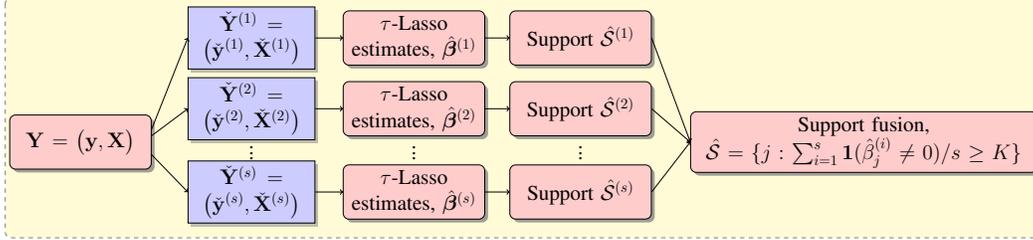


Fig. 2. The model selection is carried out in two steps: performing model selection at each node storing a distinct subset of data and combining the selected models from each node in the fusion center through a voting scheme.

*2) Fusing variable selections*

Once parameter estimation and variable selection are performed at each node, the chosen variables are communicated to the fusion center or cloud. In the fusion center, a percentage-based voting rule is used to select the relevant variables for the entire large-scale data set. The selection results of all nodes are combined according to the following rule, that is, , if a parameter is in the support within $100 \times K$ percent of subsets, it is selected to the support for the complete data set, $\hat{\mathcal{S}} = \{j : \sum_{i=1}^{s} \mathbf{1}(\hat{\beta}_j^{(i)} \neq 0)/s \geqslant K\}$. The chosen variables $\hat{\mathcal{S}}$ are broadcast into computation and storage nodes and deployed in the second stage of inference which uses the selected variables only.

Consider we skip the fusion and directly pass down the variables selected at each of $s$ nodes to the next stage. It is likely that models selected at some subsets of data are either overfitting or undefitting, which results in inaccurate inference results, thereby an appropriate rule of fusion is recommended.

### B. Stage II: Robust Inference

*1) Construction of CIs (RSOB-$\tau$)*

In this part, we perform the actual inference over the chosen variables from the model selection stage by using the RSOB-$\tau$. The derived robust $\tau$-estimation equations are used to compute bootstrap replicates instead of MM-estimation equations in [22]. Scalability is achieved by conducting inference in parallel on each distinct subset of data. Moreover, applying the RSOB-$\tau$ on only the selected variables eliminates the bias introduced by the regularization term.

In order to develop low-complexity inference procedures using bootstrap, it is required for the underlying estimator to

be expressed as a fixed-point problem. It is well-known that $\tau$-estimator fulfills the above requirement and can be represented as a solution of fixed-point equations as follows:

$$\hat{\boldsymbol{\theta}}_b = \mathbf{f}(\hat{\boldsymbol{\theta}}; \tilde{\mathbf{Y}}_b) \tag{9}$$

where $\mathbf{f} : \mathbb{R}^{|\hat{\mathcal{S}}|+1} \to \mathbb{R}^{|\hat{\mathcal{S}}|+1}$ denotes a smooth function, $\tilde{\mathbf{Y}}_b = (\check{\mathbf{y}}, \tilde{\mathbf{X}})$ with the explaining variables chosen in the model selection stage and $\hat{\boldsymbol{\theta}}_b$ is a fixed-point of $\mathbf{f}$. The terms $\hat{\boldsymbol{\theta}}_b$ and $\mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}_b)$ are defined, respectively, as follows:

$$\hat{\boldsymbol{\theta}}_b = \begin{bmatrix} \hat{\boldsymbol{\beta}}_b \\ \hat{\sigma}_b \end{bmatrix}$$

$$\mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}) = \begin{bmatrix} \left(\sum_{l=1}^{b} \hat{w}_l \tilde{\mathbf{x}}_{[l]} \tilde{\mathbf{x}}_{[l]}^T\right)^{-1} \left(\sum_{l=1}^{b} \hat{w}_l \check{y}_l \tilde{\mathbf{x}}_{[l]}\right), \\ \sum_{l=1}^{b} \hat{v}_l \hat{r}_l, \end{bmatrix} \tag{10}$$

where

$$\hat{v}_l = \frac{1}{b\delta_2} \times \frac{\rho_0(\tilde{r}_l)}{\tilde{r}_l},$$

$$\hat{w}_l = \frac{\hat{w}_\tau \rho_0'(\tilde{r}_l) + \rho_1'(\tilde{r}_l)}{\hat{r}_l},$$

$$\hat{w}_\tau = \frac{\sum_{l=1}^{b} \left[2\rho_1(\tilde{r}_l) - \rho_1'(\tilde{r}_l)\tilde{r}_l\right]}{\sum_{l=1}^{b} \rho_0'(\tilde{r}_l)\tilde{r}_l}, \tag{11}$$

$$\hat{r}_l = \check{y}_l - \tilde{\mathbf{x}}_{[l]}^T \hat{\boldsymbol{\beta}}_b,$$

$$\tilde{r}_l = \frac{\hat{r}_l}{\hat{\sigma}_b}.$$

$\hat{w}_l$ down-weights the outlying observations. In order to construct confidence intervals of the selected variables, the bootstrap replications of $\hat{\boldsymbol{\theta}}_b$ are computed as follows:

$$\hat{\boldsymbol{\theta}}_{n,b}^\star = \mathbf{f}(\hat{\boldsymbol{\theta}}_{n,b}^\star; \tilde{\mathbf{Y}}^\star), \tag{12}$$

where $\tilde{\mathbf{Y}}^\star = (\tilde{\mathbf{Y}}_b, \boldsymbol{\omega}^\star)$ denotes a bootstrap sample of size $n$ randomly drawn with replacement from the given subset of data. The multiplicity of observations is determined by the random weight vector $\boldsymbol{\omega}^\star \in \mathbb{R}^b$ drawn from a multinomial distribution $(n, (1/b)\mathbf{1}_b)$. Instead of computing a fully iterating bootstrap replicate $\hat{\boldsymbol{\theta}}_{n,b}^\star$, we can approximate it via a one-step iteration $\hat{\boldsymbol{\theta}}_{n,b}^{1\star}$ as follows:

$$\hat{\boldsymbol{\theta}}_{n,b}^{1\star} = \mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}^\star), \tag{13}$$

However, the distribution of $\hat{\boldsymbol{\theta}}_{n,b}^{1\star}$ may not exhibit the actual variability of the sampling distribution of $\hat{\boldsymbol{\theta}}_b$, mainly because all bootstrap replicates $\hat{\boldsymbol{\theta}}_{n,b}^{1\star}$ are calculated starting from the same initial value $\hat{\boldsymbol{\theta}}_b$. It has been shown that one-step iteration of bootstrap replications for many estimators could be adjusted by a correction term to provide asymptotically true estimate of bootstrap distribution [34]. In order to achieve asymptotically correct bootstrap estimates, the one-step improvement of $\hat{\boldsymbol{\theta}}_{n,b}^{1\star}$ with a linear correction term can be written as [21], [35]:

$$\hat{\boldsymbol{\theta}}_{n,b}^{R\star} = \hat{\boldsymbol{\theta}}_b + \left[\mathbf{I} - \nabla\mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}_b)\right]^{-1}\left(\hat{\boldsymbol{\theta}}_{n,b}^{1\star} - \hat{\boldsymbol{\theta}}_b\right), \tag{14}$$

where $\hat{\boldsymbol{\theta}}_{n,b}^{R\star} \in \mathbb{R}^{|\hat{\mathcal{S}}|+1}$ denotes the linearly corrected one-step bootstrap replication of $\hat{\boldsymbol{\theta}}_b$ and $\nabla\mathbf{f}(\cdot)$ is a gradient matrix with respect to $\boldsymbol{\theta}$. The linear correction term can be obtained by inverting the matrix $\left[\mathbf{I} - \nabla\mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}_b)\right]$ via the block matrix inversion lemma as described in Appendix.

Computational efficiency is attained because the correction factor and the initial estimate $\hat{\boldsymbol{\theta}}_b$ are computed only once for each subset of data. Furthermore, one-step bootstrap replications of $\hat{\boldsymbol{\theta}}_b$ are computationally inexpensive. We show in theorem 2 that $\hat{\boldsymbol{\theta}}_{n,b}^{R\star}$ would estimate the same limiting distribution as the actual bootstrap distribution $\hat{\boldsymbol{\theta}}_{n,b}^\star$ under certain regularity conditions.

The algorithm begins with generating $B$ bootstrap samples for each distinct subset of data $\tilde{\mathbf{Y}}_b^{(i)} = (\check{\mathbf{y}}^{(i)}, \tilde{\mathbf{X}}^{(i)}) \in \mathbb{R}^{b \times (|\hat{\mathcal{S}}|+1)}, i = 1, \cdots, s$. Across all computing and storage nodes, linearly corrected one-step bootstrap replicates $\hat{\boldsymbol{\beta}}_{n,b}^{R\star,(ij)}, j = 1, \cdots, B$ are computed. Then, confidence intervals associated with selected variables from stage 1 are constructed by using bootstrap percentile method for each subset of data. The estimated confidence intervals are communicated from each computing node to the fusion center for performing inference on the complete large-scale data. In the fusion center, the confidence intervals for the complete large-scale data are produced by applying coordinate-wise averaging over upper bounds and lower bounds of transmitted confidence intervals as follows:

$$\overline{\mathrm{CI}}_j = \frac{1}{s}\sum_{i=1}^{s}\overline{\mathrm{CI}}_j^{\star(i)}$$
$$\underline{\mathrm{CI}}_j = \frac{1}{s}\sum_{i=1}^{s}\underline{\mathrm{CI}}_j^{\star(i)} \tag{15}$$

where $\overline{\mathrm{CI}}_j$ and $\underline{\mathrm{CI}}_j$ denote the lower bound and upper bound of confidence interval associated with entry $j$ of the non-zero parameter.

Herein, we assumed that the proportion of outliers withing each subset of data remains below 50%, implying the estimated confidence intervals are not corrupted. Therefore, we may avoid the undesirable effects of outliers on the veracity of confidence intervals. In case more than half of observations within each subset of data are contaminated by outliers, we recommend using the adaptive trimmed mean [36] or classical trimmed mean with manually set trimming ratio [30].

---

**Algorithm 1:** The RSOB-$\tau$ procedure

**Data:** $s$ distinct subsamples
**Output: CI**

1 **for** *each subsample* **do**
2    Generate $B$ bootstrap samples of size $n$ by randomly drawing with replacement as follows: $\tilde{\mathbf{Y}}^\star = (\tilde{\mathbf{Y}}_b, \boldsymbol{\omega}^\star)$
3    Find the initial estimate $\hat{\boldsymbol{\theta}}_b$ by solving the fixed-point problem given in equation (9)
4    Calculate the one-step linearly corrected bootstrap replication corresponding to each bootstrap sample $\hat{\boldsymbol{\beta}}_{n,b}^{R\star}$ from equation (33)
5    Compute confidence intervals $\mathbf{CI}^{\star(i)}$ for each subset of data via bootstrap percentile method
6 **end**
7 Combine the confidence intervals by coordinate-wise averaging

---

## IV. STATISTICAL ROBUSTNESS AND CONSISTENCY PROPERTIES

In this section, the robustness properties of $\tau$-Lasso estimators and linearly corrected one-step replications using RSOB-$\tau$ are characterized by deriving their breakdown points. Furthermore, asymptotic normality of the linearly corrected one-step bootstrap replications is established under certain regularity conditions.

### A. Robustness Properties of $\tau$-Lasso Estimators

The finite-sample breakdown point of $\tau$-Lasso is derived. It measures the largest proportion of observations when arbitrarily replaced by outliers does not cause unbounded maximum bias or equivalently does not break down. Given a subset of data $\check{\mathbf{Y}} = (\check{\mathbf{y}}, \check{\mathbf{X}}) \in \mathbb{R}^{b \times (p+1)}$ randomly drawn from the original data $\mathbf{Y} = (\mathbf{y}, \mathbf{X}) \in \mathbb{R}^{n \times (p+1)}$, the replacement finite-sample breakdown point (FBP) $\epsilon^*(\hat{\boldsymbol{\beta}}; \check{\mathbf{Y}})$ of a regression estimator $\hat{\boldsymbol{\beta}} \in \mathbb{R}^p$ is defined as

$$\epsilon^*(\hat{\boldsymbol{\beta}}; \check{\mathbf{Y}}) = \max\{\frac{m}{b} : \sup_{\check{\mathbf{Y}}_m \in \mathcal{Y}_m} \|\hat{\boldsymbol{\beta}}(\check{\mathbf{Y}}_m)\|_{\ell_2} < \infty\} \tag{16}$$

where the set $\mathcal{Y}_m$ contains all datasets $\check{\mathbf{Y}}_m$ with $m$ $(0 < m < b)$ out of the original $b$ observations replaced by arbitrary values. The bounded *supremum* of $\ell_2$ term in the definition is equivalent to having a bounded maximum bias. The following theorem proves the results for finite-sample breakdown point of robust $\tau$-Lasso estimator by extending the theoretical result for that of robust S-PENSE estimator

demonstrated in [17] to the class of $\tau$-Lasso estimators.

**Theorem 1**: *Suppose $m(\delta)$ is the largest integer smaller than $b\min(\delta, 1-\delta)$ for a subset of data $\check{\mathbf{Y}} = (\check{\mathbf{y}}, \check{\mathbf{X}}) \in \mathbb{R}^{b \times (p+1)}$ and $\delta$ defined by equation (5). Then the finite-sample breakdown point of the $\tau$-Lasso estimator is bounded from above and below as follows*:

$$\frac{m(\delta)}{b} \leqslant \epsilon^*(\hat{\boldsymbol{\beta}}; \check{\mathbf{Y}}) \leqslant \delta \qquad (17)$$

where $\hat{\boldsymbol{\beta}}$ denotes the $\tau$-Lasso estimator.

**Proof**: The complete proof is given in Supplemental Materials.

With minor modification of the proof given in the **Supplemental Material**, we can extend the above theorem for data models including an intercept term.

### B. Asymptotic Properties of TSRD-$\tau$

In this section, we show that the proposed inference procedure will enjoy desirable asymptotic properties for fixed $p$ under certain regularity conditions if the model selection procedure is consistent. The model selection consistency in the first stage would imply that the relevant variables are selected with probability converging to 1 and the asymptotic distribution of modified bootstrap $\tau$-estimates for non-zero coefficients would be the same asymptotic normal distribution as if the true non-zero coefficients were known in advance. It seems reasonable to conjecture that the model selection procedure is consistent, extending the concept of Bolasso [27] to robust $\ell_1$-penalized estimators. We explore this possibility in simulations. The asymptotic properties of FRB for $\tau$-regression estimator were proven in Theorem 1 of [35]. However, the inaccuracy in the linear correction factor in the proof is rectified here. The asymptotic properties of BLFRB and FRB for MM-regression estimator are established in [22] and [21]. Also note that in this theorem, $\boldsymbol{\beta}_0 \in \mathbb{R}^{|\hat{\mathcal{S}}|}$ is a shorthand for $[\boldsymbol{\beta}_0]_{\hat{\mathcal{S}}}$ where parameter vector $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ as given in equation (1) and $\xrightarrow{P}$ denotes convergence in probability.

**Theorem 2** : *Let $\rho_0$ and $\rho_1$ be bounded $\rho$-functions satisfying the properties of bounded $\rho$-function defined by Maronna et al. [30] and have continuous third derivatives. Assume the model selection stage produces consistent estimates, i.e., $\mathbb{P}(\hat{\mathcal{S}} = \mathcal{S}) = o(1)$. Let $\hat{\boldsymbol{\beta}}_b$ be the $\tau$-regression estimator for a subset of data $(\check{\mathbf{y}}, \check{\mathbf{X}}) \in \mathbb{R}^{b \times (|\hat{\mathcal{S}}|+1)}$ randomly drawn from the original data $(\mathbf{y}, \underline{\mathbf{X}}) \in \mathbb{R}^{n \times (|\hat{\mathcal{S}}|+1)}$ and $\hat{\sigma}_b$ be the M-scale of residuals for the given subset of data and assume that they are consistent estimators, that is, $\hat{\boldsymbol{\beta}}_b \xrightarrow{P} \boldsymbol{\beta}_0$ and $\hat{\sigma}_b \xrightarrow{P} \sigma_0$. Given the following regularity conditions hold*:

1. *The following vectors and matrices exist and are finite*:
   1.1 $\mathbb{E}\left[\left(\bar{w}_\tau \rho_0'(r) + \rho_1'(r)\right)/r \mathbf{x}\mathbf{x}^T\right]^{-1}$
   1.2 $\mathbb{E}\left[\rho_0'(r)\mathbf{x}\right]$
   1.3 $\mathbb{E}\left[\left(\bar{w}_\tau \rho_0''(r) + \rho_1''(r)\right)\mathbf{x}\mathbf{x}^T\right]$
   1.4 $\mathbb{E}\left[\left(\bar{w}_\tau \rho_0''(r) + \rho_1''(r)\right)r\mathbf{x}\right]$
2. $\mathbb{E}[\rho_0'(r)r] \neq 0$ *and finite,*

3. $\rho_0'(u)/u$, $\rho_1'(u)/u$, $\left(\rho_0'(u) - \rho_0''(u)u\right)/u^2$ and $\left(\rho_1'(u) - \rho_1''(u)u\right)/u^2$ are continuous.

*where $\bar{w}_\tau = \left(2\mathbb{E}[\rho_1(r)] - \mathbb{E}[\rho_1'(r)r]\right)/\mathbb{E}[\rho_0'(r)r]$. Then, the distribution of $\sqrt{n}(\hat{\boldsymbol{\beta}}_{n,b}^{R\star} - \hat{\boldsymbol{\beta}}_b)$ converges weakly to the limiting distribution of $\sqrt{b}(\hat{\boldsymbol{\beta}}_b - \boldsymbol{\beta}_0)$ and consequently to the limiting distribution of $\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0)$ as $n$ and $b$ approach infinity.*

**Proof**: The complete proof is presented in Supplemental Materials.

### C. Statistical Robustness of RSOB-$\tau$

In this section, we are interested in studying the robustness properties of the proposed RSOB-$\tau$. The confidence intervals for regression parameters were constructed by using RSOB-$\tau$ quantiles and thus, the breakdown point of quantiles gives an insight into the robustness and reliability of the inferences made using the proposed method. Here, we derive some theoretical results about the breakdown point of quantile estimates of RSOB-$\tau$ and FRB using $\tau$-estimators.

Before proceeding, we define few concepts essential in comprehending the robustness results. Given $t \in (0, 1)$, $\hat{q}_t$ is defined as the $t_{th}$ upper quantile of a statistic $\hat{\boldsymbol{\beta}}_b$ such that $P[\hat{\boldsymbol{\beta}}_b > \hat{q}_t] = t$. According to [37], the upper breakdown point of a bootstrap estimate $\hat{q}_t^\star$ is defined as the minimum proportion of arbitrarily large outliers in the subset of data $\tilde{\mathbf{Y}}$ that can drive $\hat{q}_t^\star$ into infinity. In what follows theoretical results on robustness properties of estimated RSOB-$\tau$ and FRB quantiles using robust $\tau$-estimator are demonstrated. The following theorems are proved by using similar techniques as in the proofs given for BLFRB and FRB quantiles using robust MM-estimator [22] and [21]. Note that the theorems below also hold for subsets of data $\tilde{\mathbf{Y}} \in \mathbb{R}^{b \times |\hat{\mathcal{S}}|}$ having $|\hat{\mathcal{S}}|$ predictors selected from model selection stage.

**Theorem 3**: *Suppose $\mathbf{Y} = (\mathbf{y}, \mathbf{X}) \in \mathbb{R}^{n \times (p+1)}$ is a large-scale data following the linear model given in equation (1). Except, here no sparsity assumption is imposed on the parameter vector $\boldsymbol{\beta}_0$. Assume that $\mathbf{Y}$ is in general position, i.e., any subset of $p$ observations will result in a unique determination of $\boldsymbol{\beta}_0$. Let $\hat{\boldsymbol{\beta}}_n$ be a robust $\tau$-estimate of $\boldsymbol{\beta}_0$ based on the data $\mathbf{Y}$ whose breakdown point is $\epsilon$. Then, the breakdown point of the $t_{th}$ FRB quantile estimate of the regression parameters $[\beta_0]_{(l)}$, $l = 1, \cdots, p$ is determined by $\min(\epsilon_t^\star, \epsilon)$ where $\epsilon_t^\star$ is the smallest $\delta_c \in [0,1]$ that satisfies the following inequality,*

$$P\left[Binomial(n, 1-\delta_c) < p\right] \geqslant t \qquad (18)$$

Note that $\delta_c$ is different from $\delta_0$ and $\delta_1$.

