

## ***M*-ESTIMATION FOR DISCRETE DATA: ASYMPTOTIC DISTRIBUTION THEORY AND IMPLICATIONS**

BY DOUGLAS G. SIMPSON,<sup>1</sup> RAYMOND J. CARROLL<sup>2</sup> AND  
DAVID RUPPERT<sup>1</sup>

*University of Illinois, University of North Carolina at Chapel Hill  
and University of North Carolina at Chapel Hill*

The asymptotic distribution of an  $M$ -estimator is studied when the underlying distribution is discrete. Asymptotic normality is shown to hold quite generally within the assumed parametric family. When the specification of the model is inexact, however, it is demonstrated that an  $M$ -estimator whose score function is not everywhere differentiable, e.g., a Huber estimator, has a nonnormal limiting distribution at certain distributions, resulting in unstable inference in the neighborhood of such distributions. Consequently, smooth score functions are proposed for discrete data.

**1. Introduction.**  $M$ -estimation, originally proposed by Huber (1964) to estimate a location parameter robustly, has since been applied successfully to a variety of estimation problems where stability of the estimates is a concern. There is, for instance, a substantial body of literature on  $M$ -estimation for regression models; see Krasker and Welsch (1982) for a recent review. For further references on  $M$ -estimation, see Huber (1981).

Surprisingly,  $M$ -estimation for discrete data seems to have received little attention. Discrete data arise naturally in a wide variety of applications. For instance, there is a vast literature on the analysis of multinomial data, and Poisson regression models are commonly used in the analysis of medical and epidemiologic studies [see, e.g., Frome (1983)]. Discrete data are no less prone than continuous measurements to outliers or partial deviations from an otherwise reasonable model, as evidenced by data from mutation research presented in Simpson (1987). This paper investigates some aspects of  $M$ -estimation for discrete data.

A useful optimality theory has been developed by Hampel (1968, 1974) for robust  $M$ -estimation of a univariate parameter. His general prescription facilitates the construction of robust  $M$ -estimators with nearly optimum efficiency at a specified model. Proposals for robust estimation of the binomial and Poisson parameters, for instance, can be found in Hampel (1968). Hampel's univariate theory is briefly reviewed in Section 2. Extensions of this optimality theory to certain multivariate models are discussed in Krasker (1980), Krasker and Welsch (1982), Ruppert (1985) and Stefanski, Carroll and Ruppert (1986).

---

Received November 1985; revised July 1986.

<sup>1</sup>Research supported in part by NSF Grant DMS-8400602.

<sup>2</sup>Research supported by the AFOSR Contract No. F49620-85-C-0144.

AMS 1980 subject classifications. 62E20, 62F10, 62G35.

*Key words and phrases.* Robust estimation,  $M$ -estimator, discrete parametric model, smooth score function.

The score function for Hampel's optimal  $M$ -estimator is not smooth, that is, it is not everywhere differentiable. This can lead to complications in the asymptotic theory when the data are discrete. For instance, Huber (1981, page 51) considers the case where the underlying distribution is a mixture of a smooth distribution and a point mass. He observes that if the point mass is at a discontinuity of the derivative of the score function, then an  $M$ -estimate for location has a nonnormal limiting distribution. Along the same lines, Hampel (1968, page 97) notes that the optimal  $M$ -estimate for the Poisson parameter is asymptotically normal at the Poisson distribution, provided the truncation points of the score function are not integers. He conjectures that "under *any* Poisson distribution, it is asymptotically normal (with the usual variance); however, this remains to be seen."

This paper provides extensions to the asymptotic distribution theory of  $M$ -estimators especially relevant to discrete data, although Theorem 1 is somewhat broader in scope. The main results are given in Section 3 and build on results of Huber (1967). Among the applications of the theory are a more complete account of the asymptotics of the Huber  $M$ -estimate for location and a proof of Hampel's conjecture. For a related work with a different emphasis see Pollard (1985). Aside from providing a more complete asymptotic theory for  $M$ -estimation, the results have implications for choosing a score function when the data are discrete. These are discussed in the final sections. In particular, smooth score functions are proposed.

## 2. Parametric $M$ -estimation: definitions, optimality and examples.

Suppose  $X_1, X_2, \dots$  are independent observations, each thought to have distribution function (d.f.)  $F_\theta$ , where  $\theta$  belongs to a parameter set  $\Theta$ ; here  $\Theta$  is a subset of  $R^d$ ,  $d \geq 1$ . Define

$$(2.1) \quad M(t; \psi, F) = \int \psi(\cdot, t) dF,$$

where  $F$  is a d.f. on  $R^1$ ,  $\psi(\cdot, \cdot)$  is a measurable  $R^d$ -valued function on  $R^1 \times \Theta$ , and  $t \in \Theta$ . Then  $T_n$  is an  $M$ -estimator for  $\theta$ , based on a sample of size  $n$ , if it solves an equation of the form

$$(2.2) \quad M(T_n; \psi, F_n) = 0,$$

where  $F_n$  is the empirical d.f. The standard requirement

$$(2.3) \quad M(\theta; \psi, F_\theta) = 0, \quad \theta \in \Theta,$$

and additional regularity conditions ensure that  $T_n$  consistently estimates  $\theta$  when the model is correct.

