

# $\sqrt{n}$ -CONSISTENT ROBUST INTEGRATION-BASED ESTIMATION\*

Sung Jae Jun<sup>†</sup>                      Joris Pinkse<sup>‡</sup>                      Yuanyuan Wan<sup>§</sup>  
Center for the Study of Auctions, Procurements, and Competition Policy  
Department of Economics  
The Pennsylvania State University

April 2010

We propose a new robust estimator of the regression coefficients in a linear regression model. The proposed estimator is the only robust estimator based on integration rather than optimization. It allows for dependence between errors and regressors, is  $\sqrt{n}$ -consistent, and asymptotically normal. It moreover has the best-achievable *breakdown point* of regression-invariant estimators, has bounded *gross error sensitivity*, is both *affine invariant* and *regression-invariant*, and the number of operations required for its computation is linear in  $n$ . An extension would result in bounded *local shift sensitivity*, also.

---

\*This paper is based on research supported by NSF grant SES-0922127. We thank the Human Capital Foundation for their support of CAPCP. We thank Don Andrews, Roger Koenker, Peter Robinson, Neil Wallace and Haiqing Xu for helpful suggestions.

<sup>†</sup>sjun@psu.edu (corresponding author), 303 Kern Graduate Building, University Park, PA 16802

<sup>‡</sup>joris@psu.edu

<sup>§</sup>yxw162@psu.edu

## 1. INTRODUCTION

We propose a new estimator for the regression coefficients in a linear regression model, which is robust to ‘contamination.’ Our estimator is inspired by the *least median of squares* (LMS) estimator of Rousseeuw (1984) and the *Laplace estimator* of Chernozhukov and Hong (2003); see also Jun, Pinkse, and Wan (2009). Like Laplace estimators, our estimator is defined as the ratio of two integrals involving an exponential transform of (in our case) the LMS objective function, but this is where the similarity ends.

Suppose that the parameter vector of interest  $\theta_0$  is the unique minimizer of a population objective function  $\Omega$  over a compact parameter space  $\Theta$ . Laplace estimators then employ the fact that  $\theta_0$  satisfies

$$\theta_0 = \lim_{n \rightarrow \infty} \frac{\int \theta \varpi(\theta) \exp\{-\alpha_n \Omega(\theta)\} d\theta}{\int \varpi(\theta) \exp\{-\alpha_n \Omega(\theta)\} d\theta}, \quad (1)$$

where  $\varpi$  is a pseudo-prior defined on  $\Theta$  and  $\{\alpha_n\}$  is a scalar-valued deterministic sequence diverging to infinity with the sample size  $n$ . Note here that the density  $\varpi(\theta) \exp\{-\alpha_n \Omega(\theta)\} / \int \varpi(\theta) \exp\{-\alpha_n \Omega(\theta)\} d\theta$  becomes more concentrated around  $\theta_0$  as  $\alpha_n$  increases. Replacing  $\Omega$  in (1) with its sample analog  $\hat{\Omega}^1$  results in a Laplace estimator. If a quadratic expansion of  $\hat{\Omega}$  is available then the Laplace estimator is generally  $\sqrt{n}$ -consistent (Chernozhukov and Hong, 2003) and the divergence rate of  $\alpha_n$  immaterial. Absent such a quadratic expansion, as in the case of the LMS estimator, the resulting estimator is not  $\sqrt{n}$ -consistent, and the divergence rate of  $\alpha_n$  partly determines the convergence rate of the Laplace estimator (Jun, Pinkse, and Wan, 2009).

We, instead, use the fact that in our case  $\Omega$  is symmetric around  $\theta_0$ , which implies that

$$\theta_0 = \frac{\int \theta \exp\{-\Omega(\theta)\} d\theta}{\int \exp\{-\Omega(\theta)\} d\theta}, \quad (2)$$

where the integrals are taken over the entire Euclidean space. There are four fundamental differences between (1) and (2): in (2) there is no limit, there is no  $\alpha_n$ , there is no compact parameter space requirement, and there is no  $\varpi$ . Because there is no limit in (2),  $\alpha_n$  is not needed anymore. Since the symmetry of  $\Omega$  around  $\theta_0$  is used, the parameter space should not be artificially restricted and no prior can be used. Our estimator  $\hat{\theta}$  is obtained by replacing  $\Omega$  in (2) with  $\hat{\Omega}$ .

In this paper we focus our attention on the case in which  $\hat{\Omega}$  is the LMS objective function, or a close relative thereof. We show that, subject to assumptions outlined in subsequent sections,  $\hat{\theta}$  is

<sup>1</sup>We use bold face for random variables.

$\sqrt{n}$ -consistent and asymptotically normal with many robustness properties, which will be further explained below.

Instead of basing an estimator on (2), as we do in this paper, one could alternatively consider  $\hat{\theta}_L$ , the Laplace estimator using  $\hat{\Omega}$ . However, because the LMS objective function does not allow for a quadratic expansion (Kim and Pollard, 1990),  $\hat{\theta}_L$  will not be  $\sqrt{n}$ -consistent. Indeed, this scenario is similar to the one studied in Jun, Pinkse, and Wan (2009) for the objective functions of other  $\sqrt[3]{n}$ -consistent estimators.

The pioneering work of Huber (1973) has spawned an abundance of papers proposing estimators with ever more desirable robustness properties. The main differences between the estimators are their robustness properties, their asymptotic behavior absent contamination, their equivariance properties, and their degree of computational complexity. These properties are summarized in table 1. Our estimator is attractive in all four respects, as the exposition below will make apparent.

One notion of robustness is the finite sample *breakdown point* (Donoho and Huber, 1982),<sup>2</sup> which is the fraction of the sample that must be changed to push the value of an estimator arbitrarily far. The breakdown point of the least squares estimator equals  $1/n$  and the breakdown point of the least absolute deviations estimator (Koenker and Bassett, 1978) depends on the regressor distribution and can be arbitrarily close to zero in large samples (Hampel, Ronchetti, Rousseeuw, and Stahel, 1986, p.328). Most estimators, however, have a finite sample breakdown point close to 0.5 if the regressors are in *general position* (Rousseeuw, 1984). Notable exceptions are Huber (1973); Krasker (1980); Mallows (1983). Our estimator has the best achievable breakdown point of regression invariant estimators, determined in Rousseeuw (1984).

Because the requirement that regressors be in general position is strong, we provide results that are more general than that. Specifically, it can be preferable (from a breakdown point perspective) to use a quantile  $q$  other than the median. Details can be found in section 3.

Other commonly used notions of robustness are the *gross error sensitivity* (GES) and the *local shift sensitivity* (LSS), both due to Hampel (1968, 1974). The GES of an estimator is finite if its *influence function* (Hampel, 1968, 1974) is bounded. Many, but not all, robust estimators have a bounded influence function, including ours.

<sup>2</sup>An asymptotic version can be found in Hampel (1971) and a different breakdown point concept in Sakata and White (1995, 1998).

Estimation Method	acronym	BDP =0.5	GES finite	LSS finite	$\sqrt{n}$ rate	Normal	Comp. # oper.	Scale	Equivariance Affine	Regr.
Huber (1973)	HUB				✓	✓	?		✓	✓
Koenker and Bassett (1978)	LAD				✓	✓	$n^a$	✓	✓	✓
Krasker (1980)	HK		✓		✓	✓	?		✓	✓
Siegel (1982) <sup>b</sup>	RM	✓	✓		✓	✓	$n^d$	✓		✓
Mallows (1983)	MAL		✓	✓	✓	✓	?		✓	✓
Rousseeuw (1984)	LMS	✓	✓ <sup>c</sup>				$n^{d,d}$	✓	✓	✓
Rousseeuw (1984)	LTS	✓	✓		✓	✓	$n \log n$	✓	✓	✓
Rousseeuw and Yohai (1984)	SEST	✓			✓	✓	$n^2 \log n$	✓	✓	✓
Yohai (1987)	MM	✓			✓	✓	?		✓	✓
Yohai and Zamar (1988)	TAU	✓			✓	✓	?		✓	✓
Croux, Rousseeuw, and Hossjer (1994)	GS	✓			✓	✓	$n^2 \log n^e$		✓	✓
Hossjer (1994)	LTA	✓			✓	✓	$n \log n$	✓	✓	✓
Hadi and Luceno (1997)	MTLE	✓			✓	✓	?	✓	✓	✓
Chang, McKean, Naranjo, and Sheather (1999)	HBRR	✓	✓ <sup>f</sup>	?	✓	✓	?	✓	✓	✓
Zinde-Walsh (2002)	SLMS	✓			✓	✓	?	✓	✓	✓
Cizek (2008)	GTE	✓	✓		✓	✓	?	✓	✓	✓
New		✓	✓	<sup>g</sup>	✓	✓	$n$	✓	✓	✓

TABLE 1. Comparison of robust estimators of the coefficients in a linear regression model.

<sup>a</sup>With preprocessing; see Portnoy and Koenker (1997).

<sup>b</sup>Asymptotics are due to Hossjer (1994).

<sup>c</sup>If the constant is not varied, infinite if varied; see Davies (1993).

<sup>d</sup>See Croux, Rousseeuw, and Hossjer (1994).

<sup>e</sup>See Croux, Rousseeuw, and Hossjer (1994).

<sup>f</sup>If the constant is not varied.

<sup>g</sup>Can be modified to have a finite LSS.

The LSS is finite if the partial derivative of the influence function with respect to regressor and regressand-values is bounded.<sup>3</sup> We know of only one estimator, namely Mallows (1983), which is known to have a finite LSS. The proposed estimator does not have a finite LSS if the tails of the error distribution are thin. We do, however, describe a modification of our estimator which can achieve a finite LSS.

Virtually all existing robust estimators, and ours, are  $\sqrt{n}$ -consistent and asymptotically normal. The two exceptions are the two other LMS-based estimators, Rousseeuw (1984); Zinde-Walsh (2002). The original LMS estimator has been shown to be  $\sqrt[3]{n}$ -consistent and has a complicated limit distribution (see Kim and Pollard, 1990). Zinde-Walsh (2002) smoothes out the LMS objective function to obtain a better convergence rate and a limiting normal distribution, but her estimator does not achieve the desired  $\sqrt{n}$ -rate and its GES is infinite.

Like Zinde-Walsh (2002), but unlike most of the other estimators mentioned here, we do allow for dependence between errors and regressors. There are many examples in economics in which e.g. heteroskedasticity is important. Unlike Zinde-Walsh (2002), however, we do not allow for time series dependence, but follow the rest of the literature and assume independent and identically distributed (i.i.d.) data.

Almost all existing estimators, and ours, are both *affine invariant* and *regression invariant*. About half are also *scale invariant*, meaning that if the regressand is scaled, the vector of regression coefficient estimates is scaled by the same amount. Our estimator is not scale invariant, and scaling does have a material impact on its performance. Issues pertaining to scaling are discussed in detail in section 6.

Finally, there is great variation in the computational complexity of estimators, both in terms of computation time and the difficulty of writing a program. Ours is the only estimator for which the number of operations required for its computation is linear in  $n$ , albeit that the constant multiplying  $n$  can be large and increases with the number of regressors  $d$ . Because our estimator is the ratio of two integrals, it can be computed using any of a number of numerical integration techniques. For low-dimensional (small  $d$ ) problems Gaussian quadrature works well. For many regressors, (quasi) Monte Carlo techniques can be used. For the numbers produced in this paper, we use Gibbs sampling (Geman and Geman, 1984). A simple Gibbs sampling procedure is described in appendix F; a C program using a faster algorithm is available from the authors upon request.

---

<sup>3</sup>The definition of the LSS is more general in that it allows for left- and right-derivatives to be different.

The remainder of our paper is organized as follows. In section 2 we define our estimator. Its breakdown point properties are established in section 3. Section 4 contains the asymptotic results absent contamination and section 5 a discussion of its asymptotic robustness properties (GES and LSS). Finally, section 6 addresses the effects of scaling of observables and section 7 the computation of the proposed estimator.

## 2. ESTIMATOR

For some  $0 < q < 1$  to be chosen, let  $N = \lfloor qn \rfloor + 1$ , where  $\lfloor \cdot \rfloor$  denotes the largest integer no greater than its argument. Let further  $\mathcal{Q}^*(\boldsymbol{\zeta}; q^*)$  denote the  $q^*$ -quantile of the distribution of  $\boldsymbol{\zeta}$ ,  $\mathcal{Q}(\boldsymbol{\zeta}) = \mathcal{Q}^*(\boldsymbol{\zeta}; q)$ , and let  $\hat{\mathcal{Q}}(\boldsymbol{\zeta}_i)$  be the  $N$ -th order statistic of  $\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_n$  for arbitrary  $\boldsymbol{\zeta}$ 's. In case a quantile is not unique, in the sense that there are multiple values  $m$  that satisfy  $\mathbb{P}(\boldsymbol{\zeta} < m) \leq q \leq 1 - \mathbb{P}(\boldsymbol{\zeta} > m)$ ,  $\mathcal{Q}$  is taken to be any such value.<sup>4</sup>

Let  $\{(x_i, \mathbf{y}_i)\}$  be an i.i.d. sample of size  $n$  where  $x_i \in \mathbb{R}^d$ . The object of interest is the vector of regression coefficients in the linear regression model

$$\mathbf{y}_i = x_i^\top \theta_0 + \mathbf{u}_i, \quad i = 1, \dots, n.$$

Under conditions to be developed in section 4,  $\theta_0$  is unique and given by

$$\theta_0 = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \mathcal{Q}(|\mathbf{y}_i - x_i^\top \theta|^2). \quad (3)$$

Our estimator of  $\theta_0$  is

$$\hat{\theta} = \frac{\int \theta \exp\{-\hat{\mathcal{Q}}(|\mathbf{y}_i - x_i^\top \theta|^2)\} d\theta}{\int \exp\{-\hat{\mathcal{Q}}(|\mathbf{y}_i - x_i^\top \theta|^2)\} d\theta}. \quad (4)$$

The estimator  $\hat{\theta}$  resembles a Laplace estimator (Chernozhukov and Hong, 2003; Jun, Pinkse, and Wan, 2009), albeit that (as mentioned in the introduction) there is no sample-size-dependent input parameter scaling the objective function and no pseudo-prior. Indeed, in Chernozhukov and Hong (2003) and Jun, Pinkse, and Wan (2009) the objective function must be multiplied by a parameter which tends to infinity with the sample size to ensure consistency; in Chernozhukov and Hong (2003) the parameter is set to  $n$ , in Jun, Pinkse, and Wan (2009) it is chosen by the practitioner. This is not needed here because  $m(t) = \mathcal{Q}\{|\mathbf{y}_i - x_i^\top(\theta_0 + t)|\}$  happens to be symmetric in  $t$ .