**Proof**: The details of the poof are given in Appendix.

**Theorem 4**: *Suppose $\check{\mathbf{Y}} = (\check{\mathbf{y}}, \check{\mathbf{X}}) \in \mathbb{R}^{b \times (p+1)}$ is a subset of the large-scale data $\mathbf{Y} = (\mathbf{y}, \mathbf{X}) \in \mathbb{R}^{n \times (p+1)}$ formed by random resampling without replacement of the original full data set. Except, here no sparsity assumption is imposed on the parameter vector $\boldsymbol{\beta}_0$. Assume that $\check{\mathbf{Y}}$ is in general*

*position, i.e., any subset of $p$ observations will determine a bounded least-square estimate. Let $\hat{\boldsymbol{\beta}}_b$ be a robust $\tau$-estimator of $\boldsymbol{\beta}_0$ based on the bag of data $\check{\mathbf{Y}}$, $\epsilon^*(\hat{\boldsymbol{\beta}}_b, \check{\mathbf{Y}})$ denotes the finite-sample breakdown point of $\hat{\boldsymbol{\beta}}_b$ and $\delta_1 = 0.5$. Then, all the RSOB-T bootstrap quantiles estimated using the $\tau$-estimator that are constructed over $\check{\mathbf{Y}}$ will have equal asymptotic breakdown point to $\hat{\boldsymbol{\beta}}_b$.*

**Proof**: A detailed proof is given in Appendix.

## V. SIMULATIONS AND RESULTS

In this section, the performance of the proposed method is investigated in simulations considering both variable selection and inference. In particular, correct identification of sparse basis, statistical robustness of bootstrap estimates, the quality of the parameter estimates and confidence intervals are studied. The results in Theorem 2 is validated through computer simulations, indicating the distribution of bootstrap replications by using RSOB-$\tau$ asymptotically converges to the sampling distribution of $\tau$-estimator. The performance of the proposed method is assessed through computer simulations by using different proportions of outliers, large-scale data in both low-dimensional ($p < b$ or $n$) versus high-dimensional ($p \approx n$ or $\gg b$) settings. It is assumed the large-scale data follow a linear regression model where the parameter vector $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ is sparse with $k_s$ non-zero entries. The measurement noise vector, $\mathbf{v}$, is an additive white Gaussian noise with a variance (AWGN) $\sigma_v^2 = \|\mathbf{X}\boldsymbol{\beta}_0\|_{\ell_2}^2 10^{-\text{SNR}/10}/n$ (SNR in dB).

### A. Simulation Setting

Throughout simulations, the confidence intervals are reported in a nominal level of $100 \times (1 - \alpha_{cl})\% = 90\%$ with $\alpha_{cl} = 0.1$. We create a decreasing grid of 70 lambda values with logarithmic spacing of 1.1, spanning $\lambda_1$ to $\lambda_{70}$ where $\lambda_1$ is set to $\lambda_{\max}$. The maximum number of iteration in $\tau$ and MM-Lasso estimators is fixed at 30. The robustness of DPD-SIS procedure is adjusted by a tuning parameter $\alpha$ set to 0.4 (stable for a range of contamination levels) to provide robustness without significant loss in efficiency. The proposed algorithm was implemented in MATLAB except for the estimation of initial S-Lasso which was done in R using PENSE [38]. It provides a good initial estimate by constructing clean subsamples of data, potentially removing outlying observations. This is achieved by using the principal sensitivity components (PSCs) for EN estimator and removing the observations with most extreme PSCs from the subsamples in an iterative manner. A detailed description of initialization for PENSE can be found in Supplementary Materials of [39]. In addition, we used the Dual Augmented Lagrangian [40] implementation to solve the IR-LASSO, and MM- and S-estimators as in [41].

The outliers in all simulations except for **Section V-I** are introduced by randomly choosing the observations in y and replacing them with random values chosen from a standard Gaussian distribution with $\sigma_e = 250$.

### B. Data Standardization

Across all simulations for estimation problems using robust Lasso, it is assumed the linear regression model has an intercept component and all columns of the augmented regression matrix $[\mathbf{1}_{b\times 1} \check{\mathbf{X}}^{(i)}]$ except the first one are robustly standardized by centring the columns using a bisquare location estimator and scaling them using bisquare scale estimators [42]. The response vector $\check{\mathbf{y}}^{(i)}$ is centred using the bisquare location estimator [42].

### C. Calibration of Tuning Constants $c_0$ and $c_1$

In order to tune $c_0$ for the M-scale of residuals within all the robust estimators discussed in this paper, we set $\delta_i = 0.5$ for $i = 0, 1, 2$ to achieve the maximum robustness against outliers. $\delta_i$ controls the breakdown point according to Theorem 4.1 in [17] and $\delta_0$ is associated with initial S-Lasso estimator. It's worth mentioning $0.5$ is the largest value $\delta_i$ can take on. Note that $c_0$ is tuned for the desired breakdown point with the assumption $\lambda = 0$. In particular, it is recommended to have maximum robustness because bootstrapping may exacerbate the malicious behaviour of outliers in resampled data sets. The tuning constant $c_1$ of $\tau$-Lasso and MM-Lasso is respectively adjusted to $6.08$ and $4.68$ to attain $95\%$ efficiency under Gaussian errors when $\lambda = 0$. Likewise, $c_1 = 6.08$ for $\tau$-estimator and $c_1 = 4.68$ for MM-estimator are adjusted to provide $95\%$ efficiency under Gaussian errors for bootstrap replications.

### D. Choice of Fusion Parameter $K$

$K$ denotes the proportion of data subsets classified the given variable as relevant. Herein, the fusion parameter $K$ for voting rule is set to 0.5, indicating the variables selected within at least $50\%$ of subsets are regarded as relevant variables, i.e, majority voting scheme. The specific value of 0.5 is a compromise between reducing the false positive rate and maximizing the selection of true relevant variables.

### E. Choice of Batch Size $b$

The choice of batch size $b$ is a design parameter, depending on computational and storage resources available plus either model selection performance or prediction performance. On one hand, randomly partitioning the entire large-scale data of size $n$ into subsets of size $b$ and independently analyzing smaller subsets of data leads to significant computational speed-up and data storage benefits [20], [26]. On the other hand, the choice of batch size $b$ trades off between overfitting and underfitting in the variable selection procedure in additional to statistical correctness of confidence intervals. Setting the batch size $b$ to very large values results in rejecting relevant variables. This limits the invertibility of sub-design matrices as the eigenvalues are no longer strictly positive. However, one shall tune $b$ to be large enough to discard the irrelevant variables. Moreover, very small batch size $b$ equals to large number of data subsets results in wasted computation as we are occupying unnecessarily many computation and storage nodes while only a small proportion of them would be needed. The

known information-theoretic results [43] impose a condition on the batch size $b$ to ensure the support recovery for any method, $b \geqslant c_s k_s \log(p)$ for a sufficiently large constant $c_s$. On the other hand, BLB [20] requires the batch size to some values $n^{0.6} < b < n^{0.9}$ such that the statistical correctness is obtained. Combining these two constraints, one may choose the batch size in the range of $\max(n^{0.6}, c_s k_s \log(p)) \leqslant b \leqslant n^{0.9}$. We then need to take into account practical implications of computational constraints and choose an appropriate batch size $b$. Given the batch size $b$, one can obtain the number of subsamples $s$ and $\gamma$ in a straightforward manner.

We provide an example describing an appropriate batch size $b$. Let $n = 1000000$, $b = 500$ and $k_s = 50$, we then obtain that the minimum batch size from support recovery point of view is $b \geqslant 311 c_s$. Setting the constant $c_s = 30$, sufficiently large for sparse recovery, implies that $b \geqslant 9400$. On the other hand, reliable statistical correctness is attained by BLB for $b \geqslant (n^{0.6} \approx 4000)$. We shall occupy many resources over 100 computing cores. If we set $b = n^{0.75} \approx 32000$, we shall need 33 computing cores and computations over 32000 is manageable by single computing cores, which gives a good compromise between computational burden and statistical efficiency. Setting $b \approx n^{0.9}$ would mean almost 4 subsets of size 250000 which require higher computational power for each node and not very desirable.

### F. Tuning of regularization parameter

In order to calibrate the robust penalized estimation problem, $\lambda_{\max}$ is initially estimated by using the method introduced by Khan et al. [44] and then improved upon via a binary search [16]. A set of candidate lambdas in decreasing order starting from $\lambda_{max}$ is formed. Selection of $\lambda$ is carried out through a robust version of Bayesian Information Criterion (BIC) [45], defined as

$$\text{RBIC}(\lambda) = b \log(\hat{\sigma}_b^2(\lambda)) + C(b,p)\|\hat{\boldsymbol{\beta}}(\lambda)\|_{\ell_0} \qquad (19)$$

where $\hat{\sigma}_b$ denotes the robust M-scale of residuals and $C(b,p)$ is set to $\log(b)$ for settings where the dimensionality $p$ is much smaller than the sample size $n$. $C(b,p) = \log(\log(b))\log(p)$ is chosen for $\log(p)/b \to 0$ as $p \to \infty$ to attain better empirical performance [46], [47], [48]. Here, standard BIC is modified by replacing the non-robust estimate of scale with the robust M-estimate of scale to deal with outliers. Finally, the optimal tuning parameter $\lambda$ minimizes the RBIC over the pre-defined grid of lambdas

$$\lambda^* = \underset{\lambda \in \Lambda}{\arg\min}\,\text{RBIC}(\lambda). \qquad (20)$$

### G. Scenarios

We consider the following scenarios for which simulation studies are carried out.

- **Scenario 1:** We set the simulation parameters as follows: $n = 27000$, $p = 30000$, $b = 900$ ($p/b = 33.33$), $\gamma = 0.6667$, SNR = 15 dB, $B = 1000$. $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ is sparse with $k_s = 40$ non-zero entries. $[\boldsymbol{\beta}_0]_{\mathcal{S}}$ is set to $3 \times \mathbf{1}_{\mathcal{S}}$ and their positions are chosen randomly. The covariate

vectors $\mathbf{x}_{[i]}, i = 1, \cdots, n$ are drawn independently from a multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with $\Sigma_{ij} = \rho^{|i-j|}$ (Toeplitz covariance structure, $\rho = 0.5$).

- **Scenario 2:** The simulation parameters are set as follows: $n = 2000000$, $p = 80$, $b = 40000$ ($p/b = 1/500$), $\gamma = 0.730367$, SNR = 30 dB, $B = 400$. $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ is sparse with $k_s = 20$ non-zero entries. $[\boldsymbol{\beta}_0]_{\mathcal{S}}$ is set to $3 \times \mathbf{1}_{\mathcal{S}}$ and their positions are chosen randomly. Explaining variables $\mathbf{x}_{[i]}$ in the regression matrix are i.i.d, randomly drawn from a multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

- **Scenario 3:** The simulation parameters are set as follows: $n = 80000$, $p = 100$, $b = 4000$ ($p/b = 1/40$), $\gamma = 0.73466$, SNR = 15 dB, $B = 300$. The regression matrix $\mathbf{X}$ is randomly generated from mutually independent observations drawn from $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with $\Sigma_{ij} = \rho^{|i-j|}$ (Toeplitz covariance structure, $\rho = 0.5$). $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ is sparse with $k_s = 10$ non-zero entries. $[\boldsymbol{\beta}_0]_{\mathcal{S}}$ is set to $3 \times \mathbf{1}_{\mathcal{S}}$ and their positions are chosen randomly.

- **Scenario 4:** We set the simulation parameters as follows: $n = 4900$, $p = 6000$, $b = 350$ ($p/b = 17.14$), $\gamma = 0.6895$, SNR = 15 dB, $B = 1000$. $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ is sparse with $k_s = 10$ non-zero entries. $\boldsymbol{\beta}_0$ is set to

$$\boldsymbol{\beta}_0 = [2.5, 2.5, 2.5, 2, 3, 3, 3, 3.5, 3.5, 3.5, \mathbf{0}_{p-k_s}^T]^T$$

The covariate vectors $\mathbf{x}_{[i]}, i = 1, \cdots, n$ are drawn independently from a multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with $\Sigma_{ij} = \rho^{|i-j|}$ (Toeplitz covariance structure, $\rho = 0.5$).

- **Scenario 5:** We set the simulation parameters as follows: $n = 20000$, $p = 80$, SNR = 15 dB. $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ is sparse with $k_s = 15$ non-zero entries. $\boldsymbol{\beta}_0$ is set to

$$\boldsymbol{\beta}_0 = [3.5, 3.5, 3.5, 5, 5, 5, 2.5, 2.5, 2.5, 1.5, 2\times\mathbf{1}_5^T, \mathbf{0}_{p-k_s}^T]^T$$

The covariate vectors $\mathbf{x}_{[i]}, i = 1, \cdots, n$ are drawn independently from a multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with $\Sigma_{ij} = \rho^{|i-j|}$ (Toeplitz covariance structure, $\rho = 0.5$).

### H. Variable Selection Performance with Different $\ell_1$-penalized Estimators

In this subsection, we first substitute different $\ell_1$-penalized estimators with $\tau$-Lasso estimators used in the model selection stage of TSRD-$\tau$ introduced in **Section III**. We then carry out simulations to compare the variable selection performance of the resulting procedures with that of TSRD-$\tau$ and TSRD-MM. For this purpose, we use $\ell_1$-penalized estimators such as RA-Lasso [19] and Sparse-LTS [18]. The former is implemented in MATLAB by following exactly the computation algorithm described in [19] whereas we use the well-known R package **robustHD** for Sparse LTS regression [49]. We run the simulations on the synthetic data set described by **Scenario 5** except for setting SNR = 10 dB with normal

errors. In this experiment, we introduce contamination in the regression matrix $\mathbf{X}$ and the response vector $\mathbf{y}$, simultaneously. For each simulation, we consider four contamination schemes as follows:

- **Scheme 1:** 10% of observations in the response vector $\mathbf{y}$ are replaced with random values drawn from standard Gaussian with $\sigma_e = 250$. We also replace the corresponding observations in $\mathbf{X}$ with random values chosen from standard multivariate Gaussian distribution with $\mathbf{\Sigma} = \sigma_e^2 \mathbf{I}_p$. (large-variance outliers)
- **Scheme 2:** 10% of observations in the response vector $\mathbf{y}$ are replaced with random values drawn from a Gaussian distribution $\mathcal{N}(250, 1)$. We also replace the corresponding observations in $\mathbf{X}$ with random values chosen from a multivariate Gaussian distribution $\mathbf{N}(50 \times \mathbf{1}_p, \mathbf{I}_p)$. (gross outliers)
- **Scheme 3:** The additive noise is heavy-tailed Student's $t$-distribution with one degree of freedom.
- **Scheme 4:** The outlier contamination follows the same procedure as in Scheme 1 with the assumption of heavy-tailed additive errors, thereby combining Scheme 1 and 3.

Across all schemes, the same set of observations are replaced with outliers. Herein, we perform a Monte-Carlo study of 20 trials where a random realization of outlier in $\mathbf{y}$ is used at each trial. **Table I** shows the result of variable selection method with different $\ell_1$-penalized estimators for four contamination schemes, averaged over 20 trials. We observe that the proposed variable selection algorithms using $\tau$-Lasso estimator and MM-Lasso estimator perfectly recover the true relevant variables while keeping the false positive rates low. In contrast, variable selection method using Sparse-LTS results in significantly larger false positive rates for smaller subsample sizes under contamination schemes 1 and 2. Except for contamination scheme 2 (heavy-tailed errors), the variable selection with RA-Lasso results in extremely overfitted models.

*I. Robustness of Bootstrap Replications*

In this subsection, the robustness of RSOB-$\tau$ is quantified by an uncertainty measure in comparison to BLB. In particular, we calculate the standard deviation based on the bootstrap replications produced by both methods and verify the results in Theorem 4 by assessing the relative error. The bootstrap estimate of standard deviation is computed as follows:

$$\widehat{\mathrm{SD}}(\hat{\boldsymbol{\beta}}_n) = \frac{1}{|\hat{\mathcal{S}}|} \sum_{l=1}^{|\hat{\mathcal{S}}|} \left( \frac{1}{s} \sum_{i=1}^{s} \left( \sum_{j=1}^{B} \frac{\left( [\hat{\boldsymbol{\beta}}_{n,b}^{\star(ij)}]_l - [\hat{\boldsymbol{\beta}}_{n,b}^{\star(i.)}]_l \right)^2}{B-1} \right)^{1/2} \right) \tag{21}$$

where $[\hat{\boldsymbol{\beta}}_{n,b}^{\star(i.)}]_l$ is given by:

$$[\hat{\boldsymbol{\beta}}_{n,b}^{\star(i.)}]_l = \frac{1}{B} \sum_{j=1}^{B} [\hat{\boldsymbol{\beta}}_{n,b}^{\star(ij)}]_l \tag{22}$$

In order to measure how accurate the bootstrap estimate of standard deviation approximates the average standard deviation of $\hat{\boldsymbol{\beta}}_n$, we use the relative error criterion that is defined as follows:

$$\varepsilon = \frac{\widehat{\mathrm{SD}}(\hat{\boldsymbol{\beta}}_n) - \overline{\mathrm{SD}}(\hat{\boldsymbol{\beta}}_n)}{\overline{\mathrm{SD}}(\hat{\boldsymbol{\beta}}_n)} \tag{23}$$

where the average standard deviation $\overline{\mathrm{SD}}(\hat{\boldsymbol{\beta}}_n)$ is defined as $\sigma/\sqrt{n\mathcal{O}}$ based on asymptotic covariance of $\tau$-estimator [29]. Here, $\mathcal{O}$ is set to 1 for least square estimator in BLB and $\mathcal{O} = 0.95$ for $\tau$-estimator in RSOB-$\tau$ tuned to have 95% Gaussian efficiency.

First, we show that even one outlying observation could drive $\widehat{\mathrm{SD}}(\hat{\boldsymbol{\beta}}_n)$ based on bootstrap estimates obtained by BLB into infinity whereas those based on bootstrap replications obtained by the RSOB-$\tau$ remain resistant to outlier. Herein, we run the simulations on the synthetic data set described by **Scenario 2**. In this experiment, a data point within the original data set is randomly drawn and its response vector is multiplied by an extreme value $\alpha_o$ powers of 10 to imitate the situation where outlier is introduced by misplacement of decimal point. In regard to model selection, the robust and non-robust two-stage methods exhibit reliable performance in selecting the true sparse basis of the parameter vector with a TP = 1 and CER = 0. One might have expected that the non-robust inference would fail in model selection. However, only the bag of data containing the outlier yielded unreliable estimates but the voting scheme in the fusion center reduces the adversarial effect of outlier.

In the stage 2, the bootstrap estimates of $\widehat{\mathrm{SD}}(\hat{\boldsymbol{\beta}}_n)$ are computed based on the data set generated by the selected predictors from the model selection stage. As it is observed in **Fig. 3**, both two-stage algorithms TSRD-$\tau$ and TSLL perform remarkably well in terms of relative error when there are no outliers present within the data set. However, the bootstrap estimates of standard deviation obtained by BLB are severely influenced by the presence of even one outlier. As the magnitude of $\alpha_o$ increases, the relative error gets larger, implying that BLB is not robust to outliers. On contrary, the relative error of standard deviation obtained by TSRD-$\tau$ is not influenced at all by the presence of one outlier regardless of its magnitude.

