

# When does the Tukey Median work?

Banghua Zhu, Jiantao Jiao, Jacob Steinhardt\*

April 2, 2020

## Abstract

We analyze the performance of the Tukey median estimator under total variation (TV) distance corruptions. Previous results show that under Huber’s additive corruption model, the breakdown point is  $1/3$  for high-dimensional halfspace-symmetric distributions. We show that under TV corruptions, the breakdown point reduces to  $1/4$  for the same set of distributions. We also show that a certain projection algorithm can attain the optimal breakdown point of  $1/2$ . Both the Tukey median estimator and the projection algorithm achieve sample complexity linear in dimension.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Preliminaries</b>	<b>3</b>
2.1	Tukey median . . . . .	4
2.2	Halfspace-symmetric distributions . . . . .	4
2.3	Population corruption models . . . . .	5
2.4	Maximum bias and breakdown point for Tukey median . . . . .	6
<b>3</b>	<b>Population analysis of Tukey median</b>	<b>6</b>
<b>4</b>	<b>Finite sample analysis of Tukey median</b>	<b>11</b>

---

\*Banghua Zhu is with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. Jiantao Jiao is with the Department of Electrical Engineering and Computer Sciences and the Department of Statistics, University of California, Berkeley. Jacob Steinhardt is with the Department of Statistics and the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. Email: {banghua, jiantao,jsteinhardt}@berkeley.edu.

<b>5</b>	<b><math>\widetilde{TV}</math> Projection Algorithm</b>	<b>13</b>
<b>6</b>	<b>Open Problem</b>	<b>15</b>

## 1 Introduction

The Tukey median is the point(s) with largest Tukey depth (Tukey, 1975); it is a generalization of the one-dimensional median to high dimensions (see (1) for a formal definition). Its behavior is well-understood under the additive, or Huber, corruption model (Huber, 1973) in which an  $\epsilon$ -fraction of the data are arbitrary outliers. It is first shown in Donoho (1982); Donoho and Gasko (1992) that the breakdown point for Tukey median is  $1/3$  for halfspace-symmetric distributions in dimension  $d \geq 2$ , and the breakdown point is  $1/(d + 1)$  without the halfspace-symmetric assumption. Further analyses in Chen et al. (2002) quantify the influence function and maximum bias for halfspace-symmetric distributions, and the finite-sample behavior for elliptical distributions is analyzed in Chen et al. (2018).

In this paper, we consider the stronger TV corruption model, which allows both adding and deleting mass from the original distribution. We quantify the maximum bias of Tukey median and provide both upper and lower bounds for the breakdown point under TV corruptions. Interestingly, the breakdown point for halfspace-symmetric distributions in high dimensions decreases from  $1/3$  under additive corruptions to  $1/4$  under TV corruptions. We show that a different algorithm, projection under the halfspace metric, has breakdown point  $1/2$  in the same setting, which is the maximum breakdown point any translation-equivariant estimator can achieve (Rousseeuw and Leroy, 2005, Equation 1.38). We summarize the breakdown point for different algorithms in Figure 1.

We extend the population results on maximum bias and breakdown point under TV corruptions to the finite-sample case, showing that we approach the infinite-data limit within a constant factor once the number of samples  $n$  is linear in  $d$ . Our analysis holds under both the oblivious and adaptive models considered in the literature (Zhu et al., 2019).

## 2 Preliminaries

We provide definitions for the Tukey median, halfspace-symmetric distributions, the additive and TV corruption models, maximum bias, and breakdown point.

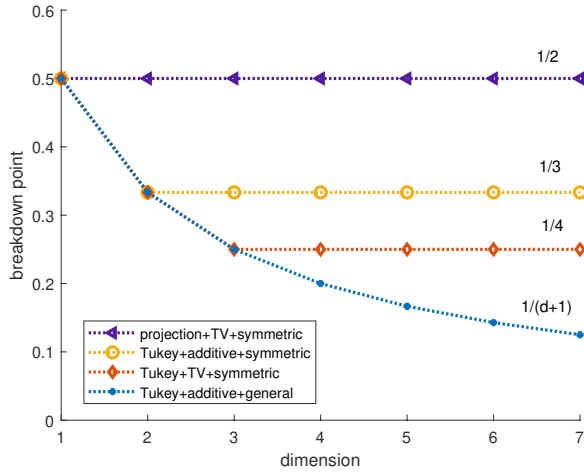


Figure 1: Summary of the breakdown point of different algorithms. Here ‘Tukey’ denotes Tukey median and ‘projection’ denotes the projection algorithm. ‘Additive’ and ‘TV’ are the two corruption models. ‘Symmetric’ denotes the family of halfspace-symmetric distribution, and ‘general’ denotes the family of all distributions.

