

# Central limit theorem for functions of weakly dependent variables

Jensen, Jens Ledet

Aarhus University, Department of Mathematical Sciences

Ny Munkegade, 8000 Aarhus C, Denmark

E-mail: jlj@imf.au.dk

## Result

We describe a central limit theorem for the sum  $S_n = \sum_{i=1}^n X_i$ ,  $X_i \in \mathbf{R}^p$ , where the  $X_i$ 's can be approximated by weakly dependent variables. The proof is for  $i \in \mathbf{Z}$ , but can be directly generalized to the case of a random field,  $i \in \mathbf{Z}^d$ . The weak dependency is formulated through a set of  $\sigma$ -algebras,  $D_j$ ,  $j \in \mathbf{Z}$ . The strong mixing coefficients of these are  $\alpha(k, l, d) = \sup |P(A_1 \cap A_2) - P(A_1)P(A_2)|$ , where the supremum is taken over sets  $A_i \in \sigma(D_j : j \in I_i)$ ,  $i = 1, 2$ , with  $|I_1| \leq k$ ,  $|I_2| \leq l$ , and the distance  $d(I_1, I_2)$  between the two sets  $I_1$  and  $I_2$  is at least  $d$ . The central limit theorem has applications in the study of nonhomogeneous hidden Markov chains (Jensen 2005).

**Theorem 1.** Assume that there exist  $\delta_0, \epsilon_0 > 0$ ,  $\delta_1, \delta_2 \geq 0$ ,  $\theta > \delta_1 + \delta_2 + \max\{(2 + \delta_0)/\delta_0, 1 + \delta_2, 2\}$  and constants  $c_0, c_1, c_2$  such that

- (1)  $EX_i = 0$ ,  $E|X_i|^{2+\delta_0} \leq c_0$ ,  $\text{Var}(a \cdot S_n) \geq \epsilon_0 n |a|^2 \quad \forall a \in \mathbf{R}^p$
- (2)  $\alpha(k, l, d) \leq c_1 k^{\delta_1} l^{\delta_2} \max\{1, d\}^{-\theta}$ ,
- (3)  $\forall m \in \mathbf{N} \exists X_j^m \in \sigma(D_k : d(k, j) \leq m) : E|X_j - X_j^m| \leq c_2 m^{-\theta}$ .

Then we have that  $\text{Var}(S_n)^{-1/2} S_n \xrightarrow{d} N_p(0, I)$  for  $n \rightarrow \infty$ .

(For the case of a random field,  $X_i$ ,  $i \in \mathbf{Z}^d$ , the lower bound on  $\theta$  is multiplied by  $\nu$ .) We divide the proof into a number of subsections. In the first two subsections we use truncation to reduce the problem to that of bounded variables. In the last section the method of Bolthausen (1982) is used for the bounded variables.

## Truncation

We use the truncation function  $T_M$  where  $T_M(x)$  equals  $x$  for  $|x| \leq M$  and equals  $Mx/|x|$  otherwise. Let  $Q_M(x) = x - T_M(x)$ . Using that  $EX_j = 0$  we write the sum  $S_n$  as  $S_n = S'_n + S''_n$ , with

$$(4) \quad S'_n = \sum_{j=1}^n [T_M(X_j) - E(T_M(X_j))] \quad \text{and} \quad S''_n = \sum_{j=1}^n [Q_M(X_j) - E(Q_M(X_j))],$$

and the idea is to prove that a CLT for  $S'_n$  for any fixed  $M$  implies a CLT for  $S_n$ .

In the lemma below we consider fixed values of  $j, k$  and fixed unit vectors  $a_j$  and  $a_k$ . We define  $U = a_j \cdot X_j$ ,  $U_M = a_j \cdot T_M(X_j)$  and  $\hat{U}_M = a_j \cdot Q_M(X_j)$ , and define  $V, V_M$  and  $\hat{V}_M$  similarly with  $j$  replaced by  $k$ .

**Lemma 2.** There exists a constant  $c_3$ , depending on  $c_1, c_2, \delta_0$  and  $\theta$  only, such that

$$\text{Cov}(U, V) \leq c_3 [c_0^{1/(2+\delta_0)} + c_0^{2/(2+\delta_0)}] \max\{1, d(j, k)\}^{-\kappa},$$

where  $\kappa = (\theta - \delta_1 - \delta_2)\delta_0/(2 + \delta_0) > 1$ .

