

CONDITIONALLY TRIMMED SUMS  
FOR INDEPENDENT RANDOM VARIABLES

Yuji KASAHARA<sup>1</sup>

Department of Information Sciences, Ochanomizu University

(Received October 10, 1995)

Abstract

This paper studies conditionally trimmed sums for triangular arrays of independent random variables and generalize the Hahn-Kuelbs result for i.i.d. cases.

1990 Mathematics Subject Classification: Primary 60F17, 60F05; Secondary 62G30.

1. Introduction

Let  $X_1, X_2, \dots$  be a sequence of i.i.d. random variables with common nondegenerate distribution function  $F(x)$ , and let  $\{X_j^{(n)}\}_{j=1}^n$  ( $n \geq 1$ ) be the order statistics based on the sample  $\{X_1, \dots, X_n\}$  in the descending order in absolute value; i.e.,  $\{X_j^{(n)}\}_{j=1}^n$  is a rearrangement of  $\{X_j\}_{j=1}^n$  so that

$$|X_1^{(n)}| \geq |X_2^{(n)}| \geq \dots \geq |X_n^{(n)}|.$$

Let  $p_n$  ( $n = 1, 2, \dots$ ) be a sequence of positive, nondecreasing integers such that  $1 \leq p_n \leq n$  and define

$$\begin{aligned} S_n^{p_n} &= \sum_{j=p_n+1}^n X_j^{(n)} \\ &= \sum_{j=1}^n X_j - \sum_{j=1}^{p_n} X_j^{(n)}, \quad n = 1, 2, \dots \end{aligned}$$

$S_n^{p_n}$ , which is often referred to as the *trimmed sum*, is the  $n$ -th partial sum with the  $p_n$  largest samples trimmed and its asymptotic distribution as  $n \rightarrow \infty$  have been studied by many authors. (e.g. Stigler [4], Mori [3], etc. See Hahn-Kuelbs [2] for more references.)

A common interest of these authors is to get asymptotic normality theorems by deleting extreme samples even in the case where the tail probability of  $X_j$  is large (and hence the CLT no longer holds for the usual partial sums), and it is one of the important philosophy that it is desirable to retain as much data as possible. From this point of view, Hahn-Kuelbs ([2]) introduced the notion of *conditionally trimmed sums*; let  $a_n$  ( $n \geq 1$ ) be positive numbers and let

$$S_n^{p_n}(a_n) = \sum_{j=1}^n X_j - \sum_{j=1}^{p_n} X_j^{(n)} I(|X_j^{(n)}| > a_n).$$

Hence  $S_n^{p_n}(a_n)$  denotes the  $n$ -th partial sum with the  $p_n$  largest samples trimmed provided that they exceed  $a_n$  in magnitude. Hahn-Kuelbs ([2]) proved that for arbitrary nondegenerate distribution function  $F(x)$ ,

<sup>1</sup>Supported in part by Grant-in-Aid for Scientific Research(No. 07640284), Ministry of Education, Science and Culture.

we can find  $p_n, a_n, b_n$  and  $c_n$  such that  $p_n/n \rightarrow 0, a_n \rightarrow \infty$  as  $n \rightarrow \infty$  and that  $(1/c_n)\{S_n^{p_n}(a_n) - b_n\}$  converges in law to a Gaussian random variable. In fact, the existence of such sequences  $p_n, a_n, b_n$  and  $c_n$  is easy itself (see Remark in Section 3). So the significance of their result is offering an explicit procedure to find  $p_n, a_n, b_n$  and  $c_n$ . Now the aim of the present paper is to consider a similar problem for sums of triangular arrays of independent random variables. The author admits that all necessary ideas are found in [2], but we believe that their idea will be clearer under our formulation.

We give the main theorem in Section 2 and the i.i.d. case will be discussed in Section 3.

## 2. Main Theorem

For every  $n = 1, 2, \dots$ , let  $\{\xi_{n,1}, \xi_{n,2}, \dots, \xi_{n,k_n}\}$  be a sequence of independent random variables ( $k_n \rightarrow \infty$ ) and let  $\{\tilde{\xi}_{n,1}, \tilde{\xi}_{n,2}, \dots, \tilde{\xi}_{n,k_n}\}$  be a rearrangement of the samples  $\{\xi_{n,1}, \xi_{n,2}, \dots, \xi_{n,k_n}\}$  so that

$$|\tilde{\xi}_{n,1}| \geq |\tilde{\xi}_{n,2}| \geq \dots \geq |\tilde{\xi}_{n,k_n}|.$$

For every  $p$  ( $1 \leq p \leq k_n$ ) and  $a (> 0)$ , we define the conditionally trimmed sum  $S_n^p(a)$  as follows:

$$S_n^p(a) = \sum_{j=1}^n \xi_{n,j} - \sum_{j=1}^p \tilde{\xi}_{n,j} I(|\tilde{\xi}_{n,j}| > a).$$

Thus  $S_n^p(a)$  denotes the sum of  $\{\xi_{n,j}\}$  with the  $p$  largest samples deleted provided that they exceed the prescribed level  $a (> 0)$  in magnitude.

