

SOME PROPERTIES OF THE LEAST SQUARES ESTIMATOR IN REGRESSION ANALYSIS WHEN THE PREDICTOR VARIABLES ARE STOCHASTIC¹

BY P. K. BHATTACHARYA

University of North Carolina

1. Introduction and summary. In the classical linear estimation set-up, we have

$$(1) \quad E\mathbf{y} = X\boldsymbol{\theta},$$

where $\mathbf{y}' = (y_1, \dots, y_n)$ is a random vector whose components are uncorrelated and have equal variance, $\boldsymbol{\theta}' = (\theta_0, \theta_1, \dots, \theta_p)$ is a vector whose elements are unknown constants and X is a matrix of n rows and $p + 1$ columns, $n \geq p + 1$, which has full rank and whose elements are known constants. Plackett [2] gives a historical note on the least squares estimator $\hat{\boldsymbol{\theta}} = (X'X)^{-1}X'\mathbf{y}$ of $\boldsymbol{\theta}$ for which the following property is well-known,

(I) Each component $\hat{\theta}_j$ of $\hat{\boldsymbol{\theta}}$ is the estimator with uniformly minimum variance among all unbiased linear estimators of the corresponding component θ_j of $\boldsymbol{\theta}$.

It can also be easily seen that

(II) For a quadratic loss function for the estimation of each component θ_j of $\boldsymbol{\theta}$, the least squares estimators have uniformly minimum risk among the class of all linear (in y 's) estimators with bounded risk.

Properties (I) and (II) hold for model (1) with uncorrelated and homoscedastic y_1, \dots, y_n . Hodges and Lehmann [1] have shown that if we do not restrict ourselves to estimators which are linear in y 's, then the least squares estimators have the following weaker property:

(III) If the loss in estimating the true vector $\boldsymbol{\theta}$ by another $\boldsymbol{\vartheta}$ is $(\boldsymbol{\theta} - \boldsymbol{\vartheta})'(\boldsymbol{\theta} - \boldsymbol{\vartheta})$, then the least squares estimator is minimax among the class of all estimators of $\boldsymbol{\theta}$ if there exists a number v such that $\text{Var } y_i \leq v, i = 1, \dots, n$, and the family of distributions \mathcal{F} of (y_1, \dots, y_n) contains the sub-family \mathcal{F}_0 of all independent normal distributions of (y_1, \dots, y_n) which satisfy (1) for some $\boldsymbol{\theta}$, and have $\text{Var } y_i = v, i = 1, \dots, n$.

In many situations, however, (y, x_1, \dots, x_p) follows a $(p + 1)$ -variate distribution on which observations are made and the method of least squares is applied to estimate the linear regression of y on x_1, \dots, x_p , regarding the x -observations to be non-stochastic. This problem differs from the classical problem of linear estimation because instead of (1) the model is $E[y | X] = X\boldsymbol{\theta}$, where the elements of the X matrix are stochastic.

For reasons given in Section 2, the loss in estimating the true regression function $\phi(x_1, \dots, x_p)$ by another function $\psi(x_1, \dots, x_p)$ is considered to be of

Received August 9, 1961; revised May 5, 1962.

¹ This research was supported by the Air Force Office of Scientific Research.

the form,

$$\int [\phi(x_1, \dots, x_p) - \psi(x_1, \dots, x_p)]^2 dF(x_1, \dots, x_p),$$

where F is the distribution function of (x_1, \dots, x_p) .

For the above loss function, it is shown under certain conditions that if the class of estimates which are linear in y 's and have bounded risk is non-empty, then the estimate obtained by the method of least squares belongs to this class and has uniformly minimum risk in this class. A necessary and sufficient condition on $F(x_1, \dots, x_p)$ is obtained for this class to be non-empty, which unfortunately is not easy to verify in particular cases. However, by a sequential modification of the sampling scheme, this condition may always be satisfied at the cost of an arbitrarily small increase in the expected sample size. It is also shown under certain further conditions on the family of admissible distributions that the least squares estimator is minimax in the class of all estimators.