*Proof.* Following Deo (1973) we expand  $\text{Cov}(U, V) = \text{Cov}(U_M + \hat{U}_M, V_M + \hat{V}_M)$  into four terms and bound each of these. Thus, using the simple bound  $|Q_M(x)| \leq (|x|/M)^{1+\delta_0} |x|$  we find

$$|\text{Cov}(U_M, \hat{V}_M)| \leq 2ME|\hat{V}_M| \leq 2Mc_0/M^{1+\delta_0} = 2c_0/M^{\delta_0}.$$

Similarly,  $|\text{Cov}(\hat{U}_M, V_M)| \leq 2c_0/M_0^\delta$ , and

$$|\text{Cov}(\hat{U}_M, \hat{V}_M)| \leq \{\text{Var}(\hat{U}_M) \text{Var}(\hat{V}_M)\}^{1/2} \leq c_0/M^{\delta_0},$$

since  $E\hat{U}_M^2 \leq E|X_j|^{2+\delta_0}/M^{\delta_0}$ . To bound  $\text{Cov}(U_M, V_M)$  we define  $U_M^m = a_j \cdot T_M(X_j^m)$  and  $V_M^m = a_k \cdot T_M(X_k^m)$ . Since  $|T_M(X_j) - T_M(X_j^m)| \leq |X_j - X_j^m|$  we find

$$|\text{Cov}(U_M - U_M^m, V_M)| \leq 2ME|U_M - U_M^m| \leq 2c_2Mm^{-\theta}, \quad |\text{Cov}(U_M^m, V_M - V_M^m)| \leq 2c_2Mm^{-\theta},$$

and therefore  $|\text{Cov}(U_M, V_M)| \leq |\text{Cov}(U_M^m, V_M^m)| + 4c_2Mm^{-\theta}$ . Finally, we use the classical bound (Ibragimov and Linnik, 1971 [17.2.1])  $|\text{Cov}(U_M^m, V_M^m)| \leq 4M^2\alpha(2m+1, 2m+1, \max\{0, d(j, k) - 2m\})$ . Putting all the terms together we obtain

$$|\text{Cov}(U, V)| \leq 5c_0/M^{\delta_0} + 4c_2Mm^{-\theta} + 4c_1M^2(2m+1)^{\delta_1+\delta_2} \max\{1, d(j, k) - 2m\}^{-\theta}.$$

Choosing  $m = \lfloor d(j, k)/3 \rfloor$  and  $M = c_0^{1/(2+\delta_0)}d(j, k)^{\kappa/\delta_0}$  we get the result of the lemma. □

**Lemma 3.** *Let  $S_n''$  be defined in (4). There exists a function  $b(M)$  with  $b(M) \rightarrow 0$  for  $M \rightarrow \infty$  such that for all unit vectors  $a$  and for all  $n$  we have  $\text{Var}(a \cdot S_n'')/n \leq b(M)$ .*

*Proof.* We use Lemma 2 with the random variable  $X$  replaced by  $Z = Q_M(X) - E(Q_M(X))$ . Let  $0 < \xi < \delta_0$  be so large that  $\kappa_1 = (\theta - \delta_1 - \delta_2)\xi/(2+\xi) > 1$ . Remembering the simple inequality  $|Q_M(x)| \leq (|x|/M)^\alpha|x|$  we have  $|ET_M(X_i)| = |EQ_M(X_i)| \leq c_0/M^{1+\delta_0}$  and  $E|Q_M(X_i)|^{2+\xi} \leq c_0/M^{\delta_0-\xi}$ . Thus, replacing  $\delta_0$  by  $\xi$  we use the bound

$$\begin{aligned} E|Z_i|^{2+\xi} &\leq 2^{1+\xi}\{E|Q_M(X_i)|^{2+\xi} + |EQ_M(X_i)|^{2+\xi}\} \\ (5) \quad &\leq 2^{1+\xi}\{(c_0/M^{\delta_0-\xi}) + (c_0/M^{1+\delta_0})^{2+\xi}\} = c_0(M, \xi). \end{aligned}$$

