

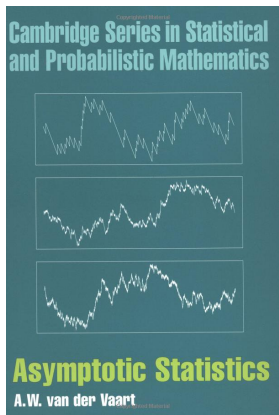
# Large Sample Theory Part I

**Song Xi Chen, Xiaojun Song**

Department of Business Statistics and Econometrics  
Center for Statistical Science  
Peking University

July 13, 2021

Van der Vaart, A. W. (1998). Asymptotic statistics. Cambridge university press.



新书推荐 | 《渐近统计》(一部介绍渐近统计的权威教材)

世界图书出版公司北京公司 1月28日

《渐近统计》是一部介绍渐近统计的权威教材,内容实用而且数学理论论述严谨。



《渐近统计》

(新书即将上市)

作者: [荷兰] 范德瓦特 (A. W. van der Vaart)

原名: Asymptotic Statistics

页数: 462

定价: 129.00元

装帧: 平装

ISBN: 9787519254025

出版社: 世界图书出版公司

书中除了介绍渐近统计的核心理论——似然推断、M估计、渐近效率、U统计和秩过程等内容,书中还涉及该领域的最新研究论题,如半参数模型、自助法、经验过程等其他应用。本书各章有习题。

## 渐近统计 中文修订版

范德瓦特 著

世界图书出版公司北京公司版权所有

张慧铭译1-5,18-19,24-25章;王鹏飞译11-15章;王然译9,20,21章;王宏浩译10,16-17章;刘云喆译23章;沈铂涵译6-7章;魏浩宇译25章;陈

February 16, 2020

### 摘要

van der Vaart教授是荷兰皇家艺术与科学院院士,也是当今具有国际影响力的统计学家。在经验过程和半参数统计方面,他做出过重要的基础性贡献。他的《弱收敛与经验过程》(Weak Convergence and Empirical Processes)与Jon Wellner合著)和《渐近统计》(Asymptotic Statistics, 1998)已成为该领域的经典文献和必读书目。近年来他的兴趣主要集中在非参数贝叶斯推断、统计遗传学、复杂网络等应用领域。由于必须以数据分析和应用为导向,统计学理论往往呈现出碎片化的倾向。他的工作却能高屋建瓴,从半参数统计的理论基础,到2000年后贝叶斯方法的频率学派解释,再到机器学习模型的研究,无不让人自然地领会到统计学理论的统一和谐之美。

Jordan教授访谈:《渐近统计》是加州大学伯克利分校的教材,书中寄托于经验过程思想的很多统计推断过程,比如M-估计,包含了最大似然估计、经验风险最小化等等。这是一本发人深省的书。Jordan教授期望学生们至少去读过本书中的一些章节而且他认为这些将要进入机器学习领域的学生需要最终读完这些书籍,更需要将它们都读上至少三遍——第一遍能够复述了,第二遍去尝试做相关的仿真实验和理论推导,第三遍看到之后会发现都是显而易见的。

## References

- **Serfling, R. (1980) Approximation Theorems in Mathematical Statistics. Wiley.** [Second Textbook]
- Ferguson, T. S. (1997). A course in large sample theory. CRC.
- Proschan, M. A., & Shaw, P. A. (2018). Essentials of probability theory for statisticians. CRC.
- White, H. (2000). Asymptotic Theory for Econometricians: Revised Edition. Emerald.
- Vaart, A. W., & Wellner, J. A. (1996). Weak convergence and empirical processes: with applications to statistics. Springer.
- Bosq, D. (2012). Nonparametric statistics for stochastic processes: estimation and prediction 2ed. Springer.
- Bijma, F., Jonker, M., & Van der Vaart, A. (2017). An introduction to mathematical statistics. Amsterdam University Press.

# Introduction: Data, Statistic and Its Distribution

Let  $Z_1, \dots, Z_n$  be a set of data, and  $T_n = T(Z_1, \dots, Z_n)$  be a Statistic.

## Task

- $T_n$  can be an estimator to a parameter  $\theta$  s.t.  $\hat{\theta}_n = T_n$ ;
- or  $T_n$  can be a test statistic for a hypothesis:  $H_0 : \theta \in \Omega_0$ .
- A key task of Inference is to derive/find the distribution of  $\hat{\theta}_n$ , say  $F_{\hat{\theta}_n}$ ,

$$F_{\hat{\theta}_n}(x) = P(\hat{\theta}_n \leq x) \quad \text{for } x \in R^d.$$

# Why Asymptotic Statistics?

## Dilemma and Benefits

- Exact **fixed sample** (non-asymptotical) analysis on statistics is HARD.
- But, letting  $n \rightarrow \infty$  simplifies things and amazingly quality approximation to  $F_{\hat{\theta}_n}(x)$  may be obtained.

The use of asymptotic approximation is two-fold.

## Van der Vaart's book:

- It can be used for **asymptotical inference** (find approximate **confidence regions** and **testing**).
- Approximations can be used theoretically to study the **quality (efficiency)** of statistical inference procedures.

# Stochastic Convergence

Let  $\{X_n\}$  be a sequence of  $\mathbb{R}^p$  random vectors and  $d(x, y)$  be a Euclidean distance in  $\mathbb{R}^p$ ,  $\{X_n, X\}$  is defined on a common  $(\Omega, \mathcal{A}, P)$ .

- **Almost-Sure Convergence** ( $X_n \xrightarrow{\text{a.s.}} X$ ): The sequence  $\{X_n\}$  is said to *converge almost surely* to  $X$

if  $d(X_n, X) \rightarrow 0$  with probability one:  $P(\lim_n d(X_n, X) = 0) = 1$ .

$X_n \xrightarrow{\text{a.s.}} X =$  "100% **sure** + 100% **accurate**."

- **Convergence in Probability** ( $X_n \xrightarrow{P} X$ ): A sequence of random variables  $\{X_n\}$  is said to *converge in probability* to  $X$  if for all  $\varepsilon > 0$ ,

$$\lim_n P(d(X_n, X) < \varepsilon) = 1.$$

$X_n \xrightarrow{P} X =$  "100% **sure** + **not 100%** **accurate**."

- **Convergence in  $r$ th mean** ( $X_n \xrightarrow{L^r} X$ ): A sequence of random variables  $\{X_n\}$  is said to *converge in  $r$ th Mean* to  $X$  if

$$\lim_n \mathbb{E}[d(X_n, X)]^r = 0.$$

# Convergence in Distribution (Weak Convergence)

Convergence in distribution is a type of weak convergence for r.v.s, it is **the most useful** stochastic convergence in statistical inferences.

**It does not require that the  $\{X_n, X\}$  are in the same  $(\Omega, \mathcal{A}, P)$ .**

**Definition 2.1** (Convergence in Distribution,  $X_n \xrightarrow{d} X$ )

Let  $\{F_n\}$  be a sequence of distribution functions for a sequence of r.v.s  $\{X_n\}$ . Then  $X_n$  is said to **converge in distribution** to a r.v.  $X$  (with distribution  $F$ ) if

$$\lim_n F_n(x) = F(x), \forall x \in \mathcal{C}_F$$

where  $\mathcal{C}_F := \{x | F(x) \text{ is continuous in } x\}$ .

## Remark 1

1. The discontinuous set  $\mathcal{C}_F^c$  is **the countable set**.
2. The spaces of  $\{X_n, X\}$  can differ as  $\xrightarrow{d}$  focus on the cdfs (free of  $(\Omega, \mathcal{A})$ ).

# A Lemma

Known in mathematical analysis: For any continuous function  $F$ ,  $F$  is u.f. continuous on  $[-M, M]$ .

## Lemma 2.2

*If  $F$  is a continuous distribution function, then  $F$  is uniformly continuous in  $\mathbb{R}$ .*

## Remark 2

*The lemma is valid for  $F$  not necessarily a cdf as long as  $F$  has limits at  $\pm\infty$ .*



## Theorem 2.3

Suppose that  $X_n \xrightarrow{d} X$  for a random variable  $X$  with  $F_n$  and  $F$  being the continuous distribution functions of  $X_n$  and  $X$ . Then  $\sup_x |F_n(x) - F(x)| \rightarrow 0$  as  $n \rightarrow \infty$ .

- The proof uses the Covering Method: divide a diverging interval  $[-M, M]$  by a partition  $\{[x_{i-1}, x_i]\}_{i=0}^{K+1}$  of equal width  $\delta$  (except the last one).
- Define  $\Delta_n = \max_{i=0}^{K+1} \{|F_n(x_i) - F(x_i)|\}$ . As  $K$  is finite for any given  $M$ ,  $\lim_{n \rightarrow \infty} \Delta_n = 0$ .
- Use the previous lemma.

## Definition 2.4

A sequence of random variables  $\{X_n\}$  is said to be asymptotically normal (AN) with “mean”  $\mu_n$  and “variance”  $\sigma_n^2$  ( $\sigma_n^2 > 0$  when  $n$  is sufficiently large) if  $\frac{X_n - \mu_n}{\sigma_n} \xrightarrow{d} N(0, 1)$ , denoted as  $X_n \sim AN(\mu_n, \sigma_n^2)$ .

- $\mu_n$  and  $\sigma_n^2$  are not necessarily to be the mean and variance of  $X_n$ . In fact, the mean and variance of  $X_n$  may not exist.
- What  $X_n \xrightarrow{d}$  is unknown !
- That  $X_n \xrightarrow{d} N(\mu_n, \sigma_n^2)$  is obviously wrong !
- Nevertheless,

$$\sup_t |P(X_n \leq t) - P(N(\mu_n, \sigma_n^2) \leq t)| \rightarrow 0$$

as  $n \rightarrow \infty$ .

## Some facts:

If  $X_n$  is AN( $\mu_n, \sigma_n^2$ ), then

- 1  $X_n$  is AN( $\bar{\mu}_n, \bar{\sigma}_n^2$ ) if and only of (iff)  $\frac{\bar{\sigma}_n}{\sigma_n} \rightarrow 1$  and  $\frac{\mu_n - \bar{\mu}_n}{\sigma_n} \rightarrow 0$ .
- 2  $a_n X_n + b_n$  is AN( $\bar{\mu}_n, \bar{\sigma}_n^2$ ) iff  $a_n \rightarrow 1$  and  $\frac{\mu_n(a_n - 1) + b_n}{\sigma_n} \rightarrow 0$ .

## Definition 2.5 (Multivariate AN)

A seq of random vectors  $\{X_n\}$  istb AN with “mean”  $\mu_n$  and “variance”  $\Sigma_n$  ( $\Sigma_n$  is positive definite when  $n$  is sufficiently large) , if

$$a'X_n \text{ is AN } (a'\mu_n, a'\Sigma_n a), \forall a \in \mathbb{R}^P.$$

# How to establish $X_n \xrightarrow{d} X$ ?

## Lemma 2.6 (Portmanteau Lemma)

For any random vectors  $X_n$  and  $X$ , the following the following statements are equivalent.

- 1  $X_n \xrightarrow{d} X$ ;
- 2  $Ef(X_n) \rightarrow Ef(X)$  for any  $f \in C_B$ ;
- 3  $Ef(X_n) \rightarrow Ef(X)$  for any  $f \in C_{B,Lip}$ ;
- 4  $\liminf Ef(X_n) \geq Ef(X)$  for all nonnegative, continuous  $f$ ;
- 5  $\liminf P(X_n \in G) \geq P(X \in G)$  for any **open set**  $G$ ;
- 6  $\limsup P(X_n \in F) \leq P(X \in F)$  for any closed set  $F$ ;
- 7  $P(X_n \in B) \rightarrow P(X \in B)$  for any Borel set  $B$  with  $P(X \in \delta B) = 0$ , where  $\delta B = \overline{B} - \overset{\circ}{B}$  is the boundary of  $B$ .
- 8 (Lévy's continuity theorem) Let  $\{X_n\}$  and  $X$  be r.v.s in  $\mathbb{R}^d$ . Then

$$X_n \xrightarrow{d} X \text{ iff } \phi_{X_n}(t) \rightarrow \phi(t) \quad \forall t \in \mathbb{R}^d.$$

# The Proof of Portmanteau Lemma

(i)  $\Rightarrow$  (ii):

*W.O.L.G.* Assume  $\sup |f(x)| \leq 1$ . First assume the df of  $X$ ,  $F_X$  is continuous. As  $X_n \xrightarrow{d} X$ , we have:

$$\lim_{n \rightarrow \infty} P(X_n \in I) = P(X \in I), \quad \text{for every rectangle } I \subset \mathbb{R}^k$$

On the other hand.  $\forall \epsilon > 0$ , choose  $I$  large enough s.t.  $P(X \in I^c) < \epsilon$ . Then we partition  $I$  into many small and non-overlapped rectangles s.t.  $I = \cup_{j=1}^K I_j$ . Choose a  $x_j \in I_j$ . Now let:

$$f_\epsilon(x) = \sum_{j=1}^K f(x_j) \mathbb{I}(x \in I_j)$$

Obviously, we can choose  $K$  large enough to have  $|f(x) - f_\epsilon(x)| < \epsilon$ ,  $x \in I$ . Thus,

# The Proof of Portmanteau Lemma

$$\begin{aligned} |Ef(X_n) - Ef_\epsilon(X_n)| &\leq E [|f(X_n) - f_\epsilon(X_n)|\mathbb{I}(X_n \in I)] \\ &\quad + E [|f(X_n) - f_\epsilon(X_n)|\mathbb{I}(X_n \in I^c)] \\ &\leq \epsilon + 2P(X_n \in I^c) \quad (\text{Recall that } \sup |f(x)| \leq 1) \end{aligned} \tag{1}$$

Similarly,

$$|Ef(X) - Ef_\epsilon(X)| \leq \epsilon + 2P(X \in I^c) < 3\epsilon \tag{2}$$

Considering,

$$\begin{aligned} |Ef_\epsilon(X_n) - Ef_\epsilon(X)| &\leq \left| \sum_{j=1}^K f(x_j) E [\mathbb{I}(X_n \in I_j) - \mathbb{I}(X \in I_j)] \right| \\ &\leq \sum_{j=1}^K |f(x_j)| |P(X_n \in I_j) - P(X \in I_j)| \rightarrow 0 \end{aligned} \tag{3}$$

due to  $X_n \xrightarrow{d} X$  and  $K$  is finite. Then considering (1), (2), and (3), we have  $\lim_{n \rightarrow \infty} Ef(X_n) = Ef(X)$ .

# The Proof of Portmanteau Lemma

If  $F_X$  is not continuous:

As  $F_X$  is right continuous and monotonous,

$$C_{F_x} := \{F_X \text{ is continuous at } x\}$$

is dense in  $\mathbb{R}^k$ . Then we could choose the vertices of rectangles  $I_j$  in  $C_{F_x}$ , then repeat the early proof.

# The Proof of Portmanteau Lemma

(iii)  $\Rightarrow$  (v)\*:

For every open set  $G$  there exists a sequence of Lipschitz functions with  $0 \leq f_m \uparrow 1_G$ . For instance  $f_m(x) = (md(x, G^c)) \wedge 1$ . For every fixed  $m$ ,

$$\liminf_{n \rightarrow \infty} P(X_n \in G) \geq \liminf_{n \rightarrow \infty} E f_m(X_n) = E f_m(X)$$

As  $m \rightarrow \infty$  the right side increases to  $P(X \in G)$  by the monotone convergence theorem.

(v)  $\Leftrightarrow$  (vi)\*:

Because a set is open if and only if its complement is closed, this follows by taking complements.



# The Proof of Portmanteau Lemma

(v) + (vi)  $\Rightarrow$  (vii)\*:

Let  $\mathring{B}$  and  $\bar{B}$  denote the interior and the closure of a set, respectively. By (iv),

$$P(X \in \mathring{B}) \leq \liminf P(X_n \in \mathring{B}) \leq \limsup P(X_n \in \bar{B}) \leq P(X \in \bar{B}),$$

by (v). If  $P(X \in \delta B) = 0$ , then left and right side are equal, whence all inequalities are equalities. The probability  $P(X \in B)$  and the limit  $\lim P(X_n \in B)$  are between the expressions on left and right and hence equal to the common value.

(vii)  $\Rightarrow$  (i)\*:

Every cell  $(-\infty, x]$  such that  $x$  is a continuity point of  $x \mapsto P(X \leq x)$  is a continuity set.

(ii)  $\Leftrightarrow$  (iv): Exercise.

## Theorem 2.7 (Mapping)

Let  $g : \mathbb{R}^k \mapsto \mathbb{R}^m$  be continuous at every point of a set  $\mathcal{C}_g$  such that  $P(X \in \mathcal{C}_g) = 1$ , then

- 1 If  $X_n \xrightarrow{d} X$ , then  $g(X_n) \xrightarrow{d} g(X)$
- 2 If  $X_n \xrightarrow{P} X$ , then  $g(X_n) \xrightarrow{P} g(X)$
- 3 If  $X_n \xrightarrow{a.s.} X$ , then  $g(X_n) \xrightarrow{a.s.} g(X)$

# The Proof of Mapping Theorem

(i). Consider using (vi) in Portmanteau Lemma: for a closed set  $F$ , define  $g^{-1}(F) = \{x_n \in g^{-1}(F)\} = \{g(X_n) \in F\}$ . We have:

$$g^{-1}(F) \subset \overline{g^{-1}(F)} \subset g^{-1}(F) \cup C_g^C \quad (4)$$

The last  $\subset$  is because of  $\forall x \in \overline{g^{-1}(F)}$ , here exist a sequence  $\{x_m\}_{m \geq 1} \subset g^{-1}(F)$  s.t.  $x_m \rightarrow x$ .

- If  $x \in C_g$ , then  $g(x_m) \rightarrow g(x) \in F$  due to  $g(x_m) \in F$  and  $F$  is a close set. Thus,  $x \in g^{-1}(F)$ .
- If  $x \notin C_g$ , (4) is evident.

# The Proof of Mapping Theorem

Then,

$$\begin{aligned}\limsup P(g(X_n) \in F) &\leq \limsup P\left(X_n \in \overline{g^{-1}(F)}\right) && \text{(by (4) left)} \\ &\leq P(X \in \overline{g^{-1}(F)}) && \text{(by Portmanteau Lemma (iv))} \\ &\leq P(X \in g^{-1}(F)) + P(X \notin C_g) && \text{(by (4) again)} \\ &= P(g(X) \in F)\end{aligned}$$

Apply Portmanteau Lemma (iv)  $\Rightarrow$  (i), we have  $g_n(X_n) \xrightarrow{d} g(X)$ .