Suppose now that  $\Theta \subset R^1$ . The influence function at  $F_\theta$  of an  $M$ -estimator for  $\theta$  has the form

$$\Omega(x, \theta) = \frac{\psi(x, \theta)}{-\int \{(d/d\theta)\psi(\cdot, \theta)\} dF_\theta},$$

provided this exists. Assume  $F_\theta$  has a density  $f_\theta$  with respect to a suitable

measure, and assume the parameterization is smooth. Letting  $l(x, \theta) = (d/d\theta)\log f_\theta(x)$ , the optimal score according to Hampel's criterion has the form

$$(2.4) \quad \psi_{c(\theta)}(l(x, \theta) - \alpha(\theta)),$$

where

$$\psi_c(u) = \begin{cases} u, & |u| \leq c, \\ c \operatorname{sign}(u), & |u| > c, \end{cases}$$

and  $\alpha$  is defined implicitly by (2.3). This estimator cannot be dominated by any  $M$ -estimator simultaneously with respect to the asymptotic variance and the bound on the influence function at  $F_\theta$ . This is assuming, of course, that the estimator is asymptotically normal at  $F_\theta$ .

The truncation point  $c(\theta)$  determines the bounds on  $\Omega(\cdot, \theta)$  and hence the robustness of the estimator to outlying data points. Observe that the maximum likelihood estimator has the form (2.4) with  $c(\theta) \equiv \infty$  and  $\alpha(\theta) \equiv 0$ .

Two examples given in Hampel (1968) will be of special interest here.

**EXAMPLE 1.** If  $F_\theta$  is the normal d.f. with mean  $\theta$  and unit variance, then  $l(x, \theta) = x - \theta$ . By symmetry  $\alpha(\theta) \equiv 0$ , and constant variance suggests setting  $c(\theta) \equiv c$ . The resulting estimator, with score  $\psi_c(x - \theta)$ , is the Huber (1964)  $M$ -estimator for location.

**EXAMPLE 2.** If  $F_\theta$  is the Poisson d.f., with density  $f_\theta(x) = e^{-\theta}\theta^x/x!$  on  $x = 0, 1, 2, \dots$ , then  $l(x, \theta) = x\theta^{-1} - 1$ . Hampel (1968, page 96) suggests taking  $c(\theta) = c\theta^{-1/2}$  on the grounds that  $l(x, \theta)$  has standard deviation  $\theta^{-1/2}$ . For this choice (2.4) is equivalent to  $\psi_c(x\theta^{-1/2} - \theta^{1/2} - \alpha(\theta))$ . The version

$$(2.5) \quad \psi_c(x\theta^{-1/2} - \beta(\theta)),$$

where  $\beta(\theta) = \theta^{1/2} + \alpha(\theta)$  is defined by (2.3), is slightly more convenient.

Figure 1 shows the influence functions associated with (2.5) for  $\theta = 1$  and  $\theta = 5$ . Observe that the influence function is centered at  $\theta^{1/2}\beta(\theta)$  rather than  $\theta$  to compensate for the bias introduced by truncating. For small  $\theta$  the score function only truncates on the right.

**3. Extended asymptotic distribution theory.** Conditions for consistency of an  $M$ -estimator can be found in Huber (1964, 1967, 1981). Since the smoothness plays no role in the consistency proofs, consistency will usually be assumed here.

Huber (1981, Theorems 3.2.4 and 6.3.1) shows under quite general conditions that if  $T_n \rightarrow \theta = T(G)$  in probability as  $n \rightarrow \infty$ , then

$$(3.1) \quad -n^{1/2}M(T_n; \psi, G) = n^{-1/2} \sum_{i=1}^n \psi(X_i, \theta) + o_p(1),$$

where  $M$  is given by (2.1). In particular,  $\psi$  need not be differentiable; monotonicity or Lipschitz continuity conditions are sufficient. Furthermore, Boos and Serfling (1980, Theorem 2.2) show that continuity of  $\psi(\cdot, \theta)$  in total variation