For the case of multivariate normal distribution of (y, x_1, \dots, x_p) , Stein [3] has considered this problem under a loss function similar to the one given above. He has shown the minimax property of the least squares estimates (which also happen to be the maximum likelihood estimates in a multivariate normal model) for the regression coefficients, and has raised many interesting questions about the admissibility of these estimates.

2. Formal statement of the problem. Let y, x_1, \dots, x_p be real-valued random variables with joint distribution function G . We assume for simplicity that G is the product of a completely specified distribution F of x_1, \dots, x_p with an unknown conditional distribution of y given x_1, \dots, x_p . Let G belong to a family of distributions \mathcal{G} . We further assume that F and \mathcal{G} satisfy the following conditions:

CONDITION (i). F is such that

$$(a)^2 \text{ for every non-null } (a_0, a_1, \dots, a_p),$$

$$\Pr [a_0 + a_1x_1 + \dots + a_px_p = 0] = 0.$$

$$(b) E(x_j^2) < \infty, j = 1, \dots, p.$$

CONDITION (iia) For every $G \in \mathcal{G}$, $E_G[y | x_1, \dots, x_p]$ is a linear function of x_1, \dots, x_p , say $\phi(x_1, \dots, x_p) = \theta_0 + \theta_1x_1 + \dots + \theta_px_p$, the row vector $\theta' = (\theta_0, \theta_1, \dots, \theta_p)$ depending on G .

CONDITION (iib) The set of all θ corresponding to all $G \in \mathcal{G}$, equals the $(p + 1)$ -dimensional Euclidean space.

CONDITION (iii) There exists a constant σ^2 such that $V_G[y | x_1, \dots, x_p] = \sigma^2$ for all $G \in \mathcal{G}$ and for all (x_1, \dots, x_p) .

In the statement of the above conditions as well as in what follows, $E_G[y | x_1, \dots, x_p]$ and $V_G[y | x_1, \dots, x_p]$ stand for the conditional expectation and the conditional variance respectively, of y given x_1, \dots, x_p under G .

² This condition is satisfied even if x has a continuous distribution and $x_j = x^j$ or if x_1, \dots, x_p are $\sin 2x, \dots, \cos x, \cos 2x, \dots$.

Now suppose $(y_i, x_{1i}, \dots, x_{pi}), i = 1, \dots, n, n \geq p + 1$, are mutually independent observations on (y, x_1, \dots, x_p) . The problem is to estimate the regression function

$$\phi(x_1, \dots, x_p) = \theta_0 + \theta_1 x_1 + \dots + \theta_p x_p,$$

or in other words the row vector $\theta' = (\theta_0, \theta_1, \dots, \theta_p)$, where the loss involved in estimating θ by β is,

$$\begin{aligned} W(\theta, \beta) &= \int [(\theta_0 + \theta_1 x_1 + \dots + \theta_p x_p) \\ &\quad - (\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)]^2 dF(x_1, \dots, x_p) \\ (2) \qquad &= \sum_{j=0}^p \sum_{j'=0}^p \mu_{jj'} (\theta_j - \beta_j) (\theta_{j'} - \beta_{j'}) \\ &= (\theta' - \beta') M (\theta - \beta), \end{aligned}$$

where

$$M = \begin{bmatrix} 1 & \mu_{01} & \dots & \mu_{0p} \\ \mu_{01} & \mu_{11} & \dots & \mu_{1p} \\ & & \dots & \\ \mu_{0p} & \mu_{1p} & \dots & \mu_{pp} \end{bmatrix}, \mu_{jj'} = E(x_j x_{j'}), j, j' = 0, 1, \dots, p, x_0 \equiv 1.$$

It follows from condition (i) that M is positive definite, and $0 < W(\theta, \beta) < \infty$ for all $\beta \neq \theta$.

