

SOME LIMIT THEOREMS IN REGRESSION THEORY

By K. R. PARTHASARATHY

and

P. K. BHATTACHARYA

Indian Statistical Institute

SUMMARY. (X, Y) follows an unknown bivariate distribution with $0 \leq X \leq 1$ and the regression of Y on X is continuous and a sequence of observations on (X, Y) are made. An estimate of the unknown regression function based on these observations and motivated by the method of Fractile Graphical Analysis has been suggested. Its large sample properties, viz., convergence in probability and almost sure uniform convergence to the true regression function have been investigated. Large sample tests for a specified regression function have also been proposed for the case when the conditional variance function of Y on X is known and for the case when it is unknown.

1. INTRODUCTION

Let us suppose that X and Y are real valued variables having a certain joint distribution function such that X takes values in the interval $[0, 1]$ and all conditional absolute moments of order up to $p \geq 3$ exist when X is fixed at any point x . If we can make a sequence of independent observations on (X, Y) the question naturally arises as to how we can construct from these observations an estimate of the unknown regression function of Y on X , possessing certain properties like convergence in probability, almost sure uniform convergence etc., to the true regression function. Another important problem is that of constructing at least a large sample test for a specified regression function. In Section 2 we shall make use of the technique of Fractile Graphical Analysis suggested by Mahalanobis (1958) to estimate the regression and analyse its large sample properties and in Section 3, construct a large sample test for the regression. The crucial part of our analysis consists of the utilization of certain results concerning the error of approximation by the central limit theorem and an upper bound of the tail probabilities in the distribution of sums of independent and bounded random variables. These results are given in the Appendix for reference.

2. ESTIMATION OF THE REGRESSION FUNCTION

Let $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be independent observations on (X, Y) , $x_{(r)}$ the r -th order statistic in the set of observed values of X and $y_{(i)} = y_j$ if $x_{(i)} = x_j$. Throughout this section we shall assume that the distribution function of X is continuous and strictly increasing so that the probability of any two observations on X being

equal is zero. Then the statistics $x_{(r)}$ are defined unambiguously for almost all samples. In order to simplify notation we write for $r = 1, 2, \dots, k$ and $s = 1, 2, \dots, n$,

$$\begin{aligned} x_{(r-1)n+s)} &= x_{rs}, & y_{(r-1)n+s)} &= y_{rs}, \\ E[Y|X=x] &= \varphi(x), & E[Y^2|X=x] &= \xi(x), \\ E[|Y-\varphi(x)|^m|X=x] &= \beta_m(x), & m &= 2, 3, \dots, p, \\ \varphi(x_{rs}) &= \varphi_{rs}, & \xi(x_{rs}) &= \xi_{rs}, & \beta_m(x_{rs}) &= \beta_{mrs}, \\ \bar{y}_r &= \sum_{s=1}^n y_{rs}/n, & \bar{y}_r^2 &= \sum_{s=1}^n y_{rs}^2/n, \\ \bar{\varphi}_r &= \sum_{s=1}^n \varphi_{rs}/n, & \bar{\varphi}_r^2 &= \sum_{s=1}^n \varphi_{rs}^2/n, \\ \bar{\beta}_{mr} &= \sum_{s=1}^n \beta_{mrs}/n, & \bar{\xi}_r &= \sum_{s=1}^n \xi_{rs}/n. \end{aligned} \quad \dots (2.1)$$

Now we define two functions $f_{nk}(x)$ and $\hat{\varphi}_{nk}(x)$ as follows :

$$\begin{aligned} f_{nk}(x) &= \bar{y}_1 & \text{if } 0 < x < x_n \\ &= \bar{y}_r & \text{if } x_{r-1n} < x < x_{rn}, \quad r = 2, \dots, k-1 \\ &= \bar{y}_k & \text{if } x_{k-1n} < x < 1. \end{aligned} \quad \dots (2.2)$$

$$\begin{aligned} \hat{\varphi}_{nk}(x) &= \bar{\varphi}_1 & \text{if } 0 < x < x_{1n} \\ &= \bar{\varphi}_r & \text{if } x_{r-1n} < x < x_{rn}, \quad r = 2 \dots k-1 \\ &= \bar{\varphi}_k & \text{if } x_{k-1n} < x < 1. \end{aligned} \quad \dots (2.3)$$

Lemma 1 : If the random variable X has a strictly increasing continuous distribution function $F(x)$, then

$$\begin{aligned} P \left[F^{-1} \left(\frac{r}{k} - \frac{1}{2k} \right) < x_{rn} < F^{-1} \left(\frac{r}{k} + \frac{1}{2k} \right), r = 1, 2, \dots, k-1 \right] \\ > 1 - 4m \cdot \exp \left[-\frac{3n}{16k} \right] - 2k \cdot \exp \left[-\frac{n}{4k} \left(1 - \frac{1}{8m} \right) \right] \end{aligned}$$

for every fixed positive integer m .