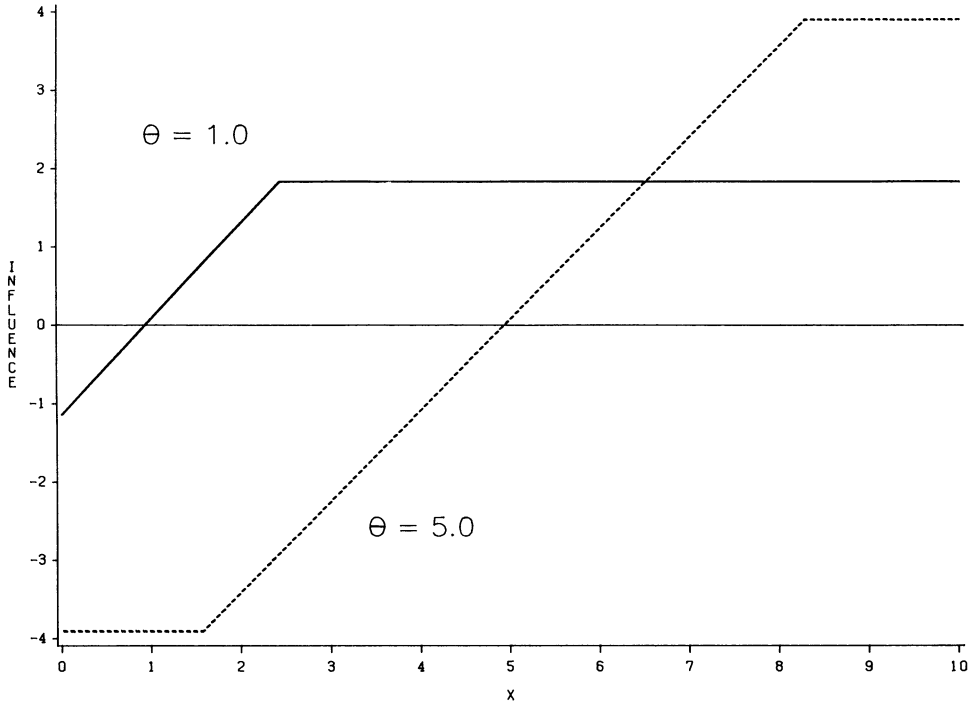


FIG. 1. The influence function of the M-estimator of the Poisson mean ( $c = 1.5$ ).

ensures that (3.1) holds, provided  $T_n \rightarrow \theta$  in probability. Continuity in total variation means

$$\lim_{t \rightarrow \theta} \|\psi(\cdot, t) - \psi(\cdot, \theta)\|_v = 0,$$

where

$$\|h\|_v = \limsup \sum_{i=1}^k |h(x_i) - h(x_{i-1})|,$$

where the supremum is over partitions  $a = x_0 < x_1 < \dots < x_k = b$  of  $[a, b]$ , and the limit is as  $a \rightarrow -\infty, b \rightarrow \infty$ . Note that the score functions of Examples 1 and 2 are continuous in total variation. For stronger almost sure representations for  $T_n$  under stronger conditions, see Carroll (1978a, 1978b).

That  $T_n$  is asymptotically normal follows immediately from (3.1) provided  $M(t; \psi, G)$  has a nonsingular derivative at  $\theta$  and  $\int \psi(\cdot, \theta) \psi(\cdot, \theta)^T dG$  is positive definite and finite. When the underlying distribution is discrete, the set of points where  $\psi$  fails to have a derivative can have positive probability for certain parameter values. In light of (3.1), it is natural to ask whether  $M$  can have a derivative at such parameter values, i.e., whether  $T_n$  can be asymptotically normal.

The following theorem addresses this question. For  $\theta \in \Theta \subset R^d$ ,  $F_\theta$  is assumed to have a density  $f_\theta = f(\cdot, \theta)$  with respect to a  $\sigma$ -finite measure  $\mu$ , and  $\psi_\theta = \psi(\cdot, \theta)$  is measurable for each  $\theta$ . Let  $\|\cdot\|$  denote sup-norm on  $R^d$ . Some regularity conditions are needed:

(A.1) There are measurable functions  $\omega_t = \omega(\cdot, t)$  and  $g_t = g(\cdot, t)$  for which  $\int \omega_t f_t d\mu$ ,  $\int \|\psi_t\| g_t d\mu$  and  $\int \omega_t g_t d\mu$  are finite and, for some  $\delta > 0$ ,

- (i)  $|f_s - f_t| \leq \|s - t\| g_t$ , and
- (ii)  $\|\psi_s\| \leq \omega_t$  almost everywhere  $[\mu]$  (a.e.) when  $\|s - t\| \leq \delta$ .

(A.2) There is a measurable  $R^d$ -valued function  $\dot{f}_t = \dot{f}(\cdot, t)$  such that

$$|f_s - f_t - \dot{f}_t^T(s - t)| = o(\|s - t\|) \quad \text{a.e.}$$

(A.3)  $\psi_s \rightarrow \psi_t$  a.e. as  $s \rightarrow t$ .

**THEOREM 1.** *If for each  $t \in \Theta$  (A.1)–(A.3) hold and*

$$(3.2) \quad M(t; F_t) = 0,$$

*then*

$$(3.3) \quad D_s M(s; F_t)|_{s=t} = - \int \psi_t \dot{f}_t^T d\mu,$$

*where  $D_s$  denotes vector differentiation, and where the dependence of  $M$  on  $\psi$  has been suppressed.*

**PROOF.** For  $s, t \in \Theta$

$$(3.4) \quad M(s; F_t) - M(t; F_t) = M(t; F_t) - M(t; F_s) - R_t(s),$$

where

$$R_t(s) = \int (\psi_s - \psi_t)(f_s - f_t) d\mu$$

and (3.2) was used. The integrand of  $R_t(s)$  is dominated in absolute value by  $2\|s - t\|\omega_t g_t$  on  $\|s - t\| \leq \delta$  because of (A.1). Hence, by (A.3) and dominated convergence

$$(3.5) \quad R_t(s) = o(\|s - t\|).$$

Similarly, (A.2) and dominated convergence imply

$$(3.6) \quad \begin{aligned} & \left| M(t; F_s) - M(t; F_t) - \int \psi_t \dot{f}_t^T d\mu(s - t) \right| \\ & \leq \int \|\psi_t\| |f_s - f_t - \dot{f}_t^T(s - t)| d\mu \\ & = o(\|s - t\|) \quad \text{as } s \rightarrow t, \end{aligned}$$

since the integrand is dominated by  $2\|s - t\|\|\psi_t\|g_t$  on  $\|s - t\| \leq \delta$ . From (3.4) to (3.6) conclude

$$\left| M(s; F_t) - M(t; F_t) + \int \psi_t \dot{f}_t^T d\mu(s - t) \right| = o(\|s - t\|).$$

Hence  $D_s M(s; F_t)$  exists at  $t$  and is given by (3.3).  $\square$

REMARKS. 1. Note that  $\psi_t$  need not be differentiable.

