

Aleš Slabý

Limit theorems for rank statistics detecting gradual changes

Commentationes Mathematicae Universitatis Carolinae, Vol. 42 (2001), No. 3, 591--600

Persistent URL: <http://dml.cz/dmlcz/119274>

Terms of use:

© Charles University in Prague, Faculty of Mathematics and Physics, 2001

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://project.dml.cz>

Limit theorems for rank statistics detecting gradual changes

ALEŠ SLABÝ

Abstract. The purpose of the paper is to investigate weak asymptotic behaviour of rank statistics proposed for detection of gradual changes, linear trends in particular. The considered statistics can be used for various test procedures. The fundamentals of the proofs are formed by results of Hušková [4] and Jarušková [5].

Keywords: gradual type of change, location model, rank statistics, weighted suprema, limit theorems

Classification: 62E20, 62G20, 62G30

1. Introduction

The basic underlying problem is testing a sequence of i. i. d. (time ordered) random observations X_1, \dots, X_n having the same common continuous distribution against an alternative that at an unknown time point there is a beginning of gradual type of change in location such that the trend after the change point is linear. Namely, we test

$$H : X_i = \theta + e_i$$

against

$$A : X_i = \theta + \delta \frac{i - m}{n} I_{\{i > m\}} + e_i \quad \text{for some } 1 \leq m < n \text{ and } \delta \neq 0$$

where e_1, \dots, e_n are i. i. d. with continuous distribution function $F(x, 0)$ provided that $F(x, \theta) = F(x - \theta) = F(x)$. However, statistics studied in the sequel can be used for testing other and more complicated alternatives, see Slabý [7].

Let R_{1n}, \dots, R_{nn} be the ranks corresponding to the observations X_1, \dots, X_n . Consider two simple linear rank statistics

$$(1) \quad S_{1k}(\mathbf{a}) = \sum_{i=1}^k (a(R_{in}) - \bar{a}_n)$$

The author gives special thanks to Marie Hušková, who was instant in granting hints and consultations. This work was also supported by grants GAČR 201/00/0769 and MSM 113200008.

and

$$(2) \quad S_{2k}(\mathbf{a}) = \sum_{i=k+1}^n (a(R_{in}) - \bar{a}_n) \frac{i-k}{n}$$

where $k = 1, \dots, n$, and $a(1), \dots, a(n)$ are scores with properties

$$(3) \quad \frac{1}{n} \sum_{i=1}^n (a(i) - \bar{a}_n)^2 \geq D_1$$

and

$$(4) \quad \frac{1}{n} \sum_{i=1}^n |a(i) - \bar{a}_n|^{2+\eta} \leq D_2$$

for some finite positive constants D_1, D_2 and η independent of n . Here

$$(5) \quad \bar{a}_n = \frac{1}{n} \sum_{i=1}^n a(i).$$

Hušková [4] studied limit behaviour of weighted suprema and L_p -functionals based on $S_{1k}(\mathbf{a})$ and, then, proposed corresponding testing procedures for abrupt change in location model setup. Whereas $S_{1k}(\mathbf{a})$ is more suitable for testing abrupt changes, analogous weighted suprema based on $S_{2k}(\mathbf{a})$ can be better employed to test gradual changes, particularly changes of linear trend in the location model.

The crucial result of Hušková [4] is that limit behaviour of the treated statistics is the same as the limit behaviour of the respective functionals of sums of certain i. i. d. random variables. Limit theorems for i. i. d. random variables, see Csörgő and Horváth [1], can be then extended to ranks. Hušková [4] shows that instead of investigating behaviour of $S_{1k}(\mathbf{a})$, it is sufficient to check behaviour of

$$(6) \quad Z_{1k}(\mathbf{a}) = \sum_{i=1}^k (a(1 + [nU_i]) - \bar{a}_n)$$

where

$$(7) \quad U_i = F(X_i), \quad i = 1, \dots, n,$$

and $[a]$ denotes the integer part of a . Notice that for $k = 1, \dots, n$ we have

$$(8) \quad ES_{1k}(\mathbf{a}) = EZ_{1k}(\mathbf{a}) = 0$$

and

$$(9) \quad \text{Var } S_{1k}(\mathbf{a}) = \frac{n}{n-1} \text{Var} \left\{ Z_{1k}(\mathbf{a}) - \frac{k}{n} Z_{1n}(\mathbf{a}) \right\} = \frac{k(n-k)}{n} \sigma_n^2(\mathbf{a})$$

where

$$(10) \quad \sigma_n^2(\mathbf{a}) = \frac{1}{n-1} \sum_{i=1}^n (a^{(i)} - \bar{a}_n)^2.$$