Furthermore, we can approximate  $Z$  by  $Q_M(X^m) - E(Q_M(X))$  with a mean error

$$\begin{aligned} E|Q_M(X) - Q_M(X^m)| &\leq E|X - X^m| + E|T_M(X) - T_M(X^m)| \\ &\leq 2E|X - X^m| \leq 2c_2m^{-\theta}. \end{aligned}$$

We can now use Lemma 2 with  $\delta_0$  replaced by  $\xi$ ,  $c_0$  replaced by  $c_0(M, \xi)$  and  $c_2$  replaced by  $2c_2$ . For some constant  $\tilde{c}_3$  and any unit vector  $a$  we then have

$$\text{Cov}(a \cdot Q_M(X_j), a \cdot Q_M(X_k)) \leq \tilde{c}_3\tilde{b}(M) \max\{1, d(j, k)\}^{-\kappa_1},$$

where  $\tilde{b}(M) = c_0(M, \xi)^{1/(2+\xi)} + c_0(M, \xi)^{2/(2+\xi)}$ . Writing  $\text{Var}(S_N'')$  as a double sum of covariances we find the result of the lemma with  $b(M) = \tilde{c}_3\tilde{b}(M)[3 + 2/(\kappa_1 - 1)]$ . From (5) we see that  $c_0(M, \xi)$  tends to zero, and therefore  $b(M)$  tends to zero as  $M$  tends to infinity. □

### Variance

Writing  $\text{Var}(a \cdot S_n)$  as a double sum of covariances we get directly from Lemma 2 the bound  $\text{Var}(a \cdot S_n) \leq c_4n$  with  $c_4 = c_3[c_0^{1/(2+\delta_0)} + c_0^{2/(2+\delta_0)}][3 + 2/(\kappa - 1)]$ .

From Lemma 3 we find

$$\begin{aligned} |\text{Var}(a \cdot S_n/\sqrt{n}) - \text{Var}(a \cdot S_n'/\sqrt{n})| &\leq 2|\text{Cov}(a \cdot S_n/\sqrt{n}, a \cdot S_n''/\sqrt{n})| + \text{Var}(a \cdot S_n''/\sqrt{n}) \\ &\leq 2\sqrt{c_4b(M)} + b(M) \rightarrow 0 \quad \text{for } M \rightarrow \infty. \end{aligned}$$

The assumption of a lower bound on the variance of  $S_n/\sqrt{n}$  in Theorem 1 therefore gives that

$$(6) \quad \text{Var}(a \cdot S_n/\sqrt{n})/\text{Var}(a \cdot S'_n/\sqrt{n}) \rightarrow 1 \quad \text{for } M \rightarrow \infty.$$

As in Ibragimov and Linnik (1971, page 346) we have from Lemma 3 and (6) that it suffices to prove a CLT for  $S'_n$  for fixed  $M$  to obtain a CLT for  $S_n$ .

**Proof of central limit theorem for  $S'_n$**

Let  $a$  be a fixed unit vector and define  $Y_i = a \cdot (T_M(X_i) - ET_M(X_i))$  and  $\tilde{S}_n = \sum_{i=1}^n Y_i = a \cdot S'_n$ .

We saw in the proof of Lemma 3 that  $|ET_M(X_i)| \leq c_0/M^{1+\delta_0}$  so that for  $M$  sufficiently large we have  $|Y_i| \leq M + 1$ . Furthermore, if we set  $Y_i^m = T_M(X_i^m) - ET_M(X_i)$  we have  $E|Y_i - Y_i^m| \leq E|X_i - X_i^m| \leq c_2m^{-\theta}$  and  $|Y_i^m| \leq M + 1$ .

We will use the method of proof from Bolthausen (1982). For this we need the following estimates.

**Lemma 4.** *There exists a constant  $c_4$  such that*

$$\text{Cov}(Y_j, Y_k) \leq c_4(M+1)^2 \max\{1, d(j, k)\}^{-\gamma}, \quad \text{Cov}(Y_j Y_k, Y_r Y_s) \leq c_4(M+1)^4 \max\{1, d(\{j, k\}, \{r, s\})\}^{-\gamma},$$

and  $\text{Cov}(Y_j, Y_k Y_r Y_s) \leq c_4(M + 1)^4 \max\{1, d(j, \{k, r, s\})\}^{-\gamma}$ , where  $\gamma = \theta - \delta_1 - \delta_2$ .