Proof : We have

$$P[x_{rn} < z] = \sum_{s=0}^{kn-rn} \binom{kn}{s} [1-F(z)]^s [F(z)]^{kn-s}. \quad \dots (2.4)$$

Applying Theorem (A 2) to the sum S_N of N independent binomial random variables each with probability for success equal to p , we have

$$P[S_N \geq N(p+\delta)] < \begin{cases} \exp \left[\frac{-N\delta}{2p} \log \left(\frac{\delta}{2q} + \sqrt{1 + \frac{\delta^2}{4q^2}} \right) \right] & \text{if } p > q \\ \exp \left[\frac{-N\delta}{2q} \log \left(\frac{\delta}{2p} + \sqrt{1 + \frac{\delta^2}{4p^2}} \right) \right] & \text{if } p < q. \end{cases}$$

which, if $\frac{\delta}{2p}$ and $\frac{\delta}{2q}$ are less than unity, becomes

$$P[S_N \geq N(p+\delta)] \leq \begin{cases} \exp \left[-N\delta^2 \left(1 - \frac{\delta}{4q} \right) \right] & \text{if } p \geq q \\ \exp \left[-N\delta^2 \left(1 - \frac{\delta}{4p} \right) \right] & \text{if } p \leq q. \end{cases} \quad \dots (2.5)$$

We have

$$\begin{aligned} P \left[F^{-1} \left(\frac{r}{k} - \frac{1}{2k} \right) \leq x_{rn} \leq F^{-1} \left(\frac{r}{k} + \frac{1}{2k} \right), r = 1, 2, \dots, k-1 \right] \\ \geq 1 - \sum_{r=1}^{k-1} P \left[x_{rn} < F^{-1} \left(\frac{r}{k} - \frac{1}{2k} \right) \right] - \sum_{r=1}^{k-1} P \left[x_{rn} > F^{-1} \left(\frac{r}{k} + \frac{1}{2k} \right) \right] \end{aligned}$$

which, after using (2.4) and an application of (2.5) for $\delta = \frac{1}{2k}$ and two binomial sums with probability for successes $\frac{r}{k} - \frac{1}{2k}$ and $1 - \frac{r}{k} - \frac{1}{2k}$ and number of summands equal to nk , becomes greater than $4m \cdot \exp \left[-\frac{3n}{16k} \right] + 2k \cdot \exp \left[-\frac{n}{4k} \left(1 - \frac{1}{8m} \right) \right]$ for every fixed integer m .

Lemma 2: Under the conditions of Lemma 1 if $\varphi(x)$, the regression function of Y on X is continuous then for any $\epsilon > 0$

$$P \left[\sup_{0 \leq x \leq 1} |\hat{\varphi}_{nk}(x) - \varphi(x)| > \epsilon \right] < c(\epsilon, \varphi) \left[m \exp \left[-\frac{3n}{16k} \right] + k \exp \left[-4k \left(1 - \frac{1}{8m} \right) \right] \right]$$

for every fixed integer m , and where $c(\epsilon, \varphi)$ is a constant which depends on ϵ and φ only.

Proof: Let $\delta[a, b] = \max_{a \leq x_1, x_2 \leq b} |\varphi(x_1) - \varphi(x_2)|$

$$\delta_2 = \max \left\{ \delta \left[0, F^{-1} \left(\frac{1}{2k} \right) \right], \delta \left[F^{-1} \left(\frac{1}{k} - \frac{1}{2k} \right), F^{-1} \left(\frac{1}{k} + \frac{1}{2k} \right) \right], \dots, \delta \left[F^{-1} \left(1 - \frac{1}{2k} \right), 1 \right] \right\}.$$

$$\text{If} \quad F^{-1} \left(\frac{r}{k} - \frac{1}{2k} \right) \leq x_m \leq F^{-1} \left(\frac{r}{k} + \frac{1}{2k} \right), r = 1, 2, \dots, k-1$$

$$\text{then} \quad d_{nk}(\varphi) = \sup_{0 \leq x \leq 1} |\hat{\varphi}_{nk}(x) - \varphi(x)| \leq \delta_2. \quad \dots (2.6)$$

Thus

$$\begin{aligned} P[d_{nk}(\varphi) > \epsilon] \leq & \left[P[d_{nk} > \epsilon] P \left[F^{-1} \left(\frac{r-1}{k} \right) \leq x_{rn} \leq F^{-1} \left(\frac{r+1}{k} \right), r = 1, 2, \dots, k-1 \right] \right. \\ & \left. + 1 - P \left[F^{-1} \left(\frac{r-1}{k} \right) \leq x_{rn} \leq F^{-1} \left(\frac{r+1}{k} \right) \right], r = 1, 2, \dots, k-1. \right] \quad \dots (2.7) \end{aligned}$$