2. When  $\psi_t = I_t = \dot{f}_t/f_t$ , (3.3) generalizes the usual information identity.

3. Huber (1981, page 51) observes a special case, namely (3.3) holds when  $\mu$  is Lebesgue measure,  $\psi(x, t) = \psi(x - t)$ , where  $\psi(\cdot)$  is skew-symmetric about zero, and  $f(x, t) = f(x - t)$ , where  $f(\cdot)$  is differentiable and symmetric about zero.

4. Equation (3.3), when it holds, also guarantees that the influence function at the model, given by

$$\{D_s M(s; F_t)|_{s=t}\}^{-1} \psi(x, t),$$

is defined for each  $t \in \Theta$ , provided that  $\int \psi_t \dot{f}_t d\mu \neq 0$ .

EXAMPLE 2 (continued). Suppose  $f(x, t) = e^{-tx}/x!$  on  $\{0, 1, 2, \dots\}$ ,  $t > 0$ . Recall that the optimal  $M$ -estimator has the score  $\psi(x, t) = \psi_c(xt^{-1/2} - \beta)$ . Theorem 2 of Huber (1967) implies the consistency of this estimator. Results of Huber (1967) also imply that this estimator is asymptotically normal at the Poisson distribution when  $t$  is in one of the open intervals where neither of the truncation points  $t^{1/2}(\beta \pm c)$  is an integer.

To show that it is asymptotically normal at every Poisson distribution, as conjectured by Hampel, use Theorem 1 with  $g(x, t) = e^{2\delta} f(x - 1, t + \delta) + \delta^{-1}(e^\delta - 1 - \delta)f(x, t)$ ,  $\omega(x, t) \equiv c$  and  $\dot{f}(x, t) = f(x - 1, t) - f(x, t)$ . Note that  $c \geq 1$  is sufficient for  $\beta$  to be continuous, and hence for (A.3). Since (3.1) holds and  $0 < \int \psi_t^2 \dot{f}_t d\mu \leq c^2$  for  $c \geq 1$ , it follows that the estimator is asymptotically normal at every Poisson distribution.

If the specification of the model is inexact (as is often suspected), no result like (3.3) is available. In certain cases, it is still possible to obtain the limiting distribution  $T_n$  from (3.1).

Assume for simplicity that  $\Theta$  is an open subset of the real line. The score functions used for robust estimation are generally at least piecewise differentiable. The one-sided derivatives of  $M(t; G)$  will then exist, in general, even when  $M$  fails to be differentiable. Write

$$m(t; G) = D_t M(t; G),$$

when the derivative exists. By a well-known result from calculus, if  $m(\theta - ; G)$  and  $m(\theta + ; G)$  exist, they are equal to the corresponding one-sided derivatives of  $M(t; G)$  at  $\theta$ ; see, e.g., Franklin (1940, page 118).

THEOREM 2. Suppose for some  $\theta$  interior to  $\Theta$  that  $M(\theta; G) = 0$ , and let  $T_n$  be a zero of  $M(t; F_n)$ ,  $n = 1, 2, \dots$ , where  $F_n$  is the empirical d.f. Assume the following:

(B.1)  $m(\theta - ; G)$  and  $m(\theta + ; G)$  exist finitely and are nonzero and of the same sign;

(B.2)  $0 < \sigma < \infty$ , where  $\sigma^2 = \int \psi_\theta^2 dG$ ;

(B.3)  $T_n \rightarrow \theta$  in probability as  $n \rightarrow \infty$ , and (3.1) holds.

Then

$$(3.7) \quad \lim_{n \rightarrow \infty} \sup_{-\infty < z < \infty} |\text{pr}\{n^{1/2}(T_n - \theta) \leq z\} - H(z)| = 0,$$

where

$$H(z) = \begin{cases} \Phi(|m(\theta + ; G)|z/\sigma), & z \geq 0, \\ \Phi(|m(\theta - ; G)|z/\sigma), & z \leq 0, \end{cases}$$

and  $\Phi$  is the standard normal d.f.

REMARKS. 1. Theorem 2 is interesting mainly when  $m(\theta - ; G) \neq m(\theta + ; G)$  and the limiting distribution is nonnormal. It is likely that this result was known previously, but it does not seem to have been published; it is included here since it will be needed later. For similar results for location estimators see Huber (1964, page 78; 1981, page 51) and Pollard (1985, Example 5).

2. The requirement that  $m(\theta \pm ; G)$  have the same sign is actually implied by the remaining conditions. If the one-sided derivatives were to have opposite signs,  $M(t; G)$  would not change signs in a neighborhood of  $\theta$  and (3.1) would not hold.

The proof of Theorem 2 is deferred to Section 7.