<sup>4</sup>For instance, the median of a binary random variable  $\boldsymbol{\zeta}$  with  $\mathbb{P}(\boldsymbol{\zeta} = 1) = 1/2$  is any value between 0 and 1, inclusive.

By substitution of  $t = \theta - \theta_0$  in (4) it follows that for  $\hat{\mathbf{m}}_n(t) = \mathcal{D}\{|\mathbf{y}_i - \mathbf{x}_i^\top(\theta_0 + t)|\}$ ,

$$\hat{\theta} - \theta_0 = \frac{\int t \exp\{-\hat{\mathbf{m}}_n^2(t)\} dt}{\int \exp\{-\hat{\mathbf{m}}_n^2(t)\} dt}. \quad (5)$$

The representation (5) will be frequently used in the remainder of the paper, especially in the proofs.

### 3. BREAKDOWN POINT

We now establish a general result concerning the breakdown properties of our estimator, which implies that the breakdown properties of our estimator are no worse than those of Rousseeuw (1984).

Let  $\mathcal{D}_n = \sup_{\|t\|=1} n^{-1} \sum_{i=1}^n I(|\mathbf{x}_i^\top t| = 0)$ ,  $\hat{\gamma} = n\mathcal{D}_n$  and  $\mathcal{D} = \sup_{\|t\|=1} \mathbb{P}(|\mathbf{x}_i^\top t| = 0)$ . The numbers  $\mathcal{D}_n, \hat{\gamma}$  represent the degree of noncollinearity in the sample and  $\mathcal{D}$  that in the population. The best breakdown point obtains when observations are in *general position* (Rousseeuw, 1984), in which case  $\hat{\gamma} = d - 1$ . However, because  $\mathcal{D} > 0$  if one (or more) of the regressors other than the constant is discrete, the general position property then occurs with probability approaching zero as  $n \rightarrow \infty$ . Our breakdown point result is hence for generic  $\hat{\gamma}$ .

**Theorem 1.** *If  $\hat{\gamma} + 1 < N < n$  then the breakdown point of  $\hat{\theta}$  satisfies  $\hat{\mathbf{b}} \geq \{\min(n - N, N - \hat{\gamma} - 1) + 1\}/n$ .*

Although theorem 1 only provides a lower bound to the breakdown point, it is straightforward to construct counter examples to the breakdown point being better than the one stated in theorem 1.

Theorem 1 has several implications. First, for  $q = 0.5$ , the breakdown point when the observations are in general position is  $(\lfloor n/2 \rfloor - d + 2)/n$  if  $d > 1$ ,<sup>5</sup> which is exactly the same as in Rousseeuw (1984, theorem 1). The best breakdown point is achieved when  $q$  is chosen to make  $N = \lfloor (n + \hat{\gamma} + 1)/2 \rfloor$ , which results in a breakdown point of  $\lfloor (n - \hat{\gamma} + 1)/2 \rfloor / n$ . If the observations are in general position then the breakdown point equals  $\lfloor (n - d)/2 \rfloor + 1$ , which is the same as that in the remark following theorem 1 of Rousseeuw (1984) and hence also as that of Siegel (1982).

<sup>5</sup>It is  $\lfloor (n + 1)/2 \rfloor / n$  if  $d = 1$ .

Asymptotically, the optimal choice of  $q$  in terms of breakdown properties is

$$q = \frac{1 + \mathcal{Y}}{2}, \quad (6)$$

resulting in a breakdown point converging to  $(1 - \mathcal{Y})/2$  as  $n \rightarrow \infty$ , which is the best achievable for any regression equivariant estimator (Rousseeuw, 1984). The rationale for the choice of  $q$  in (6) is that  $\hat{\mathcal{Y}}_n \xrightarrow{P} \mathcal{Y}$  and hence that  $\hat{\gamma} \approx n\mathcal{Y}$ , resulting in an optimal  $N$  of  $\approx n(1 + \mathcal{Y})/2$ .

#### 4. ASYMPTOTICS

We now turn to a discussion of the properties of  $\hat{\theta}$  absent contamination. Throughout we assume that  $\{(x_i, y_i)\}$  is an i.i.d. sequence of random variables and that  $0 < q < 1$ .

We start by establishing identification. Let  $m_0 = \mathcal{Q}(|y_i - x_i^\top \theta_0|) = \mathcal{Q}(|u_i|)$ .

**Assumption A.** *The conditional density  $f(\cdot|\cdot)$  of  $u_i$  given  $x_i = x$  is for any  $x$  even, continuous, positive on the entire real line, weakly decreasing at all  $u > 0$ , and strictly decreasing at  $m_0$ .*

Assumption A is strong, but for  $q = 0.5$  weaker than Kim and Pollard (1990, example 6.3) because we allow  $u_i$  and  $x_i$  to be dependent and do not assume the existence of derivatives for consistency.

Recall from section 3 that  $\mathcal{Y} = \sup_{\|t\|=1} \mathbb{P}(|x_i^\top t| = 0)$ .

**Assumption B.**  $\mathcal{Y} < 1$ .

Assumption B requires that the regressors are perfectly collinear with probability less than one. It is implied by the requirement that  $0 < \mathbb{E}(x_i x_i^\top) < \infty$  (Kim and Pollard, 1990, example 6.3), but does not assume the existence of moments for  $x_i$ . Given that for  $\mathcal{Y} > 0$  regressors are in general position with probability approaching zero (see section 3), assumption B is weak.

**Theorem 2.** *Under assumptions A and B,  $\theta_0$  defined in (3) is unique.*

We need one additional condition for consistency.

**Assumption C.**  $\mathcal{Y} < q$ .

Assumption C is the population equivalent (for  $q = 0.5$ ) of the requirement in Rousseeuw (1984) that no *vertical hyperplane* (passing through the origin) contains more than  $\lfloor n/2 \rfloor$  observations.



Assumption C can be restrictive. Indeed, with both a constant and a binary regressor it is violated when  $q \leq 0.5$ . But if  $\mathcal{Y} > q$  then with probability approaching one  $\mathcal{Y}_n > q$ , also, and the condition imposed on  $N$  in theorem 1 is violated. Consequently, none of the LMS-type estimators, Rousseeuw (1984); Zinde-Walsh (2002) and ours, will then have a breakdown point any better than the OLS estimator. So if assumption C is violated, it just means that  $q$  is chosen too small. In particular, if  $q$  is chosen according to (6) then assumption C is equivalent to assumption B.

**Theorem 3.** *Under assumptions A to C,  $\hat{\theta} \xrightarrow{P} \theta_0$ .*

We now proceed with a discussion of the asymptotic distribution of  $\hat{\theta}$ . Let

$$m(t) = \mathcal{Q}(|\mathbf{u}_i - \mathbf{x}_i^\top t|), \quad m_\infty(t) = \mathcal{Q}(|\mathbf{x}_i^\top t|). \quad (7)$$

The notation  $m_\infty$  is inspired by the fact that for any  $t \neq 0$ ,  $\lim_{\lambda \rightarrow \infty} \{m(\lambda t)/\lambda\} = m_\infty(t)$  provided that  $m_\infty(t)$  is unique.

Let  $f(\cdot)$  and  $F(\cdot)$  be the unconditional counterparts of  $f(\cdot|\cdot)$  and  $F(\cdot|\cdot)$ , and let  $\mathcal{X} = \{x : \exists \|t\| = 1, \epsilon > 0 : \|\mathbf{x}^\top t\| - m_\infty(t) < \epsilon\}$ .

**Assumption D.** (i)

$$\lim_{\eta \downarrow 0} \frac{\inf_{\|t\|=1} \mathbb{P}\{|\mathbf{x}_i^\top t| \leq m_\infty(t) + \eta\} - q}{\eta} > 0, \quad \lim_{\eta \downarrow 0} \frac{\inf_{\|t\|=1} \mathbb{P}\{|\mathbf{x}_i^\top t| \geq m_\infty(t) - \eta\} - 1 + q}{\eta} > 0,$$

where each inequality is taken to hold if the limit is infinite. Moreover, (ii) for some  $\epsilon > 0$ ,  $0 < \inf_x f(\epsilon|x) \leq \sup_x f(0|x) < \infty$  and for some  $2 < r < \infty$ , (iii)  $\lim_{s \rightarrow \infty} \{\inf_{x \in \mathcal{X}} f(s|x)/f^r(s)\} \geq 1$ , and (iv)  $\lim_{s \rightarrow \infty} [f\{s + F(-s)\}/f(s)]/f^r(s) > 0$ .

Conditions (ii) and (iii) in assumption D are automatically satisfied if  $\mathbf{u}_i$  and  $\mathbf{x}_i$  are independent and can be seen as mild conditions restricting their dependence. We have verified condition (iv) for a number of distributions satisfying assumption A, including (symmetrized versions of) the Normal, Gumbel, Laplace, and Cauchy distributions.

Finally, condition (i) is satisfied when all regressors other than the constant are continuous. For discrete distributions, (i) is satisfied for most, but not all, choices of  $q$ . Condition (i) assumes away the possibility that the  $q$ -quantile of  $|\mathbf{x}_i^\top t|$  is ambiguous for any vector  $t$  of length one. Condition (i) is unique to our paper.

Because  $q$  can be chosen to satisfy (i), condition (i) is more a nuisance than a serious obstacle for our estimator. Nevertheless, we highlight two alternatives that can be used to replace assumption D. The first solution is to assume that the tails of the conditional error density are sufficiently thick, i.e. declining more slowly than those of the density of an exponential distribution, which is not desirable. The second solution is to use a different estimator, a possibility which is discussed in section 6. We do not provide a formal justification for either solution.<sup>6</sup>

Finally, we need a condition on the derivative of  $f$ .

**Assumption E.**  $\sup_{u,x} f'(u|x) < \infty$ .

Assumption E is strong, but the assumption of the existence of the first derivative is also used in Kim and Pollard (1990); Hossjer (1994); Zinde-Walsh (2002), among others. Please note that assumption E is only used to establish asymptotic normality.

Let  $A(t, m) = \mathbb{P}(|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq m)$ ,  $\mathcal{D}(t) = \partial_m A\{t, m(t)\}$ ,<sup>7</sup>

$$H(t, s) = \text{Cov}[I\{|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq m(t)\}, I\{|\mathbf{u}_i - \mathbf{x}_i^\top s| \leq m(s)\}], \quad (8)$$

and

$$\mathcal{V} = 4 \frac{\iint ts^\top \frac{m(t)m(s)}{\mathcal{D}(t)\mathcal{D}(s)} H(t, s) \exp[-\{m^2(t) + m^2(s)\}] dt ds}{[\int \exp\{-m^2(t)\} dt]^2}. \quad (9)$$

**Theorem 4.** *Let assumptions A to E hold. Then  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N(0, \mathcal{V})$ .*

So our estimator has good breakdown properties and is both  $\sqrt{n}$ -consistent and asymptotically normal. For the sake of completeness, we now provide a consistent estimator  $\hat{\mathcal{V}}$  of  $\mathcal{V}$ . Let

$$\begin{aligned} \hat{\mathbf{H}}^*(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) &= n^{-1} \sum_{i=1}^n I\{|\mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\theta}| \leq \hat{\mathcal{D}}(|\mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\theta}|)\} I\{|\mathbf{y}_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\theta}}| \leq \hat{\mathcal{D}}(|\mathbf{y}_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\theta}}|)\} \\ &\quad - \left( n^{-1} \sum_{i=1}^n I\{|\mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\theta}| \leq \hat{\mathcal{D}}(|\mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\theta}|)\} \right) \left( n^{-1} \sum_{i=1}^n I\{|\mathbf{y}_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\theta}}| \leq \hat{\mathcal{D}}(|\mathbf{y}_i - \mathbf{x}_i^\top \tilde{\boldsymbol{\theta}}|)\} \right), \end{aligned} \quad (10)$$

and for some scalar  $h^*(\boldsymbol{\theta})$ ,

$$\hat{\mathcal{D}}^*(\boldsymbol{\theta}) = \frac{1}{2nh^*(\boldsymbol{\theta})} \sum_{i=1}^n I\{||\mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\theta}| - \hat{\mathcal{D}}(|\mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\theta}|)| \leq h^*(\boldsymbol{\theta})\}. \quad (11)$$

<sup>6</sup>Both alternatives ensure that  $t \exp\{-m^2(t)\} / \mathcal{D}(t)$  in (9) is integrable, which is needed.

<sup>7</sup> $\partial_m$  denotes the partial derivative with respect to  $m$ .

Then our estimator of the asymptotic variance is

$$\hat{\mathcal{V}} = 4 \frac{\iint (\theta - \hat{\theta})(\tilde{\theta} - \hat{\theta})^\top \frac{\mathcal{Q}(|\mathbf{y}_i - \mathbf{x}_i^\top \theta|) \mathcal{Q}(|\mathbf{y}_i - \mathbf{x}_i^\top \tilde{\theta}|)}{\hat{\mathcal{Q}}^*(\theta) \hat{\mathcal{Q}}^*(\tilde{\theta})} \hat{\mathbf{H}}^*(\theta, \tilde{\theta}) e^{-\{\mathcal{Q}(|\mathbf{y}_i - \mathbf{x}_i^\top \theta|^2) + \mathcal{Q}(|\mathbf{y}_i - \mathbf{x}_i^\top \tilde{\theta}|^2)\}} d\theta d\tilde{\theta}}{[\int \exp\{-\mathcal{Q}(|\mathbf{y}_i - \mathbf{x}_i^\top \theta|^2)\} d\theta]^2}. \quad (12)$$

The use of a uniform *kernel* in (11) is not essential but simplifies the proofs.

We need a single additional assumption, relating to the choice of *bandwidth*  $h^*$ .

**Assumption F.** *The bandwidth function  $h^*$  satisfies  $h^*(\theta) = (1 + \|\theta\|^{p^*})h_0$  for some  $2 < p^* < \infty$  and  $h_0 \prec 1 \prec n^{(1-1/p^*)/\sigma}h_0$  for some  $\sigma > 4$ .*

Because  $h_0$  and  $p^*$  are chosen by the practitioner, assumption F is not restrictive. Assumption F permits the bandwidth to converge at the ‘optimal’  $n^{-1/5}$  rate if  $p^* > 5$ . We are now in a position to state the final theorem of this section.