The loss function (2) is motivated by the following consideration. Suppose we are required to predict the value of y associated with a random observation made on (x_1, \dots, x_p) subject to a quadratic loss. If the true regression function

$$\phi(x_1, \dots, x_p) = \theta_0 + \theta_1 x_1 + \dots + \theta_p x_p$$

were known, our prediction rule would be

$$y(x_1, \dots, x_p) = \theta_0 + \theta_1 x_1 + \dots + \theta_p x_p,$$

and the risk of the procedure would be σ^2 . If however, we use the prediction rule

$$y'(x_1, \dots, x_p) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p,$$

then the risk is $\sigma^2 + W(\theta, \beta)$.

3. Optimum property of the least squares estimator in the class of linear estimators with bounded risk. In this section we shall restrict our attention only to those procedures which are linear in y 's and have bounded risk. Let us denote the class of all such procedures by \mathcal{C}_1 . Then an estimator $t \in \mathcal{C}_1$ if and only if

$$(3) \qquad t_{(p+1) \times 1} = L_{(p+1) \times n} y_{n \times 1}$$

and $\rho(\theta, t) = E[W(\theta, t)]$ is a bounded function of $\theta_0, \theta_1, \dots, \theta_p$, where $y' =$

(y_1, \dots, y_n) , and each element l_{ji} of L is a function of $x_{11}, \dots, x_{1n}, \dots, x_{p1}, \dots, x_{pn}$. Let

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{p1} \\ 1 & x_{12} & \dots & x_{p2} \\ & & \dots & \\ 1 & x_{1n} & \dots & x_{pn} \end{bmatrix}.$$

Then the estimator obtained by the method of least squares is, $t^* = (X'X)^{-1}X'y$. Since $n \geq p + 1$, it follows from condition (ia) that the matrix X has rank $p + 1$ with probability 1 and therefore t^* can be uniquely determined for almost all samples.

In what follows, $E_{y|x}[U(x, y)]$ stands for $E_y[U(x, y) | x]$.

We shall first show that if an estimator $t \in \mathcal{C}_1$, then for almost all X , the conditional expectation of each component of t given X should be equal to the corresponding component of θ , whatever θ may be. For any $t \in \mathcal{C}_1$, it follows from (3) that

$$E_{y|x}(t_j) = \sum_{i=1}^n l_{ji}(\theta_0 + \theta_1 x_{1i} + \dots + \theta_p x_{pi}) = \theta_j + \sum_{j'=0}^p h_{jj'} \theta_{j'} \quad \text{say,}$$

where $h_{jj'}$ are functions of X . Let H be a $(p + 1) \times (p + 1)$ matrix which has $h_{jj'}$ in its j th row and j' th column. To every $t \in \mathcal{C}_1$, there corresponds such a matrix H , and to show this correspondence, we shall use the notation H_t . Thus if $t \in \mathcal{C}_1$, then

$$(4) \quad E_{y|x}(t) - \theta = H_t \theta.$$

We now prove

LEMMA 1. *If $t \in \mathcal{C}_1$, then $E_{y|x}(t) \equiv \theta$ for almost all X .*

PROOF. By virtue of (4) it will be enough to show that if $P[H_t \neq 0] > 0$, then $t \notin \mathcal{C}_1$. Let t satisfy (3); then we shall show that $\rho(\theta, t)$ is unbounded. We have

$$\begin{aligned} \rho(\theta, t) &\geq E_x[(E_{y|x}(t') - \theta')M(E_{y|x}(t) - \theta)] \\ &= E_x[\theta'H_t'MH_t\theta]. \end{aligned}$$

If $P[H_t \neq 0] > 0$, some element of H_t is non-zero with positive probability. Let that element be one in the j_0 th column of H_t . Let

$$A_{j_0} = \{\theta: \theta_j = 0 \text{ for } j \neq j_0 \text{ and } \theta_{j_0} \neq 0\}.$$