## 2.1 Tukey median

For any distribution  $p$  and  $\mu \in \mathbb{R}^d$ , the Tukey depth is defined as the minimum probability density on one side of a hyperplane through  $\mu$ :

$$D_{\text{Tukey}}(\mu, p) = \inf_{v \in \mathbb{R}^d} p(v^\top (X - \mu) \geq 0). \quad (1)$$

The Tukey median of a distribution  $p$  is defined as the point(s) with largest Tukey depth:

$$T(p) = \arg \max_{\mu \in \mathbb{R}^d} D_{\text{Tukey}}(\mu, p). \quad (2)$$

When  $d = 1$ , Tukey median reduces to median. The Tukey median may not be unique even in one dimension. When the maximizer for Tukey depth is not unique, we use  $T(p)$  to denote the set for all the maximizers and refer to this set as the Tukey median for distribution  $p$ .

## 2.2 Halfspace-symmetric distributions

We adopt the definition of halfspace-symmetric distributions from (Chen et al., 2002; Zuo and Serfling, 2000). We say a distribution  $p$  is halfspace-

symmetric if there exists a point  $\mu \in \mathbb{R}^d$  such that for  $X \sim p$ ,  $(X - \mu)$  and  $-(X - \mu)$  are equal in distribution for all univariate projections, i.e.

$$\forall v \in \mathbb{R}^d, v^\top(X - \mu) \stackrel{d}{=} -v^\top(X - \mu). \quad (3)$$

Here  $\stackrel{d}{=}$  represents equal in distribution. We call the point  $\mu$  the center of the distribution  $p$ . The class of halfspace-symmetric distributions contains both the class of centro-symmetric distributions in [Donoho and Gasko \(1992\)](#) and elliptical distributions in [Chen et al. \(2018\)](#). For a halfspace-symmetric distribution  $p$ ,  $\mu$  is the mean of  $p$ . The Tukey depth satisfies  $D_{\text{Tukey}}(\mu, p) \geq 1/2$  and the Tukey median  $T(p)$  contains  $\mu$ .

### 2.3 Population corruption models

In the population level, we consider two corruption models: additive corruption model and TV corruption model [Diakonikolas et al. \(2017\)](#); [Donoho and Liu \(1988\)](#); [Zhu et al. \(2019\)](#).

**Additive corruption model** In a level- $\epsilon$  additive corruption model, given some true distribution  $p^*$ , the adversary can generate corrupted distribution  $p = (1 - \epsilon)p^* + \epsilon r$ , where  $\epsilon \in [0, 1)$  is the level of corruption, and  $r \in \mathbb{M}^d$  is an arbitrary distribution selected by adversary. We denote the set for all possible  $\epsilon$ -additive corruptions from  $p^*$  as

$$\mathcal{C}_{\text{add}}(p^*, \epsilon) = \{(1 - \epsilon)p^* + \epsilon r \mid r \in \mathbb{M}^d\}. \quad (4)$$

**Total variation distance corruption model** The total variation distance between two distributions  $p, q$  is defined as

$$\text{TV}(p, q) = \sup_A p(A) - q(A). \quad (5)$$

In a level- $\epsilon$  TV corruption model, given some true distribution  $p^*$ , the adversary can generate any corrupted distribution  $p$  with  $\text{TV}(p, p^*) \leq \epsilon$ . For any  $p \in \mathcal{C}_{\text{add}}(p^*, \epsilon)$ , it is always true that  $\text{TV}(p^*, p) = \sup_A \epsilon(p^*(A) - r(A)) \leq \epsilon$ . Thus the TV corruption model is a stronger corruption model than the additive corruptions, since TV corruptions allow not only additive corruption, but also deletion and replacement.

## 2.4 Maximum bias and breakdown point for Tukey median

Given a fixed distribution  $p^*$ , the maximum bias  $b(p^*, \epsilon)$  for Tukey median is defined as the maximum distance between  $T(p)$  and  $T(p^*)$ , where  $p$  is in the set of all possible level- $\epsilon$  corruptions:

$$b_{\text{add}}(p^*, \epsilon) = \sup_{p \in \mathcal{C}_{\text{add}}(p^*, \epsilon), x \in T(p), y \in T(p^*)} \|x - y\|, \quad (6)$$

$$b_{\text{TV}}(p^*, \epsilon) = \sup_{\text{TV}(p^*, p) \leq \epsilon, x \in T(p), y \in T(p^*)} \|x - y\|. \quad (7)$$

The corresponding breakdown point  $\epsilon^*(p^*)$  is defined as the minimum corruption level that can drive the maximum bias to infinity:

$$\epsilon_{\text{add}}^*(p^*) = \inf\{\epsilon \mid b(p^*, \epsilon) = \infty\}, \quad (8)$$

$$\epsilon_{\text{TV}}^*(p^*) = \inf\{\epsilon \mid b_{\text{TV}}(p^*, \epsilon) = \infty\}. \quad (9)$$

Based on the definition of the breakdown point for a single distribution, we define the breakdown point for a family of distribution  $\mathcal{G}$  as the worst breakdown point for any distribution inside  $\mathcal{G}$ , i.e.

$$\epsilon_{\text{add}}^*(\mathcal{G}) = \inf_{q \in \mathcal{G}} \epsilon_{\text{add}}^*(q), \quad \epsilon_{\text{TV}}^*(\mathcal{G}) = \inf_{q \in \mathcal{G}} \epsilon_{\text{TV}}^*(q). \quad (10)$$

## 3 Population analysis of Tukey median

In this section, we quantify the maximum bias and the breakdown point of the Tukey median in population level. The maximum bias of the Tukey median for halfspace-symmetric distributions under additive corruption model is determined in (Chen et al., 2002, Theorem 3.4), which shows that the worst-case perturbation is to add a single point with mass  $\epsilon$ . It is also shown in (Chen et al., 2018, Theorem 2.1) that under additive corruptions, the Tukey median achieves near optimal maximum bias for mean estimation if the true distribution  $p^*$  belongs to the family of elliptical distributions.