**Theorem 5.** *Let assumptions A to F hold. Then  $\hat{\mathcal{V}} \xrightarrow{p} \mathcal{V}$ .*

## 5. INFLUENCE FUNCTION

From the proof of theorem 4<sup>8</sup> it is apparent that the dominant asymptotic term is

$$\frac{2}{\sqrt{n}} \sum_{i=1}^n \frac{\int t \frac{m(t)}{\mathcal{Q}(t)} [I\{|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq m(t)\} - q] \exp\{-m^2(t)\} dt}{\int \exp\{-m^2(t)\} dt},$$

resulting in the influence function (Hampel, 1974)<sup>9</sup>

$$\mathcal{I}(y, x) = 2 \frac{\int t \frac{m(t)}{\mathcal{Q}(t)} [I\{|y - x^\top(\theta_0 + t)| \leq m(t)\} - q] \exp\{-m^2(t)\} dt}{\int \exp\{-m^2(t)\} dt}. \quad (13)$$

Since  $y, x$  enter (13) only through an indicator function,  $\mathcal{I}$  is uniformly bounded and the GES<sup>10</sup> of our estimator is hence finite.

The LSS is more complicated to determine. We now show that even in the constant-only case, our estimator does not have a finite LSS when the tails of the error distribution are thin.

**Theorem 6.** *Suppose that  $\mathbf{x}_i$  consists of only a constant and that  $F$  is the distribution function of a mean zero normal random variable with variance  $\zeta^2$ . Then if  $q = 0.5$ ,  $\zeta^2 > 2 \iff \sup_y |\partial_y \mathcal{I}(y)| < \infty$ .*

<sup>8</sup>See (48).

<sup>9</sup>In Hampel (1974) the influence function is defined as a functional derivative, which generally equals an element in the sum in the first order asymptotic term (Reeds, 1976; Boos and Serfling, 1980; Fernholz, 1983). We do not establish such equivalence here.

<sup>10</sup>The GES is the supremum over  $x, y$  of the norm of  $\mathcal{I}$ .

As the proof of theorem 6 illustrates, the LSS is infinite for thin-tailed error distributions because  $\exp(-m^2)$  does then not decrease fast enough as  $m$  increases. This problem can be remedied by replacing  $\exp$  in (4) by another smooth function which equals zero whenever its argument is sufficiently large negative. We do not investigate such a modification in this paper because of reasons outlined at the end of section 6.

## 6. SCALING

Like Laplace estimators (Chernozhukov and Hong, 2003), the proposed estimator is not invariant to scaling, or indeed monotonic transformations, of the objective function. In our case scaling the objective function is equivalent to scaling the data, so consider the estimator  $\hat{\theta}_\alpha$  below in which the scaling is made explicit by means of a scalar  $0 < \alpha < \infty$ .

$$\hat{\theta}_\alpha = \frac{\int \theta \exp\{-\alpha \hat{\mathcal{L}}(|\mathbf{y}_i - \mathbf{x}_i^\top \theta|^2)\} d\theta}{\int \exp\{-\alpha \hat{\mathcal{L}}(|\mathbf{y}_i - \mathbf{x}_i^\top \theta|^2)\} d\theta}. \quad (14)$$

Having  $\alpha$  be finite and nonzero is important for our results. It is apparent that  $\lim_{\alpha \rightarrow \infty} \hat{\theta}_\alpha$  (for  $q = 0.5$ ) yields Rousseeuw's LMS estimator, which is not  $\sqrt{n}$ -consistent and lacks a bounded influence function unless the supremum is only taken over the slope regressors (Davies, 1993).

To obtain the limit of  $\hat{\theta}_\alpha$  as  $\alpha \rightarrow 0$  is somewhat more complicated. We limit ourselves to the case with scalar-valued nonnegative  $x_i$  and  $q = 0.5$ , which is nonetheless instructive.<sup>11</sup> Let  $\mu$  denote the index corresponding to the median observation if the data are ordered by  $x_i$ -value and in the case of ties by  $y_i$ -value.<sup>12</sup> So if the data are arranged such that  $x_i < x_j \Rightarrow i < j$  and  $x_i \leq x_j, y_i < y_j \Rightarrow i < j$  and  $n$  is odd then  $\mu = (n + 1)/2$ .

**Theorem 7.** *Suppose that  $n$  is odd,  $d = 1$ ,  $q = 0.5$ , and that there are no ties in the  $y_i$ -values. If all  $x_i$ 's are nonnegative and  $x_\mu > 0$ , then*

$$\lim_{\alpha \rightarrow 0} \hat{\theta}_\alpha = \frac{y_\mu}{x_\mu}. \quad (15)$$

Theorem 7 has two interesting implications. First, if  $x_i$  is a constant then  $y_\mu/x_\mu$  equals the sample median, which has excellent properties. In most other cases, however,  $y_\mu/x_\mu$  is an inconsistent estimator of  $\theta_0$ . Indeed, if  $x_i$  is continuously distributed then  $y_\mu/x_\mu$  is the ratio of the  $y$  and  $x$  values of the observation corresponding to the sample median of the  $x_i$ 's.

<sup>11</sup>Nonnegativity is innocuous since  $x_i, y_i$  can be replaced with  $-x_i, -y_i$  if  $x_i$  is negative.

<sup>12</sup>We ignore the possibility of ties in the  $y_i$ 's given that they are assumed continuous throughout the paper.

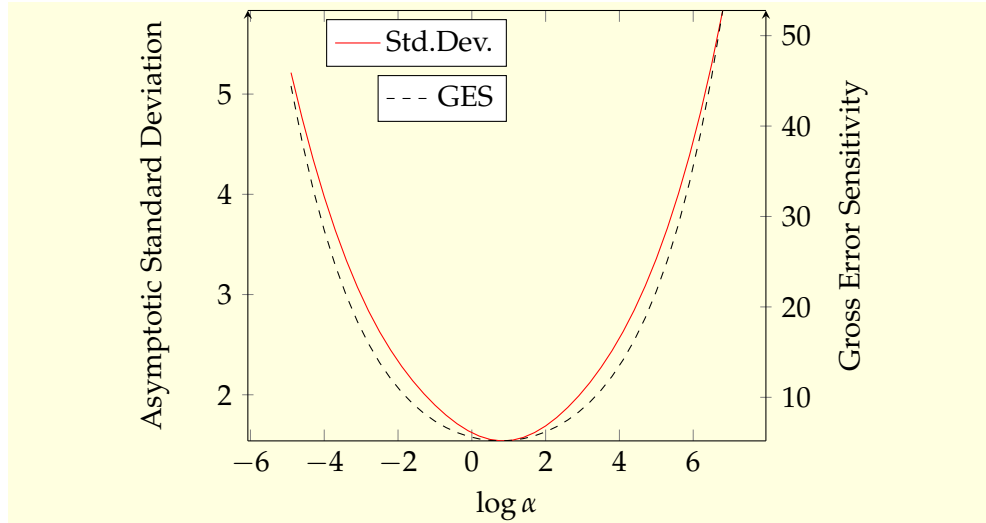


FIGURE 1. Asymptotic Variance and Gross Error Sensitivity

So the value of  $\alpha$  that minimizes the asymptotic variance is generally different from zero and infinity and the same is true for the value that minimizes the gross error sensitivity. Other estimators in this literature, including Krasker (1980); Mallows (1983), also require the choice of an input parameter but in those papers the input parameter represents a choice between efficiency (as measured by the trace of the asymptotic variance matrix) and robustness (as measured by the gross error sensitivity).<sup>13</sup>

Figure 1 demonstrates that in our case there need not be a tradeoff between efficiency and robustness. The design used to generate the graph is one in which  $u_i, x_i$  are independent standard normal random variables and there is no constant term in the model. We chose this design because it results in an explicit expression for  $m(t)$ .

The question, then, is how one should choose  $\alpha$ . One possibility is to use a first stage estimator (be it ours with a fixed  $\alpha$  or some other robust method), estimate the asymptotic variance  $\mathcal{V}_\alpha$  (or indeed the gross error sensitivity),

$$\mathcal{V}_\alpha = 4\alpha^2 \frac{\iint ts^\top \frac{m(t)m(s)}{\mathcal{D}(t)\mathcal{D}(s)} H(t,s) \exp[-\alpha\{m^2(t) + m^2(s)\}] dt ds}{[\int \exp\{-\alpha m^2(t)\} dt]^2}.$$

and choose the value of  $\alpha$  that minimizes one's estimate of (the trace of)  $\mathcal{V}_\alpha$ . We do not provide results for such a data-dependent choice of  $\alpha$ .

<sup>13</sup>See Krasker (1980).

Finally, we mentioned in section 4 that after replacing  $\hat{\mathcal{L}}(|\mathbf{y}_i - \mathbf{x}_i^\top \theta|^2)$  in (4) with a monotonic transformation thereof assumption D can be weakened. The main reason is that if one makes the exponents in (4) tend to  $-\infty$  faster as  $\|\theta\| \rightarrow \infty$  then anything multiplying the exponents can increase faster as  $\|\theta\| \rightarrow \infty$ , also. We have chosen not to take this route in this paper because  $\hat{\theta}$  defined in (4) is closer to the original LMS estimator, it is convenient, it appears to work well, and assumption D seems weaker than some other assumptions made in this paper.

## 7. COMPUTATION

There are several ways of computing our estimator; all involve numerical integration. Especially for low-dimensional ( $d$  small) problems, Gaussian quadrature works well. For high-dimensional problems, a Monte Carlo-based approach usually works better.

Because the computation of a sample median requires  $O(n)$  operations (Knuth, 1997, chapter 6), the computation of each of the integrands in (4) requires  $O(n)$  operations, also. Hence if  $\hat{\theta}$  is computed using the (classical) Monte Carlo method, with or without importance sampling (Robert and Casella, 2004, Definition 3.9), or indeed using quadrature, then the total number of operations needed is linear in  $n$ .

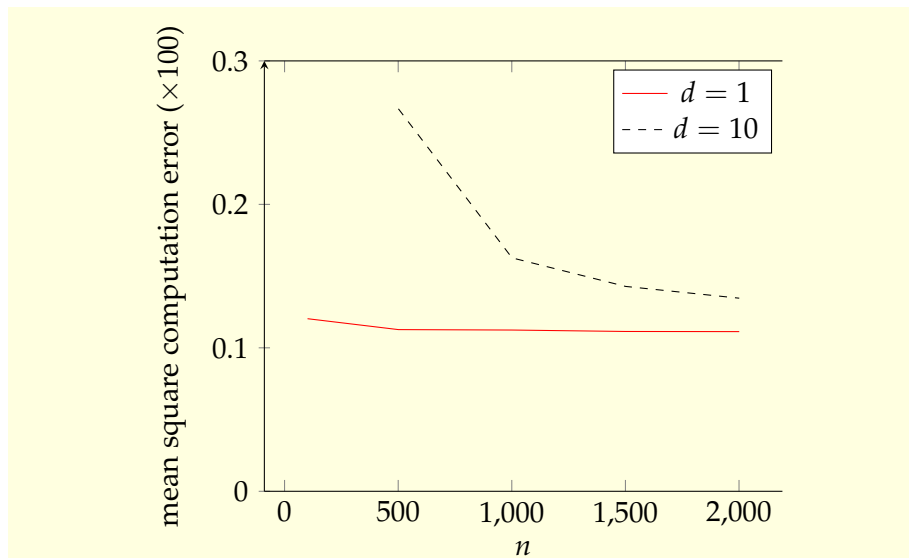


FIGURE 2. Computational Accuracy as a Function of  $n$

To illustrate, consider figure 2, for which we used the Monte Carlo method with importance sampling using a normal distribution with variance chosen to match the tails of  $\exp\{-\hat{\mathcal{L}}(|\mathbf{y}_i -$

$\mathbf{x}_i^\top \theta|^2\})$  as an instrumental distribution. For each  $(n, d)$ -combination, we constructed 1,000 samples  $s = 1, \dots, 1000$ , computed  $\hat{\theta}_{s(\infty)} = \hat{\theta}$  using 1,000,000 draws. We then computed  $\hat{\theta}_{sr}$  1000 times ( $r = 1, \dots, 1000$ ) using 1000 draws in each case. Finally, we use the average mean square deviation for each  $(n, d)$ -combination,  $\sum_{s=1}^{1000} \sum_{r=1}^{1000} (\hat{\theta}_{sr} - \hat{\theta}_{s(\infty)})^2 / 1000^2$ , as a measure of the computational accuracy of using 1000 draws.

If the number of operations needed to achieve the same level of accuracy were to increase with  $n$ , both curves in figure 2 would be increasing. The reason that they are initially decreasing is due to our choice of an instrumental distribution, which is a better match for the integrand for large  $n$  than it is for small  $n$ .

Although the results depicted in figure 2 are encouraging, some words of caution are in order. First, it is conceivable that performance is different for designs different from the one chosen here. Second, although computation is linear in  $n$ , it could be slow for any  $n$  if a large number of random draws is needed to achieve a desired level of accuracy, which arises when the instrumental distribution used is a bad match for the integrand. Indeed, the best choice of it depends on the shape of  $\exp\{-\mathcal{Q}(|\mathbf{y}_i - \mathbf{x}_i^\top \theta|^2)\}$ , in particular, on the unknown parameter vector  $\theta_0$ . Likewise, the number of draws needed to achieve the same level of accuracy need not go up linearly in  $d$ .

For these reasons, it can be preferable to use other numerical integration methods such as Gibbs sampling (Geman and Geman, 1984). A simple scheme, which requires  $O(n^2)$  operations for a draw, is described in appendix F. A faster algorithm is available from the authors.

#### REFERENCES CITED

- BOOS, D., AND R. SERFLING (1980): "A note on differentials and the CLT and LIL for statistical functions, with application to M-estimates," *The Annals of Statistics*, 8(3), 618–624.
- CHANG, W., J. MCKEAN, J. NARANJO, AND S. SHEATHER (1999): "High-Breakdown Rank Regression.," *Journal of the American Statistical Association*, 94(445), 205–219.
- CHERNOZHUKOV, V., AND H. HONG (2003): "An MCMC approach to classical estimation," *Journal of Econometrics*, 115(2), 293–346.
- CIZEK, P. (2008): "General trimmed estimation: robust approach to nonlinear and limited dependent variable models," *Econometric Theory*, 24(06), 1500–1529.
- CROUX, C., P. ROUSSEEUW, AND O. HOSSJER (1994): "Generalized S-Estimators.," *Journal of the American Statistical Association*, 89(428).