Then for $\theta \in A_{j_0}$,

$$\theta'H_t'MH_t\theta = \theta_{j_0}^2 g(X)$$

where $P[g(X) \geq 0] = 1$ and $P[g(X) > 0] > 0$, since M is positive definite. Hence there exists $\delta > 0$ such that $P(\delta) = P[g(X) > \delta] > 0$. Then for $\theta \in A_{j_0}$,

$$\rho(\theta, t) \geq E_x[\theta'H_t'MH_t\theta] \geq \theta_{j_0}^2 \delta P(\delta),$$

and for any given c , $\rho(\theta, t)$ can be made greater than c by choosing θ in A_{j_0} with $|\theta_{j_0}| > [c/\delta P(\delta)]^{\frac{1}{2}}$. This completes the proof.

The following corollary is immediate.

COROLLARY. *If $t \in \mathcal{C}_1$, then*

$$(5) \quad \rho(\theta, t) = \sigma^2 E_X \sum_{j=0}^p \sum_{j'=0}^p \mu_{jj'} \sum_{i=1}^n l_{ji} l_{j'i}$$

where l_{ji} are the elements of the matrix L through which t is defined and they satisfy

$$(6) \quad \sum_{i=1}^n l_{ji} x_{ji} = 1, \quad \sum_{i=1}^n l_{ji} x_{j'i} = 0 \quad \text{for } j \neq j' = 0, 1, \dots, p,$$

for almost all X .

It follows from the above corollary that for $t \in \mathcal{C}_1$, $\rho(\theta, t)$ does not depend on θ . Therefore, if we minimize the right side of (5) with respect to l_{ji} 's subject to the conditions in (6), the resulting matrix $\hat{L} = (X'X)^{-1}X'$ will define an estimator $\hat{t} = \hat{L}y = (X'X)^{-1}X'y$ for which $\rho(\theta, \hat{t}) \leq \rho(\theta, t)$ for all θ and for arbitrary $t \in \mathcal{C}_1$. But $\hat{t} = t^*$, a.e.

It can be easily seen that unless \mathcal{C}_1 is empty, $t^* \in \mathcal{C}_1$. We thus have

THEOREM 1. *If \mathcal{C}_1 is non-empty, then the least squares estimator $t^* \in \mathcal{C}_1$, and $\rho(\theta, t^*) \leq \rho(\theta, t)$ for all θ and for arbitrary $t \in \mathcal{C}_1$ with a strict inequality holding if $P[t = t^*] < 1$.*

If we denote by \mathcal{C}_2 the class of estimators which are linear in y 's, then the following corollary is immediate from the fact that $\rho(\theta, t^*)$ is a constant for all θ .

COROLLARY. *If \mathcal{C}_1 is non-empty, then the least squares estimator t^* is the unique minimax estimator in \mathcal{C}_2 .*

The optimum property of t^* thus depends on the non-emptiness of \mathcal{C}_1 which can be characterized in terms of $F(x_1, \dots, x_p)$. We have seen that \mathcal{C}_1 is non-empty if and only if $\rho(\theta, t^*)$ is finite. Clearly,

$$\rho(\theta, t^*) = \sigma^2 E_X \text{tr} [(X'X)^{-1}M].$$

Hence a necessary and sufficient condition for \mathcal{C}_1 being non-empty is

CONDITION (ic). $E \text{tr} [(X'X)^{-1}M] < \infty$.

This is a condition on the marginal distribution F of (x_1, \dots, x_p) , and it implies condition (ib). Thus under conditions (ia), (ic), (iia), (iib) and (iii), we can state Theorem 1 and its corollary without the qualifying clause "if \mathcal{C}_1 is non-empty".