General formulae for expectation and variance of simple linear rank statistics can be for example found in Hájek and Šidák [3], see Theorem c of Section 3.1.

A similar approach can be used in the case of statistics based on $S_{2k}(\mathbf{a})$, namely, we consider

$$(11) \quad Z_{2k}(\mathbf{a}) = \sum_{i=k+1}^n (a(1 + [nU_i]) - \bar{a}_n) \frac{i-k}{n}$$

instead of $S_{2k}(\mathbf{a})$. There is a simple but very important relationship between $Z_{1k}(\mathbf{a})$ and $Z_{2k}(\mathbf{a})$. Obviously

$$(12) \quad Z_{2k}(\mathbf{a}) = -\frac{1}{n} \sum_{i=k}^{n-1} Z_{1i}(\mathbf{a}) + \frac{n-k}{n} Z_{1n}(\mathbf{a}).$$

The same holds for $S_{1k}(\mathbf{a})$ and $S_{2k}(\mathbf{a})$, however, since $S_{1n}(\mathbf{a}) = 0$ the analog of (12) can be simplified to

$$(13) \quad S_{2k}(\mathbf{a}) = -\frac{1}{n} \sum_{i=k}^{n-1} S_{1i}(\mathbf{a}).$$

Finally note that for $k = 1, \dots, n$ we obtain

$$(14) \quad E S_{2k}(\mathbf{a}) = E Z_{2k}(\mathbf{a}) = 0$$

and

$$(15) \quad \text{Var } S_{2k}(\mathbf{a}) = \frac{n}{n-1} \text{Var} \left\{ Z_{2k}(\mathbf{a}) - \frac{(n-k)(n-k+1)}{2n^2} Z_{1n}(\mathbf{a}) \right\} \\ = v(n, k) \sigma_n^2(\mathbf{a})$$

where

$$(16) \quad v(n, k) = \frac{(n-k)(n-k+1)(2(n-k)+1)}{6n^2} - \frac{(n-k)^2(n-k+1)^2}{4n^3}.$$

2. Limit theorems

The main results are summarized below in Theorem 1. These results have been mentioned without any proof in Theorem 2 of [7] and employed there in a comparative simulation study.

Theorem 1. *Let X_1, \dots, X_n be i.i.d. random variables with common continuous distribution function F . Let assumptions (3) and (4) hold.*

(i) *As $n \rightarrow \infty$, for arbitrary $y \in \mathbb{R}$,*

$$(17) \quad P \left\{ \sqrt{2 \log \log n} \max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} \frac{1}{\sigma_n(\mathbf{a})} |S_{2k}(\mathbf{a})| \leq y + 2 \log \log n + \log \frac{\sqrt{3}}{4\pi} \right\} \rightarrow \exp(-2e^{-y})$$

where $v(n, k)$ is defined in (16).

(ii) *If moreover, as $n \rightarrow \infty$,*

$$(18) \quad \frac{n}{G} \rightarrow \infty \quad \text{and} \quad \frac{n^{2/(2+\eta)} \log n}{G} \rightarrow 0,$$

then for arbitrary $y \in \mathbb{R}$, as $n \rightarrow \infty$,

$$(19) \quad P \left\{ \sqrt{2 \log \frac{n}{G}} \max_{G < k < n-G} \frac{1}{\sqrt{w(n, G)}} \frac{1}{\sigma_n(\mathbf{a})} |S_{2, k+G}(\mathbf{a}) - 2S_{2k}(\mathbf{a}) + S_{2, k-G}(\mathbf{a})| \leq y + 2 \log \frac{n}{G} + \log \frac{\sqrt{3}}{4\pi} \right\} \rightarrow \exp(-2e^{-y})$$

where

$$(20) \quad w(n, G) = \frac{G(2G^2 + 1)}{3n^2} - \frac{G^4}{n^3}.$$

The theorem below is a result analogous to Theorem 1 but it assumes certain i.i.d. variables instead of the ranks.