# The Proof of Mapping Theorem

(ii)\*. Fix arbitrary  $\varepsilon > 0$ . For each  $\delta > 0$  let  $B_\delta$  be the set of  $x$  for which there exists  $y$  with  $|x - y| < \delta$ , but  $|g(x) - g(y)| > \varepsilon$ . If  $X \notin B_\delta$  and  $|g(X_n) - g(X)| > \varepsilon$ , then  $|X_n - X| \geq \delta$ . Consequently,

$$P(|g(X_n) - g(X)| > \varepsilon) \leq P(X \in B_\delta) + P(|X_n - X| \geq \delta)$$

The second term on the right converges to zero as  $n \rightarrow \infty$  for every fixed  $\delta > 0$ . Because  $B_\delta \cap C \downarrow \emptyset$  by continuity of  $g$ , the first term converges to zero as  $\delta \downarrow 0$ .

# Example of Mapping Theorem

- If  $X_n \xrightarrow{d.} X \sim N(0, 1)$ , then  $X_n \xrightarrow{d.} \chi_1^2$ .

- If

$$\begin{pmatrix} X_n \\ Y_n \end{pmatrix} \xrightarrow{d.} N_2 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, I_2 \right)$$

then  $X_n/Y_n \xrightarrow{d.}$  Cauchy, whose distribution has p.d.f.

$$f(x) = \frac{1}{\pi(1+x^2)}, \quad f_{\mu,\sigma}(x) = \frac{1}{\pi\sigma} \frac{1}{1 + ((x-\mu)/\sigma)^2}$$

# Example of Mapping Theorem

- $S_n^2 = n^{-1} \sum X_i^2 - \bar{X}^2$ . Let  $Y_i = (X_i, X_i^2)^\top$ . By LLN,

$$\frac{1}{n} \sum_{i=1}^n Y_i = \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n X_i \\ \frac{1}{n} \sum_{i=1}^n X_i^2 \end{pmatrix} \rightarrow \begin{pmatrix} \mu \\ \mu_2 \end{pmatrix} \quad \text{w.p.1}$$

Let  $g(x, y) = y - x^2$ . Then by Mapping Theorem:

$$S_n^2 = g\left(\bar{X}, n^{-1} \sum_{i=1}^n X_i^2\right) \xrightarrow{P.} \mu_2 - \mu^2 = \sigma^2$$

Apply Mapping Theorem again:  $S_n \xrightarrow{P.} \sigma$ .

# Example of Mapping Theorem

- If  $X_n \xrightarrow{d.} N_p(\mu, \Sigma)$ , then for any constant matrix  $C \in \mathbb{R}^{m \times p}$ ,

$$CX_n \xrightarrow{d.} N_m(C\mu, C\Sigma C^\top)$$

- If  $X_n$  is  $AN(\mu, b_n^2\Sigma)$ , then:

$$\frac{\|X_n - \mu\|}{b_n} \xrightarrow{d.} \text{a limit r.v.}$$

In fact, since  $(X_n - \mu)/b_n \xrightarrow{d.} N_p(0, \Sigma)$ , by Mapping Theorem,

$$\frac{(X_n - \mu)^\top (X_n - \mu)}{b_n^2} \xrightarrow{d.} N_p^\top(0, \Sigma) N_p(0, \Sigma)$$

Therefore,

$$\frac{\|X_n - \mu\|}{b_n} \xrightarrow{d.} \sqrt{N_p^\top(0, \Sigma) N_p(0, \Sigma)}$$



## Theorem 2.8

Let  $X_n, X$  and  $Y_n$  be random vectors. Then

- 1  $X_n \xrightarrow{\text{a.s.}} X$  implies  $X_n \xrightarrow{P} X$ ;
- 2  $X_n \xrightarrow{P} X$  implies  $X_n \xrightarrow{d} X$
- 3  $X_n \xrightarrow{P} c$  ( $c$  is a constant) if and only if  $X_n \xrightarrow{d} c$ ;
- 4 if  $X_n \xrightarrow{d} X$  and  $d(X_n, Y_n) \xrightarrow{P} 0$ , then  $Y_n \xrightarrow{d} X$ ;
- 5 if  $X_n \xrightarrow{d} X$  and  $Y_n \xrightarrow{P} c$  for a constant  $c$ , then  $(X_n, Y_n) \xrightarrow{d} (X, c)$
- 6 if  $X_n \xrightarrow{P} X$  and  $Y_n \xrightarrow{P} Y$ , then  $(X_n, Y_n) \xrightarrow{P} (X, Y)$

(i)\*. The sequence of sets  $A_n = \cup_{m \geq n} \{\|X_m - X\| > \varepsilon\}$  is decreasing for every  $\varepsilon > 0$  and decreases to the empty set if  $X_n(\omega) \rightarrow X(\omega)$  for every  $\omega$ . If  $X_n \xrightarrow{as} X$ , then  $P(\|X_n - X\| > \varepsilon) \leq P(A_n) \rightarrow 0$ .

(iv) We have:

$$A_n := |Ef(X_n) - Ef(Y_n)| \leq E\{|f(X_n) - f(Y_n)|\mathbb{I}(\|X_n - Y_n\| \leq \epsilon)\} \\ + E\{|f(X_n) - f(Y_n)|\mathbb{I}(\|X_n - Y_n\| > \epsilon)\}$$

Only for bounded Lipschitz function  $f \in C_{B,Lip}$  and  $\epsilon > 0$ ,

$$A_n \leq L\epsilon P(\|X_n - Y_n\| \leq \epsilon) + 2 \sup \|f(x)\| P(\|X_n - Y_n\| > \epsilon)$$

Thus,  $Ef(X_n) - Ef(Y_n) \rightarrow 0$ , and

$$Ef(Y_n) = Ef(X_n) + Ef(Y_n) - Ef(X_n) \rightarrow Ef(X)$$

which implies  $Y_n \xrightarrow{d} X$  due to Portmanteau Lemma (iii).

(ii)  $X_n = X + (X_n - X)$ . Since  $X \xrightarrow{d.} X$  and  $X_n - X \xrightarrow{p.} 0$ , we have  $X_n \xrightarrow{d.} X$  by using (iv).

(iii) " $\Rightarrow$ " is from (ii). For " $\Leftarrow$ ", note that:

$$\{\|X_n - c\| \geq \epsilon\} = \{X_n \in \text{ball}(c, \epsilon)^C\}$$

where  $\text{ball}(c, \epsilon) = \{x : \|x - c\| < \epsilon\}$  is open, so  $\text{ball}(c, \epsilon)^C$  is a closed set. From Portmanteau Lemma,

$$\begin{aligned} \limsup P(\|X_n - c\| \geq \epsilon) &= \limsup P(X_n \in \text{ball}(c, \epsilon)^C) \\ &\leq P(c \in \text{ball}(c, \epsilon)^C) = 0 \end{aligned}$$

Hence,  $P(\|X_n - c\| \geq \epsilon) \rightarrow 0$ , and thus  $X_n \xrightarrow{p.} c$ .

(v)  $(X_n, Y_n) = (X_n, c) + (X_n, Y_n) - (X_n, c)$ . Since

$$(X_n, Y_n) - (X_n, c) = (0, Y_n - c) \xrightarrow{P.} (0, 0)$$

From (iv), we only need to show  $(X_n, c) \xrightarrow{d.} (X, c)$ .

For any bounded continuous function  $f : (x, y) \mapsto f(x, y)$ , the marginal function  $f_m : x \mapsto f(x, c)$  is also bounded continuous. As  $X_n \xrightarrow{d.} X$ ,

$$Ef(X_n, c) = Ef_m(X_n) \rightarrow Ef_m(X) = Ef(X, c)$$

Hence  $(X_n, c) \xrightarrow{d.} (X, c)$ .

(vi) As  $\|(X_1, Y_1) - (X_2, Y_2)\| \leq \|X_1 - X_2\| + \|Y_1 - Y_2\|$ ,

$$\begin{aligned} P(\|(X_n, Y_n) - (X, Y)\| > \epsilon) &\leq P(\|X_n - X\| > \epsilon/2) + P(\|Y_n - Y\| > \epsilon/2) \\ &\rightarrow 0 \end{aligned}$$

From this theorem, we know that:

Marginal Convergence in Prob  $\Rightarrow$  Joint Convergence in Prob.

The converse is also true due to Mapping Theorem. So,

Marginal Convergence in Prob  $\Leftrightarrow$  Joint Convergence in Prob.

Nonetheless, Marginal Convergence in Dist  $\nRightarrow$  Joint Convergence in Dist, although the converse is true via Mapping Theorem.

Copula.

- If  $X_n$  is  $AN(\mu, b_n^2 \Sigma)$  with  $b_n \rightarrow 0$ , then  $X_n \xrightarrow{d} \mu$  and  $X_n \xrightarrow{P} \mu$ .

# Slutsky's theorem: algebraic operations for con. in dist.

Application of Theorem 2.8 and Continuous mapping theorem.

It is named after a Russian mathematical statistician/economist: Slutsky.

## Lemma 2.9 (Slutsky)

Let  $X_n, X$  and  $Y_n$  be random vectors or variables. If  $X_n \xrightarrow{d} X$  and  $Y_n \xrightarrow{d} c$  (or  $Y_n \xrightarrow{P} c$ ) for a constant  $c$ , then

- 1  $X_n + Y_n \xrightarrow{d} X + c$
- 2  $Y_n X_n \xrightarrow{d} cX$
- 3  $Y_n^{-1} X_n \xrightarrow{d} c^{-1}X$  provided  $c \neq 0$

The (3) is valid for matrices  $Y_n$  and  $c$  and vectors  $X_n$  provided  $c \neq 0$  is understood as  $c$  being invertible, because taking an inverse is also continuous.

# Slusky's Theorem

From (v) of Theorem 2.8, if  $X_n \xrightarrow{d.} X$  and  $Y_n \xrightarrow{p.} c$ , then

$$(X_n, Y_n) \xrightarrow{d.} (X, c)$$

Apply Mapping Theorem,

$$g(X_n, Y_n) \xrightarrow{d.} g(X, c)$$

for almost surely continuous  $g$ .

Choose  $g(x, y) = x + y$ ,  $x \times y$ ,  $x/y$  then we obtain Slusky's Theorem.

# Example

## T-statistic:

Considering  $Y_1, \dots, Y_n$  i.i.d.  $F$  with  $EY_1 = 0$  and  $EY_1^2 = \sigma^2 > 0$ . Define T-statistic:  
 $t_n := \sqrt{n}\bar{Y}/S_n$ .

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \frac{n}{n-1} \left( \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}^2 \right)$$
$$\xrightarrow{p.} 1 \cdot (EY_i^2 - (EY_i)^2) = EY_i^2$$

Hence  $S_n \xrightarrow{p.} \sqrt{EY_i^2}$ .

From the CLT for IID data.

$$\sqrt{n}\bar{Y} \xrightarrow{d.} N(0, EY_i^2)$$

Hence,

$$\frac{\sqrt{n}\bar{Y}}{S_n} \xrightarrow[\text{Slusky}]{d.} \frac{N(0, EY_i^2)}{\sqrt{EY_i^2}} \stackrel{d.}{=} N(0, 1)$$



## Definition 2.10

A sequence of random vectors  $\{X_n\}$  is said to be stochastically bounded or tight if  $\forall \varepsilon > 0$ ,  $\exists M_\varepsilon > 0$ , s.t.  $\sup_n P(\|X_n\| > M_\varepsilon) < \varepsilon$ . Denote  $X_n = O_p(1)$ .

## Theorem 2.11 (Prohorov's Theorem)

- 1 If  $X_n \xrightarrow{d} X$ , then  $X_n$  is tight.
- 2 If  $X_n$  is tight, then  $\exists$  a subsequence  $\{X_{n_i}\}$  ( $n_i > 1$ ), s.t.  $X_{n_i} \xrightarrow{d} X$  as  $n_i \rightarrow \infty$ .

## Remark 3

A single random vector is tight.

# Finite random vectors are tight.

Since the distribution function of  $\|X\|$ ,  $F(x)$ , satisfies:

$$\lim_{x \rightarrow -\infty} F(x) = 0, \quad \lim_{x \rightarrow \infty} F(x) = 1$$

So we can choose  $M_\epsilon$  s.t.  $F(M_\epsilon) - F(-M_\epsilon)$  is as close to 1 as possible, then,

$$P(\|X\| > M_\epsilon) = 1 - F(M_\epsilon) + F(-M_\epsilon)$$

can as small as possible. So we could choose  $M_\epsilon$  large enough s.t.

$$P(\|X\| > M_\epsilon) < \epsilon$$

Similarly, we have:

## Remark 4

Any finite collection of r.v.  $\{X_i\}_{i=1}^K$  is tight.

# The proof of Prohorov's Theorem

(i) only. As  $X$  is tight, one can choose  $M_\epsilon$  properly s.t.

$$P(\|X\| \geq M_\epsilon) < \epsilon, \quad \forall \epsilon > 0$$

By Portmanteau Lemma,

$$\limsup P(\|X_n\| \geq M_\epsilon) \leq P(\|X\| \geq M_\epsilon) < \epsilon$$

Hence  $\exists N$ , s.t.  $\forall n \geq N$ ,  $P(\|X_n\| \geq M_\epsilon) \leq 2\epsilon$ .

Note the  $\{X_i\}_{i=1}^{N-1}$  is tight. By modifying  $M_\epsilon$ , we can obtain

$$P(\|X_n\| \geq M_\epsilon) < \epsilon, \quad \forall n \in \mathbb{N}_+$$

So  $\{X_n\}_{n \geq 1}$  is tight.

## Remark 5

*(ii) is an extended vision of "Bounded sequences must have a convergent subsequence" in Mathematical Analysis.*

## Definition 2.12

A sequence of random vectors  $\{X_n\}$  is called

- 1 **stochastically of order  $\{a_n\}$** , with  $\{a_n\}$  being a sequence of constants,  
if  $\frac{X_n}{a_n} = O_p(1)$ .
- 2 **stochastically small**, if  $X_n \xrightarrow{P} 0$  and denoted as  $X_n = o_p(1)$ .
- 3 **stochastically of smaller order  $\{a_n\}$** , if  $\frac{X_n}{a_n} = o_p(1)$ , write as  $X_n = o_p(a_n)$ .

# Stochastic $o$ and $O$

In a previous example,  $Y_n - \mu = \bar{X} - \mu = O_p(n^{-1/2})$ , and also  $Y_n - \mu = o_p(1)$ . However,  $Y_n - \mu = O_p(n^{-1/2})$  is more accurate.

In general, considering a set of stochastic quantity  $\{T_n\}_{n \geq 1}$ ,  $n$  is typically the sample size. Let

$$\mu_n = \mathbb{E}T_n, \quad \sigma_n^2 = \text{var}(T_n)$$

if they exist. From the Chebysev inequality,

$$\mathbb{P}(\sigma_n^{-1}|T_n - \mu_n| > M) \leq \frac{\text{var}(T_n)}{(M\sigma_n)^2} = \frac{1}{M^2}$$

Then  $\forall \epsilon > 0$ , we can choose  $M_\epsilon$  s.t.  $M_\epsilon^{-2} < \epsilon$ , and

$$\mathbb{P}(\sigma_n^{-1}|T_n - \mu_n| > M_\epsilon) < \epsilon$$

which implies  $\sigma_n^{-1}(T_n - \mu_n) = O_p(1)$ , and  $T_n - \mu_n = O_p(\sigma_n)$ . This is a typical way to find the stochastic order of a quantity if we can labour out  $\sigma_n^2$ .

# Basic Rules of Stochastic $o$ and $O$

There are rules of calculus on  $o$  and  $O$  symbols, which we apply without comment. For instance,

## Some facts:

- 1  $o_P(1) + o_P(1) = o_P(1)$
- 2  $o_P(1) + O_P(1) = O_P(1)$
- 3  $O_P(1)o_P(1) = o_P(1)$
- 4  $(1 + o_P(1))^{-1} = O_P(1)$
- 5  $O_P(1) + O_P(1) = O_P(1)$
- 6  $o_P(O_P(1)) = O_P(o_P(1)) = o_P(1)$

## Remark 6

*The rules should be read from left to right.*

## Lemma 2.13

Let  $R : \mathbb{R}^k \mapsto \mathbb{R}$  be a real function with  $R(0) = 0$ . Let  $\{X_n\}$  be a sequence of r.v.s with values in  $\text{dom}(R)$  s.t.  $X_n \xrightarrow{P} 0$ . Then,  $\forall p > 0$ ,

- 1 if  $R(h) = o(\|h\|^p)$  as  $h \rightarrow 0$ , then  $R(X_n) = o_P(\|X_n\|^p)$ ;
- 2 if  $R(h) = O(\|h\|^p)$  as  $h \rightarrow 0$ , then  $R(X_n) = O_P(\|X_n\|^p)$ .

## Remark 7

The function  $R(\cdot)$  may not be continuous other than  $\{0\}$  in  $\text{dom}(R)$ .

(i). Let:

$$g(h) = \begin{cases} R(h)/\|h\|^p & , \text{ for } h \neq 0 \\ 0 & , \text{ for } h = 0 \end{cases}$$

Then  $R(X_n) = \|X_n\|^p g(X_n)$ . We will show that  $g(X_n) \xrightarrow{P} 0$ .

Note that  $g(h)$  is continuous at 0, and  $P\left(\lim_{n \rightarrow \infty} X_n = 0\right) = 1$ , by Mapping Theorem:

$$g(X_n) \xrightarrow{P} g(0) = 0$$

So  $R(X_n) = o_p(\|X_n\|^p)$ .



(ii). As  $g(h) = R(h)/\|h\|^p$  is bounded near  $x = 0$  (because  $R(h) = O(\|h\|^p)$  and  $g(0) = 0$ ),

$$\exists M, \delta > 0, \quad |g(h)| \leq M, \quad \forall |h| < \delta$$

so,

$$\{\omega : |g(X_n(\omega))| > M\} \subset \{\omega : \|X_n(\omega)\| > \delta\}$$

Thus,

$$P(|g(X_n)| > M) \leq P(\|X_n\| > \delta) \rightarrow 0$$

In fact,  $\forall \epsilon > 0, \exists N$ , s.t.  $\forall n \geq N, P(\|X_n\| > \delta) < \epsilon$ . Thus

$$P(|g(X_n)| > M) < \epsilon, \quad \forall n \geq N$$

This implies  $g(X_n) = O_p(1)$ .

# An Applied Example

Considering  $X_1, \dots, X_n$  i.i.d.  $F(\mu, \sigma^2)$  with  $EX^4 < \infty$ . Let,

$$S_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n [(X_i - \mu) - (\bar{X} - \mu)]^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X} - \mu)^2$$

Since  $EX^4 < \infty$ , by a standard CLT,

$$\frac{n^{-1} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2}{\sqrt{\text{var}((X_i - \mu)^2)/n}} \xrightarrow{d.} N(0, 1)$$

which implies,

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2 \right) \xrightarrow{d.} N(0, v^2)$$

where  $v^2 = \text{var}((X_i - \mu)^2)$ . On the other hand,  $\bar{X} - \mu = O_p(n^{-1/2})$ ,  $\sqrt{n}(\bar{X} - \mu)^2 = O_p(n^{-1/2}) = o_p(1)$ . By Slutsky Theorem,

$$\sqrt{n}(S_n - \sigma^2) = \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2 \right) - \sqrt{n}(\bar{X} - \mu)^2 \xrightarrow{d.} N(0, v^2)$$

## Chapter 2: Characteristic Functions (cf)

We need tools for derivating weak convergence.