4. Minimax property of t^* in the class of all estimators. In this section we assume F to satisfy conditions (ia) and (ic) and the family \mathcal{G} to satisfy conditions (iia) and

CONDITION (iv). There exists a number v such that $V_G[y | x_1, \dots, x_p] \leq v$ for all $G \in \mathcal{G}$ and for all x_1, \dots, x_p .

CONDITION (v). \mathcal{G} includes the class \mathcal{G}_0 of all G obtained by taking the product of the distribution F of (x_1, \dots, x_p) with a conditional distribution of y given

x_1, \dots, x_p which is normal with mean satisfying (iia) for some θ and with variance v .

Under these conditions we shall prove that t^* is a minimax estimate for θ . We shall require the following lemma to prove the minimax property.

LEMMA 2. *Let $A(z)$ be a mapping from an arbitrary space Z to the space of all $k \times k$ non-singular matrices such that $\|A(z)\|$ and $\|A(z)^{-1}\|$ are both uniformly bounded. Let U be a fixed $k \times k$ matrix and $\{b_m\}$ a sequence of real numbers converging to zero. Then,*

(a) *Det $[A(z) + b_m U]$ converges to Det $A(z)$ uniformly in z .*

(b) *There exists an integer m_0 such that $[A(z) + b_m U]^{-1}$ exists for $m \geq m_0$ and for all $z \in Z$.*

(c) *$\|[A(z) + b_m U]^{-1} - A(z)^{-1}\|$ converges to zero uniformly in z , where $[A(z) + b_m U]^{-1}$ is defined arbitrarily when $\text{Det } [A(z) + b_m U] = 0$.*

PROOF. Suppose $\max [\text{Sup}_{z \in Z} \|A(z)\|, \text{Sup}_{z \in Z} \|A(z)^{-1}\|, \|U\|] = c$. Then (a) follows from the fact that

$$|\text{Det } [A(z) + b_m U] - \text{Det } A(z)| \leq |b_m| c^k (2^k - 1) \cdot k!, \quad \text{if } |b_m| \leq 1.$$

Since $|\text{Det } A(z)^{-1}| \leq c^k \cdot k!$, $|\text{Det } A(z)| \geq 1/c^k \cdot k!$. Also, it follows from (a) that there exists an integer m_0 such that for $m \geq m_0$

$$\text{Det } A(z) - 1/2c^k \cdot k! \leq \text{Det } [A(z) + b_m U] \leq \text{Det } A(z) + 1/2c^k \cdot k!$$

for all $z \in Z$. Hence for $m \geq m_0$ and for all $z \in Z$,

$$|\text{Det } [A(z) + b_m U]| \geq 1/2c^k \cdot k!,$$

and therefore (b) follows.

(c) is proved as soon as we apply (a) to the determinants of $\{A(z) + b_m U\}$, $m = 1, 2, \dots$ and all their cofactors.

THEOREM 2. *Under conditions (ia), (ic), (iia), (iv) and (v), the least squares estimator t^* is minimax for θ in the class of all estimators.*

PROOF. We shall first prove that

$$(7) \quad \text{Sup}_{G \in \mathcal{G}_0} r(G, t^*) \leq \text{Sup}_{G \in \mathcal{G}_0} r(G, t) \quad \text{for all } t.$$

Since there is a one-one correspondence between \mathcal{G}_0 and the $(p + 1)$ -dimensional Euclidean space (in which θ takes its values), we can write $\rho(\theta, t)$ instead of $r(G, t)$ for each t where θ corresponds to G . Also since for each θ , $\rho(\theta, t) = vE \text{tr} [(X'X)^{-1}M]$, it will be enough to show that for all t ,

$$\text{Sup}_{\theta} \rho(\theta, t) \geq vE \text{tr} [(X'X)^{-1}M].$$

Suppose there exists \hat{t} such that

$$\text{Sup}_{\theta} \rho(\theta, \hat{t}) = vE \text{tr} [(X'X)^{-1}M] - \epsilon, \quad \epsilon > 0.$$