- DAVID, H. (1986): "Inequalities for ordered sums," *Annals of the Institute of Statistical Mathematics*, 38, 551–555.
- DAVIES, P. (1993): "Aspects of robust linear regression," *The Annals of Statistics*, 21(4), 1843–1899.
- DONOHU, D., AND P. HUBER (1982): "The notion of breakdown point," *A Festschrift for Erich L. Lehmann in Honor of His Sixty-fifth Birthday*, pp. 157–184.
- FERNHOLZ, L. (1983): *Von Mises calculus for statistical functionals*, vol. 19 of *Lecture Notes in Statistics*. Springer.
- GEMAN, S., AND D. GEMAN (1984): "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6, 721–741.
- GLICK, N. (1974): "Consistency conditions for probability estimators and integrals of density estimators," *Utilitas Mathematica*, 6, 61–74.
- HADI, A., AND A. LUCENO (1997): "Maximum trimmed likelihood estimators: a unified approach, examples, and algorithms," *Computational statistics & data analysis*, 25(3), 251–272.
- HAMPEL, F. (1968): "Contributions to the theory of robust estimation," Ph.D. thesis, University of California at Berkeley.
- (1971): "A general qualitative definition of robustness," *The Annals of Mathematical Statistics*, 42(6), 1887–1896.
- (1974): "The influence curve and its role in robust estimation," *Journal of the American Statistical Association*, 69(346), 383–393.
- HAMPEL, F., E. RONCHETTI, P. ROUSSEEUW, AND W. STAHEL (1986): *Robust Statistics*. Wiley New York.
- HOSSJER, O. (1994): "Rank-based estimates in the linear model with high breakdown point," *Journal of the American Statistical Association*, 89(425), 149–158.
- HUBER, P. (1973): "Robust regression: asymptotics, conjectures and Monte Carlo," *The Annals of Statistics*, 1(5), 799–821.
- JUN, S. J., J. PINKSE, AND Y. WAN (2009): "Cube root  $n$  and faster convergence, Laplace estimators, and uniform inference," Discussion paper, The Pennsylvania State University.
- KIM, J., AND D. POLLARD (1990): "Cube root asymptotics," *Annals of Statistics*, 18(1), 191–219.
- KNUTH, D. E. (1997): *The art of computer programming*, vol. 1. Addison–Wesley, 3 edn.
- KOENKER, R., AND G. BASSETT (1978): "Regression quantiles," *Econometrica*, 46(1), 33–50.



- KRASKER, W. (1980): "Estimation in linear regression models with disparate data points," *Econometrica*, 48(6), 1333–1346.
- MALLOWS, C. L. (1983): "Minimax Aspects of Bounded-Influence Regression: Comment," *Journal of the American Statistical Association*, 78(381), 77.
- PORTNOY, S., AND R. KOENKER (1997): "The Gaussian hare and the Laplacian tortoise: computability of squared-error versus absolute-error estimators, with discussion," *Statistical Science*, 12(4), 279–300.
- REEDS, J. (1976): "On the definition of von Mises functionals," Ph.D. thesis, Harvard University.
- ROBERT, C., AND G. CASELLA (2004): *Monte Carlo statistical methods*. Springer.
- ROUSSEEUW, P. (1984): "Least median of squares regression," *Journal of the American Statistical Association*, 79(388), 871–880.
- ROUSSEEUW, P., AND A. LEROY (1987): *Robust regression and outlier detection*. Wiley.
- ROUSSEEUW, P., AND V. YOHAI (1984): "Robust regression by means of S-estimators," *Robust and nonlinear time series analysis*, 26, 256–272.
- SAKATA, S., AND H. WHITE (1995): "An alternative definition of finite-sample breakdown point with applications to regression model estimators," *Journal of the American Statistical Association*, 90(431).
- (1998): "High breakdown point conditional dispersion estimation with application to S & P 500 daily returns volatility," *Econometrica*, 66(3), 529–567.
- SIEGEL, A. (1982): "Robust regression using repeated medians," *Biometrika*, 69(1), 242.
- YOHAI, V. (1987): "High breakdown-point and high efficiency robust estimates for regression," *The Annals of Statistics*, 15(2), 642–656.
- YOHAI, V., AND R. ZAMAR (1988): "High breakdown-point estimates of regression by means of the minimization of an efficient scale," *Journal of the American Statistical Association*, 83(402), 406–413.
- ZINDE-WALSH, V. (2002): "Asymptotic theory for some high breakdown point estimators," *Econometric theory*, 18(05), 1172–1196.

## APPENDIX A. BASICS

**Lemma A1.**  $\forall t : m(-t) = m(t)$ .

*Proof.* Note that

$$\begin{aligned}
q &= \mathbb{P}\{|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq m(t)\} = \mathbb{E}[F\{m(t) + \mathbf{x}_i^\top t|\mathbf{x}_i\} - F\{-m(t) + \mathbf{x}_i^\top t|\mathbf{x}_i\}] \\
&= \mathbb{E}[1 - F\{-m(t) - \mathbf{x}_i^\top t|\mathbf{x}_i\} - 1 + F\{m(t) - \mathbf{x}_i^\top t|\mathbf{x}_i\}] \\
&= \mathbb{E}[F\{m(t) + \mathbf{x}_i^\top(-t)|\mathbf{x}_i\} - F\{-m(t) + \mathbf{x}_i^\top(-t)|\mathbf{x}_i\}], \quad (16)
\end{aligned}$$

where the penultimate equality follows from the symmetry of  $F$ . Hence  $m(t) = m(-t)$ .  $\square$

## APPENDIX B. CONSISTENCY

The results in appendix B presume assumptions A and B to hold. Let for arbitrary scalar  $\lambda$  and  $t \in \mathbb{R}^d$  and any  $2 \leq p < \infty$ ,  $\mathbf{v}_i(\lambda, t) = |\mathbf{u}_i - \lambda \mathbf{x}_i^\top t| / (1 + \lambda^p)$ ,  $\mathbf{a}_i = [\mathbf{u}_i | \mathbf{x}_i^\top]^\top$ ,  $\mathcal{Q}_a = \mathcal{Q}(\|\mathbf{a}_i\|)$ ,  $\hat{\mathcal{Q}}_a = \hat{\mathcal{Q}}(\|\mathbf{a}_i\|)$ ,  $\hat{\mathbf{m}}_v(\lambda, t) = \hat{\mathcal{Q}}\{\mathbf{v}_i\}(\lambda, t)$ ,  $m_v(\lambda, t) = \mathcal{Q}\{\mathbf{v}_i(\lambda, t)\}$ . We moreover use  $m, m_\infty$  as defined in (7) and  $\mathcal{L}^\infty(\mathbb{R}^d)$  as the collection of bounded functions on  $\mathbb{R}^d$  equipped with a sup-norm.

**Lemma B1.** *For all  $\epsilon_1 > 0$  there exists a  $C_1 < \infty$  such that (i)  $\sup_{\lambda > C_1} \sup_{\|t\|=1} m_v(\lambda, t) \leq \epsilon_1$  and (ii)*

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{ \sup_{\lambda > C_1} \sup_{\|t\|=1} \hat{\mathbf{m}}_v(\lambda, t) > \epsilon_1 \right\} = 0.$$

*Proof.* We show (ii), where we establish along the way that (i) holds, also. Take  $C_1 = \max(2\mathcal{Q}_a/\epsilon_1, 1)$ .

Then  $\sup_{\lambda > C_1} \sup_{\|t\|=1} \hat{\mathbf{m}}_v(\lambda, t) \leq \hat{\mathcal{Q}}_a C_1 / (1 + C_1^p)$  and

$$\mathbb{P}\{\hat{\mathcal{Q}}_a C_1 > (1 + C_1^p)\epsilon_1\} \leq \underbrace{\mathbb{P}\{|\hat{\mathcal{Q}}_a - \mathcal{Q}_a| C_1 > (1 + C_1^p)\epsilon_1/2\}}_{<1} + \underbrace{I\{\mathcal{Q}_a C_1 > (1 + C_1^p)\epsilon_1/2\}}_{=0}. \quad \square$$

**Lemma B2.** *For any  $C_1 < \infty$  there exists a  $C_2 < \infty$  such that*

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{ \sup_{0 \leq \lambda \leq C_1} \sup_{\|t\|=1} \hat{\mathbf{m}}_v(\lambda, t) > C_2 \right\} = 0.$$

*Proof.* Take  $C_2 = 4\mathcal{Q}_a > 0$ . Noting that  $\max(1, \lambda)/(1 + \lambda^p) \leq 1$ ,

$$\mathbb{P}\left\{ \sup_{0 \leq \lambda \leq C_1} \sup_{\|t\|=1} \hat{\mathbf{m}}_v(\lambda, t) > C_2 \right\} \leq \mathbb{P}(\hat{\mathcal{Q}}_a > C_2) \leq \mathbb{P}(|\hat{\mathcal{Q}}_a - \mathcal{Q}_a| > C_2/2) < 1. \quad \square$$

Let  $\hat{A}_v(\lambda, t, m) = n^{-1} \sum_{i=1}^n I\{\mathbf{v}_i(\lambda, t) \leq m\}$  and  $A_v(\lambda, t, m) = \mathbb{P}\{\mathbf{v}_i(\lambda, t) \leq m\}$ .

**Lemma B3.**  $\sqrt{n}(\hat{A}_v - A_v) \xrightarrow{w} \mathbf{G}_v$  in  $\mathcal{L}^\infty(\mathbb{R}^{d+2})$  for a zero mean Gaussian process  $\mathbf{G}_v$  with covariance kernel  $H_v(\lambda, t, m, \tilde{\lambda}, \tilde{t}, \tilde{m}) = \text{Cov}[I\{|\mathbf{u}_i - \lambda \mathbf{x}_i^\top t| \leq (1 + \lambda^p)m\}, I\{|\mathbf{u}_i - \tilde{\lambda} \mathbf{x}_i^\top \tilde{t}| \leq (1 + \lambda^p)\tilde{m}\}]$ .

*Proof.* Let  $\mathcal{C}$  be the collection of sets of  $(u, x^\top)^\top$  indexed by  $(a, b^\top, m)^\top \in \mathbb{R}^{d+2}$  such that  $|au + x^\top b| \leq m$ . Since the collection of half spaces is a Vapnik–Chervonenkis (VC) class and  $\mathcal{C}$  is the collection of intersections of half spaces,  $\mathcal{C}$  is a VC class. Therefore,  $\mathcal{F} = \{I\{C\} : C \in \mathcal{C}\}$  is a VC subgraph class of functions that are indexed by  $(a, b^\top, m)^\top \in \mathbb{R}^{d+2}$ . Since  $\mathbb{R}^{d+2}$  is separable,  $\mathcal{F}$  is a pointwise measurable class. Therefore,  $\mathcal{F}$  is a Donsker class in  $\mathcal{L}^\infty(\mathbb{R}^{d+2})$ . Reparametrizing by  $a = 1/(1 + \lambda^p)$  and  $b = \lambda t/(1 + \lambda^p)$  does not affect the Donsker property, and therefore the weak convergence of  $\sqrt{n}(\hat{A}_v - A_v)$  in  $\mathcal{L}^\infty(\mathbb{R}^{d+2})$  follows. Apply a central limit theorem to arbitrary finite marginals and the Gaussian limit process and covariance kernel follow.  $\square$

Let  $c_{pt} = 1 + \|t\|^p$  and  $\mathbf{G}_p$  a Gaussian process with covariance kernel  $H_p^*(t, s) = H(t, s)/c_{pt}c_{ps}$ , where  $H$  is as defined in (8).

**Lemma B4.** Let  $\mathcal{S}_{n1}(t, m) = \sqrt{n}\{\hat{A}_n(t, m) - A(t, m)\}$ ,  $\mathcal{S}_{n2}(t) = \sqrt{n}[\hat{A}_n\{t, m(t)\} - A\{t, m(t)\}]$ , and  $\mathcal{S}_{n3p}(t) = \sqrt{n}[\hat{A}_n\{t, m(t)\} - A\{t, m(t)\}]/c_{pt}$ . For  $\mathbf{G}_p$  as defined above and some other Gaussian processes  $\mathcal{G}_1, \mathcal{G}_2$ , (i)  $\mathcal{S}_{n1} \xrightarrow{w} \mathcal{G}_1$  in  $\mathcal{L}^\infty(\mathbb{R}^{d+1})$ , (ii)  $\mathcal{S}_{n2} \xrightarrow{w} \mathcal{G}_2$  in  $\mathcal{L}^\infty(\mathbb{R}^d)$ , (iii)  $\mathcal{S}_{n3p} \xrightarrow{w} \mathbf{G}_p$  in  $\mathcal{L}^\infty(\mathbb{R}^d)$ .

*Proof.* First (i). Since the collection of half spaces in a Euclidean space is a VC class, the indicator functions  $\mathbb{F}^* = \{I(|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq m) : (t, m) \in \mathbb{R}^{d+1}\}$  form a VC subgraph class, because  $\mathbb{F}^*$  is generated by using a finite intersection of half spaces. Since  $\mathbb{R}^{d+1}$  is separable,  $\mathbb{F}^*$  is a pointwise measurable class and is hence Donsker.

Since the derivations for (ii) and (iii) are similar to each other, we only consider (iii). Since the Gaussianity of finite marginals follows from a central limit theorem, we focus on the stochastic equicontinuity of  $\mathcal{S}_{n3p}$ . Note that

$$\begin{aligned} |\mathcal{S}_{n3p}(t) - \mathcal{S}_{n3p}(\tilde{t})| &\leq \sup_m |\mathcal{S}_{n1}(t, m) - \mathcal{S}_{n1}(\tilde{t}, m)| + \sup_{t, m} |\mathcal{S}_{n1}(t, m)| \left| \frac{1}{c_{pt}} - \frac{1}{c_{p\tilde{t}}} \right| \\ &\quad + \sup_{t^*} |\mathcal{S}_{n1}\{t^*, m(t)\} - \mathcal{S}_{n1}\{t^*, m(\tilde{t})\}|, \end{aligned}$$

where the RHS converges in probability to 0 as  $\|t - \tilde{t}\| \rightarrow 0$ , because of (i) and since  $1/c_{pt}$  and  $m$  are continuous in  $t$ .  $\square$

**Lemma B5.**  $\sup_{\lambda, t} |A_v\{\lambda, t, \hat{m}_v(\lambda, t)\} - A_v\{\lambda, t, m_v(\lambda, t)\}| \leq 1/\sqrt{n}$ .