Theorem 4 states the upper breakdown point of RSOB-$\tau$ bootstrap quantile estimates is 0.499 for the simulation set-up described above. In order to examine the robustness of bootstrap replications, we show that the TSRD-$\tau$ bootstrap replications are robust in the face of outliers even if the data is contaminated severely by outliers. In this experiment, the outliers are introduced by randomly choosing a percentage of the observations in $\mathbf{y}$ and multiplying them with a random value $\alpha_o = 100000$. The proposed method TSRD-$\tau$ correctly recovers the true sparse basis with zero false positive rate even in the presence of severe contamination. As shown in **Fig. 4**, the bootstrap estimate of standard deviation based on the RSOB-$\tau$ is only slightly influenced by the outliers at contamination levels as high as 40%, hence verifying the results in Theorem 4. In other words, the impact of outliers on the bootstrap estimates is bounded. Note that the curves in **Fig. 3** and **Fig. 4** are obtained by averaging over 15 trials of Monte Carlo simulations

TABLE I

**COMPARISON OF VARIABLE SELECTION PERFORMANCE WITH DIFFERENT $\ell_1$- PENALIZED ESTIMATORS:** THE PROPOSED ALGORITHMS USING MM-LASSO AND $\tau$-LASSO EXHIBIT A RELIABLE PERFORMANCE IN RECOVERING THE TRUE SPARSE BASIS (TPR=1) WHILE KEEPING FALSE POSITIVE RATES LOW. VARIABLE SELECTION METHOD USING SPARSE-LTS FAILS TO SUPPRESS THE FALSE POSITIVES FOR SCHEMES 1 AND 2 (SMALLER SUBSAMPLE SIZE). VARIABLE SELECTION WITH RA-LASSO PERFORMS RELIABLY ONLY FOR SCHEME 3, HEAVY-TAILED ERRORS WITH NO OUTLIER CONTAMINATION

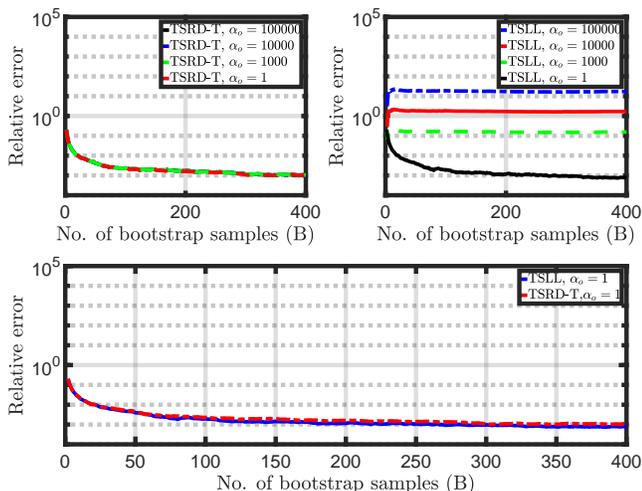| | scheme 1 | | | scheme 2 | | | scheme 3 | | | scheme 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | TP | FP | CER | TP | FP | CER | TP | FP | CER | TP | FP | CER |
| $b = 625$ | | | | | | | | | | | | |
| TSRD-$\tau$ | 1 | 0.0085 | 0.0069 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| TSRD-MM | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| Sparse-LTS | 1 | 0.14 | 0.1138 | 1 | 0.2008 | 0.1631 | 1 | 0.0008 | 0.0006 | 1 | 0.0023 | 0.0019 |
| RA-Lasso | 1 | 1 | 0.8125 | 1 | 0.1461 | 0.1187 | 1 | 0 | 0 | 1 | 1 | 0.8125 |
| $b = 800$ | | | | | | | | | | | | |
| TSRD-$\tau$ | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| TSRD-MM | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| Sparse-LTS | 1 | 0.0146 | 0.0119 | 1 | 0.05 | 0.0406 | 1 | 0 | 0 | 1 | 0 | 0 |
| RA-Lasso | 1 | 1 | 0.8125 | 1 | 0.2477 | 0.2019 | 1 | 0 | 0 | 0.9966 | 0.9938 | 0.8081 |
| $b = 1000$ | | | | | | | | | | | | |
| TSRD-$\tau$ | 1 | 0.0015 | 0.0013 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| TSRD-MM | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| Sparse-LTS | 1 | 0.0023 | 0.0019 | 1 | 0.0077 | 0.0062 | 1 | 0 | 0 | 1 | 0 | 0 |
| RA-Lasso | 0.92 | 0.9231 | 0.7650 | 1 | 0.4661 | 0.3787 | 1 | 0 | 0 | 0.3733 | 0.3569 | 0.4075 |



Fig. 3. The presence of only one extreme outlier in the data can drive the bootstrap estimate of standard deviation obtained by BLB in two-stage Lasso-LS into infinity. However, the bootstrap estimates of standard deviation obtained by RSOB-$\tau$ in TSRD-$\tau$ remains almost unaltered to the presence of one extreme outlier. The curves produced by TSRD-$\tau$ for all different values of $\alpha$ overlap, implying that one outlier has almost no effect on bootstrap estimates of standard deviation based on TSRD-$\tau$.



Fig. 4. The proposed inference method RSOB-$\tau$ at the stage 2 of TSRD-$\tau$ exhibit strong resilience to outlier even when 40% of observations are contaminated by outliers and relative error of standard deviation remains bounded.

### J. Statistical Convergence

In this subsection, the correctness of Theorem 2 is verified by computer simulations. In other words, we show the distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_{n,b}^{R\star} - \hat{\boldsymbol{\beta}}_b\big)$ converges to the limiting distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\big)$ as $n$ and $b$ approach infinity. Here, we run the simulations on the synthetic dataset described by **Scenario 2**, the same settings as in section V-I are used for the sake of con-
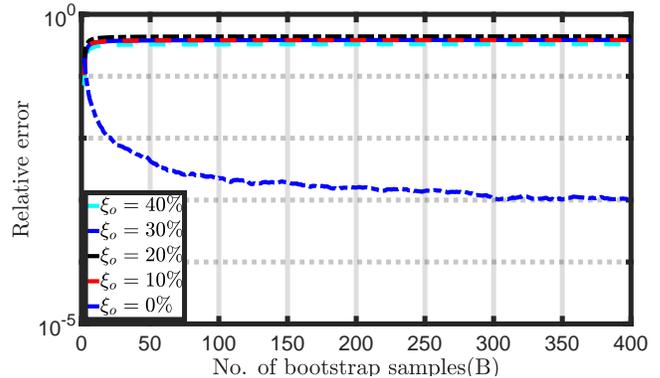
venience to study the statistical convergence. We assume that large-scale data is not contaminated by outliers. However, the number of bootstrap samples within each subset of data is set to 1000. Under the condition $\tau$-estimator is tuned for 95% normal efficiency, the limiting distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\big)$ obeys a multivariate Gaussian distribution $\mathcal{N}\big(\mathbf{0}, (\sigma^2/0.95)\mathbf{I}_{k_s}\big)$. The distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_{n,b}^{R\star} - \hat{\boldsymbol{\beta}}_b\big)$ is formed by randomly drawing a subset $\big(\check{\mathbf{y}}, \tilde{\mathbf{X}}\big)$ from the original data set $\big(\mathbf{y}, \underline{\mathbf{X}}\big) \in \mathbb{R}^{n \times |\hat{\mathcal{S}}|+1}$, computing the initial $\tau$-estimate $\hat{\boldsymbol{\beta}}_b$ and performing one-step linear correction of initial $\tau$-estimates for bootstrap samples by using the derived equations. The plot on the right-hand side of **Fig. 5** shows the empirical distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_{n,b}^{R\star} - \hat{\boldsymbol{\beta}}_b\big)$ overlaps the true limiting distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\big)$ for all

elements of $\hat{\boldsymbol{\beta}}_{n,b}^{R\star}$. On contrary, the plot on the left-hand side of **Fig. 5** shows the empirical distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_{n,b}^{1\star} - \hat{\boldsymbol{\beta}}_b\big)$ underestimates the variability of the true limiting distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\big)$ for all elements of $\hat{\boldsymbol{\beta}}_{n,b}^{1\star}$. Therefore we can conclude the distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_{n,b}^{R\star} - \hat{\boldsymbol{\beta}}_b\big)$ provides a reliable approximation of the distribution of $\sqrt{n}\big(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\big)$.



Fig. 5. The distribution of bootstrap replicates produced by linearly corrected one-step $\tau$-estimator in RSOB-$\tau$ overlaps the true limiting distribution of $\tau$-estimator whereas the distribution of bootstrap replicates produced by one-step $\tau$-estimator underestimates the variability of the true limiting distribution of $\tau$-estimator.

### K. Computational Complexity

In this subsection, we compare the computational complexity of the proposed distributed inference method, RSOB-$\tau$ at stage 2 of TSRD-$\tau$ to a robust realization of BLB method employing $\tau$-estimator for computing bootstrap replicates. We run the simulations on the synthetic data illustrated by **Scenario 3** where the proportion of outliers is set to $10\%$. The experiment was conducted in parallel on a single node of a high-performance computing cluster (Triton) where 22 computing cores and 14 GB of memory were requested and a Dell PowerEdge C4130 node was granted. The cumulative processing time is recorded after each iteration where new set of bootstrap samples are successively added to the bags. The cumulative processing time of RSOB-$\tau$ versus robustified BLB is demonstrated in **Fig. 6**. As the number of bootstrap samples increases, the proposed RSOB-$\tau$ requires significantly less processing time in comparison to robustified BLB. This implies the RSOB-$\tau$ achieves remarkably higher computational efficiency.
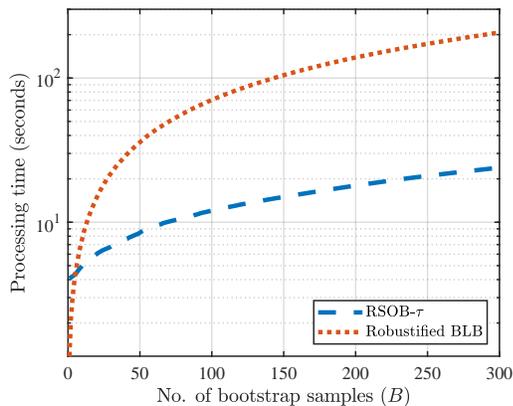


Fig. 6. The RSOB-T method employing linearly corrected one-step $\tau$-estimator is significantly faster than robustified BLB.

### L. Overfitting versus Underfitting and their implications on inference

In order to make valid inferences, we shall ensure the model selection procedure performs reliably. Herein, we examine how confidence intervals are influenced by overfitting and underfitting in the model selection. To do so, we run the simulations on the synthetic data described by **Scenario 5** where the proportion of outliers is set to $10\%$. We set the number of data partition $s = 40$. We study the validity of the statistical inference procedure TSRD-$\tau$ when model selection can not perfectly recover the true support, resulting in either overfitted models or underfitted models. We first consider a scenario where the model selection procedure fails to reject all irrelevant variables associated with zero coefficients of the parameter vector, resulting in an extremely overfitted model. We further consider the case that the model selection procedure fails to select all relevant variables and 10 out of 15 relevant variables are classified as irrelevant variables, resulting in an underfitted model.

We observe that confidence intervals for non-zero coefficients of underfitted model are much larger than that of overfitted model, as indicated by **Fig. 7**. In particular, confidence interval for one coefficient of underfitted model is very biased. Therefore, we do not only lose information about the relevant variables not chosen by the model selection in underfitted model. Also, the inference results may not reliably reflect confidence intervals for the non-zero coefficients of selected model. Moreover, we observe from **Fig. 8** that the confidence intervals for the given coefficients of parameter vector associated with false positives cover $0$. We can identify these coefficients as zero and variables corresponding to them as irrelevant. In this experiment, the confidence intervals associated with only 5 coefficients did not contain zero. Hence, the variable selection can still be improved even when one faces overfitting.
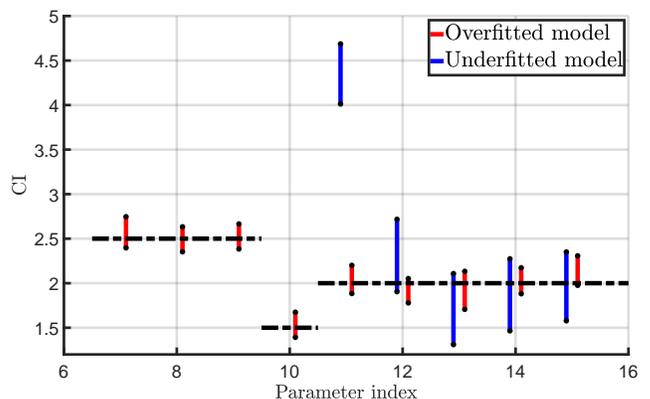


Fig. 7. Confidence intervals for non-zero coefficients of underfitted model are much larger than that of overfitted model. Besides, confidence interval for one non-zero coefficient of the underfitted model is extremely biased. (*the dash-dot lines indicate the true value of parameter vector for the corresponding entry*)
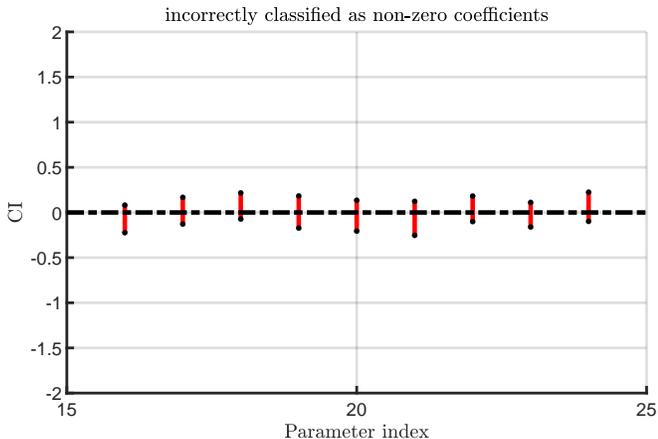
Fig. 8. In the overfitted model, we observe that confidence intervals for false positives associated with zero coefficients of the parameter vector cover zero for the given entries. This implies that we can identify the variables associated with confidence intervals covering zero as irrelevant. ( *the dash-dot lines indicate the true value of parameter vector for the corresponding entry*)

### M. Fusing the Variables Selected by S-Lasso Estimates

We now study how the variable selection algorithm performs when instead of the variables selected by $\tau$-Lasso estimates, the variables selected by the initial S-Lasso estimates across nodes are aggregated via the majority voting scheme. To do so, we carry out a series of Monte-Carlo simulations of 20 trials where a random realization of outlier is used at each trial. We run the simulations on the synthetic data set described by **Scenario 3** for low-dimensional regime and the synthetic data set described by **Scenario 4** for high-dimensional regime.

Table II shows the result of simulations. We observe that if we aggregate the variables selected by the initial S-Lasso estimates instead of those selected by $\tau$-Lasso estimates, the performance will remain almost the same. This could be explained by the fact that both $\tau$-Lasso and S-Lasso estimators promote sparsity and are robust to gross outliers. We suspect that low Gaussian efficiency has little effect on the performance of sparse recovery.

TABLE II
**FUSING THE VARIABLES SELECTED BY S-LASSO ESTIMATES:** IF WE USE THE VARIABLES SELECTED BY S-LASSO ESTIMATES INSTEAD OF THOSE SELECTED BY $\tau$-LASSO, THE MODEL SELECTION WILL REMAIN ALMOST UNCHANGED.

|  |  | TSRD-$\tau$ | | | Fusion / S-Lasso | | |
|---|---|---|---|---|---|---|---|
|  | $\xi_o$ | TP | FP | CER | TP | FP | CER |
| LD | 0.1 | 1 | 0 | 0 | 1 | 0.0011 | 0.0010 |
|  | 0.2 | 1 | 0 | 0 | 1 | 0.0017 | 0.0015 |
| HD | 0.1 | 1 | 0 | 0 | 1 | 0 | 0 |
|  | 0.2 | 1 | 0 | 0 | 1 | 0 | 0 |

## VI. CONCLUSION

This paper introduced robust and distributed inference procedures for large scale data where data exhibits an underlying low-dimensional structure and is contaminated by outliers. We propose two-stage inference procedures called TSRD-$\tau$

and TSRD-MM. The former employs the class of robust $\tau$-estimators whereas the latter employs that of MM-estimators. In the first stage, active explaining variables are selected by local variable selection employing robust Lasso estimators. The selections from each node are combined by applying a fusion rule at the fusion center or cloud. The selection is broadcast to the computational nodes. In the second stage, actual inferences on the selected variables are performed by using the robust and computationally efficient bootstrap procedures. Confidence intervals are constructed, parameter estimates are found, and standard deviations are quantified. The favorable statistical properties including consistency and robustness of the proposed method were established using analytical methods and verified in simulations. Moreover, the quantitative robustness properties of robust $\tau$-Lasso were established, in particular its finite-sample breakdown point.

Future directions of research include extending the classical asymptotic analysis to high-dimensional asymptotic analysis where $p$ grows with $n$ to infinity. It is also an open question how one can establish asymptotic results for local optima. Finally, it would be interesting to devise a rigorous proof for the model consistency of proposed TSRD-$\tau$ and TSRD-MM procedures.

## APPENDIX A
### DERIVATION OF THE LINEAR CORRECTION TERM

The linear correction may be derived by inverting a block-matrix as follows:

$$\left[\mathbf{I} - \nabla\mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}})\right]^{-1} = \begin{bmatrix} \mathcal{A} & \boldsymbol{\eta} \\ \boldsymbol{\zeta} & a \end{bmatrix}^{-1}, \qquad (24)$$

where $\mathcal{A}$, $\boldsymbol{\eta}$, $\boldsymbol{\zeta}$ and $a$ are given by

$$\mathcal{A} = \left(\hat{\mathbf{A}}_b\right)^{-1}\left[\mathcal{A}_2 - \mathcal{A}_1\right],$$

$$\boldsymbol{\eta} = \left(\hat{\mathbf{A}}_b\right)^{-1}\left[\eta_2 - \eta_1\right],$$

$$\boldsymbol{\zeta} = \frac{1}{b\delta_2}\sum_{l=1}^{b}\rho_0^{'}(\tilde{r}_l)\tilde{\mathbf{x}}_{[l]}^T,$$

$$a = \frac{1}{b\delta_2}\sum_{l=1}^{b}\rho_0^{'}(\tilde{r}_l)\tilde{r}_l. \qquad (25)$$

$\mathcal{A}_1$, $\mathcal{A}_2$, $\boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ are calculated as follows:

$$\mathcal{A}_1 = \frac{1}{b}\sum_{l=1}^{b}\tilde{\mathbf{x}}_{[l]}\nabla_{\boldsymbol{\beta}}w_\tau\rho_0^{'}(\tilde{r}_l),$$

$$\mathcal{A}_2 = \frac{1}{b\hat{\sigma}_b}\sum_{l=1}^{b}\left[\hat{w}_\tau\rho_0^{''}(\tilde{r}_l) + \rho_1^{''}(\tilde{r}_l)\right]\tilde{\mathbf{x}}_{[l]}\tilde{\mathbf{x}}_{[l]}^T,$$

$$\eta_1 = \frac{1}{b}\sum_{l=1}^{b}\nabla_\sigma w_\tau\rho_0^{'}(\tilde{r}_l)\tilde{\mathbf{x}}_{[l]}, \qquad (26)$$

$$\eta_2 = \frac{1}{b\hat{\sigma}_b}\sum_{l=1}^{b}\left[\hat{w}_\tau\rho_0^{''}(\tilde{r}) + \rho_1^{''}(\tilde{r})\right]\tilde{\mathbf{x}}_{[l]}\tilde{r}_l,$$

where $\nabla_{\boldsymbol{\beta}} w_\tau$ and $\nabla_\sigma w_\tau$ are shorthands for $\partial w_\tau(\hat{\boldsymbol{\theta}}_b)/\partial\boldsymbol{\beta}$ and $\partial w_\tau(\hat{\boldsymbol{\theta}}_b)/\partial\sigma$, respectively. $\hat{\mathbf{A}}_b$, $\nabla_{\boldsymbol{\beta}} w_\tau$ and $\nabla_\sigma w_\tau$ are given by

$$\hat{\mathbf{A}}_b = \frac{1}{b}\sum_{l=1}^{b} \hat{w}_l^{(i)} \tilde{\mathbf{x}}_{[l]}\tilde{\mathbf{x}}_{[l]}^T, \tag{27}$$

$$\nabla_{\boldsymbol{\beta}} w_\tau = \frac{\sum_{l=1}^{b}\left[\rho_1''(\tilde{r}_l)\tilde{r}_l - \rho_1'(\tilde{r}_l)\right]\tilde{\mathbf{x}}_{[l]}^T/\hat{\sigma}_b}{\sum_{l=1}^{b}\rho_0'(\tilde{r}_l)\tilde{r}_l},$$
$$+\frac{\sum_{l=1}^{b}\left[\rho_0''(\tilde{r}_l)\tilde{r}_l + \rho_0'(\tilde{r}_l)\right]\tilde{\mathbf{x}}_{[l]}^T/\hat{\sigma}_b}{\sum_{l=1}^{b}\rho_0'(\tilde{r}_l)\tilde{r}_l}\hat{w}_\tau. \tag{28}$$

$$\nabla_\sigma w_\tau = \frac{\sum_{l=1}^{b}\left[\rho_1''(\tilde{r}_l)\tilde{r}_l - \rho_1'(\tilde{r}_l)\right]\tilde{r}_l/\hat{\sigma}_b}{\sum_{l=1}^{b}\rho_0'(\tilde{r}_l)\tilde{r}_l}$$
$$+\frac{\sum_{l=1}^{b}\left[\rho_0''(\tilde{r}_l)\tilde{r}_l + \rho_0'(\tilde{r}_l)\right]\tilde{r}_l/\hat{\sigma}_b}{\sum_{l=1}^{b}\rho_0'(\tilde{r}_l)\tilde{r}_l}\hat{w}_\tau \tag{29}$$