We shall contradict this by showing that for any given $\epsilon > 0$, we can make

$$(8) \quad \rho(\xi_m, \hat{t}) > vE \text{tr} [(X'X)^{-1}M] - \epsilon$$

by choosing m sufficiently large, where $\{\xi_m\}$ is a sequence of *a priori* distributions of θ defined as follows

$$d\xi_m(\theta) = (2\pi m)^{-\frac{1}{2}(p+1)} \exp [-(1/2m)\theta'\theta] \prod_{j=0}^p d\theta_j.$$

Choose and fix $\epsilon > 0$. Since $E \operatorname{tr} [(X'X)^{-1}M]$ exists, we can find a constant $c(\epsilon)$ such that

$$\left| \int_{X \in R(\epsilon)} \operatorname{tr} [(X'X)^{-1}M] \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}) - E \operatorname{tr} [(X'X)^{-1}M] \right| < \epsilon/2v,$$

where $R(\epsilon)$ is the set $\{X: \|(X'X)\| \leq c(\epsilon), \|(X'X)^{-1}\| \leq c(\epsilon)\}$. Now,

$$\begin{aligned} \rho(\xi_m, \hat{t}) = & \int_{\theta} \left[\int_X \int_Y (\theta - \hat{t})'M(\theta - \hat{t})g_{\theta}(y|X) \right. \\ & \left. \cdot \prod_{i=1}^n dy_i \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}) \right] d\xi_m(\theta) \end{aligned}$$

where $g_{\theta}(y|X) = \operatorname{const.} \exp [-(1/2v)(y - X\theta)'(y - X\theta)]$.

Since the integral in the above expression is non-negative, we get

$$\begin{aligned} \rho(\xi_m, \hat{t}) \geq & \int_{X \in R(\epsilon)} \int_Y \left[\int_{\theta} (\theta - \hat{t})'M(\theta - t)g_{\theta}(y|X) d\xi_m(\theta) \right] \\ & \cdot \prod_{i=1}^n dy_i \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}). \end{aligned}$$

Again if $(X'X + (v/m)I)^{-1}$ exists, $g_{\theta}(y|X) d\xi_m(\theta)$ can be written as

$$\operatorname{const.} f_{1m}(\theta|X, y)f_{2m}(X, y) \prod_{j=0}^p d\theta_j,$$

where

$$\begin{aligned} f_{1m}(\theta|X, y) = & \exp [-1/(2v)\{\theta - (X'X + (v/m)I)^{-1}X'y\}'(X'X + (v/m)I) \\ & \cdot \{\theta - (X'X + (v/m)I)^{-1}X'y\}] \end{aligned}$$

and

$$f_{2m}(X, y) = \exp [-(1/2v)y'y + (1/2v)y'X(X'X + (v/m)I)^{-1}X'y].$$

But $(X'X + (v/m)I)^{-1}$ exists for all m greater than or equal to some integer m_0 , and for all $X \in R(\epsilon)$, since the hypotheses of Lemma 2 are ensured for $X \in R(\epsilon)$. Hence we have

$$\begin{aligned} \rho(\xi_m, \hat{t}) \geq & \operatorname{const.} \int_{X \in R(\epsilon)} \int_Y \left[\int_{\theta} (\theta - \hat{t})'M(\theta - \hat{t})f_{1m}(\theta|X, y) \right. \\ & \left. \cdot \prod_{j=0}^p d\theta_j \right] f_{2m}(X, y) \prod_{i=1}^n dy_i \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}) \end{aligned}$$

for $m \geq m_0$. Now since M is positive definite, for values of $m \geq m_0$ and for any given $X \in R(\epsilon)$ and \mathbf{y} ,

$$\int_{\Theta} (\boldsymbol{\theta} - \hat{\mathbf{t}})' M (\boldsymbol{\theta} - \hat{\mathbf{t}}) f_{1m}(\boldsymbol{\theta} | X, \mathbf{y}) \prod_{j=0}^p d\theta_j$$

$$\geq \int_{\Theta} \{ \boldsymbol{\theta} - (X'X + (v/m)I)^{-1} X' \mathbf{y} \}'$$

$$\cdot M \{ \boldsymbol{\theta} - (X'X + (v/m)I)^{-1} X' \mathbf{y} \} f_{1m}(\boldsymbol{\theta} | X, \mathbf{y}) \prod_{j=0}^p d\theta_j.$$