### Definition 3.1

For any random vector  $X$  with distribution function  $F$ , its cf is

$$\begin{aligned}\phi_X(t) &= \mathbb{E}[e^{itX}] = \int e^{itx} dF(x) \\ &= \mathbb{E}[\cos tX] + i\mathbb{E}[\sin tX] \quad \text{for any } t \in \mathbb{R}.\end{aligned}$$

The moment generating function (MGF) is

$$M_X(t) = \mathbb{E}[e^{tX}].$$

The cf is a frequency domain view of a distribution  $F$ , it fully characterizes  $F$ .

# Properties of Characteristic Functions

- i)  $|\phi_X(t)| \leq \phi_X(0) = 1$
- ii)  $\overline{\phi_X(t)} = \phi_X(-t)$
- iii)  $\phi_X(t)$  is uniformly continuous on  $\mathbb{R}$
- iv)  $\overline{\phi_X}$ ,  $|\phi_X|^2$  and  $\text{Re}(\phi_X)$  are c.f.s of  $-X$ ,  $X - Y$  for  $X, Y$  i.i.d.  $F$ , and  $(F_X + F_{-X})/2$  respectively.
- v) If  $\exists t_0 \neq 0$  s.t.  $|\phi_X(t_0)| = 1$ , then  $\exists a \in \mathbb{R}$  and  $a \neq 0$  s.t.  $P(X \in \{a + jh : j \in \mathbb{Z}\}) = 1$ , so  $X$  is a lattice random vector.

# Properties of Characteristic Functions (Cont.)

- (vi) If  $F$  is absolutely continuous,  $\lim_{|t| \rightarrow \infty} |\phi_X(t)| = 0$
- (vii) Two random vectors  $X$  and  $Y$  in  $\mathbb{R}^d$  are equal in distribution, denoted as  $X \stackrel{d}{=} Y$ , iff  $\phi_X(t) = \phi_Y(t) \quad \forall t \in \mathbb{R}^d$ .
- (viii) (Fourier Inversion) If  $\phi_X$  is integrable, i.e.  $\phi_X \in \mathcal{L}^1(\mathbb{R})$ ; then  $F$  is continuous with density:

$$f(x) = \frac{1}{2\pi} \int e^{-itx} \phi_X(t) dt.$$

The characteristic function (cf) of  $X$  a  $p$ -dimensional random vector is defined by

$$\phi_X(\mathbf{t}^T) = \mathbb{E}e^{it^T X} = \int_{\mathbb{R}^p} e^{it^T x} dF_X(x) \quad \text{for any } t \in \mathbb{R}^d \quad (5)$$

where  $F_X$  is the cumulative distribution function.

## *Remark 8*

*The multivariate CFs inherit the properties of univariate CFs.*

# Properties of Multivariate Characteristic Functions

- (i)  $\varphi_{\mathbf{X}}(\mathbf{t})$  exists for all  $\mathbf{t} \in \mathbb{R}^d$  and is continuous.
- (ii)  $\varphi_{\mathbf{X}}(0) = 1$  and  $|\varphi_{\mathbf{X}}(\mathbf{t})| \leq 1$  for all  $\mathbf{t} \in \mathbb{R}^d$ .
- (iii) For a scalar  $b \neq 0$ ,  $\varphi_{\mathbf{X}/b}(\mathbf{t}) = \varphi_{\mathbf{X}}(\mathbf{t}/b)$ ;  
For a vector  $\mathbf{c}$ ,  $\varphi_{\mathbf{X}+\mathbf{c}}(\mathbf{t}) = \exp\{i\mathbf{t}^T \mathbf{c}\} \varphi_{\mathbf{X}}(\mathbf{t})$ .
- (iv) For  $\mathbf{X}$  and  $\mathbf{Y}$  independent,  $\varphi_{\mathbf{X}+\mathbf{Y}}(\mathbf{t}) = \varphi_{\mathbf{X}}(\mathbf{t})\varphi_{\mathbf{Y}}(\mathbf{t})$ .
- (v) If  $E\|\mathbf{X}\| < \infty$ ,  $\dot{\varphi}_{\mathbf{X}}(\mathbf{t})$  exists and is continuous and  $\dot{\varphi}_{\mathbf{X}}(0) = i\boldsymbol{\mu}^T$ ,  
where  $\boldsymbol{\mu} = E\mathbf{X}$ .
- (vi) If  $E\|\mathbf{X}\|^2 < \infty$ ,  $\ddot{\varphi}_{\mathbf{X}}(\mathbf{t})$  exists and is continuous and  
 $\ddot{\varphi}_{\mathbf{X}}(0) = -E\mathbf{X}\mathbf{X}^T$ .
- (vii) If  $\mathbf{X}$  is  $N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\varphi_{\mathbf{X}}(\mathbf{t}) = \exp\{i\mathbf{t}^T \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}\}$ .

# Lévy-Cramér's Theorem

Lévy's continuity theorem

## Theorem 3.2

Let  $\{X_n\}$  and  $X$  be random vectors in  $\mathbb{R}^d$ . Then  $X_n \xrightarrow{d} X$  iff  $\phi_{X_n}(t) \rightarrow \phi(t) \quad \forall t \in \mathbb{R}^d$ .

**Proof:** " $\Rightarrow$ " is by Portmanteau Lemma (ii):  $Ef(X_n) \rightarrow Ef(x)$  for  $\forall f \in C_B$ , " $\Leftarrow$ " can be seen in P14 of vdv.

## Remark 9

*It provides another way for establishing weak convergence!*



# Example

Suppose  $X_1, \dots, X_n$  i.i.d.  $\text{Poisson}(\lambda)$  for fixed  $\lambda > 0$ . Then the characteristic function of  $X_i$  is:

$$\phi_X(t) = \exp\{\lambda(e^{it} - 1)\}$$

Let  $\bar{X} = n^{-1} \sum_{i=1}^n X_i$ , check the c.f.

$$\begin{aligned}\phi_{\frac{\bar{X}-\lambda}{\sqrt{\lambda/n}}}(\lambda) &= \exp\{-it\sqrt{n\lambda}\} \phi_{\bar{X}}(t/\sqrt{\lambda/n}) = \exp\{-i\sqrt{n\lambda}\} \phi_X^n(t/\sqrt{n\lambda}) \\ &= \exp\{-it\sqrt{n\lambda}\} \exp\left\{n\lambda \left(e^{\frac{it}{\sqrt{n\lambda}}} - 1\right)\right\} \\ &= \exp\left\{-it\sqrt{n\lambda} + n\lambda \left(\frac{it}{\sqrt{n\lambda}} + \frac{i^2 t^2}{2n\lambda} + o\left(\frac{1}{n\lambda}\right)\right)\right\} \\ &= \exp\{-t^2/2 + o(1)\} \rightarrow \exp\{-t^2/2\}\end{aligned}$$

Therefore,

$$\frac{\bar{X} - \lambda}{\sqrt{\lambda/n}} \xrightarrow{d.} N(0, 1)$$

# Example

## Weak Law of Large Numbers (WLLN):

Let  $Y_1, \dots, Y_n$  be IID r.v. with  $\phi_Y(t)$  being differential at  $t = 0$  and  $i\mu = \phi'(0)$ , then:

$$\bar{Y} \xrightarrow{P.} \mu$$

**Proof:** As  $\phi(0) = 1$  and  $\phi'(0)$  exists at 0,

$$\phi_Y(t) = 1 + t\phi'(0) + o(t), \quad \text{as } t \rightarrow 0$$

$$\begin{aligned}\phi_{\bar{Y}}(t) &= \phi_Y^n(t/n) = \left(1 + \frac{t}{n}\phi'(0) + o\left(\frac{t}{n}\right)\right)^n \\ &= \left(1 + \frac{it\mu}{n} + o\left(\frac{1}{n}\right)\right)^n \rightarrow e^{it\mu} = \phi_\mu(t)\end{aligned}$$

Hence,  $\bar{Y} \xrightarrow{d.} \mu$ , and  $\bar{Y} \xrightarrow{P.} \mu$ .

# Moments and Expansion of CFs

If r.v.  $X$ 's  $r$ -th moment exists, then  $\phi_X(t)$  is  $r$ -th order differentiable and,

$$\phi_X^{(r)}(t) = \int (ix)^r e^{itx} dF(x) = E\{(iX)^r e^{itX}\}$$

which implies  $\phi_X^{(r)}(0) = i^r EX^r$ .

Conversely, if  $\phi_X^{(r)}(0)$  exists for an even  $r$ , then  $X$  has finite  $r$ -th absolutely moment.

## Theorem 3.3

If  $E|X|^r < \infty$ , then

$$\phi_X(t) = \sum_{j=0}^r \frac{(it)^j}{j!} EX^j + o(|t|^r).$$

## Example: CLT

Suppose  $X_1, \dots, X_n$  i.i.d.  $F$ , with  $\mu = EX$  and  $\sigma^2 = EX^2 < \infty$ . Let  $S_n = \sum_{i=1}^n X_i$ , then

$$\frac{S_n - n\mu}{\sqrt{n\sigma^2}} \xrightarrow{d.} N(0, 1)$$

Proof:

$$\phi_{X-\mu}(t) = 1 + \frac{(it)^2\sigma^2}{2} + o(t^2)$$

$$\begin{aligned}\phi_{\frac{\bar{x}-\mu}{\sqrt{\sigma^2/n}}} &= \phi_{X-\mu}^n\left(\frac{t}{\sigma\sqrt{n}}\right) \\ &= \left(1 + \frac{1}{2}\left(\frac{t}{\sigma\sqrt{n}}\right)^2\sigma^2 + o\left(\frac{t^2}{\sigma^2n}\right)\right)^n \rightarrow e^{-t^2/2}\end{aligned}$$

Then the result comes from Lévy-Cramér's Theorem.

## Remark 10

$$\phi_{X-\mu}(t) = 1 + \frac{1}{2}(it)^2\sigma^2 + \dots + \frac{(it)^r}{r!}m_r + o(|t|^r)$$

if  $E|X|^r < \infty$  for  $r > 2$ , where  $m_j = E(X - \mu)^j$  is  $j$ -th central moment. Then,

$$\phi_{\frac{\bar{X}-\mu}{\sqrt{\sigma^2/n}}} = \left(1 - \frac{1}{2} \frac{t^2}{n} - \frac{1}{6} \frac{it^3}{n^{3/2}} \left(\frac{m_3}{\sigma}\right)^3 + \frac{1}{24} \frac{t^4}{n^2} \left(\frac{m_4}{\sigma}\right)^4 + \dots\right)^n$$

Higher order expansion of c.f.  $\Rightarrow$  **Edgeworth Expansion**.

## Remark 11

The c.f. determines all the moments of  $X$ , but  $\{m_r := E(X)^r\}_{r=1}^n$  cannot determine the law of  $X$ . This is the famous Moment problem. Carleman's condition

$$\sum_{r=1}^{\infty} m_{2r}^{-\frac{1}{2r}} = +\infty$$

gives a sufficient condition for the determinacy of  $X$ .

# Cumulants (Semi-Invariants)

## Definition 3.4

The cumulants  $\kappa_j$ 's are obtained from a power series expansion of the cumulant generating function of a r.v.  $X$ :

$$K_X(t) := \log \phi_X(t) = \sum_{j \geq 1} \frac{(it)^j}{j!} \kappa_j =: \log \left\{ 1 + \sum_{j \geq 1} \frac{1}{j!} m_j (it)^j \right\}.$$

Matching, we have  $\kappa_1 = m_1 = EX$  and

$$\kappa_2 = m_2 - m_1^2 = E(X - EX)^2 =: c_2,$$

$$\kappa_3 = m_3 - 3m_1m_2 + 2m_1^3 = E(X - EX)^3 =: c_3,$$

$$\kappa_4 = m_4 - 4m_1m_3 - 3m_2^2 + 12m_1^2m_2 - 6m_1^4 = c_4 - 3c_2^2.$$

The higher order ( $j > 3$ ) cumulants are different from central moment.

# Cumulants (Semi-Invariants)

Considering  $X_1, \dots, X_n$  i.i.d.  $F(\mu, \sigma^2)$ ,  $Y_i = (X_i - \mu)/\sigma$ . The cumulants for  $Y_i$  are:

$$\kappa_1 = 0, \quad \kappa_2 = 1, \quad \kappa_3 = \frac{E(X_i - \mu)^3}{\sigma^3}, \quad \kappa_4 = \frac{E(X_i - \mu)^4}{\sigma^4}$$

where  $\kappa_3$  is called the Skewness of  $X$ ,  $\kappa_4$  is called Kurtosis. Then,

$$\phi_Y(t) = \exp \left\{ \sum_{j \geq 1} \frac{(it)^j}{j!} \kappa_j \right\} = \exp \left\{ -\frac{t^2}{2} + \sum_{j \geq 3} \frac{(it)^j}{j!} \kappa_j \right\}$$
$$\phi_{\frac{\bar{X} - \mu}{\sqrt{n\sigma^2}}} = \phi_Y^n(t) = \exp \left\{ -\frac{t^2}{2} + \frac{(it)^3}{3! \sqrt{n}} \kappa_3 + \frac{(it)^4}{4! n} \kappa_4 + \dots \right\}$$

## Chapter 3: Central Limit Theorems (for ind. r.v.s)

Unlike the classical CLT, in this section, we will explore the general CLT when the random variables is independent but not identically distributed.

### Definition 4.1

For each  $n \geq 1$ , let  $\{X_{n1}, X_{n2}, \dots, X_{nk_n}\}$  be a collection of random vectors on a probability space  $(\Omega_n, \mathcal{F}_n, P_n)$  s.t.  $X_{n1}, \dots, X_{nk_n}$  are independent with  $k_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Then  $\{X_{nj} : 1 \leq j \leq k_n\}_{n \geq 1}$  is called a double array of independent random vectors.

Notations:

$$S_n = \sum_{j=1}^{k_n} X_{nj}, \quad \alpha_{nj} = E(X_{nj}), \quad \alpha_n = \sum_{j=1}^{k_n} E(X_{nj}) = \sum_{j=1}^{k_n} \alpha_{nj},$$

$$\sigma_{nj}^2 = \text{Var}(X_{nj}), \quad \sigma_n^2 = \sum_{j=1}^{k_n} \sigma_{nj}^2.$$



# A Useful Lemma

A useful lemma in mathematical analysis.

## Lemma 4.2

Let  $\{\theta_{nj} : 1 \leq j \leq k_n\}_{n \geq 1}$  be a double array of complex numbers satisfying as  $n \rightarrow \infty$ ,

(i)  $\max_{1 \leq j \leq k_n} |\theta_{nj}| \rightarrow 0,$

(ii)  $\sum_{j=1}^{k_n} |\theta_{nj}| \leq M < \infty$  where  $M$  is free of  $n$ ,

(iii)  $\sum_{j=1}^{k_n} \theta_{nj} \rightarrow \theta$  for a finite complex  $\theta$ , then  $\prod_{j=1}^{k_n} (1 + \theta_{nj}) \rightarrow e^\theta.$

This is a generalized formula for  $\lim_{n \rightarrow \infty} (1 + \theta/n)^n \rightarrow e^\theta$  with  $\theta_{nj} \equiv \theta/n$ .

## References for this chapter

Chung, K. L. (2001). A course in probability theory, 3rd. Academic press.

# Preliminary

For any complex  $z \neq 0$ , the complex number  $w$  satisfying  $e^w = z$  is  $\text{Log } z$ , i.e.  $\text{Log } z = w$ . Let  $w = u + vi$ , we have  $z = e^u e^{iv}$ , which means

$$|z| = e^u, \quad u = \log |z|, \quad v = \text{Arg } z = \arg z + 2k\pi, \quad \arg z \in [-\pi, \pi]$$

then:

$$\text{Log } z = \log |z| + i \text{Arg } z, \quad \log z = \log |z| + i \arg z$$

## Remark 12

For any  $|z| < 1$ ,

$$\log(1 + z) = \sum_{n=1}^{\infty} (-1)^{n-1} \frac{z^n}{n}$$

# The Proof of Lemma 4.2

$$\begin{aligned} |\log(1 + \theta_{nj}) - \theta_{nj}| &= \left| (-1)^{m-1} \frac{\theta_{nj}^m}{m} \right| \leq \frac{|\theta_{nj}|^m}{m} \\ &\leq \frac{|\theta_{nj}|^2}{2} \sum_{m=2}^{\infty} \left(\frac{1}{2}\right)^{m-2} = |\theta_{nj}|^2 < 1 \end{aligned}$$

Here the upper bound 1 is uniformly for all  $\{\theta_{nj}\}$ , hence,

$$\log(1 + \theta_{nj}) = \theta_{nj} + \Lambda_{nj} |\theta_{nj}|^2, \quad \text{where } |\Lambda_{nj}| < 1$$

Now,

$$\sum_{j=1}^{k_n} \log(1 + \theta_{nj}) = \sum_{j=1}^{k_n} \theta_j + \sum_{j=1}^{k_n} \Lambda_{nj} |\theta_{nj}|^2$$

# The Proof of Lemma 4.2

From (i) and (ii):

$$\begin{aligned} \left| \sum_{j=1}^{k_n} \Lambda_{nj} |\theta_{nj}|^2 \right| &\leq \max_{1 \leq j \leq k_n} |\theta_{nj}| \sum_{j=1}^{k_n} |\theta_{nj}| \\ &\stackrel{(ii)}{\leq} \max_{1 \leq j \leq k_n} |\theta_{nj}| M \stackrel{(i)}{\rightarrow} 0 \end{aligned}$$

And from (iii), we know that

$$\sum_{j=1}^{k_n} \log(1 + \theta_{nj}) \rightarrow \theta$$

# Liapounov's Theorem

## Theorem 4.3

For a double array  $\{X_{nj} : 1 \leq j \leq k_n\}_{n \geq 1}$ , let  $\Gamma_n = \sum_{j=1}^{k_n} \mathbb{E}|X_{nj} - \alpha_{nj}|^3$ , which is finite for every  $n$ , and if Liapounov's condition holds

$$\frac{\Gamma_n}{\sigma_n^3} = \frac{1}{\sigma_n^3} \sum_{j=1}^{k_n} \mathbb{E}|X_{nj} - \alpha_{nj}|^3 \rightarrow 0 \quad \text{as } n \rightarrow \infty, \text{ then}$$

$$\frac{S_n - \alpha_n}{\sigma_n} \xrightarrow{d} N(0, 1).$$

## Remark 13

The constant 3 can be relaxed to  $2 + \delta$  some  $\delta > 0$ .