*Proof.* By the triangle inequality and the definition of  $m_v$ ,

$$\begin{aligned} \sup_{\lambda, t} |A_v\{\lambda, t, \hat{\mathbf{m}}_v(\lambda, t)\} - A_v\{\lambda, t, m_v(\lambda, t)\}| \leq \\ \sup_{\lambda, t} |A_v\{\lambda, t, \hat{\mathbf{m}}_v(\lambda, t)\} - \hat{A}_v\{\lambda, t, \hat{\mathbf{m}}_v(\lambda, t)\}| + \sup_{\lambda, t} |\hat{A}_v\{\lambda, t, \hat{\mathbf{m}}_v(\lambda, t)\} - q|. \end{aligned} \quad (17)$$

The first right hand side term (RHS1) in (17) is  $\leq 1/\sqrt{n}$  by lemma B3 and RHS2 is  $\prec 1/\sqrt{n}$  by the definition of  $\hat{\mathbf{m}}_v$ .  $\square$

**Lemma B6.** For any  $C_1 < \infty$ ,  $\sup_{0 \leq \lambda \leq C_1} \sup_{\|t\|=1} |\hat{\mathbf{m}}_v(\lambda, t) - m_v(\lambda, t)| \leq 1/\sqrt{n}$ .

*Proof.* By lemma B5 and the mean value theorem, for some  $\hat{\mathbf{m}}_v^*(\lambda, t)$  between  $\hat{\mathbf{m}}_v(\lambda, t)$  and  $m_v(\lambda, t)$ ,

$$\sup_{0 \leq \lambda \leq C_1} \sup_{\|t\|=1} |\partial_m A_v\{\lambda, t, \hat{\mathbf{m}}_v^*(\lambda, t)\} \{\hat{\mathbf{m}}_v(\lambda, t) - m_v(\lambda, t)\}| \leq 1/\sqrt{n}.$$

It hence suffices to show that for some  $C_3 > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[ \inf_{0 \leq \lambda \leq C_1} \inf_{\|t\|=1} \partial_m A_v\{\lambda, t, \hat{\mathbf{m}}_v^*(\lambda, t)\} > C_3 \right] = 1.$$

By lemma B2 it suffices to show that

$$\inf_{0 \leq \lambda \leq C_1} \inf_{\|t\|=1} \inf_{0 \leq m \leq C_2} |\partial_m A_v(\lambda, t, m)| > 0.$$

Because  $A_v(\lambda, t, m) = \mathbb{E}[F\{\lambda \mathbf{x}_i^\top t + (1 + \lambda^p)m | \mathbf{x}_i\} - F\{\lambda \mathbf{x}_i^\top t - (1 + \lambda^p)m | \mathbf{x}_i\}]$ , it follows that for sufficiently large but finite  $C_4$ ,

$$\begin{aligned} \partial_m A_v(\lambda, t, m) &= \mathbb{E}[f\{\lambda \mathbf{x}_i^\top t + (1 + \lambda^p)m | \mathbf{x}_i\} + f\{\lambda \mathbf{x}_i^\top t - (1 + \lambda^p)m | \mathbf{x}_i\}] \\ &\geq 2 \mathbb{E}[f\{C_1 C_4 + (1 + C_1^p)C_2 | \mathbf{x}_i\} I(\|\mathbf{x}_i\| \leq C_4)] > 0, \end{aligned} \quad (18)$$

by assumption A.  $\square$

**Lemma B7.**  $\sup_{\lambda \geq 0, \|t\|=1} |\hat{\mathbf{m}}_v(\lambda, t) - m_v(\lambda, t)| \prec 1$ .

*Proof.* We show that for any  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \sup_{\lambda \geq 0, \|t\|=1} |\hat{\mathbf{m}}_v(\lambda, t) - m_v(\lambda, t)| > \epsilon \right\} = 0. \quad (19)$$

In lemma B1, take  $\epsilon_1 = \epsilon/4$  to show that for the choice of  $C_1$  given there,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \sup_{\lambda > C_1, \|t\|=1} |\hat{\mathbf{m}}_v(\lambda, t) - m_v(\lambda, t)| > \epsilon/2 \right\} = 0.$$

The case  $0 \leq \lambda \leq C_1$  is dealt with in lemma B6. □

**Lemma B8.**  $\inf_{t \in \mathbb{R}^d} \{m(t)/(1 + \|t\|)\} > 0$ .

*Proof.* We show equivalently that

$$\exists \epsilon, c > 0 : \sup_{\lambda \geq 0, \|t\|=1} \mathbb{P}\{|\mathbf{u}_i - \lambda \mathbf{x}_i^\top t| \leq c(1 + \lambda)\} \leq q - \epsilon. \quad (20)$$

By assumption C it follows that for any sufficiently small  $\epsilon > 0$ ,

$$\sup_{\|t\|=1} \mathbb{P}(|\mathbf{x}_i^\top t| \leq \epsilon) \leq q - 2\epsilon. \quad (21)$$

Now choose  $K_1 = -F^{-1}(\epsilon/2)$ ,  $K_2 = \max(2K_1/\epsilon, 1)$ ,  $c = \min(\epsilon/4, (q - \epsilon)/[2(1 + K_2) \mathbb{E}\{f(0|\mathbf{x}_i)\}])$ .

Then

$$\begin{aligned} & \sup_{\lambda > K_2, \|t\|=1} \mathbb{P}\{|\mathbf{u}_i - \lambda \mathbf{x}_i^\top t| \leq c(1 + \lambda)\} \leq \\ & \mathbb{P}(|\mathbf{u}_i| > K_1) + \sup_{\lambda > K_2, \|t\|=1} \mathbb{P}\{|\mathbf{u}_i - \lambda \mathbf{x}_i^\top t| \leq c(1 + \lambda), |\mathbf{u}_i| \leq K_1\} \\ & \leq \epsilon + \sup_{\lambda > K_2, \|t\|=1} \mathbb{P}\{\lambda |\mathbf{x}_i^\top t| \leq c(1 + \lambda) + K_1\} \leq \epsilon + \sup_{\|t\|=1} \mathbb{P}(|\mathbf{x}_i^\top t| \leq 2c + K_1/K_2) \\ & \leq \epsilon + \sup_{\|t\|=1} \mathbb{P}(|\mathbf{x}_i^\top t| \leq \epsilon) \leq q - \epsilon, \quad (22) \end{aligned}$$

by (21). Further,

$$\begin{aligned} & \sup_{\lambda \leq K_2, \|t\|=1} \mathbb{P}\{|\mathbf{u}_i - \lambda \mathbf{x}_i^\top t| \leq c(1 + \lambda)\} \\ & \leq \sup_{\lambda \leq K_2, \|t\|=1} \mathbb{E}[F\{\lambda \mathbf{x}_i^\top t + c(1 + \lambda)|\mathbf{x}_i\} - F\{\lambda \mathbf{x}_i^\top t - c(1 + \lambda)|\mathbf{x}_i\}] \\ & \leq 2c(1 + K_2) \mathbb{E}\{f(0|\mathbf{x}_i)\} \leq q - \epsilon, \quad (23) \end{aligned}$$

by assumption A and the choice of  $c$ . Combining (22) and (23) yields (20). □

**Lemma B9.**  $\inf_{t \in \mathbb{R}^d} \{\hat{\mathbf{m}}_n(t)/(1 + \|t\|)\} \succeq 1$ .

*Proof.* We show the equivalent result that for any sufficiently small  $\epsilon, c > 0$ ,

$$\mathbb{P}\left[\sup_{\lambda \geq 0, \|t\|=1} \hat{\mathbf{A}}_n\{\lambda t, c(1+\lambda)\} > q - 2\epsilon\right] \prec 1.$$

Now,

$$\begin{aligned} \mathbb{P}\left[\sup_{\lambda \geq 0, \|t\|=1} \hat{\mathbf{A}}_n\{\lambda t, c(1+\lambda)\} > q - 2\epsilon\right] &\leq \\ &I[A\{\lambda t, c(1+\lambda)\} \geq q\epsilon] + \mathbb{P}\left[\sup_{\lambda \geq 0, \|t\|=1} |\hat{\mathbf{A}}_n\{\lambda t, c(1+\lambda)\} - A\{\lambda t, c(1+\lambda)\}| > \epsilon\right]. \end{aligned} \quad (24)$$

RHS2 in (24) is  $\prec 1$  by lemma B4 and RHS1 is exactly (20).  $\square$

**Lemma B10.**  $\sup_{t \in \mathbb{R}^d} \{m(t)/(1 + \|t\|)\} < \infty$ .

*Proof.* The result follows immediately from the fact that  $\mathcal{Q}(|\mathbf{u}_i - \mathbf{x}_i^\top t|) \leq (1 + \|t\|)\mathcal{Q}_a$ , with  $\mathcal{Q}_a$  defined at the beginning of appendix B.  $\square$

**Lemma B11.**  $\sup_{t \in \mathbb{R}^d} \{\hat{m}_n(t)/(1 + \|t\|)\} \preceq 1$ .

*Proof.* The result follows immediately from the fact that  $\hat{m}_n(t) \leq (1 + \|t\|)\hat{\mathcal{Q}}_a$ , with  $\hat{\mathcal{Q}}_a$  defined at the beginning of appendix B.  $\square$

#### APPENDIX C. ASYMPTOTIC NORMALITY

The assumptions of theorem 4 are taken to hold for the lemmas below. The suprema and infima in this section are taken over  $\lambda \in [0, \infty)$  and  $\|t\| = 1$ , unless otherwise noted. Let  $f^{-1} : [0, f(0)] \rightarrow [0, \infty)$  be a function (not necessarily unique) such that  $f\{f^{-1}(t)\} = t$  for all  $t$ .

**Lemma C1.** For some  $C > 0$  and all sufficiently large  $\lambda$ ,

$$\sup_t |m(\lambda t) - \lambda m_\infty(t)| \leq f^{-1}\left(\frac{1}{C\lambda}\right) + C\lambda F\left\{-f^{-1}\left(\frac{1}{C\lambda}\right)\right\}.$$

*Proof.* For  $C$  to be chosen, set  $\epsilon = 2F\{-f^{-1}(1/C\lambda)\}$ . By David (1986, theorem 1) and assumption D for some finite  $C$  independent of  $\lambda$ ,

$$\begin{aligned} \sup_t \{m(\lambda t) - \lambda m_\infty(t)\} &\leq \mathcal{Q}^*(|\mathbf{u}_i|; 1 - \epsilon) + \lambda \sup_t \{\mathcal{Q}^*(|\mathbf{x}_i^\top t|; q + \epsilon) - m_\infty(t)\} \\ &\leq \mathcal{Q}^*(|\mathbf{u}_i|; 1 - \epsilon) + C\lambda\epsilon/2 = f^{-1}(1/C\lambda) + C\lambda F\{-f^{-1}(1/C\lambda)\}. \end{aligned}$$

The case  $\sup_t \{\lambda m_\infty(t) - m(\lambda t)\}$  is similar.  $\square$

**Lemma C2.**  $\inf_{\lambda, t} \partial_m A_v \{\lambda, t, m_v(\lambda, t)\} > 0$ .

*Proof.* For all  $\lambda$  belonging to a compact set, the result was established in (18). We now show that the result also holds for large  $\lambda$ . Noting that by the symmetry of the conditional distribution of  $\mathbf{u}_i$  given  $\mathbf{x}_i$ ,

$$\begin{aligned} A_v(\lambda, t, m) &= \mathbb{P}\{|\mathbf{u}_i - \lambda \mathbf{x}_i^\top t| \leq (1 + \lambda^p)m\} = \mathbb{E}[F\{\lambda \mathbf{x}_i^\top t + (1 + \lambda^p)m | \mathbf{x}_i\} - F\{\lambda \mathbf{x}_i^\top t - (1 + \lambda^p)m | \mathbf{x}_i\}] \\ &= \mathbb{E}[F\{\lambda |\mathbf{x}_i^\top t| + (1 + \lambda^p)m | \mathbf{x}_i\} - F\{\lambda |\mathbf{x}_i^\top t| - (1 + \lambda^p)m | \mathbf{x}_i\}]. \end{aligned}$$

Hence

$$\begin{aligned} \partial_m A_v \{\lambda, t, m_v(\lambda, t)\} &= (1 + \lambda^p) \mathbb{E}[f\{\lambda |\mathbf{x}_i^\top t| + m(\lambda t) | \mathbf{x}_i\} + f\{\lambda |\mathbf{x}_i^\top t| - m(\lambda t) | \mathbf{x}_i\}] \\ &\geq (1 + \lambda^p) \mathbb{E}[f\{\lambda |\mathbf{x}_i^\top t| - m(\lambda t) | \mathbf{x}_i\}] \\ &= (1 + \lambda^p) \mathbb{E}(f[\lambda \{|\mathbf{x}_i^\top t| - m_\infty(t)\} + \{\lambda m_\infty(t) - m(\lambda t)\} | \mathbf{x}_i]). \quad (25) \end{aligned}$$

Note first that  $\mathbb{P}\{-\epsilon/\lambda < |\mathbf{x}_i^\top t| - m_\infty(t) \leq 0\}$  and  $\mathbb{P}\{0 \leq |\mathbf{x}_i^\top t| - m_\infty(t) < \epsilon/\lambda\}$  both exceed  $C^*/\lambda$  by assumption D for any fixed  $\epsilon > 0$  and some  $C^*$  independent of  $t, \lambda$ . Hence the RHS in (25) is for sufficiently large  $\lambda$  bounded below by

$$\begin{aligned} &\lambda^p \mathbb{E}\left[f\left\{\sup_t |\lambda m_\infty(t) - m(\lambda t)| | \mathbf{x}_i\right\} I(|\mathbf{x}_i^\top t| - m_\infty(t)| \leq \epsilon/\lambda)\right] \\ &\geq C^* \lambda^{p-1} \left[f\left\{\sup_t |\lambda m_\infty(t) - m(\lambda t)|\right\}\right]^r \\ &\geq \frac{C^*}{C^{p-1}} (C\lambda)^{p-1} f\left[f^{-1}\left(\frac{1}{C\lambda}\right) + C\lambda F\left\{-f^{-1}\left(\frac{1}{C\lambda}\right)\right\}\right]^r, \quad (26) \end{aligned}$$

where the first inequality in (26) follows from assumption D and the second from lemma C1 and where  $C$  is as chosen in lemma C1. Since  $p$  can be chosen arbitrarily large, it hence suffices that for some  $p^*$ ,

$$\lim_{s \rightarrow \infty} \frac{f\{s + F(-s)/f(s)\}}{\{f(s)\}^{p^*}} > 0,$$

which was assumed in assumption D. □

For some  $\sigma > 4$ , let

$$\psi_n = n^{1/\sigma p}. \quad (27)$$

**Lemma C3.**

$$\sup_{0 \leq \lambda \leq \psi_n} \sup_{\|t\|=1} |\hat{\mathbf{m}}_v(\lambda, t) - m_v(\lambda, t)| \leq 1/\sqrt{n}.$$