Because of (2.6) the first term on the right side of (2.7) becomes less than $P(\delta_2 > \epsilon)$. Since F is continuous and strictly increasing F^{-1} is continuous. Since

$\varphi(x)$ is uniformly continuous we have $\delta_k < \epsilon$ for $k > k_0(\epsilon, \varphi)$. Thus the first term vanishes for $k > k_0(\epsilon, \varphi)$. The second term becomes less than $4m \exp \left[\frac{-3n}{16k} \right] + 2k \exp \left[\frac{-n}{4k} \left(1 - \frac{1}{8m} \right) \right]$ because of Lemma 1. Combining the two we get the required inequality.

Hereafter, for simplicity, we shall assume that $\beta_k(x) > c_1 > 0$ and $\beta_m(x) < c_2$ for all $m = 2, \dots, p$ and all x . However, from the proofs we can see that all our results hold good under more general conditions.

Lemma 3: If $\beta_k(x) > c_1 > 0$ and $\beta_m(x) < c_2$ for all $m = 1, 2, \dots, p$ and all x , then for any $\epsilon > 0$

$$P \left[\sup_{0 < x < 1} |f_{nk}x - \hat{\varphi}_{nk}(x)| > \epsilon \right] < C \cdot \frac{k}{n^{p-1}}$$

where C is a constant depending only on ϵ, p, c_1 and c_2 .

Proof: We have

$$\begin{aligned} P(\epsilon) &= P \left[\sup_{0 < x < 1} |f_{nk}(x) - \hat{\varphi}_{nk}(x)| > \epsilon |x_1, \dots, x_{kn} \right] \\ &= P \left[\sup_{r=1, 2, \dots, k} |\bar{y}_r - \bar{\varphi}_r| > \epsilon |x_1, \dots, x_{kn} \right] \\ &< \sum_{r=1}^k P \left[\frac{|\bar{y}_r - \bar{\varphi}_r|}{\sqrt{\beta_r/n}} > \frac{\epsilon\sqrt{n}}{\sqrt{c_2}} \right]. \end{aligned} \quad \dots (2.8)$$

Making use of (A7) and noting that $T_{pn} < \frac{\sqrt{n}}{4}$ we have

$$\frac{\epsilon\sqrt{n}}{\sqrt{c_2}} > \sqrt{4(p-2) \log T_{pn}} \text{ for } n > n(\epsilon).$$

Hence by Theorem A1 we obtain

$$P(\epsilon) < 2k\Phi \left(\frac{-\epsilon\sqrt{n}}{\sqrt{c_2}} \right) - \frac{\theta_p c_2}{1 + \left(\frac{\epsilon\sqrt{n}}{\sqrt{c_2}} \right)^p} \cdot \frac{k}{n^{\frac{p-2}{2}}} \quad \dots (2.9)$$

where $c_2 = \frac{4c_1^2 p}{c_1^2 k}$ and θ_p is a positive constant which depends only on p and

$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$. Since $\Phi(-x) < \frac{1}{\sqrt{2\pi} x} e^{-x^2/2}$, $x > 0$, (2.9) gives, for $n > n(\epsilon)$.

$$\begin{aligned} P(\epsilon) &< \frac{k}{\epsilon\sqrt{n}/c_2} e^{-\frac{\epsilon^2}{2c_2} \cdot n} + \theta_p \frac{c_2^{p/2}}{\epsilon^p} \frac{c_2 k}{n^{p-1}} \\ &= c_4(\epsilon, p, c_1, c_2) \left[\frac{k}{\sqrt{n}} e^{-\frac{\epsilon^2}{2c_2} \cdot n} + \frac{k}{n^{p-1}} \right] \end{aligned}$$

which completes the proof.

SOME LIMIT THEOREMS IN REGRESSION THEORY

Now we shall state and prove the main result of this section.

Theorem 1: *If the p -th conditional absolute moment of Y exists when X is fixed at any point x , $p \geq 3$, X has a strictly increasing continuous distribution function in $[0, 1]$, the regression $\phi(x)$ of Y on X is continuous and $\beta_n(x) > c_1 > 0$, $\beta_m(x) < c_2$ for all $m = 1, 2, \dots, p$ and for all x , then for $n \geq (4+\delta)k \log k$, $\delta > 0$, $d_{nk} = \sup_x |f_{nk}(x) - \phi(x)|$ converges to zero in probability as $k \rightarrow \infty$ and for $n \geq (8+\delta)k \log k$, $\delta > 0$, d_{nk} converges almost surely to zero as $k \rightarrow \infty$.*

Proof: The first part of the theorem is an immediate consequence of Lemmas 2 and 3 and the second part follows from the same lemmas and Borel-Cantelli lemma.