**Theorem 1.** Let  $\{a_n\}_{n=1}^\infty$  be a sequence of positive numbers and let

$$v_n = \sum_j \text{Var}(\xi_{n,j} I(|\xi_{n,j}| \leq a_n)).$$

If

$$(A.1) \quad \sum_j P(\varepsilon\sqrt{v_n} < |\xi_{n,j}| \leq a_n) \rightarrow 0 \quad (n \rightarrow \infty), \quad \text{for every } \varepsilon > 0,$$

then, for any  $p_n \in \mathbb{Z}_+$  ( $n \geq 1$ ) satisfying the condition

$$(A.2) \quad \frac{1}{p_n} \sum_j P(|\xi_{n,j}| > a_n) \rightarrow 0,$$

it holds

$$Z_n := \frac{1}{\sqrt{v_n}} [S_n^{p_n}(a_n) - m_n] \xrightarrow{\mathcal{L}} Z$$

where  $Z$  is an  $N(0, 1)$ -random variable and

$$m_n = \sum_j E[\xi_{n,j} I(|\xi_{n,j}| \leq a_n)].$$

(" $\xrightarrow{\mathcal{L}}$ " denotes the convergence in law.) A sufficient condition for (A.1) is

$$(A.3) \quad a_n = o(\sqrt{v_n}) \quad \text{as } n \rightarrow \infty.$$

**Proof.** Let

$$S_n^* = \sum_{j=1}^{k_n} \xi_{n,j} I(|\xi_{n,j}| \leq a_n).$$

We shall first see that

$$Z_n^* := \frac{1}{\sqrt{v_n}} \{S_n^* - m_n\} \xrightarrow{\mathcal{L}} Z.$$

This can easily be seen because  $S_1^*, S_2^*, \dots$  are usual sums of independent random variables. Indeed, it can be rewritten as follows;

$$Z_n^* = \sum_j (\zeta_{n,j} - E[\zeta_{n,j}])$$

where

$$\zeta_{n,j} = \frac{1}{\sqrt{v_n}} \xi_{n,j} I(|\xi_{n,j}| \leq a_n).$$

First, it holds

$$\begin{aligned} \sum_j P(|\zeta_{n,j}| > \varepsilon) &= \sum_j P\left(\frac{1}{\sqrt{v_n}} |\xi_{n,j}| I(|\xi_{n,j}| \leq a_n) > \varepsilon\right) \\ &= \sum_j P(|\xi_{n,j}| > \varepsilon \sqrt{v_n} \text{ and } |\xi_{n,j}| \leq a_n) \end{aligned}$$

which vanishes as  $n \rightarrow \infty$  by (A.1). Furthermore, by the definition of  $v_n$ , we also see that

$$\sum_j \text{Var}(\zeta_{n,j}) = 1.$$

Therefore, by the usual CLT, we obtain

$$(2.1) \quad Z_n^* \xrightarrow{\mathcal{L}} Z.$$

We next show that  $P(|Z_n - Z_n^*| > \varepsilon) \rightarrow 0$ , for every  $\varepsilon > 0$ :

$$\begin{aligned} P(|\tilde{\xi}_{n,p_n}| > a_n) &= P(\#\{j : |\xi_{n,j}|\} \geq p_n) = P\left[\frac{1}{p_n} \sum_j I[|\xi_{n,j}| > a_n] \geq 1\right] \\ &\leq E\left[\frac{1}{p_n} \sum_j I[|\xi_{n,j}| > a_n]\right] = \frac{1}{p_n} \sum_j P(|\xi_{n,j}| > a_n) \end{aligned}$$

which converges to 0 by (A.2). Thus we obtain

$$P(|\tilde{\xi}_{n,p_n}| > a_n) \rightarrow 0.$$

Since  $S_n^{k_n}(a_n) = S_n^*$  on the event  $\{|\tilde{\xi}_{n,p_n}| \leq a_n\}$  it holds that

$$(2.2) \quad P[Z_n = Z_n^*] \geq P[|\tilde{\xi}_{n,p_n}| \leq a_n] \rightarrow 1.$$