# The Proof of Liapounov's Theorem

Let  $\gamma_{nj} = E|X_{nj} - \alpha_{nj}|^3$ . As:

$$\sigma_{nj} = (E|X_{nj} - \alpha_{nj}|^2)^{1/2} \leq (E|X_{nj} - \alpha_{nj}|^3)^{1/3}$$

We have  $\sigma_{nj}^3 \leq \gamma_{nj}$ , thus,

$$\max_{1 \leq j \leq k_n} \sigma_{nj}^3 \leq \max_{1 \leq j \leq k_n} \gamma_{nj} < \Gamma_n \quad (6)$$

Let  $\phi_{nj}$  be the c.f. of  $(X_{nj} - \alpha_{nj})/\sigma_n$ . As  $\gamma_{nj}$  is finite, from Theorem 2.8,

$$\phi_{nj}(t) = 1 - \frac{\sigma_{nj}^2 t^2}{2\sigma_n^2} + \frac{\Lambda_{nj} \gamma_{nj} t^3}{6 \sigma_n^3}, \quad \text{where } |\Lambda_{nj}| < 1$$

$$\begin{aligned} \max_{1 \leq j \leq k_n} |\phi_{nj}(t) - 1| &\leq \frac{t^2}{2\sigma_n^2} \max_{1 \leq j \leq k_n} \sigma_{nj}^2 + \frac{t^3}{6\sigma_n^3} \max_{1 \leq j \leq k_n} \gamma_{nj} \\ &\leq \frac{t^2}{2} \left( \frac{\Gamma_n}{\sigma_n^3} \right)^{2/3} + \frac{t^3}{6\sigma_n^3} \max_{1 \leq j \leq k_n} \gamma_{nj} \xrightarrow{(6)} 0 \end{aligned} \quad (7)$$

which is the Assumption (i) of Lemma (4.2).

# The Proof of Liapounov's Theorem

In fact, note that  $\sigma_{nj}^2 = (\sigma_{nj}^3)^{2/3} \leq (\max_j \sigma_{nj}^3)^{2/3}$ ,

$$\frac{\max \sigma_{nj}^2}{\sigma_n^2} \leq \left( \frac{\max \sigma_{nj}^3}{\sigma_n^3} \right)^{2/3} \stackrel{(6)}{\leq} \left( \frac{\Gamma_n}{\sigma_n^3} \right)^{2/3} \rightarrow 0$$

and  $\max \gamma_{nj} / \sigma_n^3 \rightarrow 0$  comes from the condition directly. On the other hand,

$$\sum_{j=1}^{k_n} |\phi_{nj}(t) - 1| \leq \frac{\sum \sigma_{nj}^2 t^2}{2\sigma^2} + \frac{t^3 \Gamma_n}{6 \sigma_n^3} = \frac{t^2}{2} + \frac{t^3 \Gamma_n}{6 \sigma_n^3} \leq M(t) \quad (8)$$

which is the Assumption (ii) of Lemma (4.2).

# The Proof of Liapounov's Theorem

Finally, as

$$\left| \sum_{j=1}^{k_n} \frac{\Lambda_{nj} \gamma_{nj}}{\sigma_n^3} \right| \leq \frac{\Gamma_n}{\sigma_n^3} \rightarrow 0$$

we have,

$$\sum_{j=1}^{k_n} (\phi_{nj}(t) - 1) = -\frac{t^2}{2} + t^3 \sum_{j=1}^{k_n} \frac{\Lambda_{nj} \gamma_{nj}}{\sigma_n^3} \rightarrow -\frac{t^2}{2} \quad (9)$$

which is the Assumption (iii) of Lemma (4.2). Then from (7) to (9) and Lemma (4.2), the c.f. of  $(S_n - \alpha_n)/\sigma_n = \sum_{j=1}^{k_n} (X_{nj} - \alpha_{nj})/\sigma_n$  satisfying:

$$\prod_{j=1}^{k_n} \phi_{nj}(t) = \prod_{j=1}^{k_n} (1 + \phi_{nj}(t) - 1) \rightarrow e^{-t^2/2}$$

Apply Lévy-Cramér's Theorem., we obtain the result.



Theorem 4.3 implies the following for a single array  $\{X_n\}_{n \geq 1}$ .

## Corollary 4.4

Let  $\{X_n\}_{n \geq 1}$  be a sequence of independent random vectors,  $\alpha_j = E(X_j)$ ,  $\sigma_j^2 = \text{Var}(X_j)$  and  $\gamma_j = E|X_j - \alpha_j|^3 < \infty$ . Let  $P_n = \sum_{j=1}^n \gamma_j$ , if  $\frac{P_n}{\sigma_n^3} \rightarrow 0$ , then

$$\frac{S_n - \sum_{j=1}^n \alpha_j}{\sigma_n} \xrightarrow{d} N(0, 1).$$

It can be proved by Lindeberg's method. Here we use another approach to verify it. W.L.O.G. assuming  $\alpha_j = 0$ .

For any  $f \in C^3 := \{g \mid g^{(3)}$  is continuous in  $\mathbb{R}\}$ . Let  $Y_1, \dots, Y_n$  are independent r.v. with  $Y_j \sim N(0, \sigma_j^2)$  matching the first two moments of  $X_j$ , and  $Y_0 = \sum_{i=1}^n Y_i / \sigma_n \sim N(0, 1)$ . We want to show  $\forall f \in C^3$ ,

$$\mathbb{E}f\left(\frac{\sum_{i=1}^n X_i}{\sigma_n}\right) - \mathbb{E}f(Y_0) \rightarrow 0 \quad (10)$$

which implies  $\sum_{i=1}^n X_i / \sigma_n \xrightarrow{d} Y_0$  (Recall Portmanteau Lemma).

# Proof

Let  $Z_j = Y_1 + \dots + Y_{j-1} + X_{j+1} + \dots + X_n$  for  $2 \leq j \leq n-1$ ,  $Z_1 = X_2 + \dots + X_n$ , and  $Z_n = Y_1 + \dots + Y_{n-1}$ , then:

$$\text{Ef} \left( \frac{\sum_{i=1}^n X_i}{\sigma_n} \right) - \text{Ef} \left( \frac{\sum_{i=1}^n Y_i}{\sigma_n} \right) = \sum_{i=1}^n \left[ \text{Ef} \left( \frac{Z_i + X_i}{\sigma_n} \right) - \text{Ef} \left( \frac{Z_i + Y_i}{\sigma_n} \right) \right] \quad (11)$$

Note that:

$$\begin{aligned} f \left( \frac{Z_i + X_i}{\sigma_n} \right) &= f \left( \frac{Z_i}{\sigma_n} \right) + f' \left( \frac{Z_i}{\sigma_n} \right) \frac{X_i}{\sigma_n} + \frac{1}{2} f'' \left( \frac{Z_i}{\sigma_n} \right) \frac{X_i^2}{\sigma_n^2} + \theta_i^{(1)} \frac{X_i^3}{3! \sigma_n^3} \\ f \left( \frac{Z_i + Y_i}{\sigma_n} \right) &= f \left( \frac{Z_i}{\sigma_n} \right) + f' \left( \frac{Z_i}{\sigma_n} \right) \frac{Y_i}{\sigma_n} + \frac{1}{2} f'' \left( \frac{Z_i}{\sigma_n} \right) \frac{Y_i^2}{\sigma_n^2} + \theta_i^{(2)} \frac{Y_i^3}{3! \sigma_n^3} \end{aligned}$$

where  $|\theta_i^{(l)}| \leq \|f^{(3)}\|_\infty < \infty$ . As  $\text{E}X_i = \text{E}Y_i = 0$ ,  $\text{E}X_i^2 = \text{E}Y_i^2 = \sigma_i^2$ ,  $\text{E}Y_i^3 = \sqrt{8/\pi} \sigma_i^3$ ,

$$\begin{aligned} \left| \text{Ef} \left( \frac{Z_i + X_i}{\sigma_n} \right) - \text{Ef} \left( \frac{Z_i + Y_i}{\sigma_n} \right) \right| &\leq \frac{1}{3! \sigma_n^3} \left| \text{E} \left[ \theta_i^{(1)} X_i^3 \right] - \text{E} \left[ \theta_i^{(2)} Y_i^3 \right] \right| \\ &\leq \frac{M}{3! \sigma_n^3} \left( \gamma_i + \sqrt{\frac{8}{\pi}} \sigma_i^3 \right) \end{aligned}$$

Then from (11),

$$\begin{aligned} \left| \text{Ef} \left( \frac{\sum_{i=1}^n X_i}{\sigma_n} \right) - \text{Ef} \left( \frac{\sum_{i=1}^n Y_i}{\sigma_n} \right) \right| &\leq \sum_{i=1}^n \left| \text{Ef} \left( \frac{Z_i + X_i}{\sigma_n} \right) - \text{Ef} \left( \frac{Z_i + Y_i}{\sigma_n} \right) \right| \\ (M = \max\{|\theta_i^{(1)}|, |\theta_i^{(2)}|\}) &\leq \frac{M}{6} \left( \frac{\sum_{i=1}^n \gamma_i}{\sigma_n^3} + \sqrt{\frac{8}{\pi}} \frac{\sum_{i=1}^n \sigma_i^3}{\sigma_n^3} \right) \\ \left( \sum_{i=1}^n \sigma_i^3 \leq \Gamma_n \right) &\leq \frac{M_1}{6} \frac{\Gamma_n}{\sigma_n^3} \rightarrow 0 \end{aligned}$$

As a result, (10) is true.

### Remark 14

*The method of (11) is called Telescoping.*

An immediate consequence of Liapounov CLT is

## Corollary 4.5

For a double array (triangular array of independent variables)  $\{X_{nj}, 1 \leq j \leq k_n\}_{n \geq 1}$ , if  $|X_{nj}| \leq M_{nj}$  a.e., and  $\lim_{n \rightarrow \infty} \max_{1 \leq j \leq k_n} M_{nj} = 0$ . Let

$S_n = \sum_{j=1}^{k_n} X_{nj}$ , show that

$$\frac{S_n - E(S_n)}{\sigma_n} \xrightarrow{d} N(0, 1).$$

# Null Array

Four conditions:

- (a)  $\forall j, \lim_{n \rightarrow \infty} P(|X_{nj} - \alpha_{nj}| > \varepsilon \sigma_{nj}) = 0,$
- (b)  $\lim_{n \rightarrow \infty} \max_{1 \leq j \leq k_n} P(|X_{nj} - \alpha_{nj}| > \varepsilon \sigma_{nj}) = 0,$
- (c)  $\lim_{n \rightarrow \infty} P(\max_{1 \leq j \leq k_n} |X_{nj} - \alpha_{nj}| > \varepsilon \sigma_{nj}) = 0,$
- (d)  $\lim_{n \rightarrow \infty} \sum_{j=1}^{k_n} P(|X_{nj} - \alpha_{nj}| > \varepsilon \sigma_{nj}) = 0$

Homework: check  $(d) \Rightarrow (c) \Rightarrow (b) \Rightarrow (a)$ .

## Definition 4.6

A double array satisfying condition (b) is called a null array.

# Equivalent Form of Null Arrays

## Proposition 4.7

A double array  $\{X_{nj}, 1 \leq j \leq k_n\}_{n \geq 1}$  is a null array iff

$$\forall t \in \mathbb{R}, \lim_{n \rightarrow \infty} \max_{1 \leq j \leq k_n} |\phi_{nj}(t) - 1| = 0 \quad (e),$$

where  $\phi_{nj}$  is a c.f. of  $\frac{X_{nj} - \alpha_j}{\sigma_n}$ . Furthermore, the convergence in (e) is uniformly over any finite interval.

## Remark 15

- (i) If  $\{X_{nj}\}$  is a NA, then  $\frac{X_{nj} - \alpha_j}{\sigma_n} \xrightarrow{P} 0$ .
- (ii) Prop 4.7 implies that each element of a null array  $\left\{ \frac{X_{nj} - \alpha_j}{\sigma_n} \right\}_{j=1}^{k_n}$  uniformly degenerate at 0 on  $j$  when  $n \rightarrow \infty$ .

" $\Rightarrow$ ": WLOG assume  $\alpha_j = 0$ ,

$$\begin{aligned}
 |\phi_{nj}(t) - 1| &= \left| \mathbb{E} \left( e^{itX_{nj}/\sigma_n} - 1 \right) \right| \leq \mathbb{E} \left[ \left| e^{itX_{nj}/\sigma_n} - 1 \right| \mathbb{I}(|X_{nj}| > \epsilon\sigma_n) \right] \\
 &\quad + \mathbb{E} \left[ \left| e^{itX_{nj}/\sigma_n} - 1 \right| \mathbb{I}(|X_{nj}| \leq \epsilon\sigma_n) \right] \\
 (|e^{itu} - 1| = \sqrt{2(1 - \cos tu)} \leq |tu|) &\leq 2\mathbb{P}(|X_{nj}| > \epsilon\sigma_n) + |t| \mathbb{E} \left[ \left| \frac{X_{nj}}{\sigma_n} \right| \mathbb{I} \left( \left| \frac{X_{nj}}{\sigma_n} \right| \leq \epsilon \right) \right] \\
 &\leq 2\mathbb{P}(|X_{nj}| > \epsilon\sigma_n) + |t|\epsilon
 \end{aligned}$$

Thus,

$$\max_j |\phi_{nj}(t) - 1| \leq 2 \max_j \mathbb{P}(|X_{nj}| > \epsilon\sigma_n) + |t|\epsilon$$

(e) is veracious. In fact, for any  $|t| \leq K$ ,

$$\sup_{|t| \leq K} \max_j |\phi_{nj}(t) - 1| \leq 2 \max_j \mathbb{P}(|X_{nj}| > \epsilon\sigma_n) + K\epsilon$$

so the convergence is uniform over  $t \in [-K, K]$  as desired.



" $\Leftarrow$ ": It can be derived by using Lemma 2.3.2 ( [Homework](#)).

$$\begin{aligned} P\left(\left|\frac{X_{nj}}{\sigma_n}\right| > \frac{2}{\delta}\right) &\stackrel{(*)}{\leq} \frac{1}{\delta} \int_{|t| \leq \delta} (1 - \phi_{nj}(t)) dt = \frac{1}{\delta} \left| \int_{|t| \leq \delta} (1 - \phi_{nj}(t)) dt \right| \\ &\leq \frac{1}{\delta} \int_{|t| \leq \delta} |1 - \phi_{nj}(t)| dt \end{aligned}$$

This implies,

$$\max_j P\left(\left|\frac{X_{nj}}{\sigma_n}\right| > \frac{2}{\delta}\right) \leq \max_j \frac{1}{\delta} \int_{|t| \leq \delta} |1 - \phi_{nj}(t)| dt$$

From the BCT and (e), we know condition (h) holds.

(\*):

$$\begin{aligned} \mathbb{I}(|\delta x| > 2) &\leq 2 \left( 1 - \frac{\sin \delta x}{\delta x} \right) = \frac{2}{\delta} \int_{-\delta}^{\delta} (1 - \cos tx) dt \\ &= \frac{2}{\delta} \int_{|t| \leq \delta} (1 - \cos tx - i \sin tx) dt \end{aligned}$$

Then (\*) comes from taking expectation on both sides.

### Remark 16

*Bounded Convergence Theorem (BCT): Suppose  $f_n(t) \rightarrow f(t)$  for  $\forall t$ , and*

$$|f_n(t)| \leq g(t), \quad \int g(t) dt \text{ exists}$$

Then,

$$\int |f_n(t) - f(t)| dt \rightarrow 0, \quad \int f_n(t) dt \rightarrow \int f(t) dt$$

## Definition 4.8

A double array  $\{X_{nj}, 1 \leq j \leq k_n\}_{n \geq 1}$  is said to satisfy the Lindeberg condition, if  $\forall \varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \sigma_n^{-2} \sum_{j=1}^{k_n} \mathbb{E}\{(X_{nj} - \alpha_{nj})^2 \mathbb{I}(|X_{nj} - \alpha_{nj}| > \varepsilon \sigma_n)\} = 0,$$

where  $\alpha_{nj} = \mathbb{E}(X_{nj})$ ,  $\sigma_n^2 = \sum_{j=1}^{k_n} \text{Var}(X_{nj})$ , implicitly assumed  $\sigma_{nj}^2 = \mathbb{E}(X_{nj}^2) < \infty$  for any  $n$  and  $j$ .

# Lindeberg-Feller CLT

## Lemma 4.9

Let  $u(m, n)$  be a function of positive integers  $m$  and  $n$ , s.t.

$\forall m, \lim_{n \rightarrow \infty} u(m, n) = 0$ , then there exists a monotone increasing sequence  $\{m_n\}$ ,  $m_n \rightarrow \infty$ , s.t.  $\lim_{n \rightarrow \infty} u(m_n, n) = 0$

Lindeberg's condition is a sufficient condition (and under certain conditions also a necessary condition) for the CLT to hold for a sequence of independent random variables.

## Theorem 4.10

[Lindeberg-Feller] For a double array  $\{X_{nj}, 1 \leq j \leq k_n\}_{n \geq 1}$ , assume

$\text{Var}(X_{nj}) = \sigma_{nj}^2 < \infty$ , then

(i)  $\frac{S_n - \text{ES}_n}{\sigma_n} \xrightarrow{d} N(0, 1)$  and (ii) the double array is a null array iff the Lindeberg condition is satisfied.

# The Proof of Lemma 4.9

As  $\lim_{n \rightarrow \infty} u(m, n) = 0$  for each  $m$ .  $\exists n_m$ , s.t.

$$n \geq n_m, \quad u(m, n_m) \leq \frac{1}{m}$$

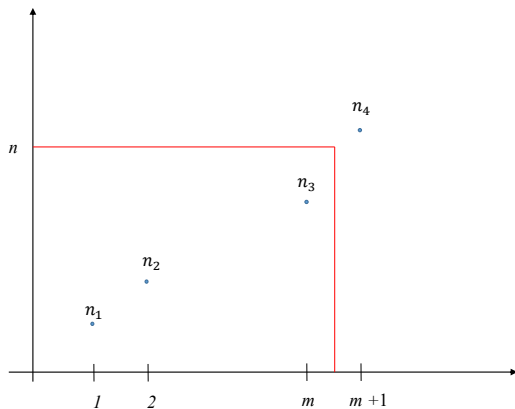
Here we obtain a sequence  $\{n_m\}_{m \geq 1}$ , and can make it strictly increase to  $\infty$ .

Now let  $m_n = m$  s.t.  $n_m \leq n \leq n_{m+1}$ . When  $n \geq n_m$ ,

$$u(m_n, n) = u(m, n) \leq \frac{1}{m}$$

As  $n_m \uparrow \infty$ ,  $m_n \uparrow \infty$  too, and  $\lim_{n \rightarrow \infty} u(m_n, n) = 0$ .