3. LARGE SAMPLE TEST OF A SPECIFIED REGRESSION FUNCTION

In this section we shall consider the problem of testing the null hypothesis H_0 that the regression function $\phi(x)$ is equal to a specified function $\mu(x)$.

Consider the statistic

$$\tau_{nk} = \sup_{r=1, \dots, k} \frac{|\bar{y}_r - \bar{\mu}_r|}{\sqrt{\beta_{nr}}}$$

where \bar{y}_r and $\bar{\beta}_{nr}$ are as defined in (2.1) and $\bar{\mu}_r = \frac{\sum_{r=1}^n \mu(x_{rr})}{n}$. If we had known the probability distribution of τ_{nk} under the null hypothesis, then for any given level of significance $0 < \alpha < 1$ we could apply the following test:

Reject H_0 if and only if $\tau_{nk} > \tau_{nk}(\alpha)$

where

$$P[\tau_{nk} > \tau_{nk}(\alpha) | H_0] = \alpha.$$

Theorem 2 enables us to apply such a test at least in the large sample.

Theorem 2: *If $\beta_n(x) > c_1 > 0$ and $\beta_m(x) < c_2 < \infty$ for all x and $m = 1, 2, \dots, p$, $p \geq 3$ then we have for $n \geq (k \log k)^{2/p-2}$*

$$\lim_{k \rightarrow \infty} P[\tau_{nk} \leq \lambda_k | H_0] = \exp \left[-\frac{1}{\sqrt{\pi}} e^{-\theta} \right]$$

where

$$\lambda_k = \sqrt{2(\theta + \log k - \frac{1}{2} \log \log k)}.$$

Proof: Under H_0 , we have

$$P(\lambda_k) = P[\tau_{nk} \leq \lambda_k | x_1, \dots, x_{kn}]$$

$$\prod_{r=1}^k P \left[\frac{|\bar{y}_r - \bar{\phi}_r|}{\sqrt{\beta_{nr}}} \leq \lambda_k | x_1, \dots, x_{kn} \right] = \prod_{r=1}^k [\Phi(\lambda_k) - \Phi(-\lambda_k) + e_{rk}]$$

where because of (A8) and the fact that $\lambda_k < \sqrt{4(p-2) \log T_{n,k}}$ for $k \geq k_0(\theta)$ we get

$$|e_{rk}| \leq c_3 \left[\frac{(1 + \lambda_k^2)}{\sqrt{n}} e^{-\lambda_k^2/2} + \frac{1}{n^{p-2}} \right]$$

where c_5 is a constant depending on p, c_1, c_2 . Thus

$$\log P(\lambda_k) = k \log [\Phi(\lambda_k) - \Phi(-\lambda_k)] + Z_k,$$

where

$$|Z_k| \leq \sum_{r=1}^k |\log [1 + \frac{e_{rk}}{\Phi(\lambda_k) - \Phi(-\lambda_k)}]|.$$

Since for sufficiently large k

$$\frac{|e_{rk}|}{\Phi(\lambda_k) - \Phi(-\lambda_k)} < \frac{1}{2}$$

and $\log(1+x) = x + vx^2, |v| < 1$ for $|x| < \frac{1}{2}$,

$$|Z_k| \leq c_4 k \left[\frac{1 + \lambda_k^2}{\sqrt{n}} e^{-\lambda_k/2} + \frac{1}{n \frac{p-2}{2}} \right] \quad \dots (3.1)$$

Substituting for λ_k in (3.1) we have

$$|Z_k| \leq c_5 \left[\frac{1 + \lambda_k^2}{\sqrt{n}} \cdot \sqrt{\log k} \cdot e^{-\theta} + \frac{k}{n \frac{p-2}{2}} \right] \quad \dots (3.2)$$

Since $n \geq (k \log k)^{2/p-1}$ it is easy to see from (3.2) that $|Z_k| \rightarrow 0$ as $k \rightarrow \infty$. This completes the proof.

The statistic τ_{nk} defined above can be used only when $\beta_{2r}(x)$ is known. The natural way of modifying this in the case of unknown variance function is to replace β_{2r} by $y_r^2 - \bar{y}^2$, where y_r^2 and \bar{y} , are as in (2.1). Since this leads to certain complications in evaluating the limiting distribution we replace it by

$$s_r^2 = \frac{1}{n} \sum_{i=1}^n [y_{ri} - \mu(x_{ri})]^2 \text{ (where } x_{ri} \text{ is as in (2.1)) and write}$$

$$t_{nk} = \sup_{r=1, 2, \dots, k} \sqrt{n} \cdot \frac{|y_r - \bar{\mu}_r|}{s_r}$$

To show that this replacement does not change τ_{nk} effectively in large samples we prove the following lemmas.