*Proof.* We use the shorthand  $\hat{A}_v\{\hat{\mathbf{m}}_v\}$  for  $\hat{A}_v\{\lambda, t, \hat{\mathbf{m}}_v(\lambda, t)\}$  and likewise for similar symbols. By the mean value theorem and assumption E for some function  $\hat{\mathbf{m}}_v^*$  between  $\hat{\mathbf{m}}_v$  and  $m_v$ ,

$$A_v(\hat{\mathbf{m}}_v) - A_v(m_v) = \partial_m A_v(m_v)(\hat{\mathbf{m}}_v - m_v) + \partial_m^2 A_v(\hat{\mathbf{m}}_v^*)(\hat{\mathbf{m}}_v - m_v)^2/2. \quad (28)$$

Further, by the triangle inequality,

$$\begin{aligned} \sup_{\lambda \geq 0} \sup_{\|t\|=1} |A_v(\hat{\mathbf{m}}_v) - A_v(m_v)| &\leq \\ &\sup_{\lambda \geq 0} \sup_{\|t\|=1} |A_v(\hat{\mathbf{m}}_v) - \hat{A}_v(\hat{\mathbf{m}}_v) - A_v(m_v) + \hat{A}_v(m_v)| + \sup_{\lambda \geq 0} \sup_{\|t\|=1} |A_v(m_v) - \hat{A}_v(m_v)| \\ &\quad + \sup_{\lambda \geq 0} \sup_{\|t\|=1} |\hat{A}_v(\hat{\mathbf{m}}_v) - q| + \sup_{\lambda \geq 0} \sup_{\|t\|=1} |q - A_v(m_v)|. \end{aligned} \quad (29)$$

RHS3 and RHS4 in (29) are  $\prec 1/\sqrt{n}$  by construction and RHS1 is  $\prec 1/\sqrt{n}$  by lemma B3. RHS2 is  $\leq 1/\sqrt{n}$ , also, by lemma B3. Combining (28) and (29) yields

$$\sup_{\lambda \geq 0} \sup_{\|t\|=1} |\partial_m A_v(m_v)(\hat{\mathbf{m}}_v - m_v) + \partial_m^2 A_v(\hat{\mathbf{m}}_v^*)(\hat{\mathbf{m}}_v - m_v)^2/2| \leq 1/\sqrt{n}. \quad (30)$$

Now let  $\kappa_n \prec 1$  be such that

$$\sup_{0 \leq \lambda \leq \kappa_n^{-2/\sigma p}} \sup_{\|t\|=1} |\hat{\mathbf{m}}_v - m_v| \leq \kappa_n. \quad (31)$$

Such  $\kappa_n$  exist by lemma B7; we will choose it later. Then for any such  $\kappa_n$  it follows from (30) that

$$\begin{aligned} \sup_{0 \leq \lambda \leq \kappa_n^{-2/\sigma p}} \sup_{\|t\|=1} |\partial_m A_v(m_v)(\hat{\mathbf{m}}_v - m_v)| &\leq \sup_{0 \leq \lambda \leq \kappa_n^{-2/\sigma p}} \sup_{\|t\|=1} |\partial_m^2 A_v(\hat{\mathbf{m}}_v^*)(\hat{\mathbf{m}}_v - m_v)^2/2| + 1/\sqrt{n} \\ &\leq \sup_{u, x} f'(u|x) \sup_{0 \leq \lambda \leq \kappa_n^{-2/\sigma p}} \sup_{\|t\|=1} |(1 + \lambda^p)^2(\hat{\mathbf{m}}_v - m_v)^2| + 1/\sqrt{n}. \end{aligned} \quad (32)$$

RHS1 in (32) is  $\prec \sup_{0 \leq \lambda \leq \kappa_n^{-2/\sigma p}} \sup_{\|t\|=1} |\hat{\mathbf{m}}_v - m_v|$  by (31) and assumption E whereas the LHS in (32) is  $\succeq \sup_{0 \leq \lambda \leq \kappa_n^{-2/\sigma p}} \sup_{\|t\|=1} |\hat{\mathbf{m}}_v - m_v|$ , by lemma C2. Hence

$$\sup_{0 \leq \lambda \leq \kappa_n^{-2/\sigma p}} \sup_{\|t\|=1} |\hat{\mathbf{m}}_v - m_v| \leq 1/\sqrt{n},$$

so we can choose  $\kappa_n = 1/\sqrt{n}$ , which corresponds to  $\psi_n = n^{1/\sigma p}$ , as defined in (27).  $\square$



**Lemma C4.**

$$\sup_{0 \leq \lambda \leq \psi_n} \sup_{\|t\|=1} |\hat{\mathbf{A}}_v\{\lambda, t, m_v(\lambda, t)\} - A_v\{\lambda, t, m_v(\lambda, t)\} - \partial_m A_v\{\lambda, t, m_v(\lambda, t)\}\{\hat{\mathbf{m}}_v(\lambda, t) - m_v(\lambda, t)\}| \prec 1/\sqrt{n}.$$

*Proof.* Using the same short hand notation as in lemma C3 we have by the triangle inequality that

$$\begin{aligned} |\hat{\mathbf{A}}_v(m_v) - A_v(m_v) + \partial_m A_v(m_v)(\hat{\mathbf{m}}_v - m_v)| &\leq |\hat{\mathbf{A}}_v(m_v) - A_v(m_v) - \hat{\mathbf{A}}_v(\hat{\mathbf{m}}_v) + A_v(\hat{\mathbf{m}}_v)| \\ &+ |\hat{\mathbf{A}}_v(\hat{\mathbf{m}}_v) - q| + |q - A_v(m_v)| + |A_v(m_v) - A_v(\hat{\mathbf{m}}_v) + \partial_m A_v(m_v)(\hat{\mathbf{m}}_v - m_v)|. \end{aligned} \quad (33)$$

RHS2 and RHS3 in (33) are  $\prec 1/\sqrt{n}$  by construction and RHS1 is  $\prec 1/\sqrt{n}$  by lemma B3, all uniformly in  $\lambda, t$ . By the mean value theorem, for some  $\hat{\mathbf{m}}_v^*$  between  $m_v, \hat{\mathbf{m}}_v$ , RHS4 in (33) is

$$\begin{aligned} \sup_{0 \leq \lambda \leq \psi_n} \sup_{\|t\|=1} |\partial_m^2 A_v(\hat{\mathbf{m}}_v^*)(\hat{\mathbf{m}}_v - m_v)^2/2| &\leq (1 + \psi_n^p)^2 \sup_{u,x} f'(u|x) \sup_{0 \leq \lambda \leq \psi_n} \sup_{\|t\|=1} (\hat{\mathbf{m}}_v - m_v)^2 \\ &\preceq n^{2/\sigma-1} \prec 1/\sqrt{n}, \end{aligned}$$

by lemma C3 and because  $\sigma > 4$ ; see (27). □

**Lemma C5.** Recalling that  $c_{pt} = 1 + \|t\|^p$ ,

$$\sup_{\|t\| \leq \psi_n} \left| \frac{\hat{\mathbf{m}}_n(t) - m(t)}{c_{pt}} + \frac{\hat{\mathbf{A}}_n\{t, m(t)\} - A\{t, m(t)\}}{c_{pt} \partial_m A\{t, m(t)\}} \right| \prec \frac{1}{\sqrt{n}}. \quad (34)$$

*Proof.* Note that for  $\lambda = \|t\|$  and  $t_0 = t/\|t\|$ ,

$$\begin{aligned} \hat{\mathbf{m}}_n(t) &= c_{pt} \hat{\mathbf{m}}_v(\lambda, t_0), & m(t) &= c_{pt} m_v(\lambda, t_0), \\ A\{t, m(t)\} &= A_v\{\lambda, t_0, m_v(\lambda, t_0)\} = q, & \hat{\mathbf{A}}_n\{t, m(t)\} &= \hat{\mathbf{A}}_v\{\lambda, t_0, m_v(\lambda, t_0)\}, \\ \partial_m A\{t, m(t)\} &= \partial_m A_v\{\lambda, t_0, m_v(\lambda, t_0)\} / c_{pt}. \end{aligned}$$

The stated result then follows from lemmas C2 and C4. □

## APPENDIX D. VARIANCE MATRIX ESTIMATION

In this appendix the assumptions of theorem 5 are used. Let  $h(t) = h^*(\theta_0 + t)$ , set  $p = p^* - 1$ , with  $p^*, \sigma$  as defined in assumption F, and let

$$\begin{aligned}\hat{\mathcal{D}}(t) &= \frac{1}{2nh(t)} \sum_{i=1}^n I\{|\|\mathbf{u}_i - \mathbf{x}_i^\top t| - \hat{\mathbf{m}}_n(t)| \leq h(t)\}, \\ \hat{\mathbf{H}}(t, s) &= n^{-1} \sum_{i=1}^n I\{|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq \hat{\mathbf{m}}_n(t)\} I\{|\mathbf{u}_i - \mathbf{x}_i^\top s| \leq \hat{\mathbf{m}}_n(s)\} \\ &\quad - \left( n^{-1} \sum_{i=1}^n I\{|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq \hat{\mathbf{m}}_n(t)\} \right) \left( n^{-1} \sum_{i=1}^n I\{|\mathbf{u}_i - \mathbf{x}_i^\top s| \leq \hat{\mathbf{m}}_n(s)\} \right).\end{aligned}\tag{35}$$

**Lemma D1.**  $\sup_{t, s \in \mathbb{R}^d} |\hat{\mathbf{H}}(t, s) - H(t, s)| \prec 1$ .

*Proof.* We will show the uniform convergence of RHS1 in the definition of  $\hat{\mathbf{H}}$  in (35), because the second term is similar. Reparametrize the first term of  $\hat{\mathbf{H}}(t, s)$  in terms of  $\lambda_0, \lambda_1, t_0, t_1$  to obtain

$$n^{-1} \sum_{i=1}^n I\{(1 + \lambda_1^p)^{-1} |\mathbf{u}_i - \mathbf{x}_i^\top \lambda_1 t_0| \leq \hat{\mathbf{m}}_v(\lambda_1, t_0)\} I\{(1 + \lambda_2^p)^{-1} |\mathbf{u}_i - \mathbf{x}_i^\top \lambda_2 s_0| \leq \hat{\mathbf{m}}_v(\lambda_2, s_0)\},$$

where  $\lambda_1, \lambda_2 \geq 0$  and  $\|t_0\| = \|s_0\| = 1$ . Since  $\hat{\mathbf{m}}_v$  converges uniformly to  $m_v$  by lemma B7, it suffices to show that for  $\mathbf{a}_i$  defined at the beginning of appendix B,

$$n^{-1} \sum_{i=1}^n I(|\mathbf{a}_i^\top t^*| \leq m^*) I(|\mathbf{a}_i^\top \tilde{t}^*| \leq \tilde{m}^*) \xrightarrow{p} \mathbb{E}\{I(|\mathbf{a}_i^\top t^*| \leq m^*) I(|\mathbf{a}_i^\top \tilde{t}^*| \leq \tilde{m}^*)\}\tag{36}$$

uniformly in  $(t^*, m^*, \tilde{t}^*, \tilde{m}^*) \in \mathbb{R}^{d+1} \times \mathbb{R}_+ \times \mathbb{R}^{d+1} \times \mathbb{R}_+$ . Since the collection of half spaces in a Euclidean space is a VC class, the collection of indicator functions  $\mathcal{F} = \{I(|\mathbf{a}_i^\top t^*| \leq m^*) I(|\mathbf{a}_i^\top \tilde{t}^*| \leq \tilde{m}^*) : (t^*, m^*, \tilde{t}^*, \tilde{m}^*) \in \mathbb{R}^{d+1} \times \mathbb{R}_+ \times \mathbb{R}^{d+1} \times \mathbb{R}_+\}$  generated by finite intersections of half spaces form a VC subgraph class. Since  $\mathbb{R}^{d+1}$  and  $\mathbb{R}_+$  are separable,  $\mathcal{F}$  is a pointwise measurable class, and it follows that  $\mathcal{F}$  is Glivenko–Cantelli.  $\square$

**Lemma D2.** For some  $\epsilon > 0$  and recalling that  $c_{pt} = 1 + \|t\|^p$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{ \inf_{t \in \mathbb{R}^d} c_{pt} \hat{\mathcal{D}}(t) < \epsilon \right\} = 0.$$

*Proof.* Note that  $\hat{\mathcal{D}}(t) = 2[\hat{\mathbf{A}}_n\{t, \hat{\mathbf{m}}_n(t) + h(t)\} - \hat{\mathbf{A}}_n\{t, \hat{\mathbf{m}}_n(t) - h(t)\}]/h(t)$ , such that by lemma B4,

$$\sup_t c_{pt} \left| 2\hat{\mathcal{D}}(t) - \frac{A\{t, \hat{\mathbf{m}}_n(t) + h(t)\} - A\{t, \hat{\mathbf{m}}_n(t) - h(t)\}}{h(t)} \right| \leq \frac{1}{\sqrt{nh_0}} \prec 1.\tag{37}$$

Using  $\hat{\mathcal{Q}}_a, \mathcal{Q}_a$  from appendix B it follows from assumption F that

$$\sup_{\|t\| > \psi_n} \frac{|\hat{\mathbf{m}}_n(t) - m(t)|}{h(t)} \leq \frac{\hat{\mathcal{Q}}_a + \mathcal{Q}_a}{h_0} \sup_{\|t\| > \psi_n} \frac{1 + \|t\|}{1 + \|t\|^p} \preceq \frac{n^{(1-p)/p\sigma}}{h_0} \prec 1. \quad (38)$$

Likewise, by lemma C3 and assumption F,

$$\sup_{\|t\| \leq \psi_n} \frac{|\hat{\mathbf{m}}_n(t) - m(t)|}{h(t)} = \sup_{\|t\| \leq \psi_n} \frac{|\hat{\mathbf{m}}_v - m_v(t)|}{h_0} \preceq \frac{1}{\sqrt{nh_0}} \prec 1. \quad (39)$$

From (37) to (39) it follows that it suffices to show that

$$\inf_t \frac{A\{t, m(t) + h(t)/2\} - A\{t, m(t) - h(t)/2\}}{h_0} > 0. \quad (40)$$