# The Proof of Lemma 4.9



$$n_3 \leq n \leq n_4, \quad \Rightarrow \quad m_n = 3$$

# The Proof of Theorem 4.10 $\Leftarrow$

WLOG, assume  $EX_{nj} = 0$  and  $\sigma_n^2 = 1$ , or we can redefine  $X_{nj} = (X_{nj} - EX_{nj})/\sigma_n$ . Now truncate  $X_{nj}$  to  $X'_{nj}$  with a  $\eta \in (0, 1)$ :

$$X'_{nj} = \begin{cases} X_{nj} & \text{if } |X_{nj}| < \eta \\ 0 & \text{o.w.} \end{cases} \quad (12)$$

Denote  $S'_n = \sum_{i=1}^{k_n} X'_{nj}$ ,  $\sigma_n'^2 = \sum_{i=1}^{k_n} \text{var}(X'_{nj}) = \text{var}(S'_n)$ .

$$|EX'_{nj}| = \left| \int_{|x| < \eta} x dF_{nj}(x) \right| \stackrel{EX_{nj}=0}{=} \left| \int_{|x| \geq \eta} x dF_{nj}(x) \right| \leq \frac{1}{\eta} \int_{|x| \geq \eta} x^2 dF_{nj}(x)$$

Hence,

$$|ES'_n| \leq \sum_{i=1}^{k_n} |EX'_{nj}| \leq \frac{1}{\eta} \sum_{j=1}^{k_n} \int_{|x| \geq \eta} x^2 dF_{nj}(x) \xrightarrow{\text{Lind Con}} 0 \quad (13)$$

# The Proof of Theorem 4.10 $\Leftarrow$

Similarly,

$$\begin{aligned}\sum_{i=1}^{k_n} \mathbb{E}X_{nj}'^2 &= \sum_{j=1}^{k_n} \int_{|x| < \eta} x^2 dF_{nj}(x) \\ &= \sum_{j=1}^{k_n} \left[ \int x^2 dF_{nj}(x) - \int_{|x| \geq \eta} x^2 dF_{nj}(x) \right] \xrightarrow{\sigma_n^2=1} 1\end{aligned}\quad (14)$$

Hence,

$$\sigma_n'^2 = \text{var}(S_n') = \sum_{j=1}^{k_n} \mathbb{E}X_{nj}'^2 - \sum_{j=1}^{k_n} (\mathbb{E}X_{nj}')^2 \xrightarrow{(14)} 1 = \sigma_n^2$$

where we use the fact:

$$\sum_{j=1}^{k_n} (\mathbb{E}X_{nj}')^2 \leq \left( \sum_{j=1}^{k_n} |\mathbb{E}X_{nj}'| \right)^2 \xrightarrow{(12)} 0$$



As

$$\frac{S'_n}{\sigma'_n} = \frac{S'_n - \mathbb{E}S'_n}{\sigma'_n} + \frac{\mathbb{E}S'_n}{\sigma'_n},$$

$\mathbb{E}S'_n \rightarrow 0$ , and  $\sigma'_n \rightarrow 1$ . From Slutsky Theorem  $S'_n/\sigma'_n$  and  $S'_n - \mathbb{E}S'_n/\sigma'_n$  would converge to the same distribution.

Next to show:

$$\frac{S'_n - \mathbb{E}S'_n}{\sigma'_n} \xrightarrow{d.} N(0, 1)$$

## The Proof of Theorem 4.10 $\Rightarrow$

From the LC, for each fixed  $m \geq 1$ ,

$$\lim_{n \rightarrow \infty} m^2 \sum_{j=1}^{k_n} \int_{|x| > 1/m} x^2 dF_{nj}(x) = 0$$

From Lemma 4.9, there exists  $\{m_n\} \uparrow \infty$  s.t.

$$\lim_{n \rightarrow \infty} \sum_{j=1}^{k_n} m_n^2 \int_{|x| > 1/m_n} x^2 dF_{nj}(x) = 0 \quad (15)$$

Let  $\eta_n = m_n^{-1} \downarrow 0$ , and use  $\eta_n$  to replace  $\eta$  in the definition of  $X'_{nj}$ , then

$$|X'_{nj}| \leq \eta_n := M_{nj}, \quad \lim_{n \rightarrow \infty} \max_{1 \leq j \leq k_n} M_{nj} = \lim_{n \rightarrow \infty} \eta_n = 0$$

Then from corollary 4.5,  $(S'_n - \mathbb{E}S'_n)/\sigma'_n \rightarrow N(0, 1)$ . So does  $S'_n/\sigma'_n$ . So  $S'_n \xrightarrow{d} N(0, 1)$  as  $\sigma'_n \rightarrow 1$ .

Since

$$\begin{aligned} P(S_n \neq S'_n) &\leq P\left(\bigcup_{j=1}^{k_n} \{X_{nj} \neq X'_{nj}\}\right) \leq \sum_{j=1}^{k_n} P(|X_{nj}| \geq \eta_n) \\ &\leq \sum_{j=1}^{k_n} \frac{1}{\eta_n^2} \int_{|x| > \eta_n} x^2 dF_{nj}(x) \xrightarrow{LC} 0 \end{aligned}$$

So we have  $S_n - S'_n \xrightarrow{P} 0$ . By Slutsky,  $S_n \xrightarrow{d} N(0, 1)$ . Here we complete the proof of sufficiency.

## The Proof of Theorem 4.10 $\Rightarrow$

Let  $\phi_{nj}$  be the c.f of  $X_{nj}$  with  $EX_{nj} = 0$  and  $\sigma_n^2 = 1$ . As  $S_n \xrightarrow{d.} N(0, 1)$  in (i),

$$\lim_{n \rightarrow \infty} \prod_{j=1}^{k_n} \phi_{nj}(t) = e^{-t^2/2}, \quad \lim_{n \rightarrow \infty} \sum_{j=1}^{k_n} \log \phi_{jn}(t) = -\frac{t^2}{2} \quad (16)$$

and it can be reinforced that it would be hold uniform over  $t \in [-K, K]$  by using covering method.

On the other hand, (ii) and the properties of NA imply,

$$\lim_{n \rightarrow \infty} \sup_{|t| < K} \max_{1 \leq j \leq k_n} |\phi_{nj}(t) - 1| = 0 \quad (17)$$

Let  $\theta_{nj} = \phi_{nj}(t) - 1$ , from the proof of Lemma 4.2, we know the following display holds.

# The Proof of Theorem 4.10 $\Rightarrow$

The following display holds.

$$\log \phi_{nj}(t) = \log(1 + \theta_{nj}) = \theta_{nj} + \Lambda_{nj}|\theta_{nj}|^2 = \phi_{nj}(t) - 1 + \Lambda_{nj}|\phi_{nj}(t) - 1|^2 \quad (18)$$

where  $|\Lambda_{nj}| < 1$  uniformly. Furthermore,

$$\sum_{j=1}^{k_n} |\phi_{nj}(t) - 1|^2 \leq \max_{1 \leq j \leq k_n} |\phi_{nj}(t) - 1| \sum_{j=1}^{k_n} |\phi_{nj}(t) - 1| \xrightarrow{(17)\&(*)} 0 \quad (19)$$

where  $(*)$  is the fact that,

$$\begin{aligned} \sum_{j=1}^{k_n} |\phi_{nj}(t) - 1| &= \sum_{j=1}^{k_n} \left| \int (e^{itx} - 1) dF_{nj}(x) \right| \\ &\stackrel{\text{Taylor Exp of } e^{itx}}{=} \sum_{j=1}^{k_n} \left| \int \left( itx + \kappa_t \frac{t^2 x^2}{2} \right) dF_{nj}(x) \right| \leq \frac{t^2}{2} < \infty \end{aligned}$$

in which  $\kappa_t \in (0, 1)$ .

# The Proof of Theorem 4.10 $\Rightarrow$

Then, from (16) and (18), we obtain:

$$\lim_{n \rightarrow \infty} \sum_{j=1}^{k_n} (\phi_{nj}(t) - 1) = \lim_{n \rightarrow \infty} \sum_{j=1}^{k_n} \log \phi_{jn}(t) = -\frac{t^2}{2}$$

By taking the real part,

$$\lim_{n \rightarrow \infty} \sum_j \int_{-\infty}^{\infty} (1 - \cos tx) dF_{nj}(x) = \frac{t^2}{2}$$

Hence for each  $\eta > 0$ , split the integral into two parts, by Chebyshev's inequality,

$$\begin{aligned} \overline{\lim}_{n \rightarrow \infty} \left| \frac{t^2}{2} - \sum_j \int_{|x| \leq \eta} (1 - \cos tx) dF_{nj}(x) \right| &= \overline{\lim}_{n \rightarrow \infty} \left| \sum_j \int_{|x| > \eta} (1 - \cos tx) dF_{nj}(x) \right| \\ &\leq \overline{\lim}_{n \rightarrow \infty} \sum_j \int_{|x| > \eta} 2 dF_{nj}(x) \\ &\leq \overline{\lim}_{n \rightarrow \infty} 2 \sum_j \frac{\sigma_{nj}^2}{\eta^2} = \frac{2}{\eta^2} \end{aligned}$$

# The Proof of Theorem 4.10 $\Rightarrow$

Note that  $0 \leq 1 - \cos \theta \leq \theta^2/2$  for every real  $\theta$ , this implies,

$$\frac{2}{\eta^2} \geq \overline{\lim}_{n \rightarrow \infty} \left\{ \frac{t^2}{2} - \sum_j \frac{t^2}{2} \int_{|x| \leq \eta} x^2 dF_{nj}(x) \right\} \geq 0$$

Therefore, for fixed  $\eta > 0$ ,

$$\overline{\lim}_{n \rightarrow \infty} \sum_{j=1}^{k_n} \mathbb{E} (|X_{nj}^2| \mathbb{I}(|X_{nj}| > \eta)) = \overline{\lim}_{n \rightarrow \infty} \left\{ 1 - \sum_{j=1}^{k_n} \int_{|x| \leq \eta} x^2 dF_{nj}(x) \right\} \leq \frac{4}{t^2 \eta^2}$$

Let  $t \rightarrow \infty$ , we have:

$$\sum_{j=1}^{k_n} \mathbb{E} (|X_{nj}^2| \mathbb{I}(|X_{nj}| > \eta)) \longrightarrow 0$$

which exactly is the LC.

# Example: Regression

$$y_j = x_j\beta + \epsilon_j, \quad \epsilon_j \text{ i.i.d. } N(0, \sigma_\epsilon^2), \quad j = 1, 2, \dots$$

where  $\{x_j\}_{j \geq 1}$  are fixed design points, s.t.

$$\max_{1 \leq j \leq n} \frac{|x_j|}{a_n} \rightarrow 0, \quad a_n^2 = \sum_{j=1}^n x_j^2$$

The ordinary least square estimate is  $\hat{\beta}_{LS} = \sum_{j=1}^n x_j y_j / a_n^2$ . We want to show that  $a_n(\hat{\beta}_{LS} - \beta) \xrightarrow{d} N(0, \sigma^2)$ .

$$a_n(\hat{\beta}_{LS} - \beta) = \frac{\sum x_j y_j - \beta \sum x_j^2}{a_n} = \frac{\sum x_j \epsilon_j}{a_n} =: \sum_{j=1}^n X_{nj}$$

where:

$$X_{nj} = \frac{x_j \epsilon_j}{\sqrt{\sum x_j^2}}, \quad \alpha_{nj} = \mathbb{E}X_{nj} = 0, \quad \sigma_{nj}^2 = \frac{x_j^2 \sigma_\epsilon^2}{a_n^2}, \quad \sigma_n^2 = \sigma_\epsilon^2$$



## Example: Regression

Here  $\{X_{nj}, 1 \leq j \leq n\}_{n \geq 1}$  is a triangular array, and

$$\begin{aligned}\sigma_n^{-2} \sum_{j=1}^n \mathbb{E} [X_{nj}^2 \mathbb{I}(|X_{nj}| > \delta \sigma_n)] &= \frac{1}{\sigma_n^2 a_n^2} \sum_{j=1}^n x_j^2 \mathbb{E} [\epsilon_j^2 \mathbb{I}(|x_j \epsilon_j / a_n| > \delta \sigma_n)] \\ (m_n = \max_j |x_j / a_n|) &\leq \frac{1}{\sigma^2 a_n^2} \sum_{j=1}^n x_j^2 \mathbb{E} [\epsilon_j^2 \mathbb{I}(|\epsilon_j| > \delta \sigma_\epsilon m_n^{-1})] \\ &= \frac{1}{\sigma_\epsilon^2} \mathbb{E} [\epsilon_j^2 \mathbb{I}(|\epsilon_j| > \delta \sigma_\epsilon m_n^{-1})] \rightarrow 0\end{aligned}$$

as  $m_n \rightarrow 0$ . From Theorem 4.10, we know that,

$$a_n(\hat{\beta}_{LS} - \beta) = \sum_{j=1}^n X_{nj} \xrightarrow{d.} N(0, \sigma_\epsilon^2)$$

# A Sufficient Condition of LC

There is a sufficient condition to verify the Lindeberg-Feller condition.

## Proposition 4.11

For a double array  $\{X_{nj}\}_{j=1}^{k_n}$  with means  $\{\mu_{nj}\}$  and variances  $\{\sigma_{nj}^2\}$ . If for some  $\nu > 2$ ,

$$\sum_{j=1}^{k_n} E|X_{nj} - \mu_{nj}|^\nu = o(\sigma_n^\nu)$$

Then, the Lindeberg condition is valid.

$$\begin{aligned}
\mathbb{E} \left[ (X_{nj} - \mu_{nj})^2 \mathbb{I}(|X_{nj} - \mu_{nj}| > \epsilon \sigma_n) \right] &= \int_{|t - \mu_{nj}| > \epsilon \sigma_n} (t - \mu_{nj})^2 dF_{nj}(t) \\
&\leq (\epsilon \sigma_n)^{2-\nu} \int_{|t - \mu_{nj}| > \epsilon \sigma_n} (t - \mu_{nj})^\nu dF_{nj}(t) \\
&\leq (\epsilon \sigma_n)^{2-\nu} \mathbb{E} |X_{nj} - \mu_{nj}|^\nu
\end{aligned}$$

Therefore,

$$\begin{aligned}
&\frac{1}{\sigma_n^2} \sum_{j=1}^{k_n} \mathbb{E} \left[ (X_{nj} - \mu_{nj})^2 \mathbb{I}(|X_{nj} - \mu_{nj}| > \epsilon \sigma_n) \right] \\
&\leq \epsilon^{2-\nu} \frac{\sum_{j=1}^{k_n} \mathbb{E} |X_{nj} - \mu_{nj}|^\nu}{\sigma_n^\nu} \rightarrow 0
\end{aligned}$$

The LC holds.

# CLT for $m$ -dependent r.v.s

## Definition 4.12

Def: A sequence of r.v.  $\{X_n\}_{n \geq 1}$  is  $m$ -dependent if  $\exists$  a positive integer  $m$  s.t. for any  $n \geq 1$  and  $j \geq m$ ,  $X_{n+j}$  is independent of  $\mathcal{F}_n = \sigma\{X_j, 1 \leq j \leq n\}$ , the  $\sigma$ -field generated by  $\{X_j\}_{j=1}^n$ .

## Theorem 4.13

Let  $\{X_n\}_{n \geq 1}$  be a sequence of  $m$ -dependent r.v.s with uniformly bounded variance s.t.  $\frac{\sigma_n}{mn^{1/3}} \stackrel{\Delta}{=} \frac{\sqrt{\text{Var}(\sum_{i=1}^n X_i)}}{mn^{1/3}} \rightarrow \infty$  as  $n \rightarrow \infty$  and  $m = o(n^{1/3})$ . Then

$$\frac{S_n - E(S_n)}{\sigma_n} \xrightarrow{d} N(0, 1).$$

A CLT for more general dependent data, the so-called mixing dependent sequence, will be discussed later.

# Proof: Use of Blocking Technique

Proof: As  $\{X_n\}_{n \geq 1}$  has uniformly bounded variance.  $\exists M$  s.t.  
 $\sup_n |Var(X_n)| \leq M.$

WLOG, assume  $E(X_j) = 0.$

We block the whole sequence by larger blocks followed by small blocks.

Let  $k = \lfloor n^{1/3} \rfloor$  be the size of large blocks,  $m$  be the size of small blocks

$p = \lfloor \frac{n}{k+m} \rfloor = O(n^{2/3})$  be the number of blocks,  $B_j = j(k + m),$

# Lay out of the Blocks

$$Y_1 = X_1 + \cdots + X_k, \quad Z_1 = X_{k+1} + \cdots + X_{k+m}$$

...

$$Y_p = X_{B_{p-1}+1} + \cdots + X_{B_{j-1}+k}, \quad Z_p = X_{B_{p-1}+k+1} + \cdots + X_{B_p}$$

$R_p = X_{B_p} + \cdots + X_n$  the residual block.

As  $k \gg m$  when  $n$  is large enough, then  $\{Y_j\}_{j=1}^p$  and  $\{Z_j\}_{j=1}^p$  are indpt rvs as they are at least  $m$ -apart.

# The Detailed Proof \*\*

Now, we have:

$$S_n = \sum_{j=1}^p Y_j + \sum_{j=1}^p Z_j + \sum_{l=1}^{n-p(k+m)} X_{B_l+l} := S'_n + S''_n + S'''_n \quad (20)$$

As  $\sup \text{var}(X_j) \leq M$ ,  $|\mathbb{E}(X_j X_l)| \leq M$ , and

$$\begin{aligned} \text{var}(S'''_n) &= \mathbb{E}(S'''_n)^2 = \left| \sum_{j,l=1}^{n-p(k+m)} \mathbb{E}(X_{p(k+m)+j} X_{p(k+m)+l}) \right| \\ &\leq (n-p(k+m))^2 M \leq (k+m)^2 M \end{aligned}$$

$$S'''_n = O_p(\sqrt{\text{var}(S'''_n)}) = O_p((n-p(k+m))) = O_p(k+m) = O_p(n^{1/3}) \quad (21)$$

Similarly,

$$EZ_j^2 = E \left( \sum_{j=1}^m X_{B_{j-1+k+j}} \right)^2 \leq m^2 M$$

And,

$$\text{var}(S_n'') \leq pm^2 M, \quad S_n'' = O_p(p^{1/2}m) = O_p(n^{1/3}m)$$

As  $\sigma_n/(mn^{1/3}) \rightarrow \infty$ , we have:

$$\frac{S_n''}{\sigma_n} = \frac{S_n''}{mn^{1/3}} \times \frac{mn^{1/3}}{\sigma_n} = o_p(1)$$

Besides, since  $k = O(n^{1/3})$ , from (21),  $S_n'''/\sigma_n = o_p(1)$ . Now,

$$\frac{S_n}{\sigma_n} = \frac{S_n'}{\sigma_n} + \frac{S_n''}{\sigma_n} + \frac{S_n'''}{\sigma_n} = \frac{\sigma_n'}{\sigma_n} \frac{S_n'}{\sigma_n'} + o_p(1) \quad (22)$$

It remains to prove  $\sigma_n'^2/\sigma_n^2 \rightarrow 1$  and  $S_n'/\sigma_n' \xrightarrow{d} N(0, 1)$ .