Here we demonstrate a gap in the breakdown point for halfspace-symmetric distributions between the additive and TV corruption models.

**Theorem 1.** *Denote  $\mathcal{G}_{\text{half}}$  as the set of all halfspace-symmetric distributions. Then the breakdown point for  $\mathcal{G}_{\text{half}}$  is*

$$\epsilon_{\text{add}}^*(\mathcal{G}_{\text{half}}) = \begin{cases} 1/2, & d = 1 \\ 1/3, & d \geq 2 \end{cases}, \quad \epsilon_{\text{TV}}^*(\mathcal{G}_{\text{half}}) = \begin{cases} 1/2, & d = 1 \\ 1/3, & d = 2 \\ 1/4, & d \geq 3 \end{cases}$$

*Proof of Theorem 1.* We first show the upper bound for both breakdown points. We defer the lower bound to Theorem 2. The construction of upper bound is summarized in Figure 1.

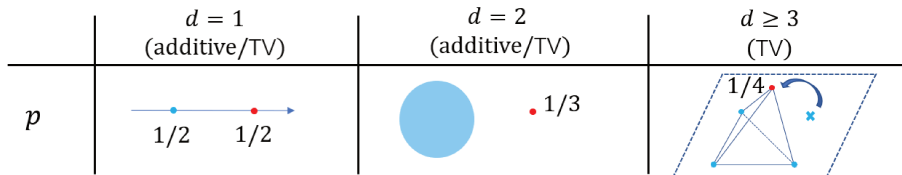


Figure 2: Illustration of worst case distributions achieving the breakdown point. Blue represents the original probability mass in  $p^*$ , blue cross represents deleted points and red represents added points by adversary. In all three cases, the red point is a Tukey median of  $p$ . Thus by driving the red point to infinity the estimator also goes to infinity.

For  $d = 1$ , by adding  $1/2$  mass onto  $z$  and letting  $z \rightarrow +\infty$ , the maximum bias can be driven to infinity. Thus  $\epsilon^*(\mathcal{G}_{\text{half}}) \leq 1/2$  under both corruption models.

For  $d \geq 2$  under additive corruption model, the upper bound of breakdown point  $\epsilon_{\text{add}}^*$  is proven in (Donoho and Gasko, 1992, Proposition 3.3). For completeness we sketch the proof here. Consider  $p^*$  as a uniform distribution supported on unit ball. The adversary adds  $1/3$  probability mass onto a point  $\mu \in \mathbb{R}^d$  outside the unit ball to get a new distribution  $p$ . Then  $D_{\text{Tukey}}(\mu, p) = 1/3$ . On the other hand, for any point  $\mu' \neq \mu$ , if  $\mu'$  is outside unit ball, there must exist a hyperplane which goes through  $\mu'$  such that the unit ball is on one side of the hyperplane. Thus  $D_{\text{Tukey}}(\mu', p) \leq 1/3$ . If  $\mu'$  is inside unit ball, consider any hyperplane that goes through 0 and  $\mu'$ . The mass of the side of hyperplane which does not contain  $\mu$  is  $1/3^1$ . Thus we also have  $D_{\text{Tukey}}(\mu', p) \leq 1/3$ . Overall  $\mu$  must be one of the Tukey median for  $p$ . By setting  $\mu \rightarrow \infty$  the proof is done.

Since TV corruption model is a stronger corruption model, the upper bound of breakdown points for  $d = 1, 2$  under TV corruptions readily follows from that under additive corruptions. Now we show the upper bound for  $\epsilon_{\text{TV}}^*$  when  $d \geq 3$ .

Consider the following example as illustrated in Figure 2: in a 3-dimensional space,  $p^*$  is a distribution with equal probability on the four nodes of a 2-dimensional square. To be precise,  $p^*(X = t) = 1/4$  for any  $t \in$

<sup>1</sup>If  $\mu$  is on the same hyperplane. One can slightly rotate the hyperplane such that  $\mu$  is not on it. This still guarantees the corresponding depth to be arbitrarily close to  $1/3$ .

$\{(-1, -1, 0), (-1, 1, 0), (1, -1, 0), (1, 1, 0)\}$ . Thus  $p^*$  is a halfspace-symmetric distribution, and  $T(p^*) = (0, 0, 0)$  gives a unique Tukey median for  $p^*$ .