*Proof.* The proof is based on successively replacing  $Y_i$  by  $Y_i^m$  in the mean of a product of  $Y$ 's. Thus, using the second inequality as an example,  $|E(Y_j Y_k Y_r Y_s) - E(Y_j^m Y_k^m Y_r^m Y_s^m)| \leq 4(M + 1)^3 c_2 m^{-\theta}$ . From inequalities of this form we obtain

$$|\text{Cov}(Y_j Y_k, Y_r Y_s) - \text{Cov}(Y_j^m Y_k^m, Y_r^m Y_s^m)| \leq 2 \cdot 4(M + 1)^3 c_2 m^{-\theta}.$$

Next, the strong mixing implies (Ibragimov and Linnik, 1971 [17.2.1])

$$|\text{Cov}(Y_j^m Y_k^m, Y_r^m Y_s^m)| \leq 4(M + 1)^4 c_1 [2(2m + 1)]^{\delta_1 + \delta_2} \max\{1, d(\{j, k\}, \{r, s\}) - 2m\}^{-\theta}.$$

Combining the two inequalities and taking  $m = \lfloor d(\{j, k\}, \{r, s\})/4 \rfloor$  we obtain the second inequality of the lemma. □

For a number  $r$  we introduce the notation

$$S_{i,n} = \sum_{j=1, d(i,j) \leq r}^n Y_j \quad \text{and} \quad \alpha_n = \sum_{i=1}^n E(Y_i S_{i,n}),$$

where eventually  $r$  will be tending to infinity with  $n$ . From Lemma 4 we find that

$$\text{Var}(\tilde{S}_n/\sqrt{n}) - \frac{\alpha_n}{n} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1, d(i,j) > r}^n \text{Cov}(Y_i, Y_j) \leq \frac{4c_4(M + 1)^2}{(\gamma - 1)r^{\gamma-1}},$$

where  $\gamma - 1 = \theta - \delta_1 - \delta_2 - 1 > 0$ . Now, consider the case  $r = n^\omega$ , with  $\omega > 0$ . Then the right hand side above tends to zero. Since also we have from (6) that  $\text{Var}(\tilde{S}_n)$  is of the same order as  $\text{Var}(S_n)$ , and we have assumed the lower bound  $\epsilon_0 n$  for the latter, we find that  $\alpha_n$  has a similar lower bound and  $\text{Var}(\tilde{S}_n)/\alpha_n \rightarrow 1$ . We must therefore show that  $\bar{S}_n = \tilde{S}_n/\sqrt{\alpha_n}$  converges to a standard normal distribution.

**Proof of CLT for  $\bar{S}_n$**  We follow the proof of Bolthausen (1982), where the definitions of  $A_1$ ,  $A_2$  and  $A_3$  below can be found. These terms depend on the argument  $t$  of the characteristic function for

$\bar{S}_n$ . One needs to show that  $E(A_1 - A_2 - A_3) \rightarrow 0$ . Using the expression in Bolthausen (1982) and Lemma 4 we have with  $J_r = \{1 \leq i, i', j, j' \leq n : d(i, j) \leq r, d(i', j') \leq r\}$ ,

$$\begin{aligned}
 E|A_1|^2 &= \frac{t^2}{\alpha_n^2} \sum_{J_r} \text{Cov}(Y_i Y_j, Y_{i'} Y_{j'}) \\
 &= \frac{t^2}{\alpha_n^2} \left\{ \sum_{J_r, d(i, i') \geq 3r} \text{Cov}(Y_i Y_j, Y_{i'} Y_{j'}) + \sum_{J_r, d(i, i') < 3r} \text{Cov}(Y_i Y_j, Y_{i'} Y_{j'}) \right\} \\
 &\leq \frac{c_4 t^2 (M+1)^4}{\alpha_n^2} \left\{ n(2r+1)^2 2 \sum_{k=3r}^{\infty} (k-2r)^{-\gamma} + n(2r+1) \cdot 8 \cdot 6r \left(1 + \sum_{j=1}^{3r} j^{-\gamma}\right) \right\} \\
 (7) \quad &= O(t^2 (M+1)^4 r^2 / n) = O(t^2 (M+1)^4 / n^{1-2\omega}),
 \end{aligned}$$