Theorem 2. *Let X_1, \dots, X_n be i.i.d. random variables with $EX_1 = 0$, $\text{Var } X_1 = 1$ and $E|X_1|^{2+\eta} < \infty$ for some $\eta > 0$. For $k = 0, \dots, n$ denote*

$$(21) \quad \widehat{S}_k = \sum_{i=k+1}^n X_i \frac{i-k}{n}.$$

(i) As $n \rightarrow \infty$, for arbitrary $y \in \mathbb{R}$,

$$(22) \quad P \left\{ \sqrt{2 \log \log n} \max_{1 \leq k \leq n-1} \frac{1}{\sqrt{\widehat{v}(n, k)}} |\widehat{S}_k| \leq y + 2 \log \log n + \log \frac{\sqrt{3}}{4\pi} \right\} \rightarrow \exp(-2e^{-y})$$

where

$$(23) \quad \widehat{v}(n, k) = \frac{(n - k)(n - k + 1)(2(n - k) + 1)}{6n^2}.$$

(ii) If moreover condition (18) holds then for arbitrary $y \in \mathbb{R}$, as $n \rightarrow \infty$,

$$(24) \quad P \left\{ \sqrt{2 \log \frac{n}{G}} \max_{G < k < n-G} \frac{1}{\sqrt{\widehat{w}(n, G)}} |\widehat{S}_{k+G} - 2\widehat{S}_k + \widehat{S}_{k-G}| \leq y + 2 \log \frac{n}{G} + \log \frac{\sqrt{3}}{4\pi} \right\} \rightarrow \exp(-2e^{-y})$$

where

$$(25) \quad \widehat{w}(n, G) = \frac{G(2G^2 + 1)}{3n^2}.$$

Corollary. Convergence (22) remains true if $\widehat{v}(n, k)$ is replaced with $v(n, k)$ defined in (16). Similarly, convergence (24) remains true in the case that $\widehat{w}(n, G)$ is replaced with $w(n, G)$ defined in (20).

The following theorem, which can be considered to be an extension of Theorem 3 in [4], poses the crucial step in the proof of Theorem 1. It is an interlink between Theorem 1 and Theorem 2.

Theorem 3. Under assumptions of Theorem 1, as $n \rightarrow \infty$,

$$(26) \quad \max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} |D_k(\mathbf{a})| = \mathcal{O}_P \left(\max \left\{ n^{-1/2}, n^{-\eta/(2+\eta)} \right\} (1 + \log n) \right)$$

where

$$D_k(\mathbf{a}) = S_{2k}(\mathbf{a}) - \left(Z_{2k}(\mathbf{a}) - \frac{(n - k)(n - k + 1)}{2n^2} Z_{1n}(\mathbf{a}) \right).$$

Moreover, if assumption (18) holds then

$$(27) \quad \max_{G < k < n-G} \frac{1}{\sqrt{w(n, G)}} |D_{k+G}(\mathbf{a}) - 2D_k(\mathbf{a}) + D_{k-G}(\mathbf{a})| = o_P \left((\log n)^{-1/2} \right).$$

Remarks.

1. Compare our Theorem 1 to Theorem 2 of [4] and notice that mutual relationship between behaviour of maxima of cumulative and moving sums is different in our case. Also, note that a printing error occurred in Theorem 2 of [4] and condition (1.17) should accord with our condition (18). The error was dragged in [7] as well. In fact condition (18) can be slightly relaxed as follows

$$\frac{n}{G} \rightarrow \infty \quad \text{and} \quad \frac{n^{2/(2+\eta)} \log(n/G)}{G} \rightarrow 0.$$

However, (27) then holds with $o_P\left((\log(n/G))^{-1/2}\right)$ on the right-hand side.