EXAMPLE 1 (continued). Recall that the Huber  $M$ -estimator for location has the score  $\psi(x, t) = \psi_c(x - t)$ . For any d.f.  $G$ ,  $M(-\infty; G) = c = -M(\infty; G)$ , and  $M(t; G)$  is continuous in  $t$  so it has a zero  $\theta$ . Assume  $\theta = 0$ . This is unique if  $G(c - ) > G(-c +)$ , in which case  $T_n \rightarrow 0$  in probability by Proposition 2.2.1 of Huber (1981). Since  $\psi_c$  is continuous in total variation, (3.1) holds. Letting  $\dot{\psi}(x, t) = D_t \psi_c(x - t) = -\psi'_c(x - t)$  if it exists, observe that  $-\dot{\psi}(x, t - ) = I(-c \leq x - t < c)$  and  $-\dot{\psi}(x, t + ) = I(-c < x - t \leq c)$ , where  $I(\cdot)$  denotes the indicator function. Bounded convergence yields  $-m(0 - ; G) = G(c - ) - G(-c - )$  and  $-m(0 + ; G) = G(c + ) - G(-c +)$ . Hence, by Theorem 2,  $n^{1/2}T_n$  is asymptotically normal if  $G(c + ) - G(c - ) = G(-c + ) - G(-c - )$ ; otherwise, it has a limiting distribution consisting of the left and right halves of two normal distributions with different variances. Pollard (1985) has previously derived the limiting distribution for this special case.

4. **A counterexample.** It is instructive to examine the extent of the nonnormality that occurs in a specific example. Consider again the optimal  $M$ -estimator for the Poisson parameter. The score function is

$$\psi(x, t) = \psi_c(xt^{-1/2} - \beta) = \begin{cases} -c, & x \leq l(t), \\ xt^{-1/2} - \beta, & l(t) < x < h(t), \\ c, & h(t) \leq x, \end{cases}$$

where  $l(t) = t^{1/2}(\beta(t) - c)$  and  $h(t) = t^{1/2}(\beta(t) + c)$  (cf. Figure 1).

Let  $G$  be the actual d.f. and let  $\theta = T(G)$ . The simplest situation is when  $\theta$  is small. Assume henceforth that  $l(\theta) < 0 < h(\theta) = 1$ . Calculation yields

$$\beta(t) = c(e^t - 1) \quad \text{for } l(t) < 0, 0 < h(t) \leq 1,$$

and

$$\beta(t) = c\{e^t(1 + t)^{-1} - 1\} + t^{1/2}(1 + t)^{-1} \quad \text{for } l(t) < 0, 1 \leq h(t) \leq 2.$$

Since  $\beta$  is continuous, equating the two expressions at  $\theta$  gives

$$(4.1) \quad \theta^{1/2}e^\theta = c^{-1}.$$

The one-sided derivatives of  $\beta$  at  $\theta$  are

$$\beta'(\theta -) = ce^\theta \quad \text{and} \quad \beta'(\theta +) = \frac{1}{2}ce^\theta(1 + \theta)^{-2},$$

where (4.1) was used. Note that  $\beta$  is strictly increasing at  $\theta$ . Since  $\psi'_c(c -) = 1$  and  $\psi'_c(c +) = 0$ ,

$$(4.2) \quad -\dot{\psi}(x, \theta -) = \begin{cases} ce^\theta, & x = 0, \\ 0, & x = 1, 2, \dots \end{cases}$$

and

$$(4.3) \quad -\dot{\psi}(x, \theta +) = \begin{cases} \frac{1}{2}ce^\theta(1 + \theta)^{-2}, & x = 0, \\ \frac{1}{2}ce^\theta\{\theta^{-1} + (1 + \theta)^{-1}\}, & x = 1, \\ 0, & x = 2, 3, \dots \end{cases}$$

Suppose  $G$  is a mixture of a Poisson distribution  $F_t$  and a point mass at an integer  $z$ , i.e.,  $G = (1 - \epsilon)F_t + \epsilon\delta_z$ . Assume  $z > h(t)$  so  $\dot{\psi}(z, \theta \pm) = 0$ . From (4.2) and (4.3)

$$(4.4) \quad \frac{m(\theta + ; G)}{m(\theta - ; G)} = \frac{1}{2} \left( \frac{t}{\theta} + \frac{1 + t}{1 + \theta} \right),$$

where  $m(\theta - ; G) = -ce^{\theta-t}(1 - \epsilon)$ . The ratio (4.4) is unity only when  $t = \theta$ , which corresponds to  $\epsilon = 0$ . By Theorem 2, the limiting distribution of  $n^{1/2}(T_n - \theta)$  consists of the right and left halves of two normal distributions. The ratio of their standard deviations is (4.4).

Solving  $0 = M(\theta; G) = c\{1 - (1 - \epsilon)e^{\theta-t}\}$  yields  $t = \theta + \log(1 - \epsilon)$ . Table 1 shows the values of  $t$  and (4.4) for several values of  $\epsilon$  when  $\theta = 0.25$  and  $c = \theta^{-1/2}e^{-\theta} = 1.5576 \dots$  [see (4.1)]. In addition, the effect on a nominal 0.05 tail probability is shown.