The LHS in (40) is by the mean value theorem for some  $0 < \omega(t) < 1$  and  $\epsilon^*(t) = \omega(t)h(t)/2$  equal to

$$\begin{aligned} c_{pt} \partial_m A\{t, m(t) + \epsilon^*(t)\} + c_{pt} \partial_m A\{t, m(t) - \epsilon^*(t)\} = \\ c_{pt} \mathbb{E} [f\{\mathbf{x}_i^\top t + m(t) + \epsilon^*(t) | \mathbf{x}_i\} + f\{\mathbf{x}_i^\top t - m(t) - \epsilon^*(t) | \mathbf{x}_i\} \\ + f\{\mathbf{x}_i^\top t + m(t) - \epsilon^*(t) | \mathbf{x}_i\} + f\{\mathbf{x}_i^\top t - m(t) + \epsilon^*(t) | \mathbf{x}_i\}], \end{aligned}$$

which by assumption A equals

$$\begin{aligned} c_{pt} \mathbb{E} [f\{|\mathbf{x}_i^\top t| + m(t) + \epsilon^*(t) | \mathbf{x}_i\} + f\{|\mathbf{x}_i^\top t| - m(t) - \epsilon^*(t) | \mathbf{x}_i\} \\ + f\{|\mathbf{x}_i^\top t| + m(t) - \epsilon^*(t) | \mathbf{x}_i\} + f\{|\mathbf{x}_i^\top t| - m(t) + \epsilon^*(t) | \mathbf{x}_i\}] \geq \\ c_{pt} \mathbb{E} [I\{|\mathbf{x}_i^\top t| - m(t) \geq 0\} f\{|\mathbf{x}_i^\top t| - m(t) - \epsilon^*(t) | \mathbf{x}_i\} + I\{|\mathbf{x}_i^\top t| - m(t) < 0\} f\{|\mathbf{x}_i^\top t| - m(t) + \epsilon^*(t) | \mathbf{x}_i\}] \\ \geq c_{pt} \mathbb{E} [I\{|\mathbf{x}_i^\top t| - m(t) \geq 0\} f\{|\mathbf{x}_i^\top t| - m(t) | \mathbf{x}_i\} + I\{|\mathbf{x}_i^\top t| - m(t) < 0\} f\{|\mathbf{x}_i^\top t| - m(t) | \mathbf{x}_i\}] \\ = c_{pt} \mathbb{E} [f\{|\mathbf{x}_i^\top t| - m(t) | \mathbf{x}_i\}]. \end{aligned}$$

Now,

$$\inf_t \left( c_{pt} \mathbb{E} [f\{|\mathbf{x}_i^\top t| - m(t) | \mathbf{x}_i\}] \right) = \inf_{\lambda \geq 0} \inf_{\|t\|=1} \left( (1 + \lambda^p) \mathbb{E} [f\{\lambda |\mathbf{x}_i^\top t| - m(\lambda t) | \mathbf{x}_i\}] \right). \quad (41)$$

The RHS in (41) is the infimum of the middle expression in (25), which is shown to be bounded away from zero, uniformly in  $\lambda, t$ , in lemma C2.  $\square$

**Lemma D3.**  $\forall t \in \mathbb{R}^d : |\hat{\mathcal{D}}(t) - \mathcal{D}(t)| \prec 1.$

*Proof.* Choose  $t$ . From (37) and assumption F it follows that

$$\left| 2\hat{\mathcal{D}}(t) - \frac{A\{t, \hat{\mathbf{m}}_n(t) + h(t)\} - A\{t, \hat{\mathbf{m}}_n(t) - h(t)\}}{h(t)} \right| \leq \frac{1}{\sqrt{nh_0}} \prec 1.$$

Take  $n$  large enough to ensure that  $\|t\| \leq \psi_n$ . We have by the mean value theorem that

$$\begin{aligned} \frac{A\{t, \hat{\mathbf{m}}_n(t) + h(t)\} - A\{t, \hat{\mathbf{m}}_n(t) - h(t)\}}{h(t)} - 2\mathcal{D}(t) = \\ \frac{\partial_m^2 A(t, \cdot)}{2h(t)} \{\hat{\mathbf{m}}_n(t) - m(t) + h(t)\}^2 - \frac{\partial_m^2 A(t, \cdot)}{2h(t)} \{\hat{\mathbf{m}}_n(t) - m(t) - h(t)\}^2. \end{aligned} \quad (42)$$

By lemma C3 and assumptions E and F both RHS terms in (42) are  $\simeq (1/n + h_0^2)/h_0 \prec 1$ .  $\square$

**Lemma D4.**

$$\forall t, s \in \mathbb{R}^d : \frac{\hat{\mathbf{H}}(t, s)}{\hat{\mathcal{D}}(t)\hat{\mathcal{D}}(s)} - \frac{H(t, s)}{\mathcal{D}(t)\mathcal{D}(s)} \prec 1.$$

*Proof.* Follows immediately from lemmas D1 to D3 and the fact that  $\hat{\mathbf{H}}$  and  $H$  are bounded.  $\square$

#### APPENDIX E. THEOREMS

**Proof of theorem 1.** The proof of this theorem is inspired by that of Rousseeuw (1984, theorem 1); see also Rousseeuw and Leroy (1987, section 3.4). Let  $\mathcal{X} = \{(x_i, y_i)\}$  be the original sample,  $\mathcal{X}^* = \{(x_i^*, y_i^*)\}$  be the contaminated sample and  $\mathcal{X}^\dagger$  the sample consisting of observations shared between the two.

Suppose that the number of observations contaminated is at most  $b^* = \min(n - N, N - \hat{\gamma} - 1) > 0$ . We show that  $\hat{\theta}$  is bounded, for which it suffices to show that (i)  $\int \|\theta\| \exp\{-\hat{\mathcal{D}}(|\mathbf{y}_i^* - \theta^\top \mathbf{x}_i^*|^2)\} d\theta < \infty$ , (ii)  $\int \exp\{-\hat{\mathcal{D}}(|\mathbf{y}_i^* - \theta^\top \mathbf{x}_i^*|^2)\} d\theta > 0$ . Let  $\bar{y} = \max_{\mathcal{X}} |y_i|$  and  $\bar{x} = \max_{\mathcal{X}} \|x_i\|$ .

First (ii). Since  $b^* > 0$  there exists for each  $\theta$  at least one observation  $(\mathbf{x}_i^\dagger(\theta), \mathbf{y}_i^\dagger(\theta)) \in \mathcal{X}^\dagger$  for which  $\hat{\mathcal{D}}(|\mathbf{y}_i^* - \theta^\top \mathbf{x}_i^*|) \leq |\mathbf{y}_i^\dagger(\theta) - \theta^\top \mathbf{x}_i^\dagger(\theta)| \leq \bar{y} + \bar{x}\|\theta\|$ . Hence the left hand side (LHS) in (ii) is bounded below by  $\exp(-2\bar{y}^2) \int \exp(-2\bar{x}^2\|\theta\|^2) d\theta > 0$ .

Now (i). Let  $\mathcal{B}(B, \rho)$  be the  $\rho$ -expansion of a  $(d - 1)$ -dimensional subspace  $B$  of  $\mathbb{R}^d$  and let  $\rho(\mathcal{X})$  be the smallest  $\rho$  for which  $\mathcal{B}(B, \rho)$  contains at least  $\hat{\gamma} + 1$  of the  $x_i$ 's. By the definition of  $\hat{\gamma}$ ,  $\rho(\mathcal{X}) > 0$ .

Since at most  $b^*$  observations are contaminated,  $\mathcal{X}^\dagger$  contains at least  $\hat{\gamma} + 1$  observations indexed  $i_1, \dots, i_{\hat{\gamma}+1}$  for which  $\hat{\mathcal{D}}(|\mathbf{y}_i^* - \theta^\top \mathbf{x}_i^*|) \geq |\mathbf{y}_{i_j} - \mathbf{x}_{i_j}^\top \theta|$ . Take  $B(\theta) = \{x \in \mathbb{R}^d : x^\top \theta = 0\}$ . Then for at least one  $j = j(\theta)$ ,  $x_{i_j} \notin \mathcal{B}\{B(\theta), \rho(\mathcal{X})/2\}$  by the definition of  $\rho(\mathcal{X})$ . For this value of  $j$ , let  $x_{i_j}^\dagger(\theta) = x_{i_j}$ ,  $\mathbf{y}_{i_j}^\dagger(\theta) = \mathbf{y}_{i_j}$ , and let  $\bar{y}$  be as defined above. Then since  $\hat{\mathcal{D}}(|\mathbf{y}_i^* - \theta^\top \mathbf{x}_i^*|) \geq |\mathbf{y}_{i_j}^\dagger -$

$\theta^\top \mathbf{x}_i^\dagger(\theta) \geq \rho(\mathcal{L})\|\theta\|/2 - \bar{\mathbf{y}}$  and  $\hat{\mathcal{L}}(|\mathbf{y}_i^* - \theta^\top \mathbf{x}_i^*|^2) \geq \rho^2(\mathcal{L})\|\theta\|^2/4 - \rho(\mathcal{L})\bar{\mathbf{y}}\|\theta\|$  for all  $\theta$ , the LHS in (i) is bounded above by  $\exp(\bar{\mathbf{y}}^2) \int \|\theta\| \exp[-\rho^2(\mathcal{L})\{\|\theta\| - 2\bar{\mathbf{y}}/\rho(\mathcal{L})\}^2/4] d\theta < \infty$ .  $\square$

**Proof of theorem 4.** We work on the numerator and denominator in (5) separately.

First the denominator. Let  $\omega(m) = \exp(-m^2)$ . Expanding  $\omega\{\hat{\mathbf{m}}_n(t)\}$  around  $m(t)$  yields for some  $\hat{\mathbf{m}}_n^*(t)$  between  $m(t)$  and  $\hat{\mathbf{m}}_n(t)$ ,

$$\int \exp\{-\hat{\mathbf{m}}_n^2(t)\} dt = \int \exp\{-m^2(t)\} dt - 2 \int \hat{\mathbf{m}}_n^*(t) \{\hat{\mathbf{m}}_n(t) - m(t)\} \exp\{-\hat{\mathbf{m}}_n^{*2}(t)\} dt. \quad (43)$$

For RHS2 in (43) we have by lemmas B7 to B11 for some  $\epsilon > 0$  that

$$\int \hat{\mathbf{m}}_n^*(t) \{\hat{\mathbf{m}}_n(t) - m(t)\} \exp\{-\hat{\mathbf{m}}_n^{*2}(t)\} dt \prec \int c_{pt}(1 + \|t\|) \exp(-\epsilon^2 \|t\|^2) dt \preceq 1.$$

Now the numerator. Let  $\psi_n$  be as defined in (27),  $\mathcal{T}_n = \{t \in \mathbb{R}^d : \|t\| \leq \psi_n\}$ , and  $\mathcal{T}_n^c = \{t \in \mathbb{R}^d : \|t\| > \psi_n\}$ . Then by lemma A1,

$$\begin{aligned} \int t \exp\{-\hat{\mathbf{m}}_n^2(t)\} dt = \\ \int_{\mathcal{T}_n} t [\exp\{-\hat{\mathbf{m}}_n^2(t)\} - \exp\{-m^2(t)\}] dt + \int_{\mathcal{T}_n^c} t \exp\{-\hat{\mathbf{m}}_n^2(t)\} dt - \int_{\mathcal{T}_n^c} t \exp\{-m^2(t)\} dt. \end{aligned} \quad (44)$$

For RHS2 in (44) note that by lemma B9 for some  $\epsilon > 0$  and since  $\psi_n$  is polynomial in  $n$ ,

$$\int_{\mathcal{T}_n^c} \|t\| \exp\{-\hat{\mathbf{m}}_n^2(t)\} dt \preceq \int_{\mathcal{T}_n^c} \|t\| \exp(-\epsilon^2 t^2) dt \prec 1/\sqrt{n}. \quad (45)$$

RHS3 in (44) can similarly be shown to be  $\prec 1/\sqrt{n}$  using lemma B8. For RHS1 in (44) we have

$$\begin{aligned} \int_{\mathcal{T}_n} t [\exp\{-\hat{\mathbf{m}}_n^2(t)\} - \exp\{-m^2(t)\}] dt = -2 \int_{\mathcal{T}_n} t m(t) \{\hat{\mathbf{m}}_n(t) - m(t)\} \exp\{-m^2(t)\} dt \\ - \int_{\mathcal{T}_n} t \{1 - 2\hat{\mathbf{m}}_n^{*2}(t)\} \{\hat{\mathbf{m}}_n(t) - m(t)\}^2 \exp\{-\hat{\mathbf{m}}_n^{*2}(t)\} dt. \end{aligned} \quad (46)$$

RHS2 in (46) is  $\prec 1/\sqrt{n}$  by lemmas B8 to B11, C2 and C3. Further, RHS1 in (46) equals

$$\begin{aligned} -2 \int_{\mathcal{T}_n} t c_{pt} m(t) \left\{ \frac{\hat{\mathbf{m}}_n(t) - m_t}{c_{pt}} + \frac{\hat{\mathbf{A}}_n\{t, m(t)\} - A\{t, m(t)\}}{c_{pt} \partial_m A\{t, m(t)\}} \right\} \exp\{-m^2(t)\} dt \\ -2 \int_{\mathcal{T}_n^c} t c_{pt} m(t) \frac{\hat{\mathbf{A}}_n\{t, m(t)\} - A\{t, m(t)\}}{c_{pt} \partial_m A\{t, m(t)\}} \exp\{-m^2(t)\} dt. \\ + 2 \int t m(t) \frac{\hat{\mathbf{A}}_n\{t, m(t)\} - A\{t, m(t)\}}{\mathcal{D}(t)} \exp\{-m^2(t)\} dt \end{aligned} \quad (47)$$

The first term in (47) is  $\prec 1/\sqrt{n}$  by lemma C5. The second term in (47) is also  $\prec 1/\sqrt{n}$ , which can be established along the lines of (45), noting that  $c_{pt}\partial_m A$  is bounded away from zero by lemma C2 and that  $tc_{pt}m(t)$  is polynomial in  $t$ .

Finally,  $\sqrt{n}$  times the last term in (47) equals

$$\frac{2}{\sqrt{n}} \sum_{i=1}^n \int t \frac{m(t)}{\mathcal{D}(t)} [I\{|\mathbf{u}_i - \mathbf{x}_i^\top t| \leq m(t)\} - q] \exp\{-m^2(t)\} dt. \quad (48)$$

Apply the Lindeberg–Levy central limit theorem.  $\square$

**Proof of theorem 5.** We show consistency of the numerator; the denominator is easier. Recall the definitions in (35) and let

$$\begin{aligned} \hat{\mathbf{R}}(t, s) &= \frac{\hat{\mathbf{m}}_n(t)\hat{\mathbf{m}}_n(s)}{\hat{\mathcal{D}}(t)\hat{\mathcal{D}}(s)} \hat{\mathbf{H}}(t, s) \exp\{-\hat{\mathbf{m}}_n^2(t) - \hat{\mathbf{m}}_n^2(s)\}, \\ R(t, s) &= \frac{m(t)m(s)}{\mathcal{D}(t)\mathcal{D}(s)} H(t, s) \exp\{-m^2(t) - m^2(s)\}. \end{aligned}$$

Substituting  $t = \theta - \theta_0$  and  $s = \tilde{\theta} - \theta_0$  in the numerator in (12) yields

$$\iint \{t + (\theta_0 - \hat{\theta})\} \{s + (\theta_0 - \hat{\theta})\}^\top \hat{\mathbf{R}}(t, s) dt ds. \quad (49)$$

To establish that (49) converges in probability to the numerator in (9), it suffices by the consistency of  $\hat{\theta}$  to show that (i)  $\iint ts^\top \{\hat{\mathbf{R}}(t, s) - R(t, s)\} dt ds \prec 1$ , (ii)  $\iint t\hat{\mathbf{R}}(t, s) dt ds \preceq 1$ , (iii)  $\iint s^\top \hat{\mathbf{R}}(t, s) dt ds \preceq 1$ , (iv)  $\iint \hat{\mathbf{R}}(t, s) dt ds \preceq 1$ . Since establishing (ii) to (iv) is similar to but easier than (i), we only establish (i) here.