On the other hand,

$$\begin{bmatrix}\mathcal{A} & \boldsymbol{\eta} \\ \boldsymbol{\zeta} & a\end{bmatrix}^{-1} = \begin{bmatrix}\mathbf{M}_b & \mathbf{d}_b \\ \mathbf{N}_b & \mathbf{q}_b\end{bmatrix}. \tag{30}$$

## APPENDIX B
### DERIVATION OF THE ONE-STEP BOOTSTRAP REPLICATES

The one-step bootstrap replicates $\hat{\boldsymbol{\beta}}_{n,b}^{1\star}$ and $\hat{\sigma}_{n,b}^{1\star}$ are calculated as follows:

$$\hat{\boldsymbol{\beta}}_{n,b}^{1\star} = \left(\sum_{l=1}^{b}\omega_l^\star w_l^\star \tilde{\mathbf{x}}_{[l]}\tilde{\mathbf{x}}_{[l]}^T\right)^{-1}\sum_{l=1}^{b}\omega_l^\star w_l^\star \tilde{y}_l\tilde{\mathbf{x}}_{[l]},$$
$$\hat{\sigma}_{n,b}^{1\star} = \sum_{l=1}^{b}\omega_l^\star \tilde{v}_l^\star\left(\tilde{y}_l - \tilde{\mathbf{x}}_{[l]}^T\hat{\boldsymbol{\beta}}_b\right), \tag{31}$$

where $v_l^\star$, $w_l^\star$ and $w_\tau^\star$ are computed as follows:

$$v_l^\star = \frac{b}{n}\hat{v}_l,$$
$$w_l^\star = \frac{w_\tau^\star \rho_0'(\tilde{r}_l) + \rho_1'(\tilde{r}_l)}{\hat{r}_l},$$
$$w_\tau^\star = \frac{\sum_{l=1}^{b}\omega_l^\star\left[2\rho_1(\tilde{r}_l) - \rho_1'(\tilde{r}_l)\tilde{r}_l\right]}{\sum_{l=1}^{b}\omega_l^\star \rho_0'(\tilde{r}_l)\tilde{r}_l}. \tag{32}$$

Therefore, the linearly corrected one-step bootstrap replications using $\tau$-estimators are calculated as follows:

$$\hat{\boldsymbol{\beta}}_{n,b}^{R\star} = \hat{\boldsymbol{\beta}}_b + \mathbf{M}_b\left(\hat{\boldsymbol{\beta}}_{n,b}^{1\star} - \hat{\boldsymbol{\beta}}_b\right) + \mathbf{d}_b\left(\hat{\sigma}_{n,b}^{1\star} - \hat{\sigma}_b\right),$$
$$\mathbf{M}_b = \left(\mathcal{A} - \boldsymbol{\eta}a^{-1}\boldsymbol{\zeta}\right)^{-1},$$
$$\mathbf{d}_b = -\mathcal{A}^{-1}\boldsymbol{\eta}\left(a - \boldsymbol{\zeta}\mathcal{A}^{-1}\boldsymbol{\eta}\right)^{-1}. \tag{33}$$

## REFERENCES

[1] K. Slavakis, G. B. Giannakis, and G. Mateos, "Modeling and optimization for big data analytics:(statistical) learning tools for our era of data deluge," *IEEE Signal Processing Magazine*, vol. 31, no. 5, pp. 18–31, 2014.

[2] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.

[3] A. Chatterjee and S. N. Lahiri, "Bootstrapping lasso estimators," *Journal of the American Statistical Association*, vol. 106, no. 494, pp. 608–625, 2011.

[4] A. Javanmard and A. Montanari, "Confidence intervals and hypothesis testing for high-dimensional regression," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 2869–2909, 2014.

[5] H. Liu, B. Yu *et al.*, "Asymptotic properties of lasso+ mls and lasso+ ridge in sparse high-dimensional linear regression," *Electronic Journal of Statistics*, vol. 7, pp. 3124–3169, 2013.

[6] A. Chatterjee, S. N. Lahiri *et al.*, "Rates of convergence of the adaptive lasso estimators to the oracle distribution and higher order refinements by the bootstrap," *The Annals of Statistics*, vol. 41, no. 3, pp. 1232–1259, 2013.

[7] L. Wasserman and K. Roeder, "High dimensional variable selection," *Annals of statistics*, vol. 37, no. 5A, p. 2178, 2009.

[8] H. Liu, X. Xu, and J. J. Li, "A bootstrap lasso+ partial ridge method to construct confidence intervals for parameters in high-dimensional sparse linear models," *arXiv preprint arXiv:1706.02150*, 2017.

[9] J. D. Lee, D. L. Sun, Y. Sun, J. E. Taylor *et al.*, "Exact post-selection inference, with application to the lasso," *Annals of Statistics*, vol. 44, no. 3, pp. 907–927, 2016.

[10] R. Berk, L. Brown, A. Buja, K. Zhang, L. Zhao *et al.*, "Valid post-selection inference," *The Annals of Statistics*, vol. 41, no. 2, pp. 802–837, 2013.

[11] C.-H. Zhang and S. S. Zhang, "Confidence intervals for low dimensional parameters in high dimensional linear models," *Journal of the Royal Statistical Society: Series B: Statistical Methodology*, pp. 217–242, 2014.

[12] R. Lockhart, J. Taylor, R. J. Tibshirani, and R. Tibshirani, "A significance test for the lasso," *Annals of statistics*, vol. 42, no. 2, p. 413, 2014.

[13] R. F. Barber, E. J. Candès *et al.*, "Controlling the false discovery rate via knockoffs," *Annals of Statistics*, vol. 43, no. 5, pp. 2055–2085, 2015.

[14] P. Bühlmann *et al.*, "Statistical significance in high-dimensional linear models," *Bernoulli*, vol. 19, no. 4, pp. 1212–1242, 2013.

[15] M. Martinez-Camara, M. Muma, B. Bejar, A. M. Zoubir, and M. Vetterli, "The regularized tau estimator: A robust and efficient solution to ill-posed linear inverse problems," *arXiv preprint arXiv:1606.00812*, 2016.

[16] E. Smucler and V. J. Yohai, "Robust and sparse estimators for linear regression models," *Computational Statistics & Data Analysis*, vol. 111, pp. 116–130, 2017.

[17] G. V. C. Freue, D. Kepplinger, M. Salibián-Barrera, and E. Smucler, "Pense: A penalized elastic net s-estimator," 2017.

[18] A. Alfons, C. Croux, and S. Gelper, "Sparse least trimmed squares regression for analyzing high-dimensional large data sets," *The Annals of Applied Statistics*, pp. 226–248, 2013.

[19] Q. Zheng, C. Gallagher, and K. Kulasekera, "Robust adaptive lasso for variable selection," *Communications in Statistics-Theory and Methods*, vol. 46, no. 9, pp. 4642–4659, 2017.

[20] A. Kleiner, A. Talwalkar, P. Sarkar, and M. I. Jordan, "A scalable bootstrap for massive data," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 76, no. 4, pp. 795–816, 2014.

[21] M. Salibian-Barrera, R. H. Zamar *et al.*, "Bootrapping robust estimates of regression," *The Annals of Statistics*, vol. 30, no. 2, pp. 556–582, 2002.

[22] S. Basiri, E. Ollila, and V. Koivunen, "Robust, scalable, and fast bootstrap method for analyzing large scale data," *IEEE Transactions on Signal Processing*, vol. 64, no. 4, pp. 1007–1017, 2015.

[23] E. Mozafari-Majd and V. Koivunen, "Robust, sparse and scalable inference using bootstrap and variable selection fusion," in *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE, 2019, pp. 271–275.

[24] ——, "Robust variable selection and distributed inference using $\tau$-based estimators for large-scale data," in *2020 28th European Signal Processing Conference (EUSIPCO)*. IEEE, 2021, pp. 2453–2457.

[25] A. Ghosh and M. Thoresen, "A robust variable screening procedure for ultra-high dimensional data," *arXiv preprint arXiv:2004.14851*, 2020.

[26] X. Chen and M.-g. Xie, "A split-and-conquer approach for analysis of extraordinarily large data," *Statistica Sinica*, pp. 1655–1684, 2014.

[27] F. R. Bach, "Bolasso: model consistent lasso estimation through the bootstrap," in *Proceedings of the 25th international conference on Machine learning*, 2008, pp. 33–40.

[28] M. Martinez-Camara, M. Muma, A. M. Zoubir, and M. Vetterli, "A new robust and efficient estimator for ill-conditioned linear inverse problems with outliers," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 3422–3426.

[29] V. J. Yohai and R. H. Zamar, "High breakdown-point estimates of regression by means of the minimization of an efficient scale," *Journal of the American statistical association*, vol. 83, no. 402, pp. 406–413, 1988.

[30] R. A. Maronna, R. D. Martin, V. J. Yohai, and M. Salibián-Barrera, *Robust statistics: theory and methods (with R)*. John Wiley & Sons, 2019.

[31] F. H. Clarke, *Optimization and nonsmooth analysis*. Siam, 1990, vol. 5.

[32] C. S. Burrus, J. Barreto, and I. W. Selesnick, "Iterative reweighted least-squares design of fir filters," *IEEE Transactions on Signal Processing*, vol. 42, no. 11, pp. 2926–2936, 1994.

[33] P. W. Holland and R. E. Welsch, "Robust regression using iteratively reweighted least-squares," *Communications in Statistics-theory and Methods*, vol. 6, no. 9, pp. 813–827, 1977.

[34] W. R. Schucany and S. Wang, "One-step bootstrapping for smooth iterative procedures," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 53, no. 3, pp. 587–596, 1991.

[35] M. Salibian-Barrera, S. Van Aelst, and V. J. Yohai, "Robust tests for linear regression models based on $\tau$-estimates," *Computational Statistics & Data Analysis*, vol. 93, pp. 436–455, 2016.

[36] R. V. Hogg, "Some observations on robust estimation," *Journal of the American Statistical Association*, vol. 62, no. 320, pp. 1179–1186, 1967.

[37] K. Singh *et al.*, "Breakdown theory for bootstrap quantiles," *The annals of Statistics*, vol. 26, no. 5, pp. 1719–1732, 1998.

[38] D. Kepplinger, M. Salibián-Barrera, and G. Cohen Freue, "https://cran.r-project.org/package=pense," 2021.

[39] G. V. C. Freue, D. Kepplinger, M. Salibián-Barrera, and E. Smucler, "Robust elastic net estimators for variable selection and identification of proteomic biomarkers," *The Annals of Applied Statistics*, vol. 13, no. 4, pp. 2065–2090, 2019.

[40] R. Tomioka, T. Suzuki, and M. Sugiyama, "Super-linear convergence of dual augmented lagrangian algorithm for sparsity regularized estimation," *Journal of Machine Learning Research*, vol. 12, no. May, pp. 1537–1586, 2011.

[41] C. Croux, G. Dhaene, and D. Hoorelbeke, "Robust standard errors for robust estimators," *CES-Discussion paper series (DPS) 03.16*, pp. 1–20, 2004.

[42] R. A. Maronna, "Robust ridge regression for high-dimensional data," *Technometrics*, vol. 53, no. 1, pp. 44–53, 2011.

[43] M. J. Wainwright, "Information-theoretic limits on sparsity recovery in the high-dimensional and noisy setting," *IEEE transactions on information theory*, vol. 55, no. 12, pp. 5728–5741, 2009.

[44] J. A. Khan, S. Van Aelst, and R. H. Zamar, "Robust linear model selection based on least angle regression," *Journal of the American Statistical Association*, vol. 102, no. 480, pp. 1289–1299, 2007.

[45] G. Schwarz *et al.*, "Estimating the dimension of a model," *Annals of statistics*, vol. 6, no. 2, pp. 461–464, 1978.

[46] Y. Fan and C. Y. Tang, "Tuning parameter selection in high dimensional penalized likelihood," *Journal of the Royal Statistical Society: SERIES B: Statistical Methodology*, pp. 531–552, 2013.

[47] Y. Kim, S. Kwon, and H. Choi, "Consistent model selection criteria on high dimensions," *The Journal of Machine Learning Research*, vol. 13, no. 1, pp. 1037–1057, 2012.

[48] A. Ghosh and S. Majumdar, "Ultrahigh-dimensional robust and efficient sparse regression using non-concave penalized density power divergence," *IEEE Transactions on Information Theory*, vol. 66, no. 12, pp. 7812–7827, 2020.

[49] A. Alfons, "robusthd: An r package for robust regression with high-dimensional data," *Journal of Open Source Software*, vol. 6, no. 67, p. 3786, 2021.

[50] Q. Li, Z. Zhu, and G. Tang, "The non-convex geometry of low-rank matrix optimization," *Information and Inference: A Journal of the IMA*, vol. 8, no. 1, pp. 51–96, 2019.

[51] Y. Chi, Y. M. Lu, and Y. Chen, "Nonconvex optimization meets low-rank matrix factorization: An overview," *IEEE Transactions on Signal Processing*, vol. 67, no. 20, pp. 5239–5269, 2019.

[52] V. J. Yohai and R. H. Zamar, "High breakdown-point estimates of regression by means of the minimization of an efficient scale," *University of Washington, Dept. of Statistics. Technical Report 84*, pp. 1–41, 1986.

[53] G. V. C. Freue, D. Kepplinger, M. Salibián-Barrera, and E. Smucler, "Proteomic biomarker study using novel robust penalized elastic net estimators," *the Annals of Applied Statistics, submitted*.

[54] J. Shao, *Mathematical statistics: exercises and solutions*. Springer Science & Business Media, 2006.

[55] P. J. Bickel and D. A. Freedman, "Some asymptotic theory for the bootstrap," *The annals of statistics*, pp. 1196–1217, 1981.

[56] P. Rousseeuw and V. Yohai, "Robust regression by means of s-estimators," in *Robust and nonlinear time series analysis*. Springer, 1984, pp. 256–272.

**Emadaldin Mozafari-Majd**

PLACE
PHOTO
HERE

**Visa Koivunen**

PLACE
PHOTO
HERE

This supplemental material contains a section on initial variable screening (Preprocessing) for very high-dimensional settings along with additional simulation results and technical proofs of the theorem discussed in the paper.

## INITIAL VARIABLE SCREENING (PREPROCESSING): DATA WITH VERY HIGH-DIMENSIONAL SUBSETS

In order to reduce the computational burden in settings with very high-dimensional subsets, an initial variable screening procedure called Density Power Divergence-SIS (DPD-SIS) [25] is employed to further reduce the model complexity, the number of variables to an order of sample size $n$. In distributed and parallel architecture, we use the variable screening only in very high-dimensional settings where $p$ is much larger than $b$. Basically, the DPD-SIS extends the Sure Independence Screening (SIS) to address robustness in the presence of outliers. The robust screening procedure assigns a certain score to each predictor and then predictors are ranked in descending order based on the calculated score. In order to compute the score, the marginal estimate of each regression coefficient is obtained via a minimum DPD estimator and then its absolute value determines the score associated with each predictor. At each node, the DPD-SIS variable screening procedure is utilized to discard a certain number of irrelevant predictors and then distinct subsets with reduced dimensionality is passed down to next step for further processing, that is, $\check{\mathbf{X}} \in \mathbb{R}^{b \times p} \to \bar{\mathbf{X}} \in \mathbb{R}^{b \times q}$ with $(q \ll p)$ an order of sample size. A reasonable choice for $q$ is the number of observations within subsets of data, that is, $q = b$. Note that variable screening procedure is dispensable in low-dimensional models, i.e. $q = p$. Subdividing a high-dimensional data into smaller subsets may be allowed as long as the batch size $b$ satisfies the requirements for sparse recovery and statistical correctness of bootstrap computations as specified in **Section V.E**, *Choice of Batch Size*, within the main body of the paper. The algorithmic details of DPD-SIS can be found in [25].

The DPD-SIS is initially applied to the distinct subsets of data and top $b$ predictors, the exact order of sample size, within each subset of data are kept based on their score and the remaining predictors are discarded. We then form a set of predictors appearing within at least half of data subsets. We now place these set of predictors on top of the ranking within each of $s$ data subsets, keep the top $b$ predictors and reject the remaining ones. The model selection proceeds with the data set of reduced dimensionality.

## ADDITIONAL SIMULATION RESULTS

### A. High-dimensional: Model Selection and Inference

In this part, we study the performance of the proposed methods, TSRD-$\tau$ and TSRD-MM, for a large-scale high-dimensional data set $(p > b)$ in terms of model selection and robustness of confidence intervals to outliers. We run the simulations on the synthtetic dataset described by **Scenario 1**. In the current high-dimensional setting, an initial variable screening procedure is employed to reduce the dimensionality prior to model selection. The DPD-SIS is initially applied to the distinct subsets of data and top $b$ covariates, the

TABLE III
**MODEL SELECTION IN HIGH-DIMENSIONAL:** THE PROPOSED TWO-STAGE ROBUST INFERENCE METHODS ACHIEVE A PERFECT RECOVERY OF TRUE SPARSE BASIS (TP=1) WITH SMALL NUMBER OF FALSE POSITIVES UNDER ZERO CONTAMINATION TO MODERATE CONTAMINATION. IN CONTRAST, THE TWO-STAGE LASSO-LS COMPLETELY FAILS AT RECOVERING THE SPARSE BASIS (TP = 0) FOR ALL SCENARIOS EXCEPT FOR OUTLIER-FREE.

| | TSRD-$\tau$ | | | TSRD-MM | | | TSLL | | |
|---|---|---|---|---|---|---|---|---|---|
| $\xi_o$ | TP | FP | CER | TP | FP | CER | TP | FP | CER |
| 0 | 1 | 0.0004 | 0.0004 | 1 | 0.0013 | 0.0013 | 1 | 0 | 0 |
| 0.1 | 1 | 0.0003 | 0.0003 | 1 | 0.0009 | 0.0009 | 0 | 0 | 0.0013 |
| 0.2 | 1 | 0.0002 | 0.0002 | 1 | 0.0002 | 0.0002 | 0 | 0 | 0.0013 |

[1] Although no truly active variables are selected by TSLL in the model selection stage, one might be mislead by low CER. This can be attributed to fact that the number of truly active variables $k_s$ are insignificant compared to $p$.

exact order of sample size, within each subset of data are kept based on their score and the remaining covariates are discarded. The model selection proceeds with the data set of reduced dimensionality. The proposed two-stage robust inference algorithms are compared to their two-stage non-robust counterpart employing Lasso in the first stage and BLB based on least square estimator in the second stage. For the sake of brevity, the non-robust two-stage inference method is regarded as *two-stage Lasso-LS* (TSLL). Before using Lasso, variable screening is carried out by the DPD-SIS to ensure the data with the reduced dimensionality contains relevant variables. In regard to model selection, the performance is quantified by using confusion matrix and CER as represented in **Table III**. The model selection algorithms in the stage 1 of TSRD-$\tau$ and TSRD-MM could perfectly identify all true non-zero parameters for different proportions of outliers at the highly underdetermined setting, $p/b = 33.33$. In contrast, the non-robust TSLL algorithm fails completely at identifying the sparse basis of the parameter vector except for outlier-free scenario.

The selected variables from the first stage are used to construct confidence intervals based on the bootstrap methods RSOB-$\tau$ and BLFRB. The confidence intervals constructed for the first 15 selected variables are shown in **Fig. 9**. The CIs formed by robust inference methods at the stage 2 of TSRD-$\tau$ and TSRD-MM remain resistant to contamination and length of CIs are slightly inflated with an increase in the proportion of outliers. In regard to outlier-free scenario, the robust bootstrap methods provide reliable estimates of the CIs constructed by bootstrap percentiles of least-square estimator. Hence, it can be concluded that the proposed two-stage inference methods, TSRD-$\tau$ and TSRD-MM, can be used to perform robust statistical inference for large-scale high-dimensional data sets. Note that no comparison is made to two-stage Lasso-LS in the presence of outlying observations due to zero true positive rate, i.e., no element of sparse basis was identified.