# The Detailed Proof \*\*

$\sigma_n'^2 / \sigma_n^2 \rightarrow 1$ :

$$E S_n^2 = E(S_n')^2 + E(S_n'')^2 + E(S_n''')^2 + 2E(S_n' S_n'') + 2E(S_n'' S_n''') + 2E(S_n''' S_n')$$

Then,

$$\begin{aligned} |E S_n^2 - E(S_n')^2| &\leq |E(S_n'')^2 + E(S_n''')^2 + 2E(S_n' S_n'') + 2E(S_n'' S_n''') + 2E(S_n''' S_n')| \\ &\leq pm^2 M^2 + (k+m)^2 M^2 + 4p(mM)^2 + 2m(k+m)M^2 \end{aligned}$$

Here we perceive that:

$$E(S_n' S_n'') = \sum_{j,l=1}^P \text{cov}(Y_j, Z_l) \stackrel{\text{indep.}}{=} \sum_{j=1}^P [\text{cov}(Y_j, Z_j) + \text{cov}(Y_j, Z_{j-1})] \leq 2p(mM)^2$$

and, similarly,

$$E(S_n'' S_n''') \leq m(k+m)M^2, \quad E(S_n''' S_n') = \text{cov}(S_n'', S_n''') = 0$$

Therefore,

$$\left| 1 - \frac{\sigma_n'^2}{\sigma_n^2} \right| \leq \frac{2pm^2 M^2 + 2(k+m)^2 M^2}{\sigma_n^2} = O\left(\frac{m^2 n^{2/3}}{\sigma_n^2}\right) \rightarrow 0$$

Hence,  $\frac{\sigma_n'^2}{\sigma_n^2} \rightarrow 1$ .

$S'_n/\sigma'_n \xrightarrow{d.} N(0, 1)$ :

We use truncation method. As  $k = [n^{1/3}]$ , let  $Y_{nj} = Y_j$ , then  $\{Y_{nj}, 1 \leq j \leq p\}$  is double array,  $|Y_{nj}| \leq km = O(n^{1/3}m) = o(\sigma'_n)$ , and

$$\begin{aligned} \frac{1}{\sigma_n'^2} \sum_{j=1}^p \mathbb{E} [Y_{nj}^2 \mathbb{I}(|Y_{nj}| \geq \eta \sigma'_n)] &\leq \frac{k^2 m^2}{\sigma_n'^2} \sum_{j=1}^p \mathbb{P}(|Y_{nj}| > \eta \sigma'_n) \\ &\leq \frac{k^2 m^2}{\sigma_n'^2} \frac{\sum_{j=1}^p \text{var}(Y_{nj})}{\eta^2 \sigma_n'^2} \leq \frac{k^2 m^2}{\eta^2 \sigma_n'^2} \rightarrow 0 \end{aligned}$$

as  $(km)/\sigma'_n \rightarrow 0$ , which implies the LC holds, then  $S'_n/\sigma'_n \xrightarrow{d.} N(0, 1)$  from Theorem 4.10.

The device allows the issue of convergence of multivariate distribution to be reduced to that of univariate ones.

## Theorem 4.14 (see Serfling (1988, p18))

*A sequence of random vectors  $\mathbf{X}_n$  in  $\mathbb{R}^d$  converges in distribution to the random vector  $\mathbf{X}$  if and only if for any linear combination of the component of  $\mathbf{X}_n$  converges in distribution to the same linear combination of the component of  $\mathbf{X}$  as  $n \rightarrow \infty$ , i.e.,*

$$\mathbf{X}_n \xrightarrow{d} \mathbf{X} \Leftrightarrow \mathbf{a}^T \mathbf{X}_n \xrightarrow{d} \mathbf{a}^T \mathbf{X}, \forall \mathbf{a} \in \mathbb{R}^d.$$

Levy's theorem implies that weak convergence of vectors is equivalent to weak convergence of linear combinations.

# Proof of Cramér-Wold

" $\Leftarrow$ " Let  $\mathbf{X}_n = (X_{n1}, \dots, X_{nd})^T$ ,  $\mathbf{X} = (X_1, \dots, X_d)^T$  have a characteristic function  $\phi_n$  and  $\phi$  respectively. As for all  $\mathbf{c} = (c_1, c_2, \dots, c_d)^T$

$$c_1 X_{n1} + \dots + c_k X_{nd} \xrightarrow{d} c_1 X_1 + \dots + c_d X_d. \quad (23)$$

The characteristic function of  $c_1 X_{n1} + \dots + c_d X_{nd}$  is

$$\phi_n(tc_1, \dots, tc_d) = E(e^{it(c_1 X_{n1} + \dots + c_d X_{nd})}).$$

The characteristic function of  $\lambda_1 X_1 + \dots + \lambda_d X_d$  is  $\phi(tc_1, \dots, tc_d)$ . Choose  $t = 1$ , then from (23)

$$\lim_{n \rightarrow \infty} \phi_n(c_1, \dots, c_d) = \phi(c_1, \dots, c_d)$$

implying  $\mathbf{X}_n \xrightarrow{d} \mathbf{X}$ .

The " $\Rightarrow$ " is obvious by mapping.

# Multivariate CLT

The Cramér-Wold device provide another approach to show Multivariate CLT.

## Theorem 4.15

Let  $\mathbf{X}_1, \mathbf{X}_2, \dots$  be i.i.d. random vectors with mean  $\boldsymbol{\mu}$  and finite covariance matrix,  $\Sigma$ . Let  $\bar{\mathbf{X}}_n = \sum_{i=1}^n \mathbf{X}_i/n$ , then

$$\sqrt{n} (\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{d} N_d(0, \Sigma).$$

by letting  $\mathbf{Y}_n := \sqrt{n} (\bar{\mathbf{X}}_n - \boldsymbol{\mu})$ , so

$$\mathbf{Y}_n \xrightarrow{d} \mathbf{Y} \quad \text{iff} \quad \mathbf{t}^T \mathbf{Y}_n \xrightarrow{d} \mathbf{t}^T \mathbf{Y} \quad \text{for all } \mathbf{t} \in \mathbb{R}^d.$$

## Chapter 4: Weakly Dependent Data

Let  $Z_1, \dots, Z_n \in \mathbb{R}^d$ , where  $d$  is the dimension of  $Z_i$ , is the equally sampled time series, i.e. daily, weekly, or yearly data.

- Strictly Stationary: for any integers  $l$  and  $m$ ,

$$(Z_{i_1}, \dots, Z_{i_m})^\top \quad \text{and} \quad (Z_{i_1+l}, \dots, Z_{i_m+l})^\top$$

have the same distribution (strong shift invariance).

- Weak Stationary (or Second Order Stationary):

$$\begin{aligned} \mathbb{E}(Z_i) &= \mathbb{E}(Z_{i+l}), & \text{var}(Z_i) &= \text{var}(Z_{i+l}), \\ \text{cov}(Z_i, Z_j) &= \text{cov}(Z_{i+l}, Z_{j+l}) \end{aligned}$$

(weak shift invariance).

Ways to make a time series stationary: difference, square root, etc.

## Definition 5.1 (ARMA models)

The sequence  $\{Z_i\}_{i \in \mathbb{Z}}$  is said to be an  $ARMA(p, q)$  if  $\{Z_i\}$  is weakly stationary and for any  $t$ ,

$$Z_t - \theta_1 Z_{t-1} - \cdots - \theta_p Z_{t-p} = \epsilon_t - \eta_1 \epsilon_{t-1} - \cdots - \eta_q \epsilon_{t-q}$$

where  $\{\epsilon_t\}$  is an independent white noise process, defined as  $WN(0, \sigma^2)$ .

Let  $\theta$  and  $\eta$  be  $p$ -th and  $q$ -th degree polynomials defined as:

$$\theta(z) = 1 - \theta_1 z - \cdots - \theta_p z^p$$

$$\eta(z) = 1 - \eta_1 z - \cdots - \eta_q z^q$$

Let  $B$  be the backward shift operator such that  $B^j Z_t = Z_{t-j}$ .

## Definition 5.2

An  $ARMA(p, q)$  process  $\{Z_t\}$  is said to be causal if there exists  $\{\psi_j\}_{j=0}^{\infty}$  such that

$$\sum_{j=0}^{\infty} |\psi_j| < \infty, \quad Z_t = \sum_{j=0}^{\infty} \psi_j \epsilon_{t-j}$$



## Theorem 5.3

Let  $\{Z_t\}$  be an ARMA( $p, q$ ) :  $\theta(B)Z_t = \eta(B)\epsilon_t$ . If  $\theta(z)$  and  $\eta(z)$  have no common zero roots, then  $\{Z_t\}$  is causal iff

$$\theta(z) \neq 0, \quad \forall z \in \mathbb{C}, |z| \leq 1$$

and the coefficients  $\{\psi_j\}$  are determined by  $\psi(z) = \eta(z)/\theta(z)$ .

Linear Process:

$$Z_t = \sum_{j=-\infty}^{\infty} \psi_j \epsilon_{t-j}, \quad \epsilon_t \text{ i.i.d. } F(0, \sigma^2)$$

ARMA process under certain conditions are special type of linear process!

# ARCH(p): Auto-Regression Conditionally Heterogeneity

Model:

$$Z_t = m(\vec{Z}_{t,p}) + \sigma(\vec{Z}_{t,p})\epsilon_t, \quad \vec{Z}_{t,p} = (Z_{t-1}, \dots, Z_{t-p})^\top$$

it generalizes  $AR(p)$ :


$$Z_t = \theta_0 + \theta_1 Z_{t-1} + \dots + \theta_p Z_{t-p} + \epsilon_t, \quad \epsilon_t \sim WN(0, \sigma^2)$$

in two aspects:

- (i) From linear to non-linear conditional mean function.
- (ii) From constant conditional variance to a function.

ARCH models (or non-linear time series models) in general are not necessarily stationary, but some conditions can guarantee they are "asymptotic stationary" (stationary after pre-burning the models for a period of "long" time). See Gouriéroux<sup>1</sup> for reference.

---

<sup>1</sup>Christian Gouriéroux. *ARCH models and financial applications*. Springer Science & Business Media, 2012. 

# Mixing Coefficients: measures of dependence

Let  $Z_1, \dots, Z_t, \dots$  be a strictly stationary process and  $\mathcal{F}_l^m$  be the sigma field generated by  $\{Z_i\}_{i=l}^m$ , where  $m \geq l$  are positive integers.

(i)  $\alpha$ -mixing (or strong mixing) coefficient:

$$\alpha(k) = \sup_{B \in \mathcal{F}_{-\infty}^t, C \in \mathcal{F}_{t+k}^\infty} |\mathbb{P}(B \cap C) - \mathbb{P}(B)\mathbb{P}(C)|, \quad k \geq 1$$

\*(ii)  $\beta$ -mixing or absolute regularity coefficient:

$$\beta(k) = \mathbb{E} \sup_{C \in \mathcal{F}_{t+k}^\infty} |\mathbb{P}(C) - \mathbb{P}(C|\mathcal{F}_{-\infty}^t)|$$

(iii)  $\phi$ -mixing:

$$\phi(k) = \sup_{B \in \mathcal{F}_{-\infty}^t, C \in \mathcal{F}_{t+k}^\infty} |\mathbb{P}(C) - \mathbb{P}(C|B)|$$

(iv)  $\rho$ -mixing:

$$\begin{aligned}\rho(k) &= \sup_{X \in L^2(\mathcal{F}_{-\infty}^t), Y \in L^2(\mathcal{F}_{t+k}^\infty)} |\text{corr}(X, Y)| \\ &= \sup_{X \in L^2(\mathcal{F}_{-\infty}^t), Y \in L^2(\mathcal{F}_{t+k}^\infty)} \left| \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}} \right|\end{aligned}$$

where  $L^2(\mathcal{F})$  is the set of all r.v.s defined on  $\mathcal{F}$  which have finite second moments, i.e.  $\forall X \in L^2(\mathcal{F}), \mathbb{E}X^2 < \infty$ .

# Relationship between these mixing coefficients

$$2\alpha(k) \leq \beta(k) \leq \phi(k), \quad 4\alpha(k) \leq \rho(k) \leq 2\phi^{1/2}(k)$$

The process  $\{Z_t\}_{t \in \mathbb{Z}}$  is said to be  $\alpha$ -mixing if  $\lim_{k \rightarrow \infty} \alpha(k) = 0$ . Similarly,  $\phi$ -mixing process can be defined as  $\lim_{k \rightarrow \infty} \phi(k) = 0$ , etc.

## Remark 17

*Mixings are different descriptions of dependence between events in two sigma fields  $(\mathcal{F}_{-\infty}^t, \mathcal{F}_{t+k}^\infty)$ . When the time gap between them goes to infinity, i.e.  $k \rightarrow \infty$ , mixings means asymptotic independence.*

# Relationship between these mixing coefficients

The inequality between different mixing coefficients means that,

$$\begin{array}{ccc} \phi\text{-mixing} & \implies & \beta\text{-mixing} \\ \Downarrow & & \Downarrow \\ \rho\text{-mixing} & \implies & \alpha\text{-mixing} \end{array}$$

So  $\alpha$ -mixing is the weakest mixing coefficient, but ironically has been called strong mixing.

# When a linear process is mixing?

For a linear causal process  $Z_t = \sum_{j=-\infty}^{\infty} \psi_j \epsilon_{t-j}$ ,  $\{\epsilon_t\}$  i.i.d.  $F(0, \sigma^2)$ . Gorodetskii<sup>2</sup> shows that  $\{Z_k\}$  is  $\alpha$ -mixing under certain conditions, and establish the rate of  $\alpha(k)$ .

Pham T. D. and Tran L. T.<sup>3</sup> show that if  $\psi_j = O(r^j)$  for  $0 < r < 1$  when  $j \rightarrow \infty$ , then the process is geometric  $\alpha$ -mixing, i.e. there exists  $C, \rho \in [0, 1)$  such that  $\alpha(k) \leq C\rho^k$ .

---

<sup>2</sup>VV Gorodetskii. "On the strong mixing properties for linear processes". In: *Theory of Probability and its Applications* 22 (1977), pp. 441–413.

<sup>3</sup>Tuan D Pham and Lanh T Tran. "Some mixing properties of time series models". In: *Stochastic processes and their applications* 19.2 (1985), pp. 297–303.

## Lemma 5.4 (Billingley's Inequality)

If  $\{Z_i\}$  is  $\alpha$ -mixing (NOT necessarily stationary),  $X \in \mathcal{F}_{-\infty}^t$  and  $Y \in \mathcal{F}_{t+k}^{\infty}$ ,  $|X| \leq C_1$ ,  $|Y| \leq C_2$ , then,

$$|\text{cov}(X, Y)| \leq 4C_1C_2\alpha(k)$$



Since,

$$\begin{aligned} |\operatorname{cov}(X, Y)| &= |\mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)| = |\mathbb{E}[X \{ \mathbb{E}(Y|\mathcal{F}_{-\infty}^t) - \mathbb{E}Y \}]| \\ &\leq C_1 \mathbb{E} |\mathbb{E}(Y|\mathcal{F}_{-\infty}^t) - \mathbb{E}Y| \\ &= C_1 \mathbb{E} [\xi (\mathbb{E}(Y|\mathcal{F}_{-\infty}^t) - \mathbb{E}Y)] \end{aligned}$$

where  $\xi = \operatorname{sgn}(\mathbb{E}(Y|\mathcal{F}_{-\infty}^t) - \mathbb{E}Y) \in \mathcal{F}_{-\infty}^t$ . Thus,

$$|\operatorname{cov}(X, Y)| \leq C_1 |\mathbb{E}(\xi Y) - \mathbb{E}\xi \mathbb{E}Y| = C_1 |\operatorname{cov}(\xi, Y)| \quad (24)$$

By using the same approach, we have,

$$|\operatorname{cov}(\xi, Y)| \leq C_2 |\mathbb{E}(\xi \eta) - \mathbb{E}\xi \mathbb{E}\eta| \quad (25)$$

where  $\eta = \operatorname{sgn}(\mathbb{E}(\xi|\mathcal{F}_{t+k}^\infty) - \mathbb{E}\xi) \in \mathcal{F}_{t+k}^\infty$ .

Let:

$$\begin{aligned} A &= \{\xi = 1\}, & B &= \{\eta = 1\} \\ A^C &= \{\xi = -1\}, & B^C &= \{\eta = -1\} \end{aligned}$$

Clearly,  $A, A^C \in \mathcal{F}_{-\infty}^t$  and  $B, B^C \in \mathcal{F}_{t+k}^\infty$ . Then,

$$\begin{aligned} |E(\xi\eta) - E\xi E\eta| &= |P(AB) + P(A^C B^C) - P(A^C B) - P(AB^C) \\ &\quad - (P(A) - P(A^C))(P(B) - P(B^C))| \leq 4\alpha(k) \end{aligned} \quad (26)$$

And the lemma is proved by combining (24), (25), and (26).

## Lemma 5.5

If  $\{Z_i\}$  is  $\alpha$ -mixing (NOT necessarily stationary),  $X \in \mathcal{F}_{-\infty}^t$  and  $Y \in \mathcal{F}_{t+k}^\infty$ ,  $E|X|^p < \infty$  for some  $p > 1$  and  $|Y| \leq C$ , then

$$|\text{cov}(X, Y)| \leq 6C\|X\|_p\alpha^{1/q}(k)$$

where  $\frac{1}{p} + \frac{1}{q} = 1$  and  $\|X\|_p = (E|X|^p)^{1/p}$ .

For some  $M > 0$ , let  $X_M = X\mathbb{I}(|X| \leq M)$ ,  $X'_M = X - X_M = X\mathbb{I}(|X| > M)$ . Then  $X = X_M + X'_M$ , and

$$|\text{cov}(X, Y)| = |\text{cov}(X_M, Y) + \text{cov}(X'_M, Y)| \leq |\text{cov}(X_M, Y)| + |\text{cov}(X'_M, Y)|$$

From Lemma 5.4,

$$|\text{cov}(X_M, Y)| \leq 4CM\alpha(k) \quad (27)$$

On the other hand, note that,

$$\mathbb{E}|X'_M| = \int_{|x|>M} |x| dF(x) \leq \int_{|x|>M} |x| \left(\frac{|x|}{M}\right)^{p-1} dF(x)$$

i.e.  $\mathbb{E}|X'_M| \leq M^{-p+1}\mathbb{E}|X|^p$ .

Thus,

$$\begin{aligned} |\text{cov}(X'_M, Y)| &= |\mathbb{E}(X'_M Y) - \mathbb{E}X'_M \mathbb{E}Y| \leq \mathbb{E}|X'_M Y| + C\mathbb{E}|X'_M| \\ &\leq 2C\mathbb{E}|X'_M| \leq 2CM^{-p+1}\mathbb{E}|X|^p \end{aligned} \quad (28)$$

Now, choose

$$M = \|X\|_p \{\alpha(k)\}^{-1/p}$$

and from (27) and (28), we prove the lemma.