Now we move one of the point  $(1, 1, 0)$  to  $(-0.5, -0.5, z)$  to get corrupted distribution  $p$ , where  $z > 0$ . Now the four points form a tetrahedron. For any point  $\mu$  that is inside the tetrahedron, the Tukey depth  $D_{\text{Tukey}}$  is always  $1/4$ . For any point that is outside the tetrahedron, the Tukey depth is always 0. Thus all the points inside the tetrahedron are a Tukey median for the corrupted distribution  $p$ . By taking  $z \rightarrow +\infty$ , the Tukey median  $T(p)$  is driven to infinity. Thus we know that  $\epsilon_{\text{TV}}^*(\mathcal{G}_{\text{half}}) \leq 1/4$  when  $d \geq 3$ .  $\square$

Without the halfspace-symmetric assumption, the breakdown point for Tukey median under additive corruption model is  $1/(d+1)$ , which is shown in (Donoho and Gasko, 1992, Proposition 2.3). This is also true under TV corruption model.

To better illustrate the behavior of Tukey median, we analyze the maximum bias for halfspace-symmetric distributions. The performance guarantee relies on the decay function of the distribution, which characterizes how much probability mass is around the center of distribution. Assume the true distribution  $p^*$  is halfspace-symmetric centered at  $\mu^*$ . We define the decay function  $h(t) : \mathbb{R}_{\geq 0} \mapsto \mathbb{R}_{\geq 0}$  as

$$h(t) = \sup_{v \in \mathbb{R}^d, \|v\|_* \leq 1} p^*(v^\top (X - \mu^*) > t), \quad (11)$$

where  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ . Note that  $h$  is a non-increasing non-negative function and  $h(0) = 1 - D_{\text{Tukey}}(\mu^*, p^*) \leq 1/2$  for halfspace-symmetric distributions.

In the next theorem, we show that the maximum bias is controlled if the distribution has enough mass around its center:

**Theorem 2.** *Assume  $p^*$  is halfspace-symmetric with center  $\mu^*$  and decay function  $h(t)$  defined in (11). Then the maximum bias satisfies:*

$$b_{\text{add}}(p^*, \epsilon) \leq \begin{cases} h^{-1} \left( \max\left(\frac{(1-\epsilon)(1-h(0))-\epsilon}{(1-\epsilon)}, \frac{1/2-\epsilon}{1-\epsilon}\right) \right), & d = 1 \\ h^{-1} \left( \max\left(\frac{(1-\epsilon)(1-h(0))-\epsilon}{(1-\epsilon)}, \frac{1/3-\epsilon}{1-\epsilon}\right) \right), & d = 2, \\ h^{-1} \left( \frac{(1-\epsilon)(1-h(0))-\epsilon}{(1-\epsilon)} \right), & d \geq 3 \end{cases} \quad (12)$$

$$b_{\text{TV}}(p^*, \epsilon) \leq \begin{cases} h^{-1}(\max(1 - h(0) - 2\epsilon, 1/2 - \epsilon)), & d = 1 \\ h^{-1}(\max(1 - h(0) - 2\epsilon, 1/3 - \epsilon)), & d = 2. \\ h^{-1}(1 - h(0) - 2\epsilon), & d \geq 3 \end{cases} \quad (13)$$



Here  $h^{-1}$  is the generalized inverse function of  $h$  defined as

$$h^{-1}(y) = \inf\{x \mid h(x) < y\}. \quad (14)$$

For Gaussian distribution with the operator norm of covariance bounded,  $\mu^*$  is the mean and  $h(t) = 1/2 - \Theta(t)$  for  $t$  small, and Theorem 2 implies that Tukey median achieves the maximum bias  $O(\epsilon)$  for robust Gaussian mean estimation, which is known to be optimal up to constant factor.

For any fixed distribution  $p^*$ , it suffices to have  $t > 0$  for  $h^{-1}(t)$  to be finite. Thus as a direct corollary of Theorem 2, it provides tight lower bound on the breakdown point of halfspace symmetric distributions in Theorem 1 via noting that  $h(0) \leq 1/2$  for halfspace-symmetric distributions.

The results in Theorem 2 can be extended beyond halfspace symmetric distributions. For any true distribution  $p^*$ , since the Tukey median of  $p^*$  may not be unique, we define the new  $h(t)$  as

$$h(t) = \sup_{v \in \mathbb{R}^d, \|v\|_* \leq 1, \mu^* \in T(p^*)} p^*(v^\top (X - \mu^*) > t). \quad (15)$$

Then following the same argument as the proof, the result in (12) and (13) still hold.

*Proof of Theorem 2.* We first show that it suffices to bound  $D_{\text{Tukey}}(T(p), p^*)$  via the following lemma:

**Lemma 1.** *Under the same condition as Theorem 2, if  $D_{\text{Tukey}}(T(p), p^*) \geq \alpha$ , we have*

$$\|T(p) - \mu^*\| \leq h^{-1}(\alpha). \quad (16)$$

*Proof of Lemma 1.* Let  $\tilde{v} = \arg \max_{\|v\|_* \leq 1} v^\top (T(p) - \mu^*)$ . Indeed, for any  $t$  such that  $h(t) < \alpha$ , if  $\|T(p) - \mu^*\| > t$ , we have

$$\begin{aligned} D_{\text{Tukey}}(T(p), p^*) &\leq p^*(\tilde{v}^\top (X - T(p)) \geq 0) \\ &= p^*(\tilde{v}^\top (X - \mu^*) \geq \|T(p) - \mu^*\|) \\ &\leq p^*(\tilde{v}^\top (X - \mu^*) > t) \\ &\leq h(t) < \alpha, \end{aligned} \quad (17)$$

resulting in a contradiction. Thus the lemma holds.  $\square$

Now it suffices to lower bound  $D_{\text{Tukey}}(T(p), p^*)$  for different dimensions and different corruption models.