Lemma 4: If X has a continuous distribution function and Y is a bounded random variable, $\beta_1(x) > c_1 > 0$ and if $n \geq (\log k)^{3+\delta}$, $\delta > 0$ then under the null hypothesis H_0

$$\lim_{k \rightarrow \infty} \log k \cdot \sup_r \left| \sqrt{\frac{n^2}{\beta_{2r}}} - 1 \right| = 0,$$

with probability one.

Proof: We have for any $\epsilon > 0$,

$$\begin{aligned} P \left[\sup_r \left| \sqrt{\frac{s_r^2}{\beta_{2r}}} - 1 \right| > \frac{\epsilon}{\log k} \mid H_0 \right] &\leq \sum_{r=1}^k P \left[\left| \frac{s_r^2}{\beta_{2r}} - 1 \right| > \frac{\epsilon}{\log k} \right] \\ &\leq \sum_{r=1}^k P \left[|s_r^2 - \beta_{2r}| > \frac{\epsilon c_1}{\log k} \right] \\ &\leq \sum_{r=1}^k P \left[\left| \bar{y}_r^2 - \bar{y}_r \right| > \frac{\epsilon c_1}{2 \log k} \right] + \sum_{r=1}^k P \left[\left| \sum_{i=1}^n y_{ri} y_{ri} - \nu_{2r}^2 \right| > \frac{n \epsilon c_1}{4 \log k} \right]. \end{aligned}$$

Then from Theorem A2 we get

$$\sum_{r=1}^k P \left[\left| \bar{y}_r^2 - \bar{t}_r \right| > \frac{ec_1}{2 \log k} |x_1, \dots, x_{kn}| \right] \leq 2k \exp \left[\frac{-b_1 n}{(\log k)^2} \right]$$

and
$$\sum_{r=1}^k P \left[\left| \sum_{s=1}^n (y_{rs} y_{rs} - \bar{y}_{rs}^2) \right| > \frac{nec_1}{4 \log k} \right] \leq 2k \exp \left[\frac{-b_2 n}{(\log k)^2} \right]$$

where b_1 and b_2 are constants depending on c_1, ϵ and the upper bound of $|Y|$. Thus whenever $n \geq (\log k)^{2+\delta}$ the series

$$\sum_{k=1}^{\infty} P \left[\sup_r \left| \sqrt{\frac{\sigma_r^2}{\beta_{2r}}} - 1 \right| > \frac{\epsilon}{\log k} \right]$$

converges. An application of Borel-Cantelli lemma completes the proof of the lemma.

Remark: It is possible to show that if $n \geq k (\log k)^{2+\delta}$, $\delta > 0$ then $\log k \sup_r \left| \sqrt{\frac{\sigma_r^2}{\beta_{2r}}} - 1 \right|$ converges to zero in probability even when Y is not a bounded random variable.

The following theorem enables us to use the statistic t_{nk} for testing a specified regression function when the conditional variance function is not known.

Theorem 3: *If X is a random variable with a continuous distribution in $[0, 1]$, Y is a bounded random variable and $\beta_k(x) > c_1 > 0$ then for $n \geq (k \log k)^{2/p-2}$ for some $p \geq 3$ we have*

$$\lim_{k \rightarrow \infty} P[t_{nk} < \lambda_k | H_0] = \exp \left[-\frac{1}{\sqrt{\pi}} e^{-\theta} \right]$$

where

$$\lambda_k = \sqrt{2(\theta + \log k - \frac{1}{2} \log \log k)}.$$

Proof: It is easily seen that

$$P[\tau_{nk} < \lambda_k(1 + V_{nk}) | H_0] \leq P[t_{nk} < \lambda_k | H_0] \leq P[\tau_{nk} < \lambda_k(1 + U_{nk}) | H_0] \dots \quad (3.3)$$

where τ_{nk} is as in Theorem 2, $U_{nk} = \sup_r \left(\sqrt{\frac{\sigma_r^2}{\beta_{2r}}} - 1 \right)$ and $V_{nk} = \inf_r \left(\sqrt{\frac{\sigma_r^2}{\beta_{2r}}} - 1 \right)$.

Further,
$$P[\tau_{nk} < \lambda_k(1 + V_{nk})] = P \left[\frac{\tau_{nk}^2 - 2 \log k + \log \log k}{2}, Z_{nk} + v_{nk} \leq \theta \right]$$

where Z_{nk} converges to unity in probability and v_{nk} converges to zero in probability because of Lemma 4. This fact together with an application of Theorem 2 leads to the result

$$\lim_{k \rightarrow \infty} P[\tau_{nk} < \lambda_k(1 + V_{nk})] = \exp \left[-\frac{1}{\sqrt{\pi}} e^{-\theta} \right]. \quad \dots \quad (3.4)$$

A similar procedure leads to the result

$$\lim_{k \rightarrow \infty} P[\tau_{nk} < \lambda_k (1 + U_{nk})] = \exp \left[-\frac{1}{\sqrt{\pi}} e^{-\theta} \right]. \quad \dots (3.5)$$

(3.3), (3.4) and (3.5) complete the proof.