For very small values of  $\epsilon$  the effect is minimal, which is in accord with the robustness of  $T_n$  in the sense of weak\* continuity [see Hampel (1971)], since it is

TABLE 1  
Effect of contaminating mass  $\epsilon$  with  $\theta = 0.25$  fixed

| $\epsilon$ | $t$   | $r = (4.4)$ | $\Phi(-1.645r)$ |
|------------|-------|-------------|-----------------|
| 0          | 0.25  | 1           | 0.05            |
| 0.01       | 0.24  | 0.976       | 0.054           |
| 0.05       | 0.199 | 0.877       | 0.074           |
| 0.10       | 0.145 | 0.748       | 0.109           |
| 0.15       | 0.087 | 0.610       | 0.158           |
| 0.20       | 0.027 | 0.465       | 0.222           |



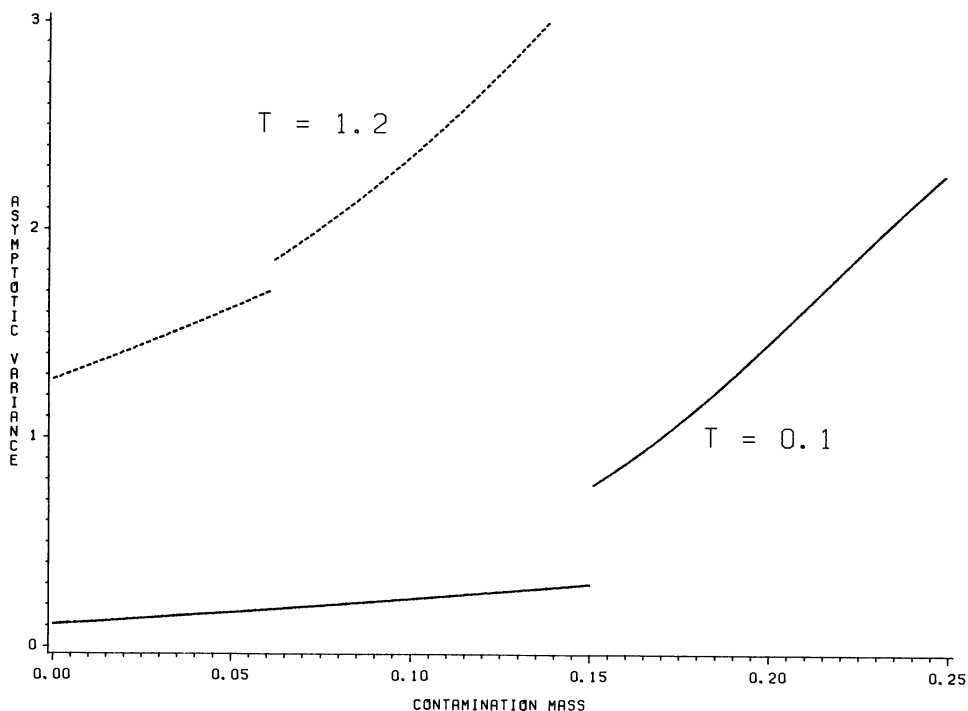


FIG. 2. Asymptotic variance as a function of contamination mass for  $M$ -estimator of Poisson mean ( $c = 1.5$ ).

asymptotically normal at the model. As  $\epsilon$  increases, however, the effect becomes more serious, and inference based on  $T_n$  can be substantially biased.

For a related work see Stigler (1973), who observes that a bias of this type can arise when the trimmed mean is used for discrete or grouped data.

**5. Smooth score functions.** In the example of the preceding section, one might argue that the parameter values where problems arise are unlikely to occur in practice, or that  $c$  can be changed slightly. It is not, however, the nonnormal limiting distribution of  $T_n$  at certain distributions that is of concern, but the instability of inference based on  $T_n$  near those distributions. This phenomenon can alternatively be interpreted as a discontinuity of the asymptotic variance functional

$$V(T(G); G) = \{m(T(G); G)\}^{-1} \int \psi_{T(G)} \psi_{T(G)}^T dG \{m(T(G); G)\}^{-1}$$

[cf. Huber (1981, page 51)]. Figure 2 shows the effect of a contaminating point mass on the asymptotic variance of the estimator of Example 2. Poisson distributions with means 0.1 and 1.2 were contaminated by a point mass located beyond the upper truncation points of the score functions. Clearly, in the

neighborhood of a distribution where  $V$  is discontinuous, inference based on  $T_n$  may be unstable.

Instability of this type can be avoided by requiring the  $M$ -estimator score function to be smooth, for example, by replacing  $\psi_c(\cdot)$  in (2.4) with a smooth approximation. A natural way to construct such a function is by rescaling a smooth distribution function.

Suppose  $F$  is an absolutely continuous d.f. with density  $f$  symmetric about zero. Then

$$(5.1) \quad \psi(x) = 2c \left( F \left( \frac{x}{2cf(0)} \right) - \frac{1}{2} \right)$$

is monotone increasing, skew-symmetric about zero and satisfies  $\psi(\infty) = c$  and  $\psi'(0) = 1$ . Observe that  $\psi_c$  is obtained from (5.1) by taking  $F$  to be the uniform distribution on  $[-\frac{1}{2}, \frac{1}{2}]$ . This can be approximated arbitrarily closely by a symmetric beta distribution with a small value for the shape parameter, i.e.,  $f(x) \propto \{(\frac{1}{2} + x)(\frac{1}{2} - x)\}^\alpha$  on  $[-\frac{1}{2}, \frac{1}{2}]$ . The resulting score function is complicated, however, and its second derivative has jump discontinuities. A more convenient choice is the logistic distribution, which leads to the smooth function

$$L_c(x) = c \tanh(x/c).$$

This has appeared previously.  $L_1(x - t)$  is the maximum likelihood score for the location of a logistic distribution with scale 1. Holland and Welsch (1977) include an  $M$ -estimator using  $L_c$  in a Monte Carlo study of robust regression estimates.

For the important special case of estimating a Poisson parameter robustly, a smooth version of the optimal  $M$ -estimator solves

$$(5.2) \quad n^{-1} \sum_{i=1}^n L_c(X_i t^{-1/2} - \beta(t)) = 0,$$

where  $\beta$  is defined in the usual manner.