Under the TV corruption model, from the definition of  $D_{\text{Tukey}}$  and TV, we have for any  $\mu \in \mathbb{R}^d$ ,

$$\begin{aligned}
& D_{\text{Tukey}}(\mu, p) - D_{\text{Tukey}}(\mu, p^*) \\
&= \inf_{v \in \mathbb{R}^d} p(v^\top(X - \mu) \geq 0) - \inf_{v \in \mathbb{R}^d} p^*(v^\top(X - \mu) \geq 0) \\
&\leq \sup_{v \in \mathbb{R}^d} p^*(v^\top(X - \mu) < 0) - p(v^\top(X - \mu) < 0) \\
&\leq \text{TV}(p, p^*) \leq \epsilon.
\end{aligned} \tag{18}$$

Under the additive corruption model, we have a tighter bound: from the definition we know that  $p(A) = (1 - \epsilon)p^*(A) + \epsilon r(A) \leq (1 - \epsilon)p^*(A) + \epsilon$  for any event  $A$ . Thus

$$\begin{aligned}
D_{\text{Tukey}}(\mu, p) &= \inf_{v \in \mathbb{R}^d} p(v^\top(X - \mu) \geq 0) \\
&\leq \inf_{v \in \mathbb{R}^d} (1 - \epsilon)p^*(v^\top(X - \mu) \geq 0) + \epsilon \\
&= (1 - \epsilon)D_{\text{Tukey}}(\mu, p^*) + \epsilon.
\end{aligned} \tag{19}$$

When  $d = 1$ , we know that  $D_{\text{Tukey}}(T(p), p) \geq 1/2$  for any distribution  $p$ . Thus  $D_{\text{Tukey}}(T(p), p^*) \geq \frac{D_{\text{Tukey}}(T(p), p) - \epsilon}{1 - \epsilon} \geq \frac{1/2 - \epsilon}{1 - \epsilon}$  under additive corruption,  $D_{\text{Tukey}}(T(p), p^*) \geq D_{\text{Tukey}}(T(p), p) - \epsilon \geq 1/2 - \epsilon$  under TV corruption.

When  $d = 2$ , we know that  $D_{\text{Tukey}}(T(p), p) \geq 1/3$  for any distribution  $p$  from (Donoho and Gasko, 1992, Proposition 2.3). Thus following the same argument as  $d = 1$ , we have  $D_{\text{Tukey}}(T(p), p^*) \geq (1/3 - \epsilon)/(1 - \epsilon)$  under additive corruption,  $D_{\text{Tukey}}(T(p), p^*) \geq 1/3 - \epsilon$  under TV corruption.

For arbitrary dimension under additive corruption model, we also have another lower bound:

$$\begin{aligned}
D_{\text{Tukey}}(T(p), p^*) &\geq (D_{\text{Tukey}}(T(p), p) - \epsilon)/(1 - \epsilon) \\
&\geq (D_{\text{Tukey}}(\mu^*, p) - \epsilon)/(1 - \epsilon) \\
&\geq ((1 - \epsilon)D_{\text{Tukey}}(\mu^*, p^*) - \epsilon)/(1 - \epsilon) \\
&= ((1 - \epsilon)(1 - h(0)) - \epsilon)/(1 - \epsilon).
\end{aligned} \tag{20}$$

Here we use that  $p(A) = (1 - \epsilon)p^*(A) + \epsilon r(A) \geq (1 - \epsilon)p^*(A)$  for any event  $A$ . For TV corruption model, we have

$$\begin{aligned}
D_{\text{Tukey}}(T(p), p^*) &\geq D_{\text{Tukey}}(T(p), p) - \epsilon \geq D_{\text{Tukey}}(\mu^*, p) - \epsilon \\
&\geq D_{\text{Tukey}}(\mu^*, p^*) - 2\epsilon = 1 - h(0) - 2\epsilon.
\end{aligned}$$

Combining the lower bounds with Lemma 1 gives the proof.  $\square$

## 4 Finite sample analysis of Tukey median

In this section, we extend the population results in the previous section to finite-sample case.

Given finite samples, there are two different corruption models: oblivious corruption and adaptive corruption (Zhu et al., 2019). In the oblivious corruption, the adversary first picks a corrupted population distribution  $p$  from  $C_{\text{TV}}(p^*, \epsilon)$ , then we take  $n$  samples from  $p$ . In the adaptive corruption, we first take  $n$  samples from  $p^*$ , then the adversary samples  $n'$  from some distribution that is stochastically dominated by a binomial distribution  $n' \sim \text{B}(n, \epsilon)$  and replace  $n'$  points in the samples by arbitrary points to get the corrupted empirical distribution  $\hat{p}_n$ . It is shown in Diakonikolas et al. (2019); Zhu et al. (2019) that adaptive corruption model is a stronger corruption model than oblivious corruptions.