Finally we shall state a theorem concerning the order of change in the level of significance if we apply the test 'reject H_0 if and only if $\tau_{nk} > \lambda_k$ ' where $[\Phi(\lambda_k) - \Phi(-\lambda_k)]^2 = 1 - \alpha$ instead of the test 'reject H_0 if and only if $\tau_{nk} < \tau_{nk}(\alpha)$ with $P[\tau_{nk} \leq \tau_{nk}(\alpha) | H_0] = 1 - \alpha$.

Theorem 4: Let λ_k satisfy $\Phi(\lambda_k) - \Phi(-\lambda_k) = 1 - \alpha$ and $\alpha < 1 - \left[1 - \sqrt{\frac{2}{e\pi}} \right]^k$.

Then under the conditions of Theorem 2 we have

$$P[\tau_{nk} \leq \lambda_k | H_0] - (1 - \alpha) \leq c \cdot \left[\frac{\left(\log \frac{k}{\alpha} \right)^2}{\sqrt{n}} + \frac{k}{n} - \frac{2}{2} \right]$$

where c is a constant depending on p , c_1 and c_2 .

4. REMARKS

In this paper we have not given any consideration to the power of the τ_{nk} and t_{nk} tests proposed in Section 3. It would be very interesting if lower bounds for the power of these tests could be given in terms of the supremum distance between the regression functions under the null hypothesis and the alternative hypothesis. Though this has not been done, we can at least easily verify that if the regression functions under null hypothesis and the alternative hypothesis are both continuous, the τ_{nk} and t_{nk} tests are consistent, i.e., for any given level of significance $0 < \alpha < 1$, the probability of rejecting H_0 tends to unity as $k \rightarrow \infty$.

For the validity of the above theorems we have imposed certain conditions on the regression function φ and the variance function β_2 separately. In some cases (e.g. Binomial or Poisson distribution) β_2 is a function of φ . This however does not affect our analysis in any way so long as the conditions on φ and β_2 remain valid separately.

For simplicity in proof we assumed that $\beta_2(x)$ and $\beta_{2n}(x)$ are bounded away from zero and bounded above respectively. However, it is easy to show that theorems 1, 2 and 4 are valid under the weaker assumption that $E[\beta_2(x)]^{-3(p-1)/2}$ and $E[\beta_2(x)]^{(p-2)/p}$ are finite.

* $\Phi(x)$ is the usual normal distribution function.

SOME LIMIT THEOREMS IN REGRESSION THEORY

Appendix

In this section we shall state and sketch the proof of some of the results which were utilised in our paper.

Let $X_1, X_2, X_3 \dots$ be a sequence of independent random variables with $EX_n = 0$ and $EX_n^2 = \sigma_n^2$ and β_{rn} the r -th absolute moment of X_n . We shall suppose that the moments of order $k > 3$ exist. We write

$$B_{rn} = \frac{1}{n} (\beta_{r1} + \dots + \beta_{rn}), \quad \rho_{rn} = \frac{B_{rn}}{B_{2n}^{r/2}}, \quad T_{kn} = \frac{\sqrt{n}}{4\rho_{kn}^{3/k}} \quad \dots \quad (A1)$$

Then
$$1 \leq \rho_{rn} \leq \rho_{kn}^{r/k}, \quad r = 2, 3, \dots, k. \quad \dots \quad (A2)$$

Let $f_n(t)$ and $F_n(x)$ be the characteristic function and distribution function of $\frac{X_1 + \dots + X_n}{\sqrt{nB_n}}$. We now state two lemmas one of which is due to Cramér (1937) and the other due to Esseen and Berry.

Lemma A1 : (Cramér). For $|t| \leq T_{kn}^{1/3}$ we have

$$\left| f_n(t) - e^{-t^2/2} - \sum_{r=1}^{k-3} \frac{P_{rn}(it)}{n^{r/2}} e^{-t^2/2} \right| \leq \frac{c_k}{T_{kn}^{k-2}} [|t|^k + |t|^{3(k-3)}] \cdot e^{-t^2/2}$$

where
$$P_{rn}(it) = \sum_{j=1}^r c_{jrn} (it)^{r+2j},$$

$c_{jrn} = c_k' \cdot \rho_{kn}^{\frac{r+2j}{k}}$ and c_2 and c_k' are constants depending only on k , c_k being positive.