### B. Effect of Substituting $\tau$-Lasso with non-regularized $\tau$- and MM-estimators on Variable Selection

We study the effect of substituting the $\tau$-Lasso estimator with $\tau$-estimator and MM-estimator under the settings de-
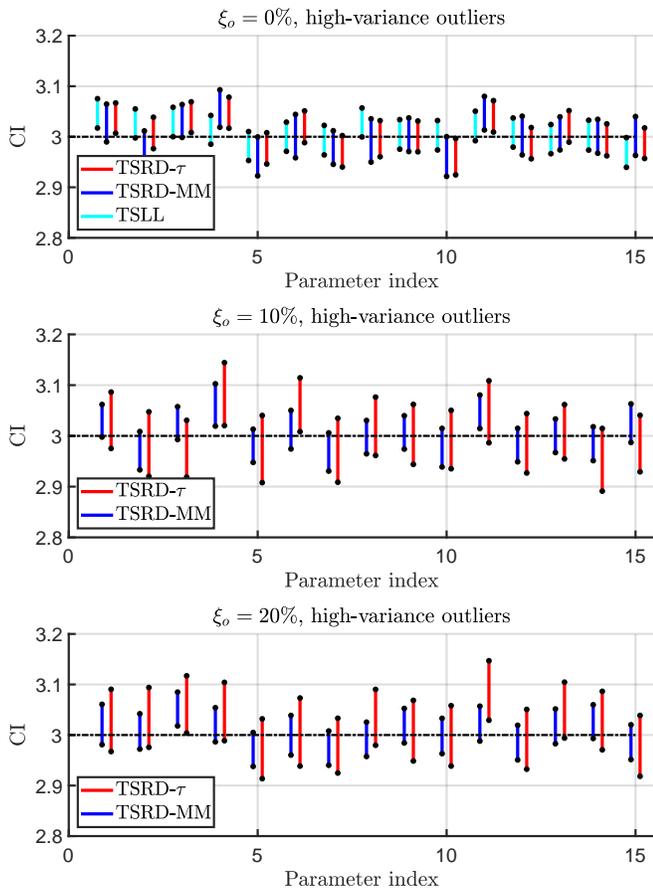
Fig. 9. When there are no outlying observations, the confidence intervals based on bootstraps methods in TSRD-$\tau$ and TSRD-MM provide reliable estimates of the CIs based on bootstrap percentiles of least-square estimator. The confidence intervals produced by TSRD-T and TSRD-MM methods exhibit robustness to outliers and their lengths are slightly affected by an increase in the proportion of outliers. Thus, reliable parameter estimates and confidence intervals are obtained even in the presence of outliers (*the dash-dot lines indicate the true value of non-zero entries of parameter vector*).

scribed by **Scenario 3** and $10\%$ outlier ratio. We conduct the above experiment with 20 trials for which a random realization of outlier is used at each trial. The $\tau$- and MM-estimators were initialized with S-Lasso estimates of the parameter vector $\beta$. We observe in our simulations that both $\tau$- and MM-estimators fail in recovering the correct support and all irrelevant variables are selected within the estimated support In contrast, TSRD-$\tau$ employing the $\tau$-Lasso estimator succeeds in recovering the correct support, indicated by the following contingency table.

TABLE IV
**Model Selection:** The Proposed TSRD-$\tau$ Method Achieves a Perfect Recovery of True Sparse Basis (TP=1) with no False Positives (FP = 0). In Contrast, all irrelevant variables are selected within the estimated support when one replaces the $\tau$-Lasso estimator with non-regularized estimators

|  | TSRD-$\tau$ | | | $\tau$ | | | MM | | |
|---|---|---|---|---|---|---|---|---|---|
| $\xi_o$ | TP | FP | CER | TP | FP | CER | TP | FP | CER |
| 0.1 | 1 | 0 | 0 | 1 | 1 | 0.9 | 1 | 1 | 0.9 |

Note that we obtained the results demonstrated in the above table by averaging over 20 trials.

### C. Effect of Initial Estimate on Variable Selection

In order to explore the influence of initial estimates on the variable selection procedure, we perform a series of simulations with high-dimensional and low-dimensional data. We run the $\tau$-Lasso estimator when randomly initialized and then when initialized with non-regularized $\tau$-estimates (applicable only to low-dimensional regime) and compare the results of model selection with recommended procedure where the $\tau$-Lasso estimator is using S-Lasso estimates as the initial point. In both high-dimensional and low-dimensional data, we carry out a Monte-Carlo study of 20 trials where a random realization of outlier is used at each trial. In case of $\tau$-Lasso estimator with random initialization, we initiate the algorithm with a randomly distributed multivariate Gaussian $\mathcal{N}(1000 \times \mathbf{1}_{p+1}, (250)^2 \times \mathbf{I}_{p+1})$ for each trial and batch of data, chosen to be far from the true coefficient $\beta_0$. We run the simulations on the synthetic data set described by **Scenario 3** for low-dimensional regime and the synthetic data set described by **Scenario 4** for high-dimensional regime.

As shown in Tables V-VI, the model selection algorithm using $\tau$-Lasso achieves the exact support recovery regardless of initialization across all trials for the low-dimensional and high-dimensional regimes. Although the model selection procedure succeeds in perfectly recovering the true support when $\tau$-Lasso estimator initialized randomly. In fact, the $\tau$-Lasso optimization problem is solved via alternating minimization where the sub-problems are non-convex themselves. Recent results show that many well-known nonconvex optimization problems possess a well behaved landscape where all second-order stationary points are global minima [50] and [51]. We conjecture that the variable succeeds in recovering the true support when the $\tau$-Lasso optimization problem is initialized randomly due to potential nice landscape of optimization problem in conjunction with collaborative nature of fusion procedure. We suspect this result may not entirely generalize to all scenarios and this topic requires further study from the optimization perspective.

### D. The Effect of Sample Size, Dimensionality , and Number of Subsamples on Computational Complexity of RSOB-$\tau$

We conduct a number of experiments to examine how the computational complexity of RSOB-$\tau$ scales with sample size,

TABLE V
**INFLUENCE OF INITIAL ESTIMATE ON MODEL SELECTION IN
LOW-DIMENSIONAL DATA:** REGARDLESS OF HOW THE $\tau$-LASSO
ESTIMATOR IS INITIALIZED, THE MODEL SELECTION METHOD
SUCCEEDS IN PERFECTLY RECOVERING THE TRUE SPARSE BASIS (TP=1)
WITH NO FALSE POSITIVES (FP = 0) UNDER CONTAMINATION.

| | TSRD-$\tau$ | | | random init. | | | init. by $\tau$-estimates | | |
|---|---|---|---|---|---|---|---|---|---|
| $\xi_o$ | TP | FP | CER | TP | FP | CER | TP | FP | CER |
| 0.1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0.2 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

TABLE VI
**INFLUENCE OF INITIAL ESTIMATE ON MODEL SELECTION IN
HIGH-DIMENSIONAL DATA:** REGARDLESS OF HOW THE $\tau$-LASSO
ESTIMATOR IS INITIALIZED, THE MODEL SELECTION METHOD
SUCCEEDS IN PERFECTLY RECOVERING THE TRUE SPARSE BASIS (TP=1)
WITH NO FALSE POSITIVES (FP = 0) UNDER CONTAMINATION.

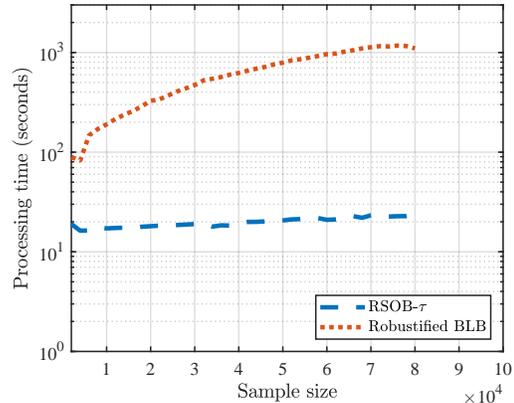| | TSRD-$\tau$ | | | random init. | | |
|---|---|---|---|---|---|---|
| $\xi_o$ | TP | FP | CER | TP | FP | CER |
| 0.1 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0.2 | 1 | 0 | 0 | 1 | 0 | 0 |



Fig. 10. The processing time of RSOB-$\tau$ method employing linearly corrected one-step $\tau$-estimator scales much better with sample size than the robustified BLB method.
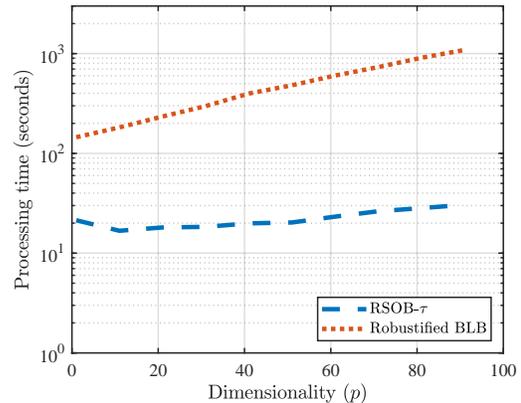


Fig. 11. The processing time of RSOB-$\tau$ method employing linearly corrected one-step $\tau$-estimator scales much better with dimensionality than the robustified BLB method.

dimensionality and the number of subsamples. To do so, we plot the processing time against one of the above parameters while keeping the remaining parameters at their defaults. This process continues until all parameters have been allowed to vary for a range of values. We run the simulations on the data set described by **Scenario 3** where $10\%$ of observations are contaminated by outliers. We then compare them to the result of simulations obtained using robust realization of BLB method based on $\tau$-estimator and report the corresponding results in **Fig. 10-12**. We conduct the experiment on a single node of a high-performance computing cluster (Triton) where 22 computing cores and 25 GB of memory were used. Note that we report the processing time of RSOB-$\tau$ versus that of robustified BLB where a distinct subsample is allocated to each computing core. The number of bootstrap samples is fixed at $B = 300$.

**Fig. 10** shows that as we increase the sample size $n$, the processing time associated with RSOB-$\tau$ grows at much slower rate than the processing time associated with the BLB method based on $\tau$-estimator, thereby achieving higher computational efficiency for larger sample size. In **Fig. 11**, we obtain similar results when plotting the processing time against dimensionality $p$. That is, RSOB-$\tau$ achieves higher computational efficiency for larger dimensions $p$.

**Fig. 12** shows that as we increase the number of subsets (data partitions), the processing time associated with both RSOB-$\tau$ and the robustified BLB decreases and simultaneously the gap between the two curves shrinks. Therefore, there would little benefit in excessively increasing the number of subsets when considering the processing time. We assume that each subsample is processed by allocating one computing core to each subsample. Setting $s$ to very large values leads to wasted computation as we are occupying so many resources while a small portion of each resource is needed. One could use $\lfloor n/n^{0.9} \rfloor < s < \lfloor n/\max(c_s k_s \log(p), n^{0.6}) \rfloor$ as a crude approximation and choose $b = n/s$ so that it uses the computational and storage capabilities of the computational nodes efficiently. As it would be difficult to allocate computing cores for large numbers, we record the processing time associated with each subsample and their maximum is considered to be the processing time for better interpretability.
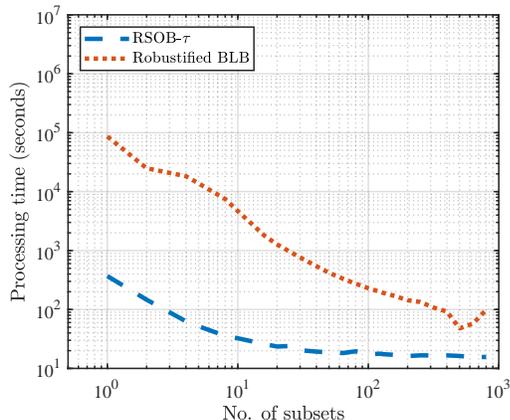
Fig. 12. The above plot indicates if one excessively increases the number of data partitions, small computational gains are made by using RSOB-$\tau$ than using the robustified BLB.

## PROOF OF THEOREM 1

The finite-sample breakdown point of $\tau$-Lasso estimator can be derived in two separate stages: first stage is dedicated to show the boundedness of finite-sample breakdown point from above and second stage is dedicated to show the boundedness of finite-sample breakdown point from below. Once the boundedness from above and below is established, the above theorem is proved immediately. Before proceeding any further, we extend the lemma 5.1 [30] for M-scale estimators to $\tau$-scale estimators as follows:

**Lemma 1**: *Consider any sequence of samples* $\left(\hat{\mathbf{Y}}^{(k)} = \left(\hat{\mathbf{y}}^{(k)}, \hat{\mathbf{X}}^{(k)}\right) \in \mathbb{R}^{b \times (p+1)}\right)_{k \in \mathbb{N}} = \left(\hat{\mathbf{Y}}^{(1)}, \hat{\mathbf{Y}}^{(2)}, \hat{\mathbf{Y}}^{(3)}, \cdots\right)$ *and corresponding residual vector* $\mathbf{r}^{(k)} = \hat{\mathbf{y}}^{(k)} - \hat{\mathbf{X}}^{(k)}\hat{\boldsymbol{\beta}}^{(k)}$ *for* $\hat{\mathbf{Y}}^{(k)}$. *Suppose* $r_l^{(k)}$ *denotes the residual for* $\left(\hat{y}_l^{(k)}, (\hat{\mathbf{x}}_{[l]}^{(k)})^T\right)$ *a row in* $\hat{\mathbf{Y}}^{(k)}$, *Then*

1) *Let* $C = \{l : |r_l^{(k)}| \to \infty\}$, *if* $\#(C) > b\delta$, *then* $\hat{\sigma}_\tau(\mathbf{r}^{(k)}) \to \infty$ *as* $k \to \infty$.
2) *Let* $D = \{l : |r_l^{(k)}| \text{ is bounded }\}$, *if* $\#(D) > b - b\delta$, *then* $\hat{\sigma}_\tau(\mathbf{r}^{(k)})$ *is bounded*.

where $\to \infty$ denotes the left-hand side of arrow tends to infinity. The part 1 of the lemma implies that if the number of unbounded entries of $\mathbf{r}^{(k)}$ exceeds $b\delta$, then $\tau$-scale estimate of $\mathbf{r}^{(k)}$ goes to infinity. On the other hand, if the number of bounded entries of $\mathbf{r}^{(k)}$ exceeds $b - b\delta$, then $\tau$-scale estimate of $\mathbf{r}^{(k)}$ remains bounded. To prove the above lemma, we use the results from lemma 3.2 in [52], the $\tau$-scale is bounded from above and below as follows:

$$\bar{c}_i \hat{\sigma}_{\mathrm{M}}(\mathbf{u}) \leqslant \hat{\sigma}_\tau(\mathbf{u}) \leqslant \sqrt{\sup_{t \in \mathbb{R}} \rho_1(t)} \hat{\sigma}_{\mathrm{M}}(\mathbf{u}) \quad \forall \mathbf{u} \in \mathbb{R}^b, \quad (34)$$

where $\bar{c}_i$ is a positive constant, $\hat{\sigma}_{\mathrm{M}}$ denotes the M-scale estimate and $\hat{\sigma}_\tau$ denotes the $\tau$-scale estimate. Now, we can take advantage of the inequality associating M-scale estimate with $\tau$-scale estimate. We extend the lemma 5.1 in p.184 of [30] for M-scale estimators and derive similar results for $\tau$-scale estimators. According to part 1 of lemma 5.1 in [30],

$$\text{if } \#(C) > b\delta \Rightarrow \hat{\sigma}_{\mathrm{M}}(\mathbf{r}^{(k)}) \to \infty,$$
$$\bar{c}_i \times \infty \leqslant \hat{\sigma}_\tau(\mathbf{r}^{(k)}) \leqslant \sqrt{\sup_{t \in \mathbb{R}} \rho_1(t)} \times \infty, \quad (35)$$

which implies both upper and lower bound of $\tau$-scale estimate goes to infinity. This proves part 1 of the lemma, if $\#(C) > b\delta$ then $\hat{\sigma}_\tau(\mathbf{r}^{(k)}) \to \infty$. To prove the second part of lemma, we know from part 2 of lemma 5.1 [30],

$$\#(D) > b - b\delta \Rightarrow \hat{\sigma}_{\mathrm{M}}(\mathbf{r}^{(k)}) \text{ is bounded},$$
$$\bar{c}_i \times \text{bounded} \leqslant \hat{\sigma}_\tau(\mathbf{r}^{(k)}) \leqslant \sqrt{\sup_{t \in \mathbb{R}} \rho_1(t)} \times \text{bounded}, \quad (36)$$

which implies both lower and upper bounds of $\tau$-scale estimates are bounded. This proves the second part of the lemma, if $\#(D) > b - b\delta$, then $\hat{\sigma}_\tau(\mathbf{r}^{(k)})$ is bounded. Now, we need to prove in two separate stages that the finite-sample breakdown point of the $\tau$-Lasso estimator is bounded from the above and the below.

### E. Bounded From Below

In order to establish the boundedness of $\hat{\boldsymbol{\beta}}$ from below, we need to show that the sequence of $\tau$-scale estimates $\left(\hat{\boldsymbol{\beta}}^{(k)} \in \mathbb{R}^p\right)_{k \in \mathbb{N}}$ is bounded for any arbitrary sequence of contaminated samples $\left(\hat{\mathbf{Y}}_m^{(k)}\right)_{k \in \mathbb{N}}$ with $m \leqslant m(\delta)$. The method of proof by contradiction can be used to establish the boundedness from below. That is, it is assumed the sequence $\left(\hat{\boldsymbol{\beta}}^{(k)}\right)_{k \in \mathbb{N}}$ is unbounded and then shown that $\hat{\boldsymbol{\beta}}^{(k)}$ violates the optimality condition where it is assumed to attain the minimum of the $\tau$-Lasso objective function. Thus, $\hat{\boldsymbol{\beta}}^{(k)}$ should be bounded to be a minimizer of the objective function.

Suppose $\tilde{\boldsymbol{\beta}} \in \mathbb{R}^p$ has a bounded $\ell_1$-norm such that $\|\tilde{\boldsymbol{\beta}}\|_{\ell_1} = K_b < \infty$. For the uncontaminated observations $(\check{y}_l, \check{\mathbf{x}}_{[l]}^T)$ within the contaminated sample $\hat{\mathbf{Y}}_m^{(k)}$, the corresponding residuals $|r_l^{(k)}(\tilde{\boldsymbol{\beta}}, \hat{\mathbf{Y}}_m^{(k)})| = |\check{y}_l - \check{\mathbf{x}}_{[l]}^T \tilde{\boldsymbol{\beta}}| < \infty$ are bounded based on the triangle inequality. Without loss of generality, it is assumed $b\delta = b\min(\delta, 1 - \delta)$. Since the inequality $m \leqslant m(\delta) \leqslant b\delta$ holds based on the theorem assumption, it can be shown the number of bounded residuals $\#(D) \geqslant b - m \geqslant b - b\delta$. Therefore, we can conclude $\hat{\sigma}_\tau\left(\mathbf{r}^{(k)}(\tilde{\boldsymbol{\beta}}, \hat{\mathbf{Y}}_m^{(k)})\right)$ can be large but still bounded based on the above lemma,

$$\sup_{k \in \mathbb{N}} \hat{\sigma}_\tau\left(\mathbf{r}^{(k)}(\tilde{\boldsymbol{\beta}}, \hat{\mathbf{Y}}_m^{(k)})\right) < \infty. \quad (37)$$

Now, let the sequence $\left(\|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}\right)_{k \in \mathbb{N}}$ be unbounded. An unbounded sequence does not converge, i.e. sequences that contain arbitrarily large numbers. Hence, there exists a sequence index $k_0$ such that $\|\hat{\boldsymbol{\beta}}^{(k_0)}\|_{\ell_1} > K_b + \frac{1}{\lambda}\sup_{k \in \mathbb{N}} \hat{\sigma}_\tau^2\left(\mathbf{r}^{(k)}(\tilde{\boldsymbol{\beta}}, \hat{\mathbf{Y}}_m^{(k)})\right)$ and as a result we can say for every $k' \geqslant k_0$,

$$\overbrace{\frac{\mathcal{L}\left(\hat{\boldsymbol{\beta}}^{(k')},\hat{\mathbf{Y}}_m^{(k')}\right)}{\hat{\sigma}_\tau^2\left(\mathbf{r}^{(k')}\left(\hat{\boldsymbol{\beta}}^{(k')},\hat{\mathbf{Y}}_m^{(k')}\right)\right)+\lambda\|\hat{\boldsymbol{\beta}}^{(k')}\|_{\ell_1}}}$$
$$> \hat{\sigma}_\tau^2\left(\mathbf{r}^{(k')}\left(\hat{\boldsymbol{\beta}}^{(k')},\hat{\mathbf{Y}}_m^{(k')}\right)\right)$$
$$+\lambda\left(K_b+\frac{1}{\lambda}\sup_{k\in\mathbb{N}}\hat{\sigma}_\tau^2\left(\mathbf{r}^{(k')}\left(\tilde{\boldsymbol{\beta}},\hat{\mathbf{Y}}_m^{(k)}\right)\right)\right) \quad (38)$$
$$\geqslant \hat{\sigma}_\tau^2\left(\mathbf{r}^{(k')}\left(\tilde{\boldsymbol{\beta}},\hat{\mathbf{Y}}_m^{(k')}\right)\right)+\lambda K_b = \mathcal{L}\left(\tilde{\boldsymbol{\beta}},\hat{\mathbf{Y}}_m^{(k')}\right) \Rightarrow$$
$$\mathcal{L}\left(\hat{\boldsymbol{\beta}}^{(k')},\hat{\mathbf{Y}}_m^{(k')}\right) \geqslant \mathcal{L}\left(\tilde{\boldsymbol{\beta}},\hat{\mathbf{Y}}_m^{(k')}\right)$$

where $\lambda$ determines the amount of regularization imposed by the $\ell_1$-norm and $\mathcal{L}(\cdot)$ is $\tau$-Lasso objective function as given in equation (38). The above result contradicts the fact that $\hat{\boldsymbol{\beta}}^{(k')}$ is the minimizer of the objective function, as the loss function is larger for $\hat{\boldsymbol{\beta}}^{(k')}$. Thus, $\hat{\boldsymbol{\beta}}^{(k)}$ should be bounded for $m \leqslant m(\delta)$ and the boundedness from below is proved.