Now we bound the maximum bias in the finite-sample case. We show that with  $d/\epsilon^2$  samples, the estimation error can be of the same order as the population error in Theorem 2:

**Theorem 3.** *Assume the true distribution  $p^*$  is halfspace-symmetric centered at  $\mu^*$  with decay function  $h(t)$  defined in (11). Denote  $\hat{p}_n$  as the corrupted empirical distribution under either oblivious or adaptive TV corruptions of level  $\epsilon$ . When  $d \geq 3$ , with probability at least  $1 - \delta$ , there exists universal constant  $C > 0$  such that for any  $\hat{\mu} \in T(\hat{p}_n)$  as the Tukey median of  $\hat{p}_n$ ,*

$$\|\hat{\mu} - \mu^*\| \leq h^{-1}(1 - h(0) - 2\tilde{\epsilon}) \quad (21)$$

when  $2\tilde{\epsilon} < 1 - h(0)$ . Here  $\tilde{\epsilon} = \epsilon + C \cdot \sqrt{\frac{d+1+\log(1/\delta)}{n}}$ ,  $h^{-1}$  is the generalized inverse function of  $h$  defined in (14).

*Proof.* It suffices to show the result for adaptive corruption model. From Lemma 1, we know that it also suffices to lower bound  $D_{\text{Tukey}}(\hat{\mu}, p^*)$ , where  $\hat{\mu} \in T(\hat{p}_n)$ .

We introduce the halfspace metric defined in Donoho and Liu (1988) as

$$\widetilde{\text{TV}}(p, q) = \sup_{v \in \mathbb{R}^d, t \in \mathbb{R}} |p(v^\top X \geq t) - q(v^\top X \geq t)|. \quad (22)$$

From the definition we have  $\widetilde{\text{TV}}(p, q) \leq \text{TV}(p, q)$  for all  $p, q$ . We first show that  $|D_{\text{Tukey}}(\mu, p) - D_{\text{Tukey}}(\mu, q)| \leq \widetilde{\text{TV}}(p, q)$  for any two distributions  $p, q$

and any  $\mu \in \mathbb{R}^d$ . To see this, note that the left hand side is

$$\begin{aligned} & |D_{\text{Tukey}}(\mu, p) - D_{\text{Tukey}}(\mu, q)| \\ &= \inf_{v \in \mathbb{R}^d} p(v^\top (X - \mu) \geq 0) - \inf_{v \in \mathbb{R}^d} q(v^\top (X - \mu) \geq 0) \\ &\leq \sup_{v \in \mathbb{R}^d} q(v^\top (X - \mu) < 0) - p(v^\top (X - \mu) < 0) \leq \widetilde{\text{TV}}(p, q). \end{aligned}$$

For Tukey median  $\hat{\mu} = T(\hat{p}_n) = \arg \max_{\mu \in \mathbb{R}^d} D_{\text{Tukey}}(\mu, \hat{p}_n)$ ,

$$\begin{aligned} D_{\text{Tukey}}(\hat{\mu}, p^*) &\geq D_{\text{Tukey}}(\hat{\mu}, \hat{p}_n) - \widetilde{\text{TV}}(\hat{p}_n, p^*) \\ &\geq D_{\text{Tukey}}(\mu^*, \hat{p}_n) - \widetilde{\text{TV}}(\hat{p}_n, p^*) \\ &\geq D_{\text{Tukey}}(\mu^*, p^*) - 2\widetilde{\text{TV}}(\hat{p}_n, p^*). \end{aligned}$$

Now let  $\hat{p}_n^*$  be the uncorrupted distribution, so that  $\hat{p}_n$  is obtained from  $\hat{p}_n^*$  by modifying part of samples as in adaptive corruption model. Then by triangle inequality of  $\widetilde{\text{TV}}$ ,

$$\begin{aligned} D_{\text{Tukey}}(\hat{\mu}, p^*) &\geq D_{\text{Tukey}}(\mu^*, p^*) - 2\widetilde{\text{TV}}(\hat{p}_n, \hat{p}_n^*) - 2\widetilde{\text{TV}}(\hat{p}_n^*, p^*) \\ &\geq 1 - h(0) - 2\text{TV}(\hat{p}_n, \hat{p}_n^*) - 2\widetilde{\text{TV}}(\hat{p}_n^*, p^*). \end{aligned}$$

where we repeatedly use the fact that for any  $p, q, \mu$ , we have  $|D_{\text{Tukey}}(\mu, p) - D_{\text{Tukey}}(\mu, q)| \leq \widetilde{\text{TV}}(p, q)$ . Here  $\hat{p}_n \mid \hat{p}_n^*$  follows adaptive corruption model. Now we upper bound the two terms  $\text{TV}(\hat{p}_n, \hat{p}_n^*)$  and  $\widetilde{\text{TV}}(\hat{p}_n^*, p^*)$ . From (Zhu et al., 2019, Lemma B.1), we know that with probability at least  $1 - \delta$ ,

$$\text{TV}(\hat{p}_n, \hat{p}_n^*) \leq (\sqrt{\epsilon} + \sqrt{\frac{\log(1/\delta)}{2n}})^2. \quad (23)$$