Lemma A2 : (Berry-Esseen). For $|t| \leq T_{3n}$

$$\left| f_n(t) - e^{-t^2/2} \right| \leq \frac{4}{T_{3n}} \cdot |t|^3 \cdot e^{-t^2/2}.$$

Let
$$g(t) = e^{-t^2/2} \left[1 + \sum_{r=1}^{k-3} \frac{P_{rn}(it)}{n^{r/2}} \right]. \quad \dots \quad (A3)$$

Hereafter we shall denote by θ_k any positive constant which depends only on k . Then by using (A1), (A2) and Lemma A1 we have

$$\begin{aligned} \left| \sum_{r=1}^{k-3} \frac{P_{rn}(it)}{n^{r/2}} \right| &\leq \theta_k \sum_{r=1}^{k-3} \sum_{j=1}^r \left| \frac{c_{jrn}}{n^{r/2}} \right| |t|^{r+2} \dots \leq \theta_k \sum_{r=1}^{k-3} \sum_{j=1}^r \frac{\rho_{kn}^{\frac{r+2j}{k}}}{n^{r/2}} \cdot |t|^{r+2j} \\ &\leq \theta_k \sum_{r=1}^{k-3} \sum_{j=1}^r \frac{\rho_{kn}^{\frac{3r}{k}}}{(\sqrt{n})^r} |t|^{r+2j} \leq \theta_k \sum_{r=1}^{k-3} \sum_{j=1}^r \frac{|t|}{T_{kn}^{3r/k}}. \quad \dots \quad (A4) \end{aligned}$$

From (A3) and (A4) it is seen that the inequality of Lemma A2 can be rewritten as

$$\left| f_n(t) - g_n(t) \right| \leq \frac{4}{T_{3n}} |t|^3 \cdot e^{-t^{1/3}} + \theta_k \cdot \sum_{r=1}^{k-3} \sum_{j=1}^r \frac{|t|^{r+3j}}{T_{kn}^r} e^{-t^{1/3}} \dots \quad (A5)$$

Lemma A3 :
$$\int_{-T_{kn}}^{+T_{kn}} \left| \frac{f_n(t) - g_n(t)}{t} \right| dt \leq \frac{\theta_k}{T_{kn}^{k-2}}.$$

Case 1 : Let $T_{kn} \leq 1$, then $T_{kn} \leq T_{kn}^{1/3}$. Hence applying Lemma A1 we get

$$\int_{-T_{kn}}^{+T_{kn}} \left| \frac{f_n(t) - g_n(t)}{t} \right| dt \leq \frac{\theta_k}{T_{kn}^{k-2}}.$$

Case 2 : Let $T_{kn} > 1$. From (A2) it is easy to see that $1 < T_{kn} \leq T_{3n}$:

Thus
$$\int_{-T_{kn}}^{+T_{kn}} \left| \frac{f_n(t) - g_n(t)}{t} \right| dt = 2 \int_{T_{kn}^{1/3}}^{T_{kn}} \left| \frac{f_n(t) - g_n(t)}{t} \right| dt + \int_{-T_{kn}^{1/3}}^{T_{kn}^{1/3}} \left| \frac{f_n(t) - g_n(t)}{t} \right| dt = 2I_1 + I_2 \text{ say.}$$

In the region of integration of I_1 Lemma A2 is applicable and for I_2 Lemma A1 is applicable. Thus from (A5) we have

$$\int_{-T_{kn}}^{+T_{kn}} \left| \frac{f_n(t) - g_n(t)}{t} \right| dt \leq \frac{\theta_k}{T_{kn}^{k-2}} + \frac{\theta_k}{T_{kn}} e^{-T_{kn}^{1/4}} \leq \theta_k / T_{kn}^{k-2}.$$

Let
$$Q(x) \geq 0, 0 \leq q(t) \leq 1, \int_{-\infty}^{+\infty} Q(x) dx = 1.$$

$$q(t) = \int_{-\infty}^{+\infty} e^{itx} Q(x) dx, q(t) \leq \frac{c}{|t|^{3(k-2)}}$$

$$\int_{-\infty}^{+\infty} |t|^{3k-7} q(t) dt < \infty. \dots \quad (A6)$$

Lemma A4 : If $q(t)$ satisfies conditions (A6) then for $|x| > 1$

$$P(T_{kn}) = \int_{-T_{kn}}^{+T_{kn}} \left| \frac{f_n(t) - g_n(t)}{t} \right| q(xt) dt < \frac{\theta}{1 + |x|^k} \cdot \frac{1}{T_{kn}^{k-2}}$$

SOME LIMIT THEOREMS IN REGRESSION THEORY

Proof: Case 1: Let $T_{kn} \leq 1$, then $T_{kn} < T_{kn}^{1/2}$. Applying Lemma A1.