2. Of course, convergence (19) and (24) also remains true if $w(n, G)$ and $\hat{w}(n, G)$ is replaced with $(2/3) \cdot (G^3/n^2)$.

3. Corollary of Theorem 2 can be further extended to

$$(28) \quad \tilde{S}_k = \sum_{i=k+1}^n (X_i - \bar{X}_n) \frac{i-k}{n}$$

where $\bar{X}_n = n^{-1} \sum_1^n X_i$. The assumption of zero mean can be relaxed in this case. The limit behaviour in question holds as well if the assumption of unit variance is relaxed and (21) or (28) is standardized by appropriate estimate of variance. See [5] for the discussion. Note in this context that the problem of estimation of mean and variance does not arise in the case of ranks because (5) and (10) are known.

4. The use of the above functionals of $S_{2k}(\mathbf{a})$ for testing changes in the location model as well as applicability of the limit theorems is — besides others — discussed in [7].

3. Proofs

Proof of Theorem 3. Using (12) and (13) we can write

$$D_k(\mathbf{a}) = -\frac{1}{n} \sum_{i=k}^{n-1} \Delta_i(\mathbf{a})$$

where

$$\Delta_i(\mathbf{a}) = S_{1i}(\mathbf{a}) - \left(Z_{1i}(\mathbf{a}) - \frac{i}{n} Z_{1n}(\mathbf{a}) \right).$$

It follows that

$$\begin{aligned}
 (29) \quad \max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} |D_k(\mathbf{a})| &= \frac{1}{n} \max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} \left| \sum_{i=k}^{n-1} \Delta_i(\mathbf{a}) \right| \\
 &\leq \frac{1}{n} \max_{1 \leq k \leq n-1} \sqrt{\frac{(n-k)^3}{v(n, k)}} \max_{1 \leq i \leq n-1} \sqrt{\frac{n}{i(n-i)}} |\Delta_i(\mathbf{a})|
 \end{aligned}$$

and

$$\begin{aligned}
 (30) \quad \max_{G < k < n-G} \frac{1}{\sqrt{w(n, G)}} |D_{k+G}(\mathbf{a}) - 2D_k(\mathbf{a}) + D_{k-G}(\mathbf{a})| \\
 &= \max_{G < k < n-G} \frac{1}{\sqrt{w(n, G)}} \frac{1}{n} \left| \sum_{i=k-G}^{k-1} \Delta_i(\mathbf{a}) - \sum_{i=k}^{k+G-1} \Delta_i(\mathbf{a}) \right| \\
 &\leq \frac{1}{\sqrt{w(n, G)}} \frac{2G}{\sqrt{n}} \max_{1 \leq i \leq n-1} \frac{1}{\sqrt{n}} |\Delta_i(\mathbf{a})|.
 \end{aligned}$$

Theorem 3 in [4] implies that

$$(31) \quad \max_{1 \leq i \leq n-1} \sqrt{\frac{n}{i(n-i)}} |\Delta_i(\mathbf{a})| = \mathcal{O}_P \left(\max \left\{ n^{-1/2}, n^{-\eta/(2+\eta)} \right\} (1 + \log n) \right)$$

and

$$(32) \quad \max_{1 \leq i \leq n-1} \frac{1}{\sqrt{n}} |\Delta_i(\mathbf{a})| = \mathcal{O}_P \left(\max \left\{ n^{-1/2}, n^{-\eta/(2+\eta)} \right\} \right).$$

Since

$$\frac{v(n, k)}{(n-k)^3} = \frac{1}{n^2} \left(1 + \frac{1}{n-k} \right) \left(\frac{1}{3} + \frac{1}{6(n-k)} - \frac{n-k+1}{4n} \right)$$

we can easily see that $v(n, k)/(n-k)^3$ is increasing in k independently of n . It follows that $v(n, k)/(n-k)^3 = \mathcal{O}(n^{-2})$ for arbitrary $1 \leq k \leq n-1$ and hence

$$\max_{1 \leq k \leq n-1} \sqrt{\frac{(n-k)^3}{v(n, k)}} = \mathcal{O}(n).$$

This along with (29) and (31) yields behaviour (26).