### F. Bounded From Above

The boundedness from above can be established by showing that the estimator breaks down for $m > b\delta$. In order to prove such a property, proof by contradiction is used in this work. That is, it is assumed that the sequence of estimates $\left(\hat{\boldsymbol{\beta}}^{(k)}\right)_{k\in\mathbb{N}}$ minimizing the $\tau$-Lasso objective function over the contaminated sample $\hat{\mathbf{Y}}_m^{(k)}$ is bounded. and then shown the $\tau$-Lasso objective function evaluated at $\hat{\boldsymbol{\beta}}^{(k)}$ achieves larger value than the $\tau$-Lasso objective function evaluated at a given unbounded $\tilde{\boldsymbol{\beta}}^{(k)}$. This contradicts the assumption that $\hat{\boldsymbol{\beta}}^{(k)}$ is the minimum of $\tau$-Lasso objective function. Thus, it can be concluded $\hat{\boldsymbol{\beta}}^{(k)}$ must be unbounded for $m > b\delta$ establishing the boundedness from above.

As the $\tau$-Lasso objective function is comprised of three primary components, its evaluation is carried out in three separate steps as follows:

1) *M-scale of residuals, component 1*
2) $\frac{1}{b}\sum_{l=1}^b \rho_1\left(\frac{\hat{r}_l^{(k)}}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}\right)$, *component 2*
3) $\lambda\|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}$, *component 3*

Once these components are evaluated, the $\tau$-Lasso objective function is formed by a simple addition and multiplication operation over these components. In order to proceed with the proof, the $\tau$-Lasso objective function is evaluated for the bounded sequences $\hat{\boldsymbol{\beta}}^{(k)}$ and unbounded sequences $\tilde{\boldsymbol{\beta}}^{(k)}$ and then compared with each other to arrive at a contradiction.

#### 1) Evaluation of $\tau$-Lasso Objective Function for $\hat{\boldsymbol{\beta}}^{(k)}$
##### a) Component 1
Suppose the set $C \subset \{1,\cdots,b\}$ denotes the indices of the observations within the original contaminated-free sample $\check{\mathbf{Y}}$ replaced by outliers to construct the contaminated sample $\hat{\mathbf{Y}}_m^{(k)}$ with $m = \#(C)$. To simplify the proof without loss of the generality, we choose an arbitrary $\mathbf{x}_0 \in \mathbb{R}^p$ with unit $\ell_2$-norm,

$\|\mathbf{x}_0\|_{\ell_2} = 1$. Now, we construct a contaminated sequence of samples with $m$ outliers given by

$$\left(\hat{y}_l^{(k)},\hat{\mathbf{x}}_{[l]}^{(k)}\right) = \begin{cases} \left(k^{\nu+1},\mathbf{x}_0 k\right) & l \in C \\ \left(\check{y}_l,\check{\mathbf{x}}_{[l]}\right) & l \notin C \end{cases}, \quad (39)$$

where $0 < \nu \leqslant 1$ and $\left(k^{\nu+1},\mathbf{x}_0 k\right)$ are chosen to account for outlying observations. The sequence for outlying observations diverges as $k$ goes to infinity.

First, we assume that $\hat{\boldsymbol{\beta}}^{(k)}$ is bounded in norm and consequently, we have $|r_l^{(k)}(\hat{\boldsymbol{\beta}}^k,\hat{\mathbf{Y}}_m^{(k)})| = |\hat{y}_l - (\hat{\mathbf{x}}_{[l]}^{(k)})^T\hat{\boldsymbol{\beta}}^{(k)}| < \infty$ is bounded for $l \notin C$ and $k \in \mathbb{N}$. $\hat{r}_l^{(k)}$ is a shorthad for $r_l^{(k)}(\hat{\boldsymbol{\beta}}^{(k)},\hat{\mathbf{Y}}_m^{(k)})$. On the other hand, the residuals corresponding to contaminated observations, $l \in C$ are lower bounded by

$$|\hat{r}_l^{(k)}| = |k^{\nu+1} - k\mathbf{x}_0^T\hat{\boldsymbol{\beta}}^{(k)}| = k|k^\nu - \mathbf{x}_0^T\hat{\boldsymbol{\beta}}^{(k)}| \text{ and } \mathbf{x}_0^T\hat{\boldsymbol{\beta}}^{(k)}$$
$$\leqslant \|\mathbf{x}_0\|_{\ell_1}\|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1} \Rightarrow |k^\nu - \mathbf{x}_0^T\hat{\boldsymbol{\beta}}^{(k)}| \geqslant |k^\nu - \|\mathbf{x}_0\|_{\ell_1}\|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}|$$
$$\Rightarrow |\hat{r}_l^{(k)}| \geqslant k|k^\nu - \|\mathbf{x}_0\|_{\ell_1}\|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}| \quad (40)$$

In addition, the right-hand side of inequality goes to infinity as $k$ approaches infinity which implies the residuals $\hat{r}_l^{(k)}$ go to infinity for $l \in C$ as well. Based on the above lemma and lemma 5.1 in p.184 of [30] , we conclude that both $\hat{\sigma}_\tau(\hat{\mathbf{r}}^{(k)})$ and $\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})$ go to infinity for $\#(C) > b\delta$. As a result, we can decompose the M-estimation of scale equation as follows:

$$\sum_{l\notin C} \rho_0\left(\frac{\hat{r}_l^{(k)}}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}\right) + \sum_{l\in C} \rho_0\left(\frac{\hat{r}_l^{(k)}}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}\right) = b\delta \quad (41)$$

Recalling the proof of Theorem 4.1 in [53], it follows

$$\rho_0\left(\frac{1}{\gamma}\right) = \frac{b\delta}{m} \quad (42)$$

where

$$\gamma = \lim_{k\to\infty} \frac{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}{k^{\nu+1}} \quad (43)$$

##### b) Component 2
On the other hand, we have

$$\frac{1}{b}\sum_{l=1}^b \rho_1\left(\frac{\hat{r}_l^{(k)}}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}\right) = \frac{m}{b}\rho_1\left(\frac{1-\mathbf{x}_0^T\hat{\boldsymbol{\beta}}^{(k)}/k^\nu}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})/k^{\nu+1}}\right) \Rightarrow$$
$$\lim_{k\to\infty} \frac{m}{b}\rho_1\left(\frac{1-\mathbf{x}_0^T\hat{\boldsymbol{\beta}}^{(k)}/k^\nu}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})/k^{\nu+1}}\right) = \frac{m}{b}\rho_1\left(\frac{1}{\gamma}\right) \quad (44)$$

##### c) Component 3
$\lim_{k\to\infty}\|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}$ remains bounded according to the assumption of the proof.

*d) Deriving $\tau$-Lasso Objective Function*

Now, we can evaluate the $\tau$-Lasso loss function according to equation (38) as follows:

$$\lim_{k\to\infty} \frac{\hat{\sigma}_\tau^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} = \lim_{k\to\infty} \frac{\hat{\sigma}_M^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} \times \frac{m}{b}\rho_1\left(\frac{1-\mathbf{x}_0^T\hat{\boldsymbol{\beta}}^{(k)}/k^\nu}{\hat{\sigma}_M^2(\hat{\mathbf{r}}^{(k)})/k^{\nu+1}}\right)$$

$$\Rightarrow \lim_{k\to\infty} \frac{\hat{\sigma}_\tau^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} = \frac{m\gamma^2}{b}\rho_1\left(\frac{1}{\gamma}\right) \Rightarrow \lim_{k\to\infty} \frac{\mathcal{L}(\hat{\boldsymbol{\beta}}^{(k)}, \hat{\mathbf{Y}}_m^{(k)})}{k^{2\nu+2}}$$

$$= \lim_{k\to\infty}\left(\frac{\hat{\sigma}_\tau^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} + \lambda\frac{\|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}}{k^{2\nu+2}}\right) = \frac{m\gamma^2}{b}\rho_1\left(\frac{1}{\gamma}\right) \tag{45}$$

where $\lim_{k\to\infty} \|\hat{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}/k^{2\nu+2} = 0$ due to the bounded norm of $\hat{\boldsymbol{\beta}}^{(k)}$.

*2) Evaluation of $\tau$-Lasso Objective Function for $\tilde{\boldsymbol{\beta}}^{(k)}$*

Now, we evaluate the $\tau$-Lasso objective function for an unbounded sequence $\tilde{\boldsymbol{\beta}}^{(k)}$ in a three-step procedure as follows:

*a) Component 1*

The unbounded sequence is assumed to be $\tilde{\boldsymbol{\beta}}^{(k)} = \frac{k^\nu}{2}\mathbf{x}_0$. The residuals for this sequence become

$$\hat{r}_l^{(k)} = \begin{cases} k^{\nu+1} - \frac{k^{\nu+1}}{2}\mathbf{x}_0^T\mathbf{x}_0 & l \in C \\ \check{y}_l - \frac{k^\nu}{2}\mathbf{x}_0^T\check{\mathbf{x}}_{[l]} & l \notin C \end{cases} = \begin{cases} \frac{k^{\nu+1}}{2} & l \in C \\ \check{y}_l - \frac{k^\nu}{2}\mathbf{x}_0^T\check{\mathbf{x}}_{[l]} & l \notin C \end{cases} \tag{46}$$

where $\mathbf{x}_0^T\mathbf{x}_0 = \|\mathbf{x}_0\|_{\ell_2}^2 = 1$ and all residuals go to infinity as $k \to \infty$. Hence, both $\hat{\sigma}_\tau(\hat{\mathbf{r}}^{(k)})$ and $\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})$ tend to infinity. The decomposition of M-estimation of scale equation yields,

$$\sum_{l\notin C}\rho_0\left(\frac{\check{y}_l - \frac{k^\nu}{2}\mathbf{x}_0^T\check{\mathbf{x}}_{[l]}}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}\right) + \sum_{l\in C}\rho_0\left(\frac{k^{\nu+1}/2}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}\right) = b\delta \tag{47}$$

Using the proof of Theorem 4.1 in [53] , it can be inferred

$$\rho_0\left(\frac{1}{\lim_{k\to\infty}\frac{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}{k^{\nu+1}/2}}\right) = \frac{b\delta}{m} \tag{48}$$

where

$$\lim_{k\to\infty}\frac{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}{k^{\nu+1}/2} = \gamma \tag{49}$$

*b) Component 2*

On the other hand, we have

$$\frac{1}{b}\sum_{l=1}^{b}\rho_1\left(\frac{\hat{r}_l^{(k)}}{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}\right) = \frac{m}{b}\rho_1\left(\frac{1}{\frac{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}{k^{\nu+1}/2}}\right) \Rightarrow$$

$$\lim_{k\to\infty}\frac{m}{b}\rho_1\left(\frac{1}{\frac{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}{k^{\nu+1}/2}}\right) = \frac{m}{b}\rho_1\left(\frac{1}{\gamma}\right) \tag{50}$$

*c) Component 3*

$\lim_{k\to\infty} \|\tilde{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}$ diverges as $k$ goes to infinity.

*d) Deriving $\tau$-Lasso objective function*

Now, we can evaluate the $\tau$-Lasso loss function according to equation (38) as follows:

$$\lim_{k\to\infty} \frac{\hat{\sigma}_\tau^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} = \lim_{k\to\infty} \frac{\hat{\sigma}_M^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} \times \frac{m}{b}\rho_1\left(\frac{1}{\frac{\hat{\sigma}_M(\hat{\mathbf{r}}^{(k)})}{k^{\nu+1}/2}}\right)$$

$$\Rightarrow \lim_{k\to\infty} \frac{\hat{\sigma}_\tau^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} = \frac{m\gamma^2}{4b}\rho_1\left(\frac{1}{\gamma}\right) \Rightarrow \lim_{k\to\infty} \frac{\mathcal{L}(\tilde{\boldsymbol{\beta}}^{(k)}, \hat{\mathbf{Y}}_m^{(k)})}{k^{2\nu+2}}$$

$$= \lim_{k\to\infty}\left(\frac{\hat{\sigma}_\tau^2(\hat{\mathbf{r}}^{(k)})}{k^{2\nu+2}} + \lambda\frac{\|\tilde{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}}{k^{2\nu+2}}\right) = \frac{m\gamma^2}{4b}\rho_1\left(\frac{1}{\gamma}\right) \tag{51}$$

where

$$\lim_{k\to\infty} \|\tilde{\boldsymbol{\beta}}^{(k)}\|_{\ell_1}/k^{2\nu+2} = \lim_{k\to\infty} \|k^\nu\mathbf{x}_0/2\|_{\ell_1}/k^{2\nu+2}$$
$$= \lim_{k\to\infty} \|\mathbf{x}_0/2\|_{\ell_1}/k^{\nu+2} = 0 \tag{52}$$

*3) Comparison*

Now, we can compare the $\tau$-Lasso objective function for the given bounded and unbounded sequences and conclude that for large enough $k_0$,

$$\frac{\mathcal{L}(\tilde{\boldsymbol{\beta}}^{(k)}, \hat{\mathbf{Y}}_m^{(k)})}{k^{2\nu+2}} < \frac{\mathcal{L}(\hat{\boldsymbol{\beta}}^{(k)}, \hat{\mathbf{Y}}_m^{(k)})}{k^{2\nu+2}}, \quad \forall k \geqslant k_0 \tag{53}$$

The above results contradict the fact that the bounded $\hat{\boldsymbol{\beta}}^{(k)}$ is the minimum of the $\tau$-Lasso objective function for the contaminated sample with $m > b\delta$. Because the $\tau$-Lasso objective function for the unbounded $\tilde{\boldsymbol{\beta}}^{(k)}$ is smaller than that of $\hat{\boldsymbol{\beta}}^{(k)}$. This implies the $\hat{\boldsymbol{\beta}}^{(k)}$ have to be unbounded and thus, the robust $\tau$-Lasso estimator breaks down for $m > b\delta$.

### PROOF OF THEOREM 2

To begin with the proof, it follows from the first-order condition that $\tau$-estimates of regression parameter and scale for the given subset of data, $\hat{\boldsymbol{\beta}}_b$ and $\hat{\sigma}_b$ must satisfy the following equations [29]:

$$\frac{1}{b}\sum_{l=1}^{b}\left[\hat{w}_\tau\rho_0'\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right) + \rho_1'\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right)\right]\check{\mathbf{x}}_{[l]} = \mathbf{0},$$
$$\frac{1}{b}\sum_{l=1}^{b}\rho_0\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right) = \delta, \tag{54}$$

where $\hat{r}_l = \check{y}_l - \check{\mathbf{x}}_{[l]}^T\hat{\boldsymbol{\beta}}_b$ and $\hat{w}_\tau$ is given by

$$\hat{w}_\tau = \frac{\sum_{l=1}^{b}\left[2\rho_1\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right) - \rho_1'\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right)\frac{\hat{r}_l}{\hat{\sigma}_b}\right]}{\sum_{l=1}^{b}\rho_0'\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right)\frac{\hat{r}_l}{\hat{\sigma}_b}}. \tag{55}$$

Therefore, we can obtain $\tau$-estimates of regression $\hat{\boldsymbol{\beta}}_b$ and scale $\hat{\sigma}_b$ for the subset of data with observations $(\check{\mathbf{y}}, \check{\mathbf{X}})$ as follows:

$$\hat{\boldsymbol{\beta}}_b = \mathbf{A}_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b)^{-1}\mathbf{v}_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b),$$
$$\hat{\sigma}_b = \hat{\sigma}_b u_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b), \tag{56}$$

where

$$\mathbf{A}_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b) = \frac{1}{b} \sum_{l=1}^{b} \hat{w}_l \tilde{\mathbf{x}}_{[l]} \tilde{\mathbf{x}}_{[l]}^T,$$

$$\mathbf{v}_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b) = \frac{1}{b} \sum_{l=1}^{b} \hat{w}_l \breve{y}_l \tilde{\mathbf{x}}_{[l]},$$

$$u_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b) = \frac{1}{b\delta} \sum_{l=1}^{b} \rho_0\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right),$$

$$\hat{w}_l = \frac{\hat{w}_\tau \rho_0'\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right) + \rho_1'\left(\frac{\hat{r}_l}{\hat{\sigma}_b}\right)}{\hat{r}_l}.$$

(57)

Alternatively, the robust $\tau$-estimates of regression parameters and scale can be written as the solution of a fixed-point problem as follows:

$$\hat{\boldsymbol{\theta}}_b = \mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}), \qquad (58)$$

where $\mathbf{f} : \mathbb{R}^{(|\hat{\mathcal{S}}|+1)} \to \mathbb{R}^{(|\hat{\mathcal{S}}|+1)}$, $\hat{\boldsymbol{\theta}}_b \in \mathbb{R}^{(|\hat{\mathcal{S}}|+1)}$ and $\mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}})$ are given by

$$\hat{\boldsymbol{\theta}}_b = \begin{bmatrix} \hat{\boldsymbol{\beta}}_b \\ \hat{\sigma}_b \end{bmatrix},$$

$$\mathbf{f}(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}) = \begin{bmatrix} \mathbf{A}_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b)^{-1} \mathbf{v}_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b) \\ \hat{\sigma}_b u_b(\hat{\boldsymbol{\beta}}_b, \hat{\sigma}_b) \end{bmatrix}.$$

(59)

Conditioned on $\hat{\mathcal{S}} = \mathcal{S}$, $|\hat{\mathcal{S}}|$ can be replaced with $k_s$. Given $\rho_0$ and $\rho_1$ are differentiable functions, we can expand $\mathbf{f}$ by using Taylor expansion around the limiting values, $\boldsymbol{\theta}_0 = [\boldsymbol{\beta}_0, \sigma_0]^T$, as follows:

$$\hat{\boldsymbol{\theta}}_b = \mathbf{f}(\boldsymbol{\theta}_0) + \nabla \mathbf{f}(\boldsymbol{\theta}_0)(\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0)$$
$$+ \underbrace{\frac{1}{2}\left[\mathbf{I}_{k_s} \otimes (\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0)^T\right] \nabla^2 \mathbf{f}(\bar{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0)}_{R_b}, \qquad (60)$$

where $\mathbf{f}(\boldsymbol{\theta}_0)$ is a short-hand for $\mathbf{f}(\boldsymbol{\theta}_0; \tilde{\mathbf{Y}})$, $\nabla \mathbf{f}(\cdot) \in \mathbb{R}^{(k_s+1) \times (k_s+1)}$ is the matrix of partial derivatives, $\nabla^2 \mathbf{f}(\cdot) \in \mathbb{R}^{(k_s+1)^2 \times (k_s+1)}$ is the Hessian matrix of $\mathbf{f}(\cdot)$, $\otimes$ denotes the Kronecker product and $\bar{\boldsymbol{\theta}}$ lies on the line segment between $\boldsymbol{\theta}_0$ and $\hat{\boldsymbol{\theta}}_b$. The term $\left[\mathbf{I}_{k_s} \otimes (\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0)^T\right] \nabla^2 \mathbf{f}(\bar{\boldsymbol{\theta}})$ in the remainder is a $(k_s+1)^2 \times (k_s+1)$ matrix whose $(i, j)$-entry is given by

$$M_{ij} = (\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0)^T \left[\frac{\partial^2 \mathbf{f}(\bar{\boldsymbol{\theta}})}{\partial \theta_i \partial \theta_j}\right]. \qquad (61)$$