$$P(T_{kn}) \leq \frac{\theta_k}{T_{kn}^{k-2}} \int_{-\infty}^{+\infty} [|t|^k + |t|^{2(k-2)}] e^{-t^2} \cdot \frac{q(xt)dt}{|t|}$$

$$\leq \frac{\theta_k}{T_{kn}^{k-2}} \cdot \frac{1}{|x|^k} \cdot \int_{-\infty}^{+\infty} [|t|^{k-1} + |t|^{2k-7}] e^{-t^2} q(t)dt \leq \frac{\theta_k}{(1+|x|^k)} \frac{1}{T_{kn}^{k-2}}$$

Case 2: Let $T_{kn} > 1$. Then as before $1 < T_{kn} \leq T_{3n}$.

$$P(T_{kn}) = 2 \int_{T_{kn}^{1/2}}^{T_{kn}} + \int_{-T_{kn}^{1/2}}^{-T_{kn}} \left| \frac{f_n(t) - g_n(t)}{t} \right| q(xt)dt = 2I_1 + I_2 \text{ say.}$$

As before by applying Lemma A1 for I_2 and Lemma A2 for I_1 and proceeding as in Case 1 we obtain

$$P(T_{kn}) \leq \frac{\theta_k}{1+|x|^k} \cdot \frac{1}{T_{kn}^{k-2}}$$

Utilising Lemmas A3 and A4 and proceeding along the same lines as Esseen (1944) we can prove the following results.

Theorem A1: Let X_1, X_2, \dots, X_n be a sequence of random variables with mean zero and finite absolute moments of order $k(k \geq 3)$. Then for $x \geq \lambda_n$

$$\left| F_n(x) - \Phi(x) - \sum_{r=1}^{k-3} \frac{P_{rn}(-\Phi)}{n^{r/2}} \right| \leq \frac{\theta_k}{1+|x|^k} \cdot \frac{1}{T_{kn}^{k-2}} \text{ and for } x < \lambda_n$$

$$\left| F_n(x) - \Phi(x) - \sum_{r=1}^{k-3} \frac{P_{rn}(-\Phi)}{n^{r/2}} \right| \leq \theta'_k \left[\frac{(1+|x|^3)e^{-x^2}}{T_{kn}} + \frac{1}{T_{kn}^{k-2}} \right]$$

where $\lambda_n = 1 + \delta$ if $T_{kn} \leq 1$ and $\max[1 + \delta, \sqrt{2 \cdot (1 + \delta)(k-2) \log T_{kn}}]$ if $T_{kn} > 1$.

θ_k and θ'_k are constants which depends only on k and δ , $P_{rn}(-\Phi) = \sum_{j=1}^r (-1)^{r+j} c_{jrn} \cdot \Phi^{(r+2j)}(x)$ and c_{jrn} are as in Lemma A 1.

Remark: From the above theorem it can be easily deduced that

$$\left| F_n(x) - \Phi(x) \right| \leq \frac{\theta_k}{1+|x|^k} \cdot \frac{1}{T_{kn}^{k-2}} \text{ for } x > \lambda_n \quad \dots \quad (\text{A } 7)$$

$$\left| F_n(x) - \Phi(x) \right| \leq \theta'_k \left[\frac{(1+|x|^3)e^{-x^2}}{T_{kn}} + \frac{1}{T_{kn}^{k-2}} \right] \text{ for } x \leq \lambda_n. \quad (\text{A } 8)$$

SANKHYĀ : THE INDIAN JOURNAL OF STATISTICS : SERIES A

Let $\xi_1, \xi_2, \dots, \xi_n$ be n independently distributed random variables with $E(\xi_i) = 0$, $|\xi_i| \leq c$, $i = 1, 2, \dots, n$, $E(\xi_1 + \dots + \xi_n)^2 = \sigma^2$ and $\xi = \xi_1 + \dots + \xi_n$. Then we have the following theorem due to Prohorov (1950).

Theorem A2: Under the conditions stated above, for $x > 0$ we have

$$P\{\xi > x\} \leq \exp \left[-\frac{x}{2c} \sinh^{-1} \frac{xc}{2\sigma^2} \right].$$

REFERENCES

- CRAMÉR, H. (1937): Random variables and probability distributions. *Cambridge Tracts in Mathematics*, No. 36.
- ESSEEN, C. G. (1944): Fourier analysis of distribution functions. A mathematical study of the Laplace-Gaussian Law. *Acta Math.*, 77, 1.
- MARALANANDIS, P. C. (1958): Lectures in Japan: Fractile Graphical Analysis. Indian Statistical Institute.
- PROHOROV, YU. V. (1950): An extremal problem in probability theory. *Theoriai Veroyatnosti i ee Primeneniya* 4, 211.

Paper received: May, 1960.

ACKNOWLEDGEMENT

Thanks are due to the following for their kind help in the editing of the papers published in this issue.

Shri S. John

Professor D. B. Lahiri

Dr. Sujit Kumar Mitra

Dr. C. R. Rao

Shri R. Ranga Rao

Shri J. Sethuraman

Editor,

Sankhyā : The Indian Journal of Statistics.