For the second term  $\widetilde{\text{TV}}(\hat{p}_n^*, p^*)$ , from the VC inequality (Devroye and Lugosi, 2012, Chap 2, Chapter 4.3) and the fact that the family of sets  $\{\{x \mid v^\top x \geq t\} \mid \|v\| = 1, t \in \mathbb{R}, v \in \mathbb{R}^d\}$  has VC dimension  $d + 1$ , there exists some universal constant  $C^{\text{vc}}$  such that with probability at least  $1 - \delta$ :

$$\widetilde{\text{TV}}(p^*, \hat{p}_n^*) \leq C^{\text{vc}} \cdot \sqrt{\frac{d + 1 + \log(1/\delta)}{n}}. \quad (24)$$

Denote  $\tilde{\epsilon} = (\sqrt{\epsilon} + \sqrt{\frac{\log(1/\delta)}{2n}})^2 + C^{\text{vc}} \cdot \sqrt{\frac{d+1+\log(1/\delta)}{n}}$ . Combining the two lemmata together, we know that with probability at least  $1 - 2\delta$ ,  $D_{\text{Tukey}}(\hat{\mu}, p^*) \geq 1 - h(0) - 2\tilde{\epsilon}$ . The proof is completed by combining the result with Lemma 1.  $\square$

As a direct corollary of the finite sample result, we can show that for Gaussian distribution the estimation error is  $O(\epsilon)$  with sample complexity  $O(d/\epsilon^2)$ . We remark that with the same proof, the population results in Theorem 2 for  $d = 1, 2$  and additive corruptions can all be extended to finite-sample results with sample complexity  $O(d/\epsilon^2)$ . Similarly the halfspace-symmetric assumption can be discarded.

## 5 $\widetilde{\text{TV}}$ Projection Algorithm

In the previous two sections, we show that Tukey median can achieve breakdown point  $1/4$  for halfspace symmetric distributions under TV corruptions and the sample complexity is linear in dimension. In this section, we show that projection under halfspace metric  $\widetilde{\text{TV}}$ , as defined in (22), is able to improve the breakdown point to  $1/2$  under the same conditions. The  $\widetilde{\text{TV}}$  projection algorithm is first proposed in Donoho and Liu (1988) for robust mean estimation, and later generalized in Zhu et al. (2019) for general robust inference problems.

Denote  $\mathcal{G}(h)$  as the set of halfspace-symmetric distributions with controlled cumulative density function around its center:

$$\mathcal{G}(h) = \{p \mid X \sim p \text{ is halfspace-symmetric around } \mu \text{ and} \\ \sup_{v \in \mathbb{R}^d, \|v\|_* \leq 1} p(v^\top(X - \mu) > t) \leq h(t)\}. \quad (25)$$

The  $\widetilde{\text{TV}}$  projection algorithm projects the corrupted empirical distribution  $p$  onto the set  $\mathcal{G}(h)$  under  $\widetilde{\text{TV}}$  distance, i.e. the output is

$$\hat{\mu}(p) = \mathbb{E}_q[X], \text{ where } q = \arg \min_{q \in \mathcal{G}(h)} \widetilde{\text{TV}}(q, p). \quad (26)$$

Note that the  $\widetilde{\text{TV}}$  projection algorithm requires the knowledge of the set  $\mathcal{G}(h)$ , while the Tukey median is agnostic to the distributional assumption on  $p^*$ . In return, the  $\widetilde{\text{TV}}$  projection algorithm achieves a breakdown point of  $1/2$  and better maximum bias than the Tukey median, as shown in the following theorem:

**Theorem 4.** *Assume the true distribution  $p^*$  is halfspace-symmetric centered at  $\mu^*$  with decay function  $h(t)$  defined in (11). Then for any  $p$  with  $\text{TV}(p^*, p) \leq \epsilon$ , the projection estimator  $\hat{\mu}(p)$  in (26) satisfies*

$$\|\hat{\mu}(p) - \mu^*\| \leq 2h^{-1}(1/2 - \epsilon) \quad (27)$$

when  $\epsilon < 1/2$ . Here  $h^{-1}$  is the generalized inverse function of  $h$  defined in (14).