In addition, $\nabla \mathbf{f}(\cdot)$ is defined as follows:

$$\nabla \mathbf{f}(\boldsymbol{\theta}) = \begin{bmatrix} \partial[(\mathbf{A}_b)^{-1} \mathbf{v}_b]/\partial \boldsymbol{\beta} & \partial[(\mathbf{A}_b)^{-1} \mathbf{v}_b]/\partial \sigma \\ \partial[\sigma u_b]/\partial \boldsymbol{\beta} & \partial[\sigma u_b]/\partial \sigma \end{bmatrix}. \qquad (62)$$

Here, $u_b$, $\mathbf{A}_b$ and $\mathbf{v}_b$ are short-hands for $u_b(\boldsymbol{\theta})$, $\mathbf{A}_b(\boldsymbol{\theta})$ and $\mathbf{v}_b(\boldsymbol{\theta})$, respectively. Tedious but straightforward calculations show that the second-order terms $\partial^2 \mathbf{f}(\bar{\boldsymbol{\theta}})/(\partial \theta_i \partial \theta_j)$ are a combination of sample mean products. Taking into account the convergence of the sample mean products to their corresponding population mean according to Lemma 2 in [21], (an extension of law of large numbers) and continuity of derivatives of $\rho_0$ and $\rho_1$, we can guarantee

$\partial^2 \mathbf{f}(\bar{\boldsymbol{\theta}})/(\partial \theta_i \partial \theta_j) = O_p(1)$ for $i, j = 1, \cdots, k_s + 1$. On the other hand, it follows from the root-$n$ consistency of estimators [54] that $\|\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0\|_{\ell_2} = O_p(1/\sqrt{b})$. Hence, the $(i, j)$-entry $M_{ij} = O_p(1/\sqrt{b})$. Noting that $i_{th}$ entry in the remainder is a linear combination $\sum_{j=1}^{k_s+1} M_{ij}([\hat{\theta}_b]_j - [\theta_0]_j) = o_p(1/\sqrt{b})$, implying the remainder term $R_b = o_p(1/\sqrt{b})$. Therefore, we can re-express the Taylor expansion given in equation (60) as follows:

$$(\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0) = (\mathbf{f}(\boldsymbol{\theta}_0) - \boldsymbol{\theta}_0) + \nabla \mathbf{f}(\boldsymbol{\theta}_0)(\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0) + o_p(1/\sqrt{b}) \Rightarrow$$
$$[\mathbf{I} - \nabla \mathbf{f}(\boldsymbol{\theta}_0)](\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0) = (\mathbf{f}(\boldsymbol{\theta}_0) - \boldsymbol{\theta}_0) + o_p(1/\sqrt{b}) \Rightarrow$$
$$\sqrt{b}(\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0) = [\mathbf{I} - \nabla \mathbf{f}(\boldsymbol{\theta}_0)]^{-1} \sqrt{b}[\mathbf{f}(\boldsymbol{\theta}_0) - \boldsymbol{\theta}_0] + o_p(1).$$

(63)

On the other hand, we know that $o_p(1)$ term converges to $\mathbf{0}_{k_s+1}$ in probability as $b$ tends to infinity. As a result, both sides of the following converge to the same limiting distribution.

$$\sqrt{b}(\hat{\boldsymbol{\theta}}_b - \boldsymbol{\theta}_0) \sim [\mathbf{I} - \nabla \mathbf{f}(\boldsymbol{\theta}_0)]^{-1} \sqrt{b}[\mathbf{f}(\boldsymbol{\theta}_0) - \boldsymbol{\theta}_0], \qquad (64)$$

where the notation $\sim$ stands for weak convergence of both sides to the same limiting distribution. Consider $n$ bootstrap samples are drawn from the given subset of data, we approximate the actual bootstrap estimates $\hat{\boldsymbol{\theta}}_{n,b}^\star$ for the given subset of data using linearly corrected one-step bootstrap estimates $\hat{\boldsymbol{\theta}}_{n,b}^{R\star}$ as follows:

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_{n,b}^{R\star} - \hat{\boldsymbol{\theta}}_b) \sim [\mathbf{I} - \nabla \mathbf{f}(\hat{\boldsymbol{\theta}}_b)]^{-1} \sqrt{n}[\mathbf{f}^\star(\hat{\boldsymbol{\theta}}_b) - \hat{\boldsymbol{\theta}}_b], \quad (65)$$

where $n$ denotes the number of observations in the complete data set and $\mathbf{f}^\star(\hat{\boldsymbol{\theta}}_b)$ is one-step bootstrap estimate and shorthand for $\mathbf{f}^\star(\hat{\boldsymbol{\theta}}_b; \tilde{\mathbf{Y}}^\star)$ where $\tilde{\mathbf{Y}}^\star = \left(\breve{\mathbf{y}}, \tilde{\mathbf{X}}; \omega^\star\right) \in \mathbb{R}^{n \times (k_s+1)}$. Since $[\mathbf{I} - \nabla \mathbf{f}(\hat{\boldsymbol{\theta}}_b)]^{-1}$ is a consistent estimator of $[\mathbf{I} - \nabla \mathbf{f}(\boldsymbol{\theta}_0)]^{-1}$, we only need to show that $\sqrt{n}[\mathbf{f}^\star(\hat{\boldsymbol{\theta}}_b) - \hat{\boldsymbol{\theta}}_b]$ converges to the same limiting distribution as $\sqrt{b}[\mathbf{f}(\boldsymbol{\theta}_0) - \boldsymbol{\theta}_0]$. To proceed with estimation of the correction matrix $[\mathbf{I} - \nabla \mathbf{f}(\hat{\boldsymbol{\theta}}_b)]^{-1}$, we compute the gradient matrix $\mathbf{I} - \nabla \mathbf{f}(\hat{\boldsymbol{\theta}}_b)$ as follows:

$$\mathbf{I} - \nabla \mathbf{f}(\hat{\boldsymbol{\theta}}_b) = \left[\begin{array}{c|c} \mathcal{A} & \boldsymbol{\eta} \\ \hline \boldsymbol{\zeta} & a \end{array}\right]. \qquad (66)$$

where

$$\mathcal{A} = \mathbf{I} - \frac{\partial[(\mathbf{A}_b)^{-1} \mathbf{v}_b]}{\partial \boldsymbol{\beta}}\Big|_{\hat{\boldsymbol{\theta}}_b}, \quad \boldsymbol{\eta} = -\frac{\partial[(\mathbf{A}_b)^{-1} \mathbf{v}_b]}{\partial \sigma}\Big|_{\hat{\boldsymbol{\theta}}_b}$$
$$, a = 1 - \frac{\partial[\sigma u_b]}{\partial \sigma}\Big|_{\hat{\boldsymbol{\theta}}_b}, \boldsymbol{\zeta} = -\frac{\partial[\sigma u_b]}{\partial \boldsymbol{\beta}}\Big|_{\hat{\boldsymbol{\theta}}_b}.$$

(67)

Let's begin with calculating $\boldsymbol{\zeta}$,

$$\boldsymbol{\zeta} = -\frac{\partial[\sigma u_b]}{\partial \boldsymbol{\beta}}\Big|_{\hat{\boldsymbol{\theta}}_b} \Rightarrow -\frac{\partial[\sigma u_b]}{\partial \boldsymbol{\beta}} = -\left(\frac{\partial \sigma}{\partial \boldsymbol{\beta}} u_b + \sigma \frac{\partial u_b}{\partial \boldsymbol{\beta}}\right) \Rightarrow$$
$$\boldsymbol{\zeta} = \frac{1}{b\delta} \sum_{l=1}^{b} \rho_0'(\breve{r}_l) \tilde{\mathbf{x}}_{[l]}^T.$$

(68)

where $\check{r}_l$ denotes a shorthand for $\check{r}_l(\boldsymbol{\beta}) = \check{y}_l - \tilde{\mathbf{x}}_{[l]}^T\boldsymbol{\beta}$ and $\check{r}_l = \hat{r}_l/\hat{\sigma}_b$. In order to calculate $a$, we need to derive $\partial[\sigma u_b]/\partial\sigma$,

$$\frac{\partial[\sigma u_b]}{\partial\sigma} = \left(\frac{\partial\sigma}{\partial\sigma}u_b + \sigma\frac{\partial u_b}{\partial\sigma}\right) = \frac{1}{b\delta}\left[\sum_{l=1}^{b}\rho_0\left(\frac{\check{r}_l}{\sigma}\right) - \sum_{l=1}^{b}\rho_0'\left(\frac{\check{r}_l}{\sigma}\right)\frac{\check{r}_l}{\sigma}\right].$$
(69)

Therefore, $a$ can be derived as follows:

$$a = 1 - \frac{\partial[\sigma u_b]}{\partial\sigma}\Big|_{\hat{\boldsymbol{\theta}}_b} = \frac{1}{b\delta}\sum_{l=1}^{b}\rho_0'(\check{r}_l)\check{r}_l.$$
(70)

Finding $\mathcal{A}$ requires differentiating $(\mathbf{A}_b)^{-1}\mathbf{v}_b$ with respect to $\boldsymbol{\beta}$. To do so, $\boldsymbol{\alpha}_b = (\mathbf{A}_b)^{-1}\mathbf{v}_b$ is defined to simplify the derivations as follows:

$$\mathbf{v}_b = \mathbf{A}_b\boldsymbol{\alpha}_b \Rightarrow \frac{\partial}{\partial\boldsymbol{\beta}}[\mathbf{A}_b\boldsymbol{\alpha}_b] = \frac{\partial\mathbf{v}_b}{\partial\boldsymbol{\beta}}.$$
(71)

To avoid confusion, the subscripts are dropped from $\boldsymbol{\alpha}_b$, $\mathbf{A}_b$ and $\mathbf{v}_b$. Hence, we can express $\frac{\partial}{\partial\boldsymbol{\beta}}[\mathbf{A}\boldsymbol{\alpha}]$ as

$$\frac{\partial}{\partial\boldsymbol{\beta}}[\mathbf{A}\boldsymbol{\alpha}] = \mathbf{A}\frac{\partial\boldsymbol{\alpha}}{\partial\boldsymbol{\beta}} + \begin{bmatrix} | & \vdots & \vdots & | \\ \frac{\partial[\mathbf{A}]}{\partial\beta_1}\boldsymbol{\alpha} & \vdots & \vdots & \frac{\partial[\mathbf{A}]}{\partial\beta_{k_s}}\boldsymbol{\alpha} \\ | & \vdots & \vdots & | \end{bmatrix} \Rightarrow$$

$$\frac{\partial\boldsymbol{\alpha}}{\partial\boldsymbol{\beta}} = \mathbf{A}^{-1}\left[\frac{\partial\mathbf{v}}{\partial\boldsymbol{\beta}} - \begin{bmatrix} | & \vdots & \vdots & | \\ \frac{\partial[\mathbf{A}]}{\partial\beta_1}\boldsymbol{\alpha} & \vdots & \vdots & \frac{\partial[\mathbf{A}]}{\partial\beta_{k_s}}\boldsymbol{\alpha} \\ | & \vdots & \vdots & | \end{bmatrix}\right].$$
(72)

Next, we find the expression for $\frac{\partial\mathbf{v}}{\partial\boldsymbol{\beta}}$

$$\frac{\partial\mathbf{v}}{\partial\boldsymbol{\beta}} = \frac{1}{b}\sum_{l=1}^{b}\frac{\check{r}_l\tilde{\mathbf{x}}_{[l]}\partial w_\tau/\partial\boldsymbol{\beta}\rho_0'(\check{r}_l/\sigma)}{\check{r}_l^2}\check{y}_l$$

$$+\frac{1}{b}\sum_{l=1}^{b}\left[\frac{\check{r}_l\left(-w_\tau\rho_0''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}/\sigma - \rho_1''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}/\sigma\right)}{\check{r}_l^2}\right.$$
(73)

$$\left.+\frac{\tilde{\mathbf{x}}_{[l]}\left(w_\tau\rho_0'(\check{r}_l/\sigma) + \rho_1'(\check{r}_l/\sigma)\right)}{\check{r}_l^2}\right]\check{y}_l\tilde{\mathbf{x}}_{[l]}^T,$$

where $w_\tau$ is given by

$$w_\tau = \frac{\sum_{l=1}^{b}\left[2\rho_1\left(\frac{\check{r}_l}{\sigma}\right) - \rho_1'\left(\frac{\check{r}_l}{\sigma}\right)\frac{\check{r}_l}{\sigma}\right]}{\sum_{l=1}^{b}\rho_0'\left(\frac{\check{r}_l}{\sigma}\right)\frac{\check{r}_l}{\sigma}}.$$
(74)

On the other hand, $\frac{\partial w_\tau}{\partial\boldsymbol{\beta}}$ can be calculated as follows:

$$\frac{\partial w_\tau}{\partial\boldsymbol{\beta}} = \frac{\sum_{l=1}^{b}\left[\rho_1''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}^T\check{r}_l/\sigma^2 - \rho_1'(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}^T/\sigma\right]}{\sum_{l=1}^{b}\rho_0'(\check{r}_l/\sigma)\check{r}_l/\sigma}$$

$$+\frac{\sum_{l=1}^{b}\left[\rho_0''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}^T\check{r}_l/\sigma^2 + \rho_0'(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}^T/\sigma\right]}{\sum_{l=1}^{b}\rho_0'(\check{r}_l/\sigma)\check{r}_l/\sigma} \times w_\tau.$$
(75)

Now, we need to compute $\partial[\mathbf{A}]/\partial\beta_j$,

$$\frac{\partial[\mathbf{A}]}{\partial\beta_j}\boldsymbol{\alpha} = \frac{1}{b}\sum_{l=1}^{b}\left[\frac{\partial w_\tau/\partial\beta_j\rho_0'(\check{r}_l/\sigma) - w_\tau\rho_0''(\check{r}_l/\sigma)\tilde{x}_{lj}/\sigma}{\check{r}_l}\right.$$

$$\left.-\frac{\rho_1'(\check{r}_l/\sigma)\tilde{x}_{lj}/\sigma}{\check{r}_l} + \frac{\left(w_\tau\rho_0'(\check{r}_l/\sigma) + \rho_1'(\check{r}_l/\sigma)\right)\tilde{x}_{lj}}{\check{r}_l^2}\right]\tilde{\mathbf{x}}_{[l]}\tilde{\mathbf{x}}_{[l]}^T\boldsymbol{\alpha}$$
(76)

Therefore, we have

$$\begin{bmatrix} | & \vdots & \vdots & | \\ \frac{\partial[\mathbf{A}]}{\partial\beta_1}\boldsymbol{\alpha} & \vdots & \vdots & \frac{\partial[\mathbf{A}]}{\partial\beta_{k_s}}\boldsymbol{\alpha} \\ | & \vdots & \vdots & | \end{bmatrix} = \frac{1}{b}\sum_{l=1}^{b}\frac{\check{r}_l\tilde{\mathbf{x}}_{[l]}\partial w_\tau/\partial\boldsymbol{\beta}\rho_0'(\check{r}_l/\sigma)}{\check{r}_l^2}\tilde{\mathbf{x}}_{[l]}^T\boldsymbol{\alpha}$$

$$+\frac{1}{b}\sum_{l=1}^{b}\left[\frac{\check{r}_l\left(-w_\tau\rho_0''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}/\sigma - \rho_1''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}/\sigma\right)}{\check{r}_l^2}\right.$$

$$\left.+\frac{\tilde{\mathbf{x}}_{[l]}\left(w_\tau\rho_0'(\check{r}_l/\sigma) + \rho_1'(\check{r}_l/\sigma)\right)}{\check{r}_l^2}\right]\tilde{\mathbf{x}}_{[l]}^T\tilde{\mathbf{x}}_{[l]}^T\boldsymbol{\alpha}$$
(77)

Now, we can calculate $\frac{\partial\boldsymbol{\alpha}}{\partial\boldsymbol{\beta}}$ as follows:

$$\frac{\partial\boldsymbol{\alpha}}{\partial\boldsymbol{\beta}} = \mathbf{A}^{-1}\left[\frac{1}{b}\sum_{l=1}^{b}\frac{\tilde{\mathbf{x}}_{[l]}\partial w_\tau/\partial\boldsymbol{\beta}\rho_0'(\check{r}_l/\sigma)}{\check{r}_l}(\check{y}_l - \tilde{\mathbf{x}}_{[l]}^T\boldsymbol{\alpha})\right.$$

$$+\frac{1}{b}\sum_{l=1}^{b}\left[\frac{-w_\tau\rho_0''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}/\sigma - \rho_1''(\check{r}_l/\sigma)\tilde{\mathbf{x}}_{[l]}/\sigma}{\check{r}_l}\right.$$
(78)

$$\left.\left.+\frac{\tilde{\mathbf{x}}_{[l]}\left(w_\tau\rho_0'(\check{r}_l/\sigma) + \rho_1'(\check{r}_l/\sigma)\right)}{\check{r}_l^2}\right]\tilde{\mathbf{x}}_{[l]}^T(\check{y}_l - \tilde{\mathbf{x}}_{[l]}^T\boldsymbol{\alpha})\right]$$

Plugging in $\hat{\boldsymbol{\theta}}_b$ into $\frac{\partial\boldsymbol{\alpha}}{\partial\boldsymbol{\beta}}$, we can proceed with $\mathcal{A}$,

$$\mathcal{A} = \mathbf{I} - \frac{\partial[\mathbf{A}^{-1}\mathbf{v}]}{\partial\boldsymbol{\beta}}\Big|_{\hat{\boldsymbol{\theta}}_b} = (\hat{\mathbf{A}}_b)^{-1}\left[-\frac{1}{b}\sum_{l=1}^{b}\tilde{\mathbf{x}}_{[l]}\nabla_{\boldsymbol{\beta}}w_\tau\rho_0'(\check{r}_l)\right.$$

$$\left.+\frac{1}{b}\sum_{l=1}^{b}\left[\hat{w}_\tau\rho_0''(\check{r}_l) + \rho_1''(\check{r}_l)\right]\tilde{\mathbf{x}}_{[l]}\tilde{\mathbf{x}}_{[l]}^T/\hat{\sigma}_b\right]$$
(79)

where $\boldsymbol{\alpha}\big|_{\hat{\boldsymbol{\theta}}_b} = \hat{\boldsymbol{\beta}}_b$ is given by the fixed-point assumption. Now, we turn our attention to deriving the last missing expression, $\frac{\partial[\mathbf{A}^{-1}\mathbf{v}]}{\partial\sigma}$,

$$\frac{\partial\mathbf{v}}{\partial\sigma} = \mathbf{A}\frac{\partial\boldsymbol{\alpha}}{\partial\sigma} + \frac{\partial\mathbf{A}}{\partial\sigma}\boldsymbol{\alpha} \Rightarrow \mathbf{A}\frac{\partial\boldsymbol{\alpha}}{\partial\sigma} = \frac{\partial\mathbf{v}}{\partial\sigma} - \frac{\partial\mathbf{A}}{\partial\sigma}\boldsymbol{\alpha} \Rightarrow$$

$$\frac{\partial\boldsymbol{\alpha}}{\partial\sigma} = \mathbf{A}^{-1}\left(\frac{\partial\mathbf{v}}{\partial\sigma} - \frac{\partial\mathbf{A}}{\partial\sigma}\boldsymbol{\alpha}\right)$$
(80)

Subsequently, $\frac{\partial\boldsymbol{\alpha}}{\partial\sigma}$ can be derived by following analogous procedures to $\frac{\partial\boldsymbol{\alpha}}{\partial\boldsymbol{\beta}}$ as follows:

$$\frac{\partial\boldsymbol{\alpha}}{\partial\sigma} = \mathbf{A}^{-1}\left[\frac{1}{b}\sum_{l=1}^{b}\left[\frac{\partial w_\tau/\partial\sigma\rho_0'(\check{r}_l/\sigma) - w_\tau\rho_0''(\check{r}_l/\sigma)\check{r}_l/\sigma^2}{\check{r}_l}\right.\right.$$

$$\left.\left.-\frac{\rho_1''(\check{r}_l/\sigma)\check{r}_l/\sigma^2}{\check{r}_l}\right]\tilde{\mathbf{x}}_{[l]}(\check{y}_l - \tilde{\mathbf{x}}_{[l]}^T\boldsymbol{\alpha})\right]$$
(81)

where $\frac{\partial w_\tau}{\partial \sigma}$ is given as follows:

$$\frac{\partial w_\tau}{\partial \sigma} = \frac{\sum_{l=1}^b \left[\rho_1''(\check r_l/\sigma)\check r_l^2/\sigma^3 - \rho_1'(\check r_l/\sigma)\check r_l/\sigma^2\right]}{\sum_{l=1}^b \rho_0'(\check r_l/\sigma)\check r_l/\sigma}$$
$$+ \frac{\sum_{l=1}^b \left[\rho_0''(\check r_l/\sigma)\check r_l^2/\sigma^3 + \rho_0'(\check r_l/\sigma)\check r_l/\sigma^2\right]}{\sum_{l=1}^b \rho_0'(\check r_l/\sigma)\check r_l/\sigma} \times w_\tau \tag{82}$$