We use the weak version of Glick (1974, theorem A).<sup>14</sup> Pointwise convergence in probability of  $\hat{\mathbf{R}}$  to  $R$  follows from lemmas B7 and D4 and Slutsky.

We now only need to find a convergent upper bound to  $\|t\| \|s\| |\hat{\mathbf{R}}(t, s)|$  whose limit is integrable. Because  $\sup_m \{|m| \exp(-|m|)\} = 1/e$ ,  $\hat{\mathbf{H}}$  is bounded, it suffices by lemma D2 to show that for any fixed order polynomial  $P$ ,  $\int |P(t)| \exp\{\hat{\mathbf{m}}_n(t) - \hat{\mathbf{m}}_n^2(t)\} dt$  is convergent, which follows from lemma B9.  $\square$

**Proof of theorem 6.** Let  $u = y - \theta_0$ . It follows from (13) that we need to determine when

$$\partial_u \int t \frac{m(t)}{\mathcal{D}(t)} I\{|u - t| \leq m(t)\} \exp\{-m^2(t)\} dt \quad (50)$$

<sup>14</sup>See the comment on page 67 in Glick (1974).

is uniformly bounded in  $u$ .<sup>15</sup> Let  $t^- = t^-(u), t^+ = t^+(u)$  be such that  $u - t^- - m(t^-) = 0$ ,  $u - t^+ + m(t^+) = 0$  if such  $t^-, t^+$  exist; set  $t^-$  and/or  $t^+$  equal to  $-\infty, \infty$ , respectively if no solution exists. So  $|u - t| \leq m(t) \iff t^-(u) \leq t \leq t^+(u)$ .

From the definition of  $m$  it follows that if solutions for  $t^+, t^-$  exist they solve

$$F(2t^+ - u) - F(u) = 1/2, \quad F(u) - F(2t^- - u) = 1/2.$$

Hence  $t^+ = \infty$  if  $u \geq 0$  and  $t^- = -\infty$  if  $u \leq 0$ . We concentrate on the case  $u < 0$ ; the case  $u > 0$  is similar and the left derivative at  $u = 0$  (if it is finite) equals the limit as  $u \uparrow 0$ . The expression in (50) then equals

$$\partial_u \int_{-\infty}^{t^+(u)} t \frac{m(t)}{\mathcal{D}(t)} \exp\{-m^2(t)\} dt = t^+ \frac{m(t^+)}{\mathcal{D}(t^+)} \exp\{-m^2(t^+)\} \partial_u t^+. \quad (51)$$

By the implicit function theorem  $\partial_u t^+ = \mathcal{D}(t^+)/2f(2t^+ - u)$ , such that twice the RHS in (51) equals

$$t^+ \frac{m(t^+)}{f(2t^+ - u)} \exp\{-m^2(t^+)\} = \frac{t^+(t^+ - u)}{f(2t^+ - u)} \exp\{-(t^+ - u)^2\} = \\ \zeta \sqrt{2\pi} t^+(t^+ - u) \exp\{-(1 - 2/\zeta^2)t^{+2} + 2(1 - 1/\zeta^2)ut^+ - (1 - 1/2\zeta^2)u^2\}.$$

Take the limit as  $u \uparrow 0$  noting that for fixed negative  $u$ ,  $t^+$  is finite and that (tedious but simple derivations yield)  $\lim_{u \uparrow 0} \{ut^+(u)\} = 0$ . □

**Proof of theorem 7.** Note that  $\hat{\mathcal{Q}}(|\mathbf{y}_i - \mathbf{x}_i\theta|)$  is a piecewise linear function of  $\theta$ . Indeed, on each such segment  $[\underline{\theta}, \bar{\theta}]$ ,  $\hat{\mathcal{Q}}(|\mathbf{y}_i - \mathbf{x}_i\theta|) = |\mathbf{y}_j - \mathbf{x}_j\theta|$  for some  $j = 1, \dots, n$ . If  $\underline{\theta}, \bar{\theta}$  are both finite then for  $s = 0, 1$ ,

$$\lim_{\alpha \rightarrow 0} \int_{\underline{\theta}}^{\bar{\theta}} \theta^s \exp\{-\alpha \hat{\mathcal{Q}}(|\mathbf{y}_i - \mathbf{x}_i\theta|^2)\} d\theta = \lim_{\alpha \rightarrow 0} \int_{\underline{\theta}}^{\bar{\theta}} \theta^s \exp\{-\alpha (\mathbf{y}_j - \mathbf{x}_j\theta)^2\} d\theta = (\bar{\theta}^{s+1} - \underline{\theta}^{s+1})/(s+1),$$

which is finite. The limit is hence determined by the terms for which  $\bar{\theta}$  or  $\underline{\theta}$  is infinite. Note that  $\hat{\mathcal{Q}}(|\mathbf{y}_i - \mathbf{x}_i\theta|) = |\mathbf{y}_\mu - \mathbf{x}_\mu\theta|$  for any sufficiently large  $|\theta|$ . Now, for the denominator,

$$\lim_{\alpha \rightarrow 0} \left( \sqrt{\alpha/\pi} \int_{-\infty}^{\bar{\theta}} \exp\{-\alpha (\mathbf{y}_\mu - \mathbf{x}_\mu\theta)^2\} d\theta \right) = 1/2x_\mu. \quad (52)$$

<sup>15</sup>There is a slight abuse of notation here since the left and right derivatives at  $u = 0$  can differ. If that is the case, let for the remainder of this proof  $\partial_u$  denote the greater in absolute value of the left and right derivatives.

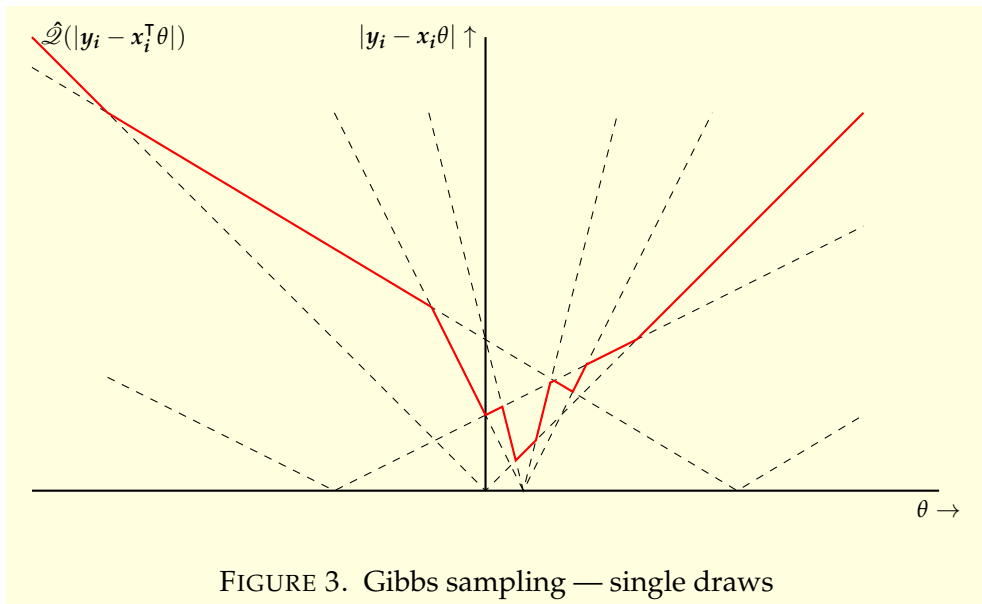


FIGURE 3. Gibbs sampling — single draws

By symmetry, the same limit applies to  $\int_{\underline{\theta}}^{\infty}$ , so the denominator converges to  $1/x_\mu$ . For the numerator and arbitrary  $\underline{\theta}, \bar{\theta}$  by substitution of  $t = \sqrt{2\alpha}(x_\mu\theta - y_\mu)$ ,

$$\begin{aligned} \sqrt{\frac{\alpha}{\pi}} \int_{\underline{\theta}}^{\bar{\theta}} \theta \exp\{-\alpha(y_\mu - x_\mu\theta)^2\} d\theta &= \frac{1}{x_\mu^2} \int_{\sqrt{2\alpha}(x_\mu\underline{\theta} - y_\mu)}^{\sqrt{2\alpha}(x_\mu\bar{\theta} - y_\mu)} \left(\frac{t}{\sqrt{2\alpha}} + y_\mu\right) \phi(t) dt \\ &= \frac{1}{\sqrt{2\alpha}x_\mu^2} [\phi\{\sqrt{2\alpha}(x_\mu\bar{\theta} - y_\mu)\} - \phi\{\sqrt{2\alpha}(x_\mu\underline{\theta} - y_\mu)\}] \\ &\quad + \frac{y_\mu}{x_\mu^2} [\Phi\{\sqrt{2\alpha}(x_\mu\bar{\theta} - y_\mu)\} - \Phi\{\sqrt{2\alpha}(x_\mu\underline{\theta} - y_\mu)\}]. \end{aligned}$$

Hence adding the terms for  $\underline{\theta} = -\infty$  with arbitrary  $\bar{\theta}$  and  $\bar{\theta} = \infty$  with arbitrary  $\underline{\theta}$  and taking  $\alpha \rightarrow 0$ , we obtain

$$\frac{y_\mu}{x_\mu^2} + \lim_{\alpha \rightarrow 0} \frac{\phi\{\sqrt{2\alpha}(x_\mu\bar{\theta} - y_\mu)\} - \phi\{\sqrt{2\alpha}(x_\mu\underline{\theta} - y_\mu)\}}{\sqrt{2\alpha}} = \frac{y_\mu}{x_\mu^2}, \quad (53)$$

since  $\phi$  has derivative zero at zero. Hence  $\lim_{\alpha \rightarrow 0} \hat{\theta}_\alpha = (y_\mu/x_\mu^2)/(1/x_\mu) = y_\mu/x_\mu$ .  $\square$

#### APPENDIX F. COMPUTATION

The method we describe here is Gibbs sampling (Geman and Geman, 1984); for other possibilities see section 7. Because  $\hat{\theta}$  can be thought of as the mean of a distribution with density function  $\psi(\theta) \propto \exp\{-\mathcal{Q}(|y_i - x_i^T\theta|^2)\}$ , all we need is a method to draw random numbers from that distribution. The idea is to draw random numbers from the conditional distribution of each element  $\theta^*$  of  $\theta$  given the remaining elements.



Since the linear combination of remaining regressors and corresponding coefficients can be absorbed into  $\mathbf{y}_i$ , the remainder of our discussion presumes  $d = 1$  and  $n$  odd.<sup>16</sup> To simplify the discussion further, we will presume that all regressors are positive-valued and that there are no two regressors taking the same value. Zeroes and ties require minor adjustments to the procedure and negative values can be accommodated by replacing  $(\mathbf{x}_i, \mathbf{y}_i)$  with  $(-\mathbf{x}_i, -\mathbf{y}_i)$ .

Figure 3 represents the way  $\hat{\mathcal{Q}}(|\mathbf{y}_i - \mathbf{x}_i\theta|)$  varies with the value of  $\theta$ . Since  $\hat{\mathcal{Q}}(|\mathbf{y}_i - \mathbf{x}_i\theta|)$  is piecewise linear in  $\theta$ ,  $\psi$  is a different normal density in each of a number of intervals  $(\theta_{(j)}, \theta_{(j+1)})$ ,  $j = 0, \dots, J$ . It can be shown that  $J \leq 2n + 1$ . Thus, since for some observation  $i_j$ ,  $\hat{\mathcal{Q}}(|\mathbf{y}_i - \mathbf{x}_i\theta|) = |\mathbf{y}_{i_j} - \mathbf{x}_{i_j}\theta|$  for all  $\theta \in (\theta_{(j)}, \theta_{(j+1)})$  we get for  $\Phi_j(\theta) = \Phi\{\sqrt{2}(\mathbf{x}_{i_j}\theta_{(j+1)} - \mathbf{y}_{i_j})\}$  that

$$\Psi(\theta) \propto \sum_{j=0}^J \left\{ I(\theta_{(j+1)} \leq \theta) \left( \frac{\Phi_j(\theta_{(j+1)}) - \Phi_j(\theta_{(j)})}{\mathbf{x}_{i_j}} \right) + I(\theta_{(j)} \leq \theta < \theta_{(j+1)}) \left( \frac{\Phi_j(\theta) - \Phi_j(\theta_{(j)})}{\mathbf{x}_{i_j}} \right) \right\}.$$

So one only needs to find the boundary points  $\theta_{(j)}$  and corresponding observation  $i_j$  in order to compute  $\Psi$ . Once the boundaries  $\theta_{(j)}$  are known,  $\Psi^{-1}(\zeta)$  can be computed for any value  $\zeta \in (0, 1)$  by identifying the value of  $j$  for which  $\Psi(\theta_{(j)}) \leq \zeta < \Psi(\theta_{(j+1)})$  and then computing the inverse of a univariate normal distribution function.

What remains to be done, therefore, is to find the corner points  $\theta_{(j)}$ . The simplest way of achieving this is to compute all intersection points by brute force, which takes  $O(n^2)$  operations. Noting that both  $i_0$  and  $i_j$  equal the index of the median of  $\mathbf{x}_i$ , one can start on the left and look at all intersection points of the downward sloping line corresponding to  $i_0$  with the remaining observations. The observation whose line intersects first will be  $i_1$ , and so forth. Such a brute force approach suffices for most applications.

For the simulations in section 6 we used a more complicated algorithm, which is significantly faster for large  $n$ . A C program is available from the authors.

---

<sup>16</sup>To accommodate even  $n$ , just add an observation with both  $\mathbf{x}_i, \mathbf{y}_i$  set to zero.