*Proof.* By triangle inequality and the property of projection,

$$\begin{aligned}\widetilde{\text{TV}}(p^*, q) &\leq \widetilde{\text{TV}}(p^*, p) + \widetilde{\text{TV}}(p, q) \\ &\leq \widetilde{\text{TV}}(p^*, p) + \widetilde{\text{TV}}(p, p^*) \\ &= 2\widetilde{\text{TV}}(p^*, p) \leq 2\text{TV}(p^*, p) \leq 2\epsilon.\end{aligned}\tag{28}$$

We also know that  $p^*, q \in \mathcal{G}(h)$ . Let  $\tilde{v} = \arg \max_{\|v\|_* \leq 1} v^\top (\hat{\mu} - \mu^*)$ . We have

$$\begin{aligned}q(\tilde{v}^\top (X - \frac{\mu^* + \hat{\mu}}{2}) < 0) \\ = q(\tilde{v}^\top (X - \hat{\mu}) < -\frac{\|\hat{\mu} - \mu^*\|}{2}) \leq h(\frac{\|\mu^* - \hat{\mu}\|}{2}).\end{aligned}\tag{29}$$

We show that it implies for any  $\epsilon < 1/2$ ,  $\|\hat{\mu} - \mu^*\| \leq 2h^{-1}(1/2 - \epsilon)$ . For any  $t$  such that  $h(t) < 1/2 - \epsilon$ , if  $\|\hat{\mu} - \mu^*\| > 2t$ ,

$$\begin{aligned}p^*(\tilde{v}^\top (X - \frac{\mu^* + \hat{\mu}}{2}) < 0) &= 1 - p^*(\tilde{v}^\top (X - \mu^*) \geq \frac{\|\hat{\mu} - \mu^*\|}{2}) \\ &\geq 1 - p^*(\tilde{v}^\top (X - \mu^*) > t) \geq 1 - h(t) > 1/2 + \epsilon.\end{aligned}\tag{30}$$

On the other hand, from  $\widetilde{\text{TV}}(p^*, q) \leq 2\epsilon$ , we know that

$$\begin{aligned}p^*(\tilde{v}^\top (X - \frac{\mu^* + \hat{\mu}}{2}) < 0) &\leq q(\tilde{v}^\top (X - \frac{\mu^* + \hat{\mu}}{2}) < 0) + 2\epsilon \\ &\leq h(\frac{\|\mu^* - \hat{\mu}\|}{2}) + 2\epsilon < 1/2 + \epsilon,\end{aligned}$$

resulting in a contradiction.  $\square$

The population result can also be extended to finite-sample case by projecting  $\hat{p}_n$  instead of  $p$  under  $\widetilde{\text{TV}}$ . The proof follows the same technique in Theorem 3. The key to the success of  $\widetilde{\text{TV}}$  projection is that it allows us to check the halfspace that goes through the middle of  $\mu^*$  and  $\hat{\mu}$ , while Tukey median is only allowed to check the halfspace that goes through  $\mu^*$  and  $\hat{\mu}$ . Although projection under  $\text{TV}$  would also give the same population rate, the finite sample error can be huge since  $\text{TV}(\hat{p}_n, p) = 1$ . For both Theorem 3 and 4, the results can be extended to a more general perturbation model of corruptions under  $\widetilde{\text{TV}}$  distance.

## 6 Open Problem

Considering the TV corruption model, Tukey median is an affine-equivariant estimator with breakdown point  $1/4$  in high dimensions and good finite sample error for halfspace-symmetric distributions. The  $\widetilde{\text{TV}}$  projection algorithm is not affine-equivariant, but achieves breakdown point  $1/2$  and good finite sample error in the same set of distributions. Both algorithms may not be efficiently solvable.

It is an open problem to find an estimator that is affine-equivariant, with breakdown point  $1/2$  and good finite sample error for halfspace-symmetric distributions without considering computational efficiency.

## References

- Mengjie Chen, Chao Gao, and Zhao Ren. Robust covariance and scatter matrix estimation under hubers contamination model. *The Annals of Statistics*, 46(5):1932–1960, 2018.
- Zhiqiang Chen, David E Tyler, et al. The influence function and maximum bias of tukey’s median. *The Annals of Statistics*, 30(6):1737–1759, 2002.
- Luc Devroye and Gábor Lugosi. *Combinatorial methods in density estimation*. Springer Science & Business Media, 2012.
- Ilias Diakonikolas, Gautam Kamath, Daniel M Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Being robust (in high dimensions) can be practical. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 999–1008. JMLR. org, 2017.
- Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high-dimensions without the computational intractability. *SIAM Journal on Computing*, 48(2):742–864, 2019.
- David L Donoho. Breakdown properties of multivariate location estimators. Technical report, Technical report, Harvard University, Boston, 1982.
- David L Donoho and Miriam Gasko. Breakdown properties of location estimates based on halfspace depth and projected outlyingness. *The Annals of Statistics*, 20(4):1803–1827, 1992.

- David L Donoho and Richard C Liu. The “automatic” robustness of minimum distance functionals. *The Annals of Statistics*, 16(2):552–586, 1988.
- Peter J Huber. Robust regression: asymptotics, conjectures and monte carlo. *The Annals of Statistics*, 1(5):799–821, 1973.
- Peter J Rousseeuw and Annick M Leroy. *Robust regression and outlier detection*, volume 589. John wiley & sons, 2005.
- John W Tukey. Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians, Vancouver, 1975*, volume 2, pages 523–531, 1975.
- Banghua Zhu, Jiantao Jiao, and Jacob Steinhardt. Generalized resilience and robust statistics. *arXiv preprint arXiv:1909.08755*, 2019.
- Yijun Zuo and Robert Serfling. General notions of statistical depth function. *Annals of statistics*, pages 461–482, 2000.