Hence,

$$\rho(\xi_m, \hat{\mathbf{t}}) \geq \text{const.} \int_{X \in R(\epsilon)} \int_{\mathbf{y}} \left[\int_{\Theta} \{ \boldsymbol{\theta} - (X'X + (v/m)I)^{-1} X' \mathbf{y} \}'$$

$$\cdot M \{ \boldsymbol{\theta} - (X'X + (v/m)I)^{-1} X' \mathbf{y} \} f_{1m}(\boldsymbol{\theta} | X, \mathbf{y}) \prod_{j=0}^p d\theta_j \right]$$

$$\cdot f_{2m}(X, \mathbf{y}) \prod_{i=1}^n dy_i \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}).$$

After some simplifications, the right side of the last inequality reduces to

$$v \int_{X \in R(\epsilon)} \text{tr} [(X'X + (v/m)I)^{-1} M] \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}) = K_m(\epsilon) \text{ say.}$$

It follows from Lemma 2 that

$$\lim_{m \rightarrow \infty} K_m(\epsilon) = v \int_{X \in R(\epsilon)} \text{tr} [(X'X)^{-1} M] \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}).$$

Therefore, for sufficiently large m ,

$$K_m(\epsilon) > v \int_{X \in R(\epsilon)} \text{tr} [(X'X)^{-1} M] \prod_{i=1}^n dF(x_{1i}, \dots, x_{pi}) - \epsilon/2$$

$$> vE \text{tr} [(X'X)^{-1} M] - \epsilon.$$

Hence, for such large values of m , (8) holds. This proves (7). To complete the proof of the theorem, we have only to note that since $\mathcal{G}_0 \subset \mathcal{G}$,

$$\text{Sup}_{\sigma \in \mathcal{G}_0} r(G, \mathbf{t}) \leq \text{Sup}_{\sigma \in \mathcal{G}} r(G, \mathbf{t})$$

for arbitrary \mathbf{t} , whereas

$$\text{Sup}_{\sigma \in \mathcal{G}_0} r(G, \mathbf{t}^*) = vE \text{tr} [(X'X)^{-1} M] = \text{Sup}_{\sigma \in \mathcal{G}} r(G, \mathbf{t}^*).$$

5. Remarks.

(a) The results in the previous sections are proved under the assumption that the marginal distribution F of (x_1, \dots, x_p) is completely specified. Now suppose F belongs to a family \mathcal{F} of p -variate distribution functions. Then the family

\mathcal{G} of distributions of (y, x_1, \dots, x_p) can be expressed as the union of disjoint sub-families $\mathcal{G}(F)$, $F \in \mathcal{F}$, where each $G \in \mathcal{G}(F)$ has F as the marginal distribution of (x_1, \dots, x_p) . It can be easily seen that if conditions (ia) and (ic) are satisfied by each $F \in \mathcal{F}$ and if conditions (iia), (iib) and (iii) are satisfied by \mathcal{G} , then Theorem 1 and its corollary hold. Also, Theorem 2 remains valid if condition (ic) is replaced by

CONDITION (id). $\text{Sup}_{F \in \mathcal{F}} E_F \text{tr} [(X'X)^{-1}M_F] < \infty$, where M_F is the matrix defined in Section 2 in which $\mu_{jj'} = E_F(x_j x_{j'})$, and if condition (v) holds for each $\mathcal{G}(F)$, $F \in \mathcal{F}$.