Combining (2.1) and (2.2) we obtain the assertion of the theorem.  $\square$

### 3. The Case of IID Random Variables

In this section we shall study the case where the triangular array  $\{\xi_{n,j}\}$  comes from an i.i.d. sequence: Let  $X_1, X_2, \dots$  be a nondegenerate i.i.d. random variables with the common distribution function  $F$  as in Introduction and let  $k_n = n$ ,  $\xi_{n,j} = X_j$ . Then, as a special case of Theorem 1 in the previous section we have the following theorem, which is a modification of Hahn-Kuelbs ([2]).

**Theorem 2.** Let  $a_n$  be a sequence tending to infinity and one of the following two conditions is satisfied.

$$(B.1a) \quad \frac{n}{a_n^2} \int_{|x| \leq a_n} x^2 dF(x) \longrightarrow \infty \quad (n \rightarrow \infty)$$

$$(B.1b) \quad n \int_{|x| > \varepsilon \sqrt{n}} dF(x) \longrightarrow 0 \quad (n \rightarrow \infty), \quad \text{for every } \varepsilon > 0.$$

Then, for any  $p_n$  ( $1 \leq p_n \leq n$ ) such that

$$(B.2) \quad \frac{n}{p_n} \int_{|x| > a_n} dF(x) \longrightarrow 0 \quad (n \rightarrow \infty),$$

it holds that

$$(2.1) \quad \frac{1}{b_n} \{S_n^{p_n}(a_n) - c_n\} \xrightarrow{\mathcal{L}} Z \quad (n \rightarrow \infty)$$

where

$$b_n = \sqrt{n \int_{|x| \leq a_n} x^2 dF(x)},$$

$$c_n = n \int_{|x| \leq a_n} x dF(x)$$

and  $Z$  is an  $N(0, 1)$ -random variable.

**Proof.** Since  $v_n$  in Theorem 1 is equal to  $n \int_{|x| < a_n} x^2 dF(x)$ , (B.1a) is equivalent to (A.3). If (B.1b) is satisfied instead, then it is easy to see that  $nP[|X_1| > \varepsilon \sqrt{v_n}] \rightarrow 0$ , which implies (A.1). Furthermore, (A.2) can be rewritten as (B.2). Thus we have the assertion.  $\square$

**Remarks.** (i) We can always find  $a_n$  ( $\uparrow \infty$ ) satisfying (B.1a) unless  $P[X_1 = 0] = 1$ . Also we can choose  $p_n$  so that  $p_n = o(n)$  in (B.2).

(ii) A sufficient condition for (B.1b) is  $E[X_1^2] < \infty$  since, in general,

$$nP[|X/\sqrt{n}| > \varepsilon] \leq \varepsilon^{-2} E[X^2 : |X| > \varepsilon \sqrt{n}].$$

(iii) As we mentioned in Introduction, if we are interested only in the existence of  $\{p_n, a_n, b_n, c_n\}$  satisfying (2.1), the proof is easier: For every fixed  $a > 0$ , choose  $\alpha > 0$  small enough so that  $|X_n^{[\alpha n]}| > a$  for all sufficiently large  $n$  (a.s.), which is possible due to the law of large numbers. Then, for all sufficiently large  $n$ ,  $S_n^{[\alpha n]}(a)$  is equal to  $Y_1 + \cdots + Y_n$  where  $Y_k = X_k I[|X_k| \leq a]$ . Therefore, under a suitable linear normalization,  $S_n^{[\alpha n]}(a)$  converges in law to an  $N(0, 1)$ -random variables. Now by a standard argument, choosing  $\alpha_n$  tending slowly enough, we can find  $a_n$  ( $\uparrow \infty$ ),  $b_n$  and  $c_n$  satisfying (2.1).

## References

- [1] S. Csörgő, E. Haeusler and D. M. Mason, The asymptotic distribution of trimmed sums, *Ann. Probab.* **16**(1988) 672-699.
- [2] M. G. Hahn and J. Kuelbs, Universal asymptotic normality for conditionally trimmed sums, *Statis. Probab. Letters* **7**(1989) 9-15.
- [3] T. Mori, On the limit distribution of lightly trimmed sums, *Math. Proc. Cambridge Philos. Soc.* **96**(1984) 507-516.
- [4] S. M. Stigler, The asymptotic distribution of the trimmed mean. *Ann. Statist.* **1**(1973) 472-477.