## Lemma 5.6 (Rio's Inequality)

Let  $X$  and  $Y$  be two integrable<sup>a</sup> real-valued r.v.s and let  $Q_X(u) = \inf\{t : P(|X| > t) \leq u\}$  be the quantile function of  $|X|$ . Then if  $Q_X Q_Y$  is integrable over  $(0, 1)$ . We have:

$$|\text{cov}(X, Y)| \leq 2 \int_0^{2\alpha} Q_X(u) Q_Y(u) du$$

where  $\alpha = \alpha(\sigma(X), \sigma(Y))$  is the  $\alpha$ -mixing coefficient between sigma fields  $\sigma(X)$  and  $\sigma(Y)$ .

---

<sup>a</sup>  $\lim_{c \rightarrow \infty} E\{|X| \mathbb{I}(|X| > c)\} = 0$

Denote  $X^+ = 0 \vee X$  and  $X^- = 0 \vee (-X)$ , then

$$\begin{aligned} \operatorname{cov}(X, Y) &= \operatorname{cov}(X^+, Y^+) + \operatorname{cov}(X^-, Y^-) \\ &\quad - \operatorname{cov}(X^-, Y^+) - \operatorname{cov}(X^+, Y^-) \end{aligned} \quad (29)$$

Note that:

$$\operatorname{cov}(X^+, Y^+) = \int_{\mathbb{R}_+^2} [P(X > u, Y > v) - P(X > u)P(Y > v)]dudv$$

Recall the definition of  $\alpha$ , we have

$$\operatorname{cov}(X^+, Y^+) \leq \int_{\mathbb{R}_+^2} \inf(\alpha, P(X > u), P(Y > v))dudv \quad (30)$$

Now apply (29), (30), and the elementary inequality

$$\begin{aligned} & (\alpha \wedge a \wedge c) + (\alpha \wedge a \wedge d) + (\alpha \wedge b \wedge c) + (\alpha \wedge b \wedge d) \\ & \leq 2[(2\alpha) \wedge (a + b) \wedge (c + d)] \end{aligned}$$

to  $a = P(X > u)$ ,  $b = P(-X > u)$ ,  $c = P(Y > v)$ ,  $d = P(-Y > v)$ , we get

$$|\text{Cov}(X, Y)| \leq 2 \int_{\mathbb{R}_+^2} \inf(2\alpha, P(|X| > u), P(|Y| > v)) dudv =: I$$

Then we only need to show,

$$I = 2 \int_0^{2\alpha} Q_X(u) Q_Y(u) du \quad (31)$$



Note the form of  $I$ , define a bivariate r.v.  $(Z, T)$ ,

$$(Z, T) = (0, 0)\mathbb{I}(U \geq 2\alpha) + (Q_X(U), Q_Y(U))\mathbb{I}(U < 2\alpha)$$

where  $U$  is a uniform distributed r.v. over  $[0, 1]$ . So for any  $u, v > 0$ ,

$$\{Z > u, T > v\} = \{U < 2\alpha, U < P(|X| > u), U < P(|Y| > v)\}$$

by calculating the true integral value,

$$\begin{aligned} \int_0^{2\alpha} Q_X(u)Q_Y(u)du &= E(ZT) = \int_{\mathbb{R}_+^2} P(Z > u, T > v)dudv \\ &= \int_{\mathbb{R}_+^2} \inf(2\alpha, P(|X| > u), P(|Y| > v))dudv \end{aligned}$$

which entails (31) and the proof is thus complete.

## Lemma 5.7 (Davydov's Inequality)

Let  $X$  and  $Y$  be two real r.v.s such that  $X \in L^q(\mathcal{F}_{-\infty}^t)$ ,  $Y \in L^r(\mathcal{F}_{t+k}^\infty)$  where  $q > 1, r > 1$  and  $\frac{1}{q} + \frac{1}{r} = 1 - \frac{1}{p}$ , then

$$|\text{cov}(X, Y)| \leq 2p(2\alpha(k))^{1/p} \|X\|_q \|Y\|_r$$

(i) Suppose first that  $q$  and  $r$  are finite. Then Markov's inequality yields,

$$P\left(|X| > \frac{\|X\|_q}{u^{1/q}}\right) \leq \frac{E|X|^q}{(\|X\|_q/u^{1/q})^q} = u, \quad 0 < u \leq 1$$

which implies,

$$Q_X(u) \leq \frac{\|X\|_q}{u^{1/q}}, \quad 0 < u \leq 1$$

similarly,  $Q_Y(u) \leq \|Y\|_r/u^{1/r}$ . Using Rio's inequality,

$$\begin{aligned} |\text{cov}(X, Y)| &\leq 2 \int_0^{2\alpha(k)} \frac{\|X\|_q}{u^{1/q}} \frac{\|Y\|_r}{u^{1/r}} du \\ &= 2\|X\|_q\|Y\|_r \int_0^{2\alpha(k)} u^{\frac{1}{p}-1} du \quad \left(\frac{1}{p} + \frac{1}{q} + \frac{1}{r} = 1\right) \\ &= 2p(2\alpha(k))^{1/p}\|X\|_q\|Y\|_r \end{aligned}$$

(ii) If  $r = +\infty$ ,  $q$  is finite, then  $\frac{1}{q} + \frac{1}{p} = 1$ . Note that  $Q_Y(u) \leq Q_Y(0) = \|Y\|_\infty := \sup |Y|$ . This is because,

$$\begin{aligned} \{t \mid P(|Y| > t) = 0\} &\subset \{t \mid P(|Y| > t) \leq u\}, \quad \forall u \in [0, 1] \\ \implies \inf\{t \mid P(|Y| > t) = 0\} &\geq \inf\{t \mid P(|Y| > t) \leq u\} \\ \implies Q_Y(u) &\leq Q_Y(0) \end{aligned}$$

Also, clearly  $\inf\{t \mid P(|Y| > t) = 0\} = \|Y\|_\infty$ . Use Rio's inequality again,

$$\begin{aligned} |\text{cov}(X, Y)| &\leq 2 \int_0^{2\alpha(k)} \frac{\|X\|_q}{u^{1/q}} \|Y\|_\infty du \\ &= 2 \|X\|_q \|Y\|_\infty \int_0^{2\alpha(k)} u^{\frac{1}{p}-1} du \\ &= 2p(2\alpha(k))^{1/p} \|X\|_q \|Y\|_\infty \end{aligned}$$

which is similar to, but not exactly the same as Lemma 5.5.

(iii) If  $r = +\infty$  and  $q = +\infty$ , then  $p = 1$ . Similarly, use Rio's inequality again,

$$|\text{cov}(X, Y)| \leq 2 \times 2 \int_0^{2\alpha(k)} \|X\|_\infty \|Y\|_\infty du = 4 \|X\|_\infty \|Y\|_\infty \alpha(k)$$

which is the same as Lemma 5.4. Now, we have completed all the proof of this lemma.

# Weakly Dependent Stationary Process

Suppose  $\{X_i\}$  be a weakly stationary process with finite second moments. Let  $\gamma(j) = \text{cov}(X_i, X_{i+j})$ .

## Definition 5.8

The process is said to be weakly dependent if  $\sum_{k=0}^{\infty} |\gamma(k)| < \infty$ , or it would be said to be a long memory process.

# Weakly Dependent Stationary Process

Let  $\alpha(k)$  be the strong-mixing coefficient defined on the sigma fields generated by  $\{X_i\}_{i \in \mathbb{Z}}$ .

From Lemma 5.7 ( $r = q$ ), if  $E|X_i|^q < \infty$  for  $q > 2$ , and  $\sum_{k=0}^{\infty} \alpha^{1/p}(k) < \infty$  for

$p = \frac{q}{q-2}$ , then

$$\sum_{k=0}^{\infty} |\gamma(k)| \leq 2p \|X\|_q^2 \sum_{k=0}^{\infty} \alpha^{1/p}(k) < \infty$$

Thus, this process is weakly dependent (short-memory). In particular, if  $\alpha(k) \leq C\rho^k$ , i.e. geometric strong mixing, then

$$\sum_{k=0}^{\infty} \alpha^{1/p}(k) \leq C \sum_{k=0}^{\infty} \rho^{k/p} = \frac{C}{1 - \rho^{1/p}} < \infty.$$

## Remark 18

In general, to ensure  $\sum_{k=0}^{\infty} \alpha^{1/p}(k) < \infty$ , we require  $\alpha^{1/p}(k) \sim k^{-(1+\eta)}$ , i.e.  $\alpha(k) \sim k^{-p(1+\eta)}$  for  $\eta > 0$  when  $k$  is sufficiently large, which implies  $\alpha(k) \rightarrow 0$  as  $k \rightarrow \infty$  sufficiently fast.

## Remark 19

Note that geometric strong mixing (GSM) means  $\alpha(k) \leq C\rho^k = Ce^{-\beta k}$ , which entails  $\alpha(k) \rightarrow 0$  at exponential rate.



# Spectral Density

Define:

$$\gamma(h) = \int_{-\pi}^{\pi} e^{ih\nu} dF(\nu) = \int_{-\pi}^{\pi} e^{ih\nu} f(\nu) d\nu$$

Then by Laplace transformation:

$$f(\lambda) = \frac{1}{2\pi} \sum_{n=-\infty}^{\infty} e^{-in\lambda} \gamma(n)$$

## Theorem 5.9

If  $\sum_{k=-\infty}^{\infty} |\gamma(k)| < \infty$ , then  $\{X_n\}$  has spectral density  $f$ , and we have  $\sum_{h=-\infty}^{\infty} \gamma(h) = 2\pi f(0)$ .

## Lemma 5.10

Let  $\{X_t\}_{t \in \mathbb{Z}}$  be a zero-mean real-valued weakly stationary process such that for some  $r > 2$ ,

$$\sup_{t \in \mathbb{Z}} \mathbb{E} |X_t|^r < \infty, \quad \sum_{k \geq 1} \alpha(k)^{1 - \frac{2}{r}} < +\infty$$

then the series  $\sum_{k \in \mathbb{Z}} \gamma(k)$  is absolutely convergent, has a nonnegative sum  $\sigma^2$  and,

$$n \operatorname{var}(S_n/n) \longrightarrow \sigma^2 \quad (32)$$

where  $\gamma(k) = \operatorname{cov}(X_0, X_k)$ .

# Proof

First we study the series  $\sum_{k \in \mathbb{Z}} \gamma(k)$ , by using Lemma 5.7 with  $q = r$  and  $\frac{1}{p} = 1 - \frac{2}{r}$ , we get

$$|\gamma(k)| \leq \frac{2r}{r-2} (\mathbb{E}|X_0|^r)^{2/r} (2\alpha(k))^{1-2/r}$$

which proves the absolute convergence of the series since  $\sum_{k \geq 1} \alpha(k)^{1-2/r} < +\infty$ .  
Now clearly,

$$n \operatorname{var} \left( \frac{S_n}{n} \right) = n^{-1} \sum_{0 \leq s, t \leq n-1} \operatorname{cov}(X_s, X_t) = \sum_{k=-(n-1)}^{n-1} \left( 1 - \frac{|k|}{n} \right) \gamma(k)$$

due to  $\{X_t\}$  being weakly stationary, thus,

$$\lim_{n \rightarrow \infty} n \operatorname{var} \left( \frac{S_n}{n} \right) = \sigma^2 \geq 0$$

and the theorem is thus established.

## Theorem 5.11

Let  $\{X_t\}_{t \in \mathbb{Z}}$  be a zero-mean real-valued strictly stationary process such that for some  $r > 2$  and some  $\beta > 0$ ,

$$\mathbb{E}|X_t|^r < \infty, \quad \alpha(k) \leq ak^{-\beta}$$

where  $a$  is a positive constant and  $\beta > r/(r-2)$ . Then if  $\sigma^2 = \sum_{k=-\infty}^{\infty} \gamma(k) > 0$ , we have

$$\frac{S_n}{\sigma\sqrt{n}} \xrightarrow{d} N(0, 1)$$

The proof can be seen in Theorem 1.7 of D. Bosq's book<sup>4</sup>.

---

<sup>4</sup>Denis Bosq. *Nonparametric statistics for stochastic processes: estimation and prediction*. Vol. 110. Springer Science & Business Media, 2012.

## Chapter 5: Delta Method

Suppose we have a sequence of estimators  $\{T_n\}_{n \geq 1}$  on  $\mathbb{R}^k$  for a parameter  $\theta \in \mathbb{R}^k$ .

- For the  $\phi(\theta)$ , the parameter of interest, considering convergence of  $\phi(T_n)$  to  $\phi(\theta)$ , where  $\phi : \mathbb{R}^k \rightarrow \mathbb{R}^m$ .
- From Mapping,  $T_n \xrightarrow{P} \theta \Rightarrow \phi(T_n) \xrightarrow{P} \phi(\theta)$ , if  $\phi$  was continuous at  $\theta$ .
- Does  $\sqrt{n}(\phi(T_n) - \phi(\theta))$  have an asymptotic distribution, if suppose  $\sqrt{n}(T_n - \theta) \xrightarrow{d} T$ ?
- Is the convergence in distribution persevered under smooth transformation?

# The Derivative of Vector-valued Functions

Recall that  $\phi(\cdot)$  is differentiable at  $\theta$  if there exists a linear map (matrix)  $\phi'_\theta : \mathbb{R}^k \mapsto \mathbb{R}^m$  such that

$$\begin{aligned}\phi(\theta + h) - \phi(\theta) &:= \phi'(\theta)h + R(h) \\ &= \phi'(\theta)h + o(\|h\|), \quad h \rightarrow 0.\end{aligned}$$

Define  $\phi'_\theta(h) := \phi'(\theta)h$ .

## Derivative map (Jacobian matrix)

The derivative map  $h \mapsto \phi'_\theta(h)$  is matrix multiplication by the matrix

$$\phi'(\theta) \triangleq \phi'_\theta = \begin{pmatrix} \frac{\partial \phi_1}{\partial \theta_1}(\theta) & \cdots & \frac{\partial \phi_1}{\partial \theta_k}(\theta) \\ \vdots & & \vdots \\ \frac{\partial \phi_m}{\partial \theta_1}(\theta) & \cdots & \frac{\partial \phi_m}{\partial \theta_k}(\theta) \end{pmatrix} = \left( \frac{\partial \phi_i(\theta)}{\partial \theta_j} \right)_{m \times k}.$$

## Remark 20

If  $m = 1, k > 1$ , the derivative map is called the gradient of the function.

# Main Result of Delta Method

## Theorem 6.1 (First Order Delta Method)

If  $\phi$  differentiable at  $\theta$ ,  $\phi'(\theta) \neq 0$ , and  $r_n(T_n - \theta) \xrightarrow{d} T$  for a deterministic sequence of  $\{r_n\}$ , satisfied  $r_n \rightarrow \infty$ , then:

- (i)  $r_n(\phi(T_n) - \phi(\theta)) - \phi'(\theta)(r_n(T_n - \theta)) \xrightarrow{P} 0$ ;
- (ii)  $r_n(\phi(T_n) - \phi(\theta)) \xrightarrow{d} \phi'(\theta)T$ .

## The proof: (i)

- Since  $r_n(T_n - \theta) \xrightarrow{d} T$ , then  $T_n - \theta \xrightarrow{P} 0$  by stochastic boundedness.
- By the differentiability of  $\phi$  at  $\theta$ , then

$$\phi(\theta + h) - \phi(\theta) - \phi'(\theta)h := R(h) = o(\|h\|).$$

- Replace the  $h$  by  $T_n - \theta$ , multiply  $r_n$  to get

$$r_n[\phi(T_n) - \phi(\theta) - \phi'_{\theta}(T_n - \theta)] = o_P(r_n \|T_n - \theta\|) = o_P(1).$$

by Lemma of Stochastic Plug-in (since  $T_n - \theta \xrightarrow{P} 0$ ).

## The proof: (ii)

- Matrix multiplication is continuous, so  $\phi'_\theta(r_n(T_n - \theta)) \xrightarrow{d} \phi'_\theta(T)$  by the continuous-mapping theorem.
- Apply Slutsky's lemma to conclude that

$$r_n(\phi(T_n) - \phi(\theta)) \xrightarrow{d} \phi'(\theta)T$$

which has the same weak limit as  $\phi'_\theta(r_n(T_n - \theta))$ .

## Example: Normal delta method

Let  $T_n$  be a sequence of statistics such that

$$\sqrt{n}(T_n - \theta) \xrightarrow{d} N(0, \sigma^2(\theta))$$

Let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be once differentiable at  $\theta$  with  $g'(\theta) \neq 0$ . Then

$$\sqrt{n}[g(T_n) - g(\theta)] \xrightarrow{d} N(0, [g'(\theta)]^2 \sigma^2(\theta))$$



# High Order Delta Method: what if $\phi'(\theta) = 0$ ?

First Order Delta Method is largely based on Taylor's expansions with  $\phi'(\theta) \neq 0$ . If  $\phi'(\theta) = 0$  but  $\phi''(\theta) \neq 0$ , we have

$$\phi(T_n) = \phi(\theta) + \frac{1}{2}\phi''(\theta)(T_n - \theta)^2 + \dots$$

Then

$$n(\phi(T_n) - \phi(\theta)) = \frac{1}{2}\phi''(\theta)[\sqrt{n}(T_n - \theta)]^2 \xrightarrow{d} \dots$$

If  $\sqrt{n}\bar{X}_n \xrightarrow{d} N(0, 1)$ , then  $n\phi(\bar{X}_n) \xrightarrow{d} \frac{1}{2}\phi''(\theta)\chi_1^2$ .

## Theorem 6.2

Suppose  $\phi$  be a univ.  $m$  times differentiable at  $\theta$  with  $\phi^{(m)}(\theta) \neq 0$ ,  $\phi^{(j)}(\theta) = 0$ ,  $j < m$ , then:

$$\frac{r_n^m (\phi(T_n) - \phi(\theta))}{\frac{1}{m!}\phi^{(m)}(\theta)} \xrightarrow{d} T^m.$$

A multivariable version of this theorem is available in Serfling P124.

## Examples: 2ed Delta Method

- ① Suppose  $X_1, \dots, X_n$  are iid with mean  $\mu$  and known variance  $\sigma^2$ , and we want to test  $H_0 : \mu = 0$ . Under the null hypothesis  $H_0 : \mu = 0$ , the following statistic

$$T(\mathbf{X}) := n\bar{X}_n^2/\sigma^2 \xrightarrow{d} [N(0, 1)]^2 = \chi_1^2.$$

- ② Suppose that  $\sqrt{n}\bar{X}_n$  converges in law to a standard normal distribution. Now consider the limiting behavior of  $\cos(\bar{X}_n)$ .
- Because the derivative of  $\cos(x)$  is zero at  $x = 0$ , we still use the proof of First Order Delta Method. It yields that

$$\sqrt{n}(\cos(\bar{X}_n) - 1) \xrightarrow{d} \delta_0$$

which implies that  $\sqrt{n}(\cos(\bar{X}_n) - 1) \xrightarrow{P} \delta_0$ .