Having $\frac{\partial \alpha}{\partial \sigma}$, we can compute $\eta$ as follows:

$$\eta = -\frac{\partial[(\mathbf{A}^{-1}\mathbf{v})}{\partial \sigma}|_{\hat\theta_b} = (\hat{\mathbf{A}}_b)^{-1}\left[\frac{1}{b}\sum_{l=1}^b\left[-\nabla_\sigma w_\tau \rho_0'(\check r_l)\right.\right.$$
$$\left.\left.+\hat w_\tau \rho_0''(\check r_l)\check r_l/\hat\sigma_b + \rho_1'(\check r_l)\check r_l/\hat\sigma_b\right]\tilde{\mathbf{x}}_{[l]}\right] \tag{83}$$

Consequently, we can exploit the block matrix inversion lemma to compute $[\mathbf{I} - \nabla\mathbf{f}(\hat\theta_b)]^{-1}$ as follows:

$$\left[\mathbf{I} - \nabla\mathbf{f}(\hat\theta_b)\right]^{-1} = \begin{bmatrix}\mathcal{A} & \eta \\ \zeta & a\end{bmatrix}^{-1} = \begin{bmatrix}\mathbf{M}_b & \mathbf{d}_b \\ \mathbf{N}_b & \mathbf{q}_b\end{bmatrix},$$
$$\mathbf{M}_b = \left(\mathcal{A} - \eta a^{-1}\zeta\right)^{-1},$$
$$\mathbf{d}_b = -\mathcal{A}^{-1}\eta\left(a - \zeta\mathcal{A}^{-1}\eta\right)^{-1}, \tag{84}$$
$$\mathbf{N}_b = -\left(a - \zeta\mathcal{A}^{-1}\eta\right)^{-1}\zeta\mathcal{A}^{-1},$$
$$\mathbf{q}_b = \left(a - \zeta\mathcal{A}^{-1}\eta\right)^{-1}.$$

Now, we only need to prove both $\sqrt{n}[\mathbf{f}^\star(\hat\theta_b) - \hat\theta_b]$ and $\sqrt{b}[\mathbf{f}(\theta_0) - \theta_0]$ are convergent to the same limiting distribution. We show that $\mathbf{f}(\theta_0) - \theta_0$ can be expressed as smooth function of means. Therefore, we can exploit the results on central limit theorem, and its extension to bootstrapping of smooth functions of means [55] and show that $\sqrt{n}[\mathbf{f}^\star(\hat\theta_b)-\hat\theta_b]$ and $\sqrt{b}[\mathbf{f}(\theta_0) - \theta_0]$ are convergent to the same limiting distribution. Let's define $\mathbf{Q}(\theta_0)$ and its expected value $\mu(\theta_0)$ as follows:

$$\mathbf{Q}(\theta_0)$$
$$= \left(\frac{\rho_0'(r/\sigma_0)}{r}\mathbf{x}\mathbf{x}^T, \frac{\rho_1'(r/\sigma_0)}{r}\mathbf{x}\mathbf{x}^T, \frac{\rho_0'(r/\sigma_0)}{r}y\mathbf{x}, \frac{\rho_1'(r/\sigma_0)}{r}y\mathbf{x},\right.$$
$$\left.2\rho_1(r/\sigma_0) - \rho_1'(r/\sigma_0)r/\sigma_0, \rho_0'(r/\sigma_0)r/\sigma_0, \sigma_0\frac{\rho_0(r/\sigma_0)}{\delta}\right),$$
$$\mu(\theta_0) = \mathbb{E}[\mathbf{Q}(\theta_0)] = \left(\mathbf{Z}_1, \mathbf{Z}_2, \mathbf{z}_3, \mathbf{z}_4, z_5, z_6, \sigma_0\right) \tag{85}$$

where $\mathbf{Q}(\theta_0)$ and $\mu(\theta_0) \in \mathbb{R}^{k_s \times k_s} \times \mathbb{R}^{k_s \times k_s} \times \mathbb{R}^{k_s} \times \mathbb{R}^{k_s} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}$ and

$$r = y - \mathbf{x}^T\beta_0$$
$$\beta_0 = \left(\frac{z_5}{z_6} \times \mathbf{Z}_1 + \mathbf{Z}_2\right)^{-1} \times \left(\frac{z_5}{z_6} \times \mathbf{z}_3 + \mathbf{z}_4\right) \tag{86}$$

Given $b$ observations of $\left(\check y_l, \tilde{\mathbf{x}}_{[l]}\right)$, the sample mean $\bar{\mathbf{Q}}_b(\theta_0)$ is given by

$$\bar{\mathbf{Q}}_b(\theta_0)$$
$$= \left(\frac{1}{b}\sum_{l=1}^b \frac{\rho_0'(\check r_l(\beta_0)/\sigma_0)}{\check r_l(\beta_0)}\tilde{\mathbf{x}}_{[l]}\tilde{\mathbf{x}}_{[l]}^T, \frac{1}{b}\sum_{l=1}^b \frac{\rho_1'(\check r_l(\beta_0)/\sigma_0)}{\check r_l(\beta_0)}\tilde{\mathbf{x}}_{[l]}\tilde{\mathbf{x}}_{[l]}^T\right.$$
$$, \frac{1}{b}\sum_{l=1}^b \frac{\rho_0'(\check r_l(\beta_0)/\sigma_0)}{\check r_l(\beta_0)}\check y_l\tilde{\mathbf{x}}_{[l]}, \frac{1}{b}\sum_{l=1}^b \frac{\rho_1'(\check r_l(\beta_0)/\sigma_0)}{\check r_l(\beta_0)}\check y_l\tilde{\mathbf{x}}_{[l]},$$
$$\frac{1}{b}\sum_{l=1}^b \left[2\rho_1(\check r_l(\beta_0)/\sigma_0) - \rho_1'(\check r_l(\beta_0)/\sigma_0)\check r_l(\beta_0)/\sigma_0\right],$$
$$\left.\frac{1}{b}\sum_{l=1}^b \rho_0'(\check r_l(\beta_0)/\sigma_0)\check r_l(\beta_0)/\sigma_0, \frac{\sigma_0}{b\delta}\sum_{l=1}^b \rho_0(\check r_l(\beta_0)/\sigma_0)\right). \tag{87}$$

where $\check r_l(\beta_0) = \check y_l - \check{\mathbf{x}}_{[l]}^T\beta_0$. Then, consider the function $\mathbf{g} : \mathbb{R}^{k_s \times k_s} \times \mathbb{R}^{k_s \times k_s} \times \mathbb{R}^{k_s} \times \mathbb{R}^{k_s} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R} \to \mathbb{R}^{k_s} \times \mathbb{R}$ defined as follows:

$$\mathbf{g}(\bar{\mathbf{Z}}_1, \bar{\mathbf{Z}}_2, \bar{\mathbf{z}}_3, \bar{\mathbf{z}}_4, \bar z_5, \bar z_6, \bar z_7)$$
$$= \begin{bmatrix}\left(\frac{\bar z_5}{\bar z_6} \times \bar{\mathbf{Z}}_1 + \bar{\mathbf{Z}}_2\right)^{-1} \times \left(\frac{\bar z_5}{\bar z_6} \times \bar{\mathbf{z}}_3 + \bar{\mathbf{z}}_4\right), \\ \bar z_7\end{bmatrix}. \tag{88}$$

which is a composition of differentiable functions, yielding a smooth function. Now, we can express $\mathbf{f}(\theta_0) = \mathbf{g}(\bar{\mathbf{Q}}_b(\theta_0))$, $\theta_0 = \mathbf{g}(\mu(\theta_0))$, $\mathbf{f}(\hat\theta_b) = \mathbf{g}(\bar{\mathbf{Q}}_b(\hat\theta_b))$, $\mathbf{f}^\star(\hat\theta_b) = \mathbf{g}^\star(\bar{\mathbf{Q}}_{n,b}(\hat\theta_b))$ and $\mathbf{f}(\hat\theta_b) = \hat\theta_b$ (based on the fixed-point property) as smooth function of means and consequently

$$\sqrt{b}[\mathbf{f}(\theta_0) - \theta_0] = \sqrt{b}[\mathbf{g}(\bar{\mathbf{Q}}_b(\theta_0)) - \mathbf{g}(\mu(\theta_0))]. \tag{89}$$

Based on Theorem 2.2 in [55], we have

$$\sqrt{b}[\bar{\mathbf{Q}}_b(\theta_0) - \mu(\theta_0)] \sim \sqrt{b}[\bar{\mathbf{Q}}_b^\star(\theta_0) - \bar{\mathbf{Q}}_b(\theta_0)],$$
$$\sqrt{b}[\bar{\mathbf{Q}}_b^\star(\theta_0) - \bar{\mathbf{Q}}_b(\theta_0)] \sim \sqrt{n}[\bar{\mathbf{Q}}_{n,b}^\star(\theta_0) - \bar{\mathbf{Q}}_b(\theta_0)] \tag{90}$$

and since the estimator $\hat\theta_b$ is consistent, we can show that

$$\sqrt{n}[\bar{\mathbf{Q}}_{n,b}^\star(\theta_0) - \bar{\mathbf{Q}}_b(\theta_0)] \sim \sqrt{n}[\bar{\mathbf{Q}}_{n,b}^\star(\hat\theta_b) - \bar{\mathbf{Q}}_b(\hat\theta_b)]. \tag{91}$$

Given $\mathbf{g}$ is a smooth function, the following holds using Lemma 8.10 given in [55],

$$\sqrt{b}[\mathbf{g}(\bar{\mathbf{Q}}_b(\theta_0)) - \mathbf{g}(\mu(\theta_0))]$$
$$\sim \nabla\mathbf{g}(\mu(\theta_0))\sqrt{b}[\bar{\mathbf{Q}}_b(\theta_0) - \mu(\theta_0)],$$
$$\sqrt{n}[\mathbf{g}(\bar{\mathbf{Q}}_{n,b}^\star(\hat\theta_b)) - \mathbf{g}(\bar{\mathbf{Q}}_b(\hat\theta_b))]$$
$$\sim \nabla\mathbf{g}(\mu(\hat\theta_b))\sqrt{n}[\bar{\mathbf{Q}}_{n,b}^\star(\hat\theta_b) - \bar{\mathbf{Q}}_b(\hat\theta_b)]. \tag{92}$$

Therefore, we have

$$\sqrt{b}[\mathbf{g}(\bar{\mathbf{Q}}_b(\theta_0)) - \mathbf{g}(\mu(\theta_0))]$$
$$\sim \sqrt{n}[\mathbf{g}(\bar{\mathbf{Q}}_{n,b}^\star(\hat\theta_b)) - \mathbf{g}(\bar{\mathbf{Q}}_b(\hat\theta_b))] \tag{93}$$

which basically proves

$$\sqrt{b}[\mathbf{f}(\theta_0) - \theta_0] \sim \sqrt{n}[\mathbf{f}^\star(\hat\theta_b) - \hat\theta_b] \tag{94}$$

As discussed earlier, $[\mathbf{I} - \nabla\mathbf{f}(\hat\theta_b)]^{-1}$ is a consistent estimate of $[\mathbf{I} - \nabla\mathbf{f}(\theta_0)]^{-1}$. Therefore, we have

$$\sqrt{b}(\hat\theta_b - \theta_0) \sim \sqrt{n}(\hat\theta_{n,b}^{R\star} - \hat\theta_b)$$
$$\sqrt{b}(\hat\theta_b - \theta_0) \sim \sqrt{n}(\hat\theta_n - \theta_0) \tag{95}$$
$$\sqrt{n}(\hat\theta_n - \theta_0) \sim \sqrt{n}(\hat\theta_{n,b}^{R\star} - \hat\theta_b)$$

where they will converge to the same limiting distribution as $n$ and $b$ tend to infinity.

## PROOF OF THEOREM 3

To begin with the proof, it is assumed a certain proportion of observations in $\mathbf{Y}$ are contaminated by outliers that no longer comply with the linear regression model given in equation (1). Basically, we will show that if there are at least $p$ non-outlying observations in a bootstrap sample, it is guaranteed the FRB $\hat{\boldsymbol{\beta}}_n^{R\star}$ remains bounded. Hence, we will determine under what conditions, the FRB $\hat{\boldsymbol{\beta}}_n^{R\star}$ becomes unbounded or equivalently the maximum bias goes to infinity. The FRB $\hat{\boldsymbol{\beta}}_n^{R\star}$ is given by

$$\hat{\boldsymbol{\beta}}_n^{R\star} = \mathbf{M}_n(\hat{\boldsymbol{\beta}}_n^{1\star} - \hat{\boldsymbol{\beta}}_n) + \mathbf{d}_n(\hat{\sigma}_n^{1\star} - \hat{\sigma}_n) \tag{96}$$

where $\mathbf{M}_n$ and $\mathbf{d}_n$ can be computed by using equation (33) by changing $b$ to $n$ and $(\check{\mathbf{y}}, \check{\mathbf{X}})$ to $(\mathbf{y}, \mathbf{X})$. It is easy to show that correction factors $\mathbf{M}_n$ and $\mathbf{d}_n$ depend on the original data set $\mathbf{Y}$ rather than bootstrap sample and stay bounded if the original $\tau$-estimator $\hat{\boldsymbol{\beta}}_n$ does not break down. Next, we need to discuss how $\hat{\boldsymbol{\beta}}_n^{1\star}$ and $\hat{\sigma}_n^{1\star}$ are influenced by bootstrapping. Recall from equation (31) that

$$\hat{\sigma}_n^{1\star} = \frac{\hat{\sigma}_n}{n\delta} \sum_{l=1}^{n} \rho_0 \left( \frac{y_l^\star - \mathbf{x}_{[l]}^{\star T} \hat{\boldsymbol{\beta}}_n}{\hat{\sigma}_n} \right) \tag{97}$$

which implies $\hat{\sigma}_n^{1\star}$ remains bounded for any bootstrap sample due to boundedness of $\rho_0(\cdot)$. Subsequently, we will study under what conditions one-step bootstrap $\tau$-estimates $\hat{\boldsymbol{\beta}}_n^{1\star}$ can break down. $\hat{\boldsymbol{\beta}}_n^{1\star}$ is given by

$$\hat{\boldsymbol{\beta}}_n^{1\star} = \left[ \sum_{l=1}^{n} \tilde{\mathbf{x}}_{[l]}^\star \tilde{\mathbf{x}}_{[l]}^{\star T} \right]^{-1} \left[ \sum_{l=1}^{n} \tilde{y}_l^\star \tilde{\mathbf{x}}_{[l]}^\star \right], \tag{98}$$

where $\tilde{\mathbf{x}}_{[l]}^\star = \sqrt{\hat{w}_l^\star} \mathbf{x}_{[l]}^\star$, $\tilde{y}_l^\star = \sqrt{\hat{w}_l^\star} y_l^\star$ and $\hat{w}_l^\star$ is given by

$$\hat{w}_l^\star = \frac{\hat{w}_\tau^\star \rho_0'\left(\frac{\hat{r}_l^\star}{\hat{\sigma}_n}\right) + \rho_1'\left(\frac{\hat{r}_l^\star}{\hat{\sigma}_n}\right)}{\hat{r}_l^\star},$$

$$\hat{w}_\tau^\star = \frac{\sum_{l=1}^{n} \left[ 2\rho_1\left(\frac{\hat{r}_l^\star}{\hat{\sigma}_n}\right) - \rho_1'\left(\frac{\hat{r}_l^\star}{\hat{\sigma}_n}\right) \frac{\hat{r}_l^\star}{\hat{\sigma}_n} \right]}{\sum_{l=1}^{n} \rho_0'\left(\frac{\hat{r}_l^\star}{\hat{\sigma}_n}\right) \frac{\hat{r}_l^\star}{\hat{\sigma}_n}}, \tag{99}$$

$$\hat{r}_l^\star = y_l^\star - \mathbf{x}_{[l]}^{\star T} \hat{\boldsymbol{\beta}}_n.$$

Thus, $\hat{\boldsymbol{\beta}}_n^{1\star}$ can be expressed as the solution of a least-square problem with the observations $(\tilde{\mathbf{y}}^\star, \tilde{\mathbf{X}}^\star)$. It can be inferred from the equation (99) that the weights $\hat{w}_l^\star$ will remain bounded as $\hat{r}_l^\star$ approaches infinity. Therefore, one needs to verify that as long as there are at least $p$ good, non-outlying observations within the bootstrap sample, the corresponding one-step bootstrap estimate $\hat{\boldsymbol{\beta}}_n^{1\star}$ will remain bounded. In other words, it would suffice to show that contamination by outliers will influence one-step bootstrap estimate $\hat{\boldsymbol{\beta}}_n^{1\star}$ by a finite amount whose value is independent of outliers.

The remaining of the proof follows exactly that of Theorem 2 in [21] with an exception $c_1$ is assumed to be greater than or equal to $c_0$ or equivalently $c_1 = \max(c_1, c_0)$ without loss of generality.

## PROOF OF THEOREM 4

According to Theorem 3, the qunatile estimates $\hat{q}_t^\star$ obtained by FRB employing $\tau$-estimator can breakdown under two scenarios as follows:

- If $\hat{\boldsymbol{\beta}}_b$ is an unreliable estimate of $\boldsymbol{\beta}_0$, which may be attributed to the higher proportion of outliers than the finite-sample breakdown point of the estimator in $\check{\mathbf{Y}}$.
- If the number of bootstrap samples containing less than $p$ good, non-outlying observations constitutes at least $t\%$ of the total number of bootstrap samples, $B$.

Unreliable implies the estimate does not remain bounded any longer. In regard to RSOB-T replicates, all the bootstrap qunatiles $\hat{q}_t^*$, $t \in (0, 1)$, will be driven above any bound if $\hat{\boldsymbol{\beta}}_b$ is already unreliable. By contrast, we can show all the bootstrap quantile estimates $\hat{q}_t^*$ will remain bounded under the given assumptions of the theorem with high probability approaching to one in large-scale datasets, $n \to \infty$, as long as $\hat{\boldsymbol{\beta}}_b$ is reliable, i.e. the proportion of outliers in $\check{\mathbf{Y}}$ is less than the finite-sample breakdown point of the estimator. This implies all bootstrap samples formed according to RSOB-T scheme will contain at least $p$ good observations.

Following a similar approach given in Theorem 1 and the relationship 34 between $\tau$-scale and M-scale, it is easy to show that finite-sample breakdown of $\tau$-estimator is the same as the finite-sample breakdown point of S-estimator. Under the given assumptions, the finite-sample breakdown point of S-estimator is given by P. Rousseuw in theorem 1 of [56]. Then, the finite-sample breakdown point of $\tau$-estimator is equal to that of S-estimator as follows:

$$\epsilon^*(\hat{\boldsymbol{\beta}}_b, \check{\mathbf{Y}}) = \frac{\lfloor b/2 \rfloor - |\hat{\mathcal{S}}| + 1}{b} \tag{100}$$

where the initial estimate of $\boldsymbol{\beta}_0$ with high breakdown is an S-estimate of $\boldsymbol{\beta}_0$. Considering $\hat{\boldsymbol{\beta}}_b$ is a bounded and reliable estimate of $\boldsymbol{\beta}_0$, there exists at least $h = b - (\lfloor b/2 \rfloor - |\hat{\mathcal{S}}| + 1)$ non-outlying observations. Using lemma 1 in [22], all observations within $\check{\mathbf{Y}}$ will be drawn at least once in the bootstrap sample with high probability converging to 1 as $n \to \infty$ (Big Data). Given the assumption of general position and the fact that $h > |\hat{\mathcal{S}}|$ implies there exists at least more than $|\hat{\mathcal{S}}|$ non-outlying observations in the bootstrap sample, we can conclude all bootstrap quantiles $\hat{q}_t^\star$ are bounded and reliable with high probability converging to 1 as $n \to \infty$ (Big Data).