Table 2 gives asymptotic variances  $V_\theta$  and bounds  $\gamma_\theta$  on influence functions for the estimator defined by (5.2), labeled  $L_c$ , and the optimal estimator, labeled  $\psi_c$ . In each case  $c = 1.5$ . The calculations are at the Poisson model, and  $V_\theta$  and  $\gamma_\theta$  are stabilized by dividing by  $\theta$  and  $\theta^{1/2}$ , respectively.

Note that  $V_\theta/\theta$  is the asymptotic relative efficiency of the maximum likelihood estimator (sample mean) with respect to the corresponding  $M$ -estimator. The asymptotic variances for the logistic score are slightly smaller than those for the "optimal" score. This is possible because the bounds on the influence function of  $L_c$  are slightly higher for  $\psi_c$ . In terms of performance at the model, there appears to be little difference between  $L_c$  and  $\psi_c$ .

**6. Further remarks.** The need for smooth score functions is most clear when the data consist of counts. In this case every deviation from the model involves point masses.

An important consequence of Theorem 1 is that Hampel's optimal estimator (2.4) is indeed optimal as claimed when the model distribution is discrete. It would be disturbing if the theory were to break down at a countable number of

TABLE 2  
Asymptotic variances and influence function bounds at the Poisson model

| Mean<br>$\theta$ | $\psi_c$          |                              | $L_c$             |                              |
|------------------|-------------------|------------------------------|-------------------|------------------------------|
|                  | $V_\theta/\theta$ | $\gamma_\theta/\theta^{1/2}$ | $V_\theta/\theta$ | $\gamma_\theta/\theta^{1/2}$ |
| 0.1              | 1.052             | 3.16                         | 1.048             | 3.27                         |
| 0.2              | 1.107             | 2.24                         | 1.081             | 2.53                         |
| 0.3              | 1.138             | 1.98                         | 1.094             | 2.29                         |
| 0.4              | 1.114             | 2.00                         | 1.095             | 2.19                         |
| 0.5              | 1.092             | 1.98                         | 1.083             | 2.14                         |
| 1.0              | 1.071             | 1.84                         | 1.059             | 2.07                         |
| 2.0              | 1.057             | 1.74                         | 1.045             | 2.04                         |
| 5.0              | 1.043             | 1.75                         | 1.038             | 2.02                         |
| 10.0             | 1.040             | 1.74                         | 1.035             | 2.02                         |
| 100.0            | 1.037             | 1.73                         | 1.033             | 2.01                         |

parameter values. Moreover, the smooth versions discussed in Section 5, which provide more stable inference, are justified for every parameter value as being nearly optimal.

Although the discussion has focused on the score functions arising from Hampel's optimality theory, it is not limited to that context. For instance, a score based on Hampel's three part redescending  $\psi$  [see Huber (1981, page 102)] will be prone to the same difficulties, and a smooth version will be more stable.

**7. Proof of Theorem 2.** Since the d.f.  $H$  is continuous, uniform convergence in (3.7) will follow from pointwise convergence via Pólya's theorem [Serfling (1980, page 18)].

Write  $M(t)$  for  $M(t; G)$  and  $m(t)$  for  $m(t; G)$ . Denote by  $U(\delta)$  the set  $\{t: 0 < |t - \theta| < \delta\}$ . By (B.1),  $m$  is defined on  $U(\delta)$  if  $\delta$  is sufficiently small. Moreover, given  $\epsilon > 0$ , there is a  $\delta$  for which  $t \in U(\delta)$  implies

$$|m(t) - m(\theta -)| < \epsilon, \text{ if } t < \theta$$

and

$$|m(t) - m(\theta +)| < \epsilon, \text{ if } t > \theta.$$

Choosing  $\epsilon < \min\{|m(\theta -)|, |m(\theta +)|\}$  then guarantees that  $|m(t)|$  is bounded away from zero on  $U(\delta)$ . Fix such a  $\delta$ .

Since  $M(\theta) = 0$ ,  $t \in U(\delta)$  implies

$$(7.1) \quad M(t) = m(\tau)(t - \theta),$$

for some  $\tau$  strictly between  $t$  and  $\theta$ , by the mean value theorem (which only requires one-sided derivatives at the endpoints of the interval on which it is applied). Since  $m$  is bounded away from zero on  $U(\delta)$ , (7.1) shows

$$|t - \theta| = O(|M(t)|),$$

as  $t \rightarrow \theta$ . The right-hand side of (7.1) equals

$$(7.2) \quad D(t)(t - \theta) + R(t),$$

where

$$D(t) = m(\theta +)I(t > \theta) + m(\theta -)I(t < \theta),$$

$$R(t) = [\{m(\tau) - m(\theta +)\}I(t > \theta) + \{m(\tau) - m(\theta -)\}I(t < \theta)](t - \theta),$$

and  $I(A)$  is the indicator for  $A$ . Note that (7.2) also holds if  $t = \theta$ . Since  $R(t) = o(|t - \theta|) = o(|M(t)|)$ , (7.1) and (7.2) yield

$$(7.3) \quad D(T_n)n^{1/2}(T_n - \theta) = n^{1/2}M(T_n) + o(|n^{1/2}M(T_n)|).$$

Because of (B.2), (B.3) and the Lindeberg–Lévy central limit theorem, the right-hand side of (7.3) converges in distribution to a  $N(0, \sigma^2)$  random variable, and, hence, so does the left-hand side.