By condition (18) we have

$$\frac{1}{\sqrt{w(n, G)}} \frac{2G}{\sqrt{n}} = \mathcal{O} \left((n/G)^{1/2} \right) = o \left(n^{1/2-1/(2+\eta)} (\log n)^{-1/2} \right).$$

This along with (30) and (32) yields (27).

Proof of Theorem 2. Assertion (i) of Theorem 2 coincides with Theorem 1 in [5]. The corollary presents a result which is half a step to Theorem 2 of the paper and is proven ibidem. Theorem 2 of [5] shows the convergence for \tilde{S}_k defined in (28) and standardized by $\sqrt{v(n, k)}$.

To prove assertion (ii) assume at first that X_1, \dots, X_n are normally distributed and define a zero-mean standardized Gaussian process

$$\xi_n(t) = \frac{1}{\sqrt{\hat{w}(n, G)}} Y_{[G(t+1)]}, \quad 0 \leq t \leq \frac{n}{G} - 2,$$

where

$$(33) \quad Y_k = \hat{S}_{k+G} - 2\hat{S}_k + \hat{S}_{k-G} = \sum_{i=k-G+1}^k X_i \frac{i-k+G}{n} + \sum_{i=k+1}^{k+G-1} X_i \frac{k+G-i}{n}.$$

For $k \geq G$ and $k+l \leq n-G$ we have

$$\text{Cov}(Y_k, Y_{k+l}) = \text{Cov}(Y_G, Y_{G+l}) = R_{n,G}(l).$$

Hence, if $n/G \rightarrow \infty$ then $\xi_n(t)$ converges to a zero-mean standardized stationary Gaussian process $\{\xi(t), t \geq 0\}$ with autocovariance function

$$(34) \quad \rho(t) = \lim_{n/G \rightarrow \infty} \frac{R_{n,G}([Gt])}{\hat{w}(n, G)}.$$

Now investigate properties of $\rho(t)$. For $0 \leq l \leq G-2$ we get

$$\begin{aligned} R_{n,G}(l) &= \sum_{i=l+1}^G \frac{i}{n} \cdot \frac{i-l}{n} + \sum_{i=G+1}^{G+l} \frac{2G-i}{n} \cdot \frac{i-l}{n} + \sum_{i=G+l+1}^{2G-1} \frac{2G-i}{n} \cdot \frac{2G+l-i}{n} \\ &= \frac{1}{n^2} \left(\sum_{i=1}^{G-l} (i+l)i + \sum_{i=1}^l (G-i)(i+G-l) + \sum_{i=1}^{G-l-1} (G-l-i)(G-i) \right). \end{aligned}$$

Plugged in (34) it implies

$$\begin{aligned} \rho(t) &= \frac{3}{2} \left(\int_0^{1-t} (s+t) s \, ds + \int_0^t (1-s)(1+s-t) \, ds + \int_0^{1-t} (1-s)(1-s-t) \, ds \right) \\ &= 1 - \frac{3}{2}t^2 + \frac{3}{4}t^3 \end{aligned}$$

for $t \in \langle 0, 1 \rangle$. Further if $l \geq 2G$ then $R_{n,G}(l) = 0$ and hence $\rho(t) = 0$ for $t > 2$. Realize that

$$\max_{t \geq 0} \{\xi_n(t) - \xi(t)\} = o_P(1)$$

and apply Lemma 1 of [5] or the original Theorem 12.3.5 of [6] for $\xi(t)$ to obtain convergence (24) for normally distributed X_1, \dots, X_n .