(b) Condition (ic) is not satisfied in general and no simple way of verifying this condition is known. When this condition is not satisfied, t^* is not only inadmissible but is the worst possible estimate. In such cases, or even when the condition cannot be verified due to difficulties in analysis, it is very dangerous to use the method of least squares. It is not known whether there exists an estimator of θ with bounded risk when condition (ic) is not satisfied. Even if such an estimator exists, it has to be non-linear in y 's.

(c) Under the following sequential modification of the sampling scheme, condition (ic) is always satisfied. The procedure given below is very crude but if the loss function is the sum of two components, one given by (2) and another proportional to the sample size, it has a bounded risk function.

Suppose (y, x_1, \dots, x_p) follows a $(p + 1)$ -variate distribution. Let us choose and fix a constant c , however large. We then say the independent observations $(y_i, x_{1i}, \dots, x_{pi})$, $i = 1, \dots, n \geq p + 1$, on (y, x_1, \dots, x_p) have risk of order c if $\|(X'X)^{-1}\| \leq c$. Then our sampling scheme is as follows:

Sampling Scheme. Choose and fix a positive constant c . Make n independent observations on (y, x_1, \dots, x_p) . If the x -observations have risk of order c , stop sampling; if not, reject the observations and repeat the procedure till a set of observations having risk of order c is obtained, which is called the set of effective observations up to a risk of order c .

For any c , the effective observations up to a risk of order c can be considered to be observations on a process $(y'_i, x'_{1i}, \dots, x'_{pi})$, $i = 1, \dots, n$ for which y'_1, \dots, y'_n given x'_{1i}, \dots, x'_{pi} , $i = 1, \dots, n$, are mutually independent, the regression function of y' on x'_1, \dots, x'_p and the conditional variance of y' given x'_1, \dots, x'_p are the same as those for (y, x_1, \dots, x_p) , while the marginal distribution of (x'_1, \dots, x'_p) satisfies condition (ic). Hence if (y, x_1, \dots, x_p) satisfy conditions (iia) and (iib), then (y', x'_1, \dots, x'_p) also satisfies conditions (iia) and (iib), and similarly for condition (iii) or (iv) or (v). It can also be noticed that we have never made use of the independence of (x_{1i}, \dots, x_{pi}) , $i = 1, \dots, n$ in the course of our analysis; all that we required was the independence of (y_1, \dots, y_n) given (x_{1i}, \dots, x_{pi}) , $i = 1, \dots, n$ and this property is preserved in the process $(y'_i, x'_{1i}, \dots, x'_{pi})$, $i = 1, \dots, n$. Thus we see that under conditions (ia), (ib), (iia), (iib) and (iii), the least squares estimator t^* obtained from a set of effective observations up to a risk of some order c belongs to \mathcal{C}_1 , and is the unique estimator for θ having uniformly minimum risk among

all members of \mathcal{C}_1 obtained from the same set of observations. Also, under conditions (ia), (ib), (iia), (iv) and (v), the above estimator is minimax for θ in the class of all estimators obtained from the same set of observations. Under this sampling scheme, the sample size becomes a random variable with expectation greater than n but since $(X'X)^{-1}$ exists with probability 1 by virtue of condition (ia), the increase in the expected sample size over n can be made arbitrarily small by taking c sufficiently large.

Acknowledgments. The author wishes to thank Professor Wassily Hoeffding and Professor S. N. Roy for making some helpful suggestions.

REFERENCES

- [1] HODGES, J. L., JR. and LEHMANN, E. L. (1950). Some problems in minimax point estimation. *Ann. Math. Statist.* **21** 182-197.
- [2] PLACKETT, R. L. (1949). A historical note on the method of least squares. *Biometrika* **36** 458-460. Stanford University Press. Stanford, Calif.
- [3] STEIN, CHARLES (1960). Multiple regression. *Contributions to Probability and Statistics (Essays in Honor of Harold Hotelling)*. 424-443.