To obtain the limiting distribution of  $T_n$ , partition its range and consider cases. If  $z < 0$ , then

$$\text{pr}\{n^{1/2}(T_n - \theta) \leq z, T_n > \theta\} = 0,$$

while

$$\text{pr}\{n^{1/2}(T_n - \theta) \leq z, T_n < \theta\} = \text{pr}\{|D(T_n)|n^{1/2}(T_n - \theta) \leq |D(T_n)|z\}.$$

Since  $D(T_n) = m(\theta -)$  when  $T_n < \theta$ , and  $D(t)$  does not change sign on  $(\theta - \delta, \theta + \delta)$  by (B.1), (7.3) implies that this last probability converges to  $\Phi(|m(\theta -)|z/\sigma)$  as  $n \rightarrow \infty$ . Similar arguments establish that, for  $z > 0$ ,

$$\text{pr}\{n^{1/2}(T_n - \theta) \leq z, T_n < \theta\} = \text{pr}\{|m(\theta -)|n^{1/2}(T_n - \theta) < 0\} \rightarrow \frac{1}{2}$$

and

$$\begin{aligned} &\text{pr}\{n^{1/2}(T_n - \theta) \leq z, T_n > \theta\} \\ &= \text{pr}\{0 < |m(\theta +)|n^{1/2}(T_n - \theta) \leq z|m(\theta +)|\} \rightarrow \Phi(|m(\theta +)|z/\sigma) - \frac{1}{2} \end{aligned}$$

and, finally,

$$\text{pr}\{n^{1/2}(T_n - \theta) \leq 0\} = 1 - \text{pr}\{|m(\theta +)|n^{1/2}(T_n - \theta) > 0\} \rightarrow \frac{1}{2},$$

as  $n \rightarrow \infty$ . The result follows by collecting terms.  $\square$

**Acknowledgments.** The authors thank Barry H. Margolin for stimulating this research. The authors also thank the referees for helpful comments, including the reference to Pollard (1985).

### REFERENCES

BOOS, D. D. and SERFLING, R. J. (1980). A note on differentials and the CLT and LIL for statistical functions, with applications to  $M$ -estimates. *Ann. Statist.* **8** 618–624.  
 CARROLL, R. J. (1978a). On almost sure expansions for  $M$ -estimates. *Ann. Statist.* **6** 314–318.  
 CARROLL, R. J. (1978b). On the asymptotic distribution of multivariate  $M$ -estimates. *J. Multivariate Anal.* **8** 361–371.  
 FRANKLIN, P. (1940). *A Treatise on Advanced Calculus*. Wiley, New York.  
 FROME, E. L. (1983). The analysis of rates using Poisson regression models. *Biometrics* **39** 665–674.  
 HAMPEL, F. (1968). Contributions to the theory of robust estimation. Ph.D. thesis, Univ. California, Berkeley.

- HAMPEL, F. (1971). A general qualitative definition of robustness. *Ann. Math. Statist.* **42** 1887–1896.
- HAMPEL, F. (1974). The influence curve and its role in robust estimation. *J. Amer. Statist. Assoc.* **69** 383–393.
- HOLLAND, P. W. and WELSCH, R. E. (1977). Robust regression using iterativity reweighted least-squares. *Comm. Statist. A—Theory Methods* **6** 813–827.
- HUBER, P. J. (1964). Robust estimation of a location parameter. *Ann. Math. Statist.* **35** 73–101.
- HUBER, P. J. (1967). The behavior of maximum likelihood estimates under nonstandard conditions. *Proc. Fifth Berkeley Symp. Math. Statist. Probab.* **1** 221–233. Univ. California Press.
- HUBER, P. J. (1981). *Robust Statistics*. Wiley, New York.
- KRASKER, W. S. (1980). Estimation in linear regression models with disparate data points. *Econometrica* **48** 1333–1346.
- KRASKER, W. S. and WELSCH, R. E. (1982). Efficient bounded-influence regression estimation. *J. Amer. Statist. Assoc.* **77** 595–604.
- POLLARD, D. (1985). New ways to prove central limit theorems. *Econometric Theory* **1** 295–314.
- RUPPERT, D. (1985). On the bounded influence regression estimator of Krasker and Welsch. *J. Amer. Statist. Assoc.* **80** 205–208.
- SERFLING, R. J. (1980). *Approximation Theorems of Mathematical Statistics*. Wiley, New York.
- SIMPSON, D. G. (1987). Minimum Hellinger distance estimation for the analysis of count data. *J. Amer. Statist. Assoc.* To appear.
- STEFANSKI, L. A., CARROLL, R. J. and RUPPERT, D. (1986). Optimally bounded score functions for generalized linear models with applications to logistic regression. *Biometrika* **73** 413–424.
- STIGLER, S. M. (1973). The asymptotic distribution of the trimmed mean. *Ann. Statist.* **1** 472–477.

DOUGLAS G. SIMPSON  
DEPARTMENT OF STATISTICS  
UNIVERSITY OF ILLINOIS  
CHAMPAIGN, ILLINOIS 61820

RAYMOND J. CARROLL  
DAVID RUPPERT  
DEPARTMENT OF STATISTICS  
UNIVERSITY OF NORTH CAROLINA  
CHAPEL HILL, NORTH CAROLINA 27514