where the last expression follows upon taking  $r = n^\omega$ . For the  $A_2$  term we have from Bolthausen (1982) and Lemma 4

$$\begin{aligned}
 E|A_2| &\leq \frac{nt^2(M+1)}{\alpha_n^{3/2}} \max_i \sum_{j,k=1, d(i,j) \leq r, d(i,k) \leq r}^n \text{Cov}(Y_j, Y_k) \\
 (8) \quad &\leq \frac{c_4 nt^2(M+1)^3}{\alpha_n^{3/2}} (2r+1) \left\{ 1 + 2 \sum_{j=1}^{2r+1} j^{-\gamma} \right\} = O(t^2 (M+1)^3 / n^{\frac{1}{2}-\omega}).
 \end{aligned}$$

Finally, we need to consider  $A_3 = \alpha_n^{-1/2} \sum_{i=1}^n Y_i \exp\{it(\bar{S}_n - \bar{S}_{i,n})\}$ , where  $\bar{S}_{i,n} = S_{i,n} / \sqrt{\alpha_n}$ . Considering the exponential part, and replacing all the variables by the approximating variables  $Y_j^m$ , we see that

$$E|\exp\{it(\bar{S}_n - \bar{S}_{i,n})\} - \exp\{it(\bar{S}_n^m - \bar{S}_{i,n}^m)\}| \leq \frac{|t|}{\sqrt{\alpha_n}} E|(\tilde{S}_n - S_{i,n}) - (\tilde{S}_n^m - S_{i,n}^m)| \leq \frac{|t|}{\sqrt{\alpha_n}} n c_2 m^{-\theta}.$$

Using this we obtain

$$|\text{Cov}(Y_i, \exp\{it(\bar{S}_n - \bar{S}_{i,n})\}) - \text{Cov}(Y_i^m, \exp\{it(\bar{S}_n^m - \bar{S}_{i,n}^m)\})| \leq c_2 m^{-\theta} + (M+1) \frac{|t|}{\sqrt{\alpha_n}} n c_2 m^{-\theta}.$$

From the mixing we have (Ibragimov and Linnik, 1971 [17.2.1])

$$|\text{Cov}(Y_i^m, \exp\{it(\bar{S}_n^m - \bar{S}_{i,n}^m)\})| \leq 4(M+1)c_1(2m+1)^{\delta_1} n^{\delta_2} (r-2m)^{-\theta}.$$

Taking  $m = r/3$  and combining the two bounds we find for some constant  $c_5$

$$|\text{Cov}(Y_i, \exp\{it(\bar{S}_n - \bar{S}_{i,n})\})| \leq c_5 \{ (M+1)n^{\delta_2} r^{-\theta+\delta_1} + r^{-\theta} [1 + (M+1)|t|n/\sqrt{\alpha_n}] \}.$$

Since  $EY_i = 0$  the terms in  $A_3$  are covariances, and taking  $r = n^\omega$  we therefore have the bound

$$(9) \quad |EA_3| = O((M+1)[n^{-\omega\theta+1} + n^{-\omega(\theta-\delta_1)+\delta_2+\frac{1}{2}}]).$$

Thus, if we choose  $\omega$  such that  $\omega < \frac{1}{2}$ ,  $\omega\theta - 1 > 0$  and  $\omega(\theta - \delta_1) - \delta_2 - \frac{1}{2} > 0$ , which is possible from the assumption on  $\theta$  in Theorem 1, we see that all of (7), (8), and (9) tend to zero.

**REFERENCES**

Bolthausen, E. (1982): On the central limit theorem for stationary mixing random variables. *Ann. Probab.*, **10**, 1047–1050.  
 Deo, C. M. (1973): A note on empirical processes of strong-mixing sequences. *Ann. Probab.*, **1**, 870–875.

Ibragimov, I. A. and Linnik, Yu. V. (1971): *Independent and Stationary Sequences of Random Variables*. Wolters-Noordhoff Publishing, Groningen.

Jensen, J.L. (2005): Context dependent DNA evolutionary models. Research Report, No. 458, Department of Mathematical Sciences, University of Aarhus.