By virtue of results of Einmahl [2] there are i. i. d. random variables $X_k^{(N)}$ with standard normal distribution such that

$$\max_{1 \leq k \leq n} k^{-1/(2+\eta)} \left| \sum_{i=1}^k (X_i - X_i^{(N)}) \right| = \mathcal{O}_P(1).$$

Hence, by (33) and (25) and according to condition (18) we have

(35)

$$\begin{aligned} \max_{G < k < n-G} \frac{1}{\sqrt{\widehat{w}(n, G)}} |Y_k - Y_k^{(N)}| &\leq \frac{1}{\sqrt{\widehat{w}(n, G)}} \frac{2G}{n} \max_{1 \leq k \leq n} \left| \sum_{i=1}^k (X_k - X_k^{(N)}) \right| \\ &= \mathcal{O}_P \left(\frac{n^{1/(2+\eta)}}{\sqrt{G}} \right) = o_P \left(\left(\log \frac{n}{G} \right)^{-1/2} \right) \end{aligned}$$

where $Y_k^{(N)}$ are obtained by replacing X_i with $X_i^{(N)}$ in (33). It concludes the proof of assertion (ii). Its extension in the corollary is straightforward since (34) and (35) hold true for $w(n, G)$ as well.

Proof of Theorem 1. According to (3), (4) and (7)

$$(36) \quad \frac{a(1 + [nU_i]) - \bar{a}_n}{\sigma_n(\mathbf{a})}, \quad i = 1, \dots, n,$$

are i. i. d. random variables with zero mean, with essentially unit variance, and with finite absolute moment of order $2+\eta$, $\eta > 0$. Thus, by Corollary to Theorem 2 we have

$$(37) \quad \max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} \frac{1}{\sigma_n(\mathbf{a})} |Z_{2k}(\mathbf{a})| = \mathcal{O}_P \left(\sqrt{\log \log n} \right)$$

and the rate in (37) cannot be improved.

It can be easily seen that $4n^4(n-k)^{-2}(n-k+1)^{-2}v(n, k)$ is decreasing for $0 < k < n$ and

$$\max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} \frac{(n-k)(n-k+1)}{2n^2} = \mathcal{O} \left(n^{-1/2} \right)$$

that is attained for $k = 1$. Thus, by central limit theorem we have

$$\max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} \frac{1}{\sigma_n(\mathbf{a})} \frac{(n-k)(n-k+1)}{2n^2} |Z_{1n}(\mathbf{a})| = \mathcal{O}_P(1).$$

This along with (37) implies that

$$(38) \quad \max_{1 \leq k \leq n-1} \frac{1}{\sqrt{v(n, k)}} \frac{1}{\sigma_n(\mathbf{a})} \left| Z_{2k}(\mathbf{a}) - \frac{(n-k)(n-k+1)}{2n^2} Z_{1n}(\mathbf{a}) \right| \\ = \mathcal{O}_P \left(\sqrt{\log \log n} \right)$$

and properly standardized l. h. s. in (37) has the same limit distribution as the standardized l. h. s. in (38). Now apply Theorem 3 to show convergence (17). Proof of convergence (19) is analogous.

REFERENCES

- [1] Csörgő M., Horváth L., *Limit Theorems in Change-Point Analysis*, John Wiley & Sons, 1997.
- [2] Einmahl U., *A useful estimate in the multidimensional invariance principle*, Probab. Theory Related Fields **76** (1987), 81–101.
- [3] Hájek J., Šidák Z., *Theory of Rank Tests*, Academia, Praha, 1967.
- [4] Hušková M., *Limit theorems for rank statistics*, Statist. Probab. Lett. **32** (1997), 45–55.
- [5] Jarušková D., *Testing appearance of linear trend*, J. Statist. Plann. Inference **70** (1998), 263–276.
- [6] Leadbetter M.R., Lindgren G., Rootzén H., *Extremes and Related Properties of Random Sequences and Processes*, Springer-Verlag, Heidelberg, 1983.
- [7] Slabý A., *Behaviour of Some Rank Statistics for Detecting Changes*, in Measuring Risk in Complex Stochastic Systems (proc. conf.), Berlin, 1999 (Franke J., Härdle W., Stahl G., Eds.), Springer-Verlag, Berlin, 2000, pp. 161–172.

DEPARTMENT OF PROBABILITY AND MATHEMATICAL STATISTICS, FACULTY OF MATHEMATICS AND PHYSICS, CHARLES UNIVERSITY, SOKOLOVSKÁ 83, 186 75 PRAGUE 8, CZECH REPUBLIC

(Received May 5, 2000, revised January 19, 2001)