- Thus, it should be concluded that  $\sqrt{n}$  is not the right norming rate for the random sequence  $\cos(\bar{X}_n) - 1$ . 2ed Order Delta Method

$$\cos \bar{X} - \cos 0 = (\bar{X} - 0)0 + \frac{1}{2}(\bar{X} - 0)^2(\cos x)''|_{x=0} + \dots$$

implies

$$-2n(\cos \bar{X} - 1) \xrightarrow{d} \chi_1^2.$$

## Example: Variance

Let  $X_1, \dots, X_n$  i.i.d.  $F$  with finite 4-th moments. Let  $\alpha_i = EX_1^i$  for  $i = 1, 2, 3, 4$  and  $m_{n1} = n^{-1} \sum_{j=1}^n X_j$ . Then,

$$S_n = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2 = n^{-1} \sum_{i=1}^n X_i^2 - \bar{X}^2 = \phi(m_{n1}, m_{n2})$$

where  $\phi(x_1, x_2) = x_2 - x_1^2$ . From MCLT,

$$\sqrt{n} \left[ \begin{pmatrix} m_{n1} \\ m_{n2} \end{pmatrix} - \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} \right] \xrightarrow{d.} N \left( 0, \text{var} \begin{pmatrix} X_1 \\ X_1^2 \end{pmatrix} \right)$$

and  $\phi'(\alpha_1, \alpha_2) = (-2\alpha_1, 1)$ . Hence,

$$\begin{aligned} \sqrt{n}(S_n - \sigma^2) &= \sqrt{n}(\phi(m_{n1}, m_{n2}) - \phi(\alpha_1, \alpha_2)) \\ &\xrightarrow{d.} (-2\alpha_1, 1) N \left( 0, \text{var} \begin{pmatrix} X_1 \\ X_1^2 \end{pmatrix} \right) = N(0, c_4 - c_2^2) \end{aligned}$$

where we use the fact that:

$$(-2\alpha_1, 1) \text{var} \begin{pmatrix} X_1 \\ X_1^2 \end{pmatrix} \begin{pmatrix} -2\alpha_1 \\ 1 \end{pmatrix} = E(X_1 - \alpha_1)^4 - (E(X_1 - \alpha_1)^2)^2 = c_4 - c_2^2$$

## Example: Standard Deviation

Considering the unbiased estimator:

$$S_{n-1} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} S_n$$

So,

$$\begin{aligned} \sqrt{n}(S_{n-1} - \sigma^2) &= \sqrt{n}\left(\frac{n}{n-1} S_n - \sigma^2\right) \\ &= \sqrt{n}\left(S_n - \sigma^2 + \left(\frac{n}{n-1} - 1\right) S_n\right) \\ &= \sqrt{n}(S_n - \sigma^2) + \sqrt{n}\left(\frac{n}{n-1} - 1\right) S_n \\ &\xrightarrow{d.} N(0, c_4 - c_2^2) \quad \left(\sqrt{n}\left(\frac{n}{n-1} - 1\right) S_n = o_p(1)\right) \end{aligned}$$

Furthermore,  $S_n^{1/2} = \sqrt{S_n} = \phi(S_n)$ ,  $\phi(x) = \sqrt{x}$ , and  $\phi'(x) = \frac{1}{2}x^{-1/2}$ ,

$$\sqrt{n}(S_n^{1/2} - \sigma) \xrightarrow{d.} N\left(0, \frac{c_4 - c_2^2}{4\sigma^2}\right)$$

# More Examples

If  $X_n$  is  $AN(\mu, \sigma_n^2)$  and  $\sigma_n \rightarrow 0$ . Then,

- (i)  $X_n^2$  is  $AN(\mu^2, 4\mu^2\sigma_n^2)$  for  $\mu \neq 0$ .
- (ii)  $\frac{1}{X_n}$  is  $AN\left(\mu^{-1}, \frac{\sigma_n^2}{\mu^4}\right)$  for  $\mu \neq 0$ .
- (iii)  $e^{X_n}$  is  $AN(e^\mu, e^{2\mu}\sigma_n^2)$  for any  $\mu$ .
- (iv)  $\log |X_n|$  is  $AN(\log |\mu|, \mu^{-2}\sigma_n^2)$  if  $\mu \neq 0$ ;  $\log |\sigma_n^{-1}X_n| \xrightarrow{d.} \log |N(0, 1)|$  for  $\mu = 0$ .
- (v) Suppose  $X_1, \dots, X_n$  i.i.d.  $F$  on  $\mathbb{R}^p$  with  $(\mu, \Sigma)$ . Let

$$\theta = \mu^\top \mu, \quad \hat{\theta} = \bar{X}^\top \bar{X} = \phi(\bar{X})$$

If  $\mu \neq 0$ ,  $\phi'(\mu) = 2\mu^\top$ ,  $\phi''(\mu) = 2I_p$ . As  $\sqrt{n}(\bar{X} - \mu) \xrightarrow{d.} N_p(0, \Sigma)$ .  
So,

$$\sqrt{n} \left( \bar{X}^\top \bar{X} - \mu^\top \mu \right) \xrightarrow{d.} 2\mu^\top N_p(0, \Sigma) = N(0, 4\mu^\top \Sigma \mu)$$

## Example: weighted $\chi_1^2$ distribution

If  $\mu = 0$ ,  $\mu^\top \Sigma \mu = 0$ , the above display  $\xrightarrow{d.} 0$  is not useful. In fact, as  $\sqrt{n}\bar{X} \xrightarrow{d.} N_p(0, \Sigma)$ ,

$$n\bar{X}^\top \bar{X} \xrightarrow{d.} N_p^\top(0, \Sigma) N_p(0, \Sigma) \stackrel{d.}{=} Z^\top \Sigma^{1/2} \Sigma^{1/2} Z = Z^\top \Sigma Z$$

where  $Z \sim N(0, I_p)$ . Suppose

$$\Sigma = U^\top \text{diag}(\lambda_1, \dots, \lambda_p) U, \quad \tilde{Z} = UZ \sim N_p(0, I_p)$$

Then,

$$Z^\top \Sigma Z \stackrel{d.}{=} \sum_{i=1}^p \lambda_i \tilde{Z}_i^2 \stackrel{d.}{=} \sum_{i=1}^p \lambda_i \chi_{1i}^2$$

where  $\{\chi_{1i}^2\}_{i=1}^p$  i.i.d.  $\chi_1^2$ .

So  $n\bar{X}^\top \bar{X}$  converges to a weighted  $\chi_1^2$  distribution.

Thus  $n\bar{X}^\top \bar{X} = O_p(1)$  if  $\mu = 0$ , and  $\bar{X}^\top \bar{X} - \mu^\top \mu = O_p\left(\frac{1}{\sqrt{n}}\right)$  if  $\mu \neq 0$ .

# Chi-Square Test for Variance

Suppose  $X_1, \dots, X_n$  i.i.d.  $F$  with  $EX_1^4 < \infty$ .

$$H_0 : \sigma^2 \leq 1 \quad \text{vs} \quad H_1 : \sigma^2 > 1$$

Denote  $\Theta_0 = \{0 < \sigma^2 \leq 1\}$  and  $\Theta_1 = \{\sigma^2 > 1\}$ . If  $F = N(\mu, \sigma^2)$ ,  $\frac{nS_n}{\sigma^2} \sim \chi_{n-1}^2$ .

Test statistic for  $H_0 : \sigma^2 \leq 1$  is  $nS_n$  by setting  $\sigma^2 = 1$ , and we reject  $H_0$  if  $nS_n > \chi_{n-1, \alpha}^2$ . The size of this test is

$$\begin{aligned} P_{\Theta_0} (nS_n > \chi_{n-1, \alpha}^2 \mid \sigma^2 \in \Theta_0) &= P_{\Theta_0} \left( \frac{nS_n}{\sigma^2} > \frac{1}{\sigma^2} \chi_{n-1, \alpha}^2 \mid \sigma^2 \leq 1 \right) \\ &\leq P(\chi_{n-1}^2 > \chi_{n-1, \alpha}^2) = \alpha \end{aligned}$$

So the size  $\leq \alpha$  with the maximum size at  $\sigma^2 = 1$  equals to  $\alpha$ .

# Chi-Square Test for Variance

If  $F \neq$  Normal, the excessive kurtosis:

$$\kappa = \frac{E(X - \mu)^4}{\sigma^4} - 3 \neq 0$$

From CLT and the fact that  $\chi_{n-1}^2 = \sum_{i=1}^{n-1} Z_i^2$  for  $\{Z_i\}_{i=1}^{n-1}$  i.i.d.  $N(0, 1)$ . Then,

$$\frac{\chi_{n-1}^2 - (n-1)}{\sqrt{2(n-1)}} \xrightarrow{d.} N(0, 1) \quad (33)$$

From the previous example, we know that:

$$\sqrt{n} \left( \frac{S_n}{\sigma^2} - 1 \right) \xrightarrow{d.} N(0, \kappa - 1) \neq N(0, 2)$$

And from (33),

$$P \left( \frac{\chi_{n-1}^2 - (n-1)}{\sqrt{2(n-1)}} \geq Z_\alpha \right) \rightarrow P(N(0, 1) \geq Z_\alpha) = \alpha$$



# Chi-Square Test for Variance

As  $P(\chi_{n-1}^2 > \chi_{n-1,\alpha}^2) = \alpha$ , we have  $\chi_{n-1,\alpha}^2 \approx (n-1) + Z_\alpha \sqrt{2(n-1)}$ . i.e.

$$\lim_{n \rightarrow \infty} \frac{\chi_{n-1,\alpha}^2 - (n-1)}{\sqrt{2(n-1)}} = Z_\alpha$$

Consequently, the level of the Chi-Square Test is:

$$\begin{aligned} P_{\sigma^2=1} \left( \frac{nS_n}{\sigma^2} > \chi_{n-1,\alpha}^2 \right) &= P_{\sigma^2=1} \left( \sqrt{n} \left( \frac{S_n}{\sigma^2} - 1 \right) > \frac{\chi_{n-1,\alpha}^2 - n}{\sqrt{n}} \right) \\ &\approx P_{\sigma^2=1} \left( \sqrt{n} \left( \frac{S_n}{\sigma^2} - 1 \right) > \frac{(n-1) + Z_\alpha \sqrt{2(n-1)} - n}{\sqrt{n}} \right) \\ &\rightarrow P(N(0, \kappa + 2) > \sqrt{2}Z_\alpha) = 1 - \Phi \left( \frac{\sqrt{2}Z_\alpha}{\sqrt{\kappa + 2}} \right) \end{aligned}$$

For heavy-tail  $F$ ,  $\kappa > 0$ , so that,

$$1 - \Phi \left( \frac{\sqrt{2}Z_\alpha}{\sqrt{\kappa + 2}} \right) > 1 - \Phi(Z_\alpha) = \alpha$$

# Chi-Square Test for Variance

Power of this test:

$$\begin{aligned} P_{\sigma^2 > 1} \left( \frac{nS_n}{\sigma^2} > \chi_{n-1, \alpha}^2 \right) &= P_{\sigma^2 > 1} \left( \sqrt{n} \left( \frac{S_n}{\sigma^2} - 1 \right) > \frac{\sigma^{-2} \chi_{n-1, \alpha}^2 - n}{\sqrt{n}} \right) \\ &\rightarrow 1 - P \left( N(0, \kappa + 2) > \frac{\sigma^{-2} \{ (n-1) + Z_\alpha \sqrt{2(n-1)} \} - n}{\sqrt{n}} \right) \\ &= 1 - \Phi \left( \frac{(\sigma^{-2} - 1)\sqrt{n}}{\sqrt{\kappa + 2}} - \frac{\sigma^{-2}}{\sqrt{n(\kappa + 2)}} + \frac{\sqrt{2}Z_\alpha}{\sqrt{\kappa + 2}} \sqrt{\frac{n-1}{n}} \right) \\ &\approx 1 - \Phi \left( \frac{(\sigma^{-2} - 1)\sqrt{n}}{\sqrt{\kappa + 2}} + \frac{\sqrt{2}}{\sqrt{\kappa + 2}} Z_\alpha \right) \rightarrow 1 \end{aligned}$$

So the power  $\rightarrow 1$  as  $n \rightarrow \infty$ , the test is consistent.

# Multinomial Vectors and $\chi^2$ statistic

Let  $(n_1, \dots, n_K)$  be multinomial  $(n; p_1, \dots, p_K)$  with each  $p_i > 0$ . Then

$$X_n = \sqrt{n} \left( \frac{n_1}{n} - p_1, \dots, \frac{n_K}{n} - p_K \right) := (X_{n1}, \dots, X_{nK}) \xrightarrow{d.} N(0, \Sigma)$$

where  $\Sigma = (\sigma_{ij})$  with

$$\sigma_{ij} = \begin{cases} p_i(1 - p_j) & i = j \\ -p_i p_j & i \neq j \end{cases}$$

A test statistic for goodness-of-fit is:

$$\begin{aligned} T_n &= \sum_{i=1}^K \frac{(n_i - np_i)^2}{np_i} = n \sum_{i=1}^K \frac{1}{p_i} \left( \frac{n_i}{n} - p_i \right)^2 \\ &= \sum_{i=1}^K \frac{1}{p_i} X_{ni}^2 = X_n^\top C X_n \end{aligned}$$

where  $C = \text{diag}(p_1^{-1}, \dots, p_K^{-1})$ .

# Multinomial Vectors and $\chi^2$ statistic

By mapping theorem,

$$T_n \xrightarrow{d.} Z^\top \Sigma^{1/2} C \Sigma^{1/2} Z \stackrel{d.}{=} \chi_{n-1}^2$$

In fact, we can show that  $A = \Sigma^{1/2} C \Sigma^{1/2}$  is an idempotent:

$$\begin{aligned} \sigma_{ij} &= p_i(\delta_{ij} - p_j), & C\Sigma &= (p_i^{-1}\sigma_{ij}) = (\delta_{ij} - p_j) \\ (C\Sigma)^2 &= \left( \sum_{l=1}^K (\delta_{il} - p_l)(\delta_{lj} - p_j) \right) = (\delta_{ij} - p_j) = C\Sigma \end{aligned}$$

As a result,  $C\Sigma$  is an idempotent, so is  $A$ , which entails:

$$\text{tr} \left( \Sigma^{1/2} C \Sigma^{1/2} \right) = \text{tr}(C\Sigma) = n - 1$$

# Wald Test

Suppose  $X_1, \dots, X_n$  i.i.d.  $F$  on  $\mathbb{R}^p$  with  $\mu$  and  $\Sigma > 0$  ( $p$  fixed).

$$H_0 : \mu = \mu_0 \quad \text{vs} \quad H_1 : \mu \neq \mu_0$$

Wald statistic is:

$$W_n = n(\bar{X} - \mu_0)^\top S_n^{-1}(\bar{X} - \mu_0), \quad S_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^\top$$

From LLN,  $S_n \xrightarrow{P} \Sigma$ ,  $S_n^{-1} \xrightarrow{P} \Sigma^{-1}$ . And we note that:

$$\bar{X} - \mu = O_p(n^{-1/2})$$

Hence,

$$\begin{aligned}W_n &= n(\bar{X} - \mu_0)^\top (\Sigma^{-1} + S_n^{-1} - \Sigma^{-1}) (\bar{X} - \mu_0) \\&= \sqrt{n}(\bar{X} - \mu_0)^\top \Sigma^{-1} \sqrt{n}(\bar{X} - \mu_0) \\&\quad + \sqrt{n}(\bar{X} - \mu_0)^\top (S_n^{-1} - \Sigma^{-1}) \sqrt{n}(\bar{X} - \mu_0) \\&= \sqrt{n}(\bar{X} - \mu_0)^\top \Sigma^{-1} \sqrt{n}(\bar{X} - \mu_0) + o_p(1) \xrightarrow{d.} \chi_p^2\end{aligned}$$

as  $\sqrt{n}\Sigma^{-1/2}(\bar{X} - \mu_0) \xrightarrow{d.} N(0, I_p)$ .

So Wald test requires  $H_0$  if  $W_n > \chi_{p,1-\alpha}^2$ .

# Variance Stabilizing Transform (VST)

Typically, we have various means to obtain,

$$\sqrt{n}(T_n - \theta) \xrightarrow{d.} N(0, \sigma^2(\theta))$$

where  $\sigma^2(\theta)$  is the asymptotic variance depending on  $\theta$ . The asymptotic CI for  $\theta$  are:

$$T_n \pm Z_{1-\alpha/2} \frac{\sigma(\hat{\theta})}{\sqrt{n}}$$

So the width of the CIs varies with respect to  $\sigma(\theta)$ .

The purpose of VST is to transform  $T_n$  to  $\phi(T_n)$  such that

$$\sqrt{n}(\phi(T_n) - \phi(\theta)) \xrightarrow{d.} N(0, c^2)$$

where  $c > 0$  is a constant.

# Variance Stabilizing Transform (VST)

From earlier result,

$$\sqrt{n}(\phi(T_n) - \phi(\theta)) \xrightarrow{d.} \phi'(\theta)N(0, \sigma^2(\theta)) \stackrel{d.}{=} N(0, (\phi'(\theta))^2 \sigma^2(\theta))$$

So  $\phi'(\theta)\sigma(\theta) = c$ , which implies:

$$\phi'(\theta) = \frac{c}{\sigma(\theta)}, \quad \phi(\theta) = \int \frac{d\theta}{\sigma(\theta)}$$

is the VST.



# Tukey's Hanging Rootogram

Let  $X_1, \dots, X_n$  i.i.d. the pdf  $f$ . The Kernel Density Estimator is:

$$\hat{f}_{nh}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

where  $K$  is a symmetric pdf,  $N(0, 1)$ 's pdf, for instance. It can be shown (Serfling<sup>5</sup>, P114) that:

$$\sqrt{nh}(\hat{f}_{nh}(x) - f(x)) \xrightarrow{d} N(0, f(x))$$

provided  $nh^5 \rightarrow 0$  and  $nh \rightarrow \infty$  as  $n \rightarrow \infty$ . So  $\hat{f}_{nh}(x)$  is  $AN(f(x), \frac{f(x)}{nh})$ . By Delta method and VST:

$$\phi(f) = \int \frac{df}{\sqrt{f}} = f^{1/2}$$

So we do "root-gram":  $\hat{f}_{nh}^{1/2}(x)$  is:

$$AN(f^{1/2}(x), \frac{1}{4hn})$$

<sup>5</sup>Robert J Serfling. *Approximation theorems of mathematical statistics*. Vol. 162. John Wiley & Sons, 2009. 

# Asymptotically Uniformly Integrable

## Definition 6.3 (Uniform Integrability)

A sequence of random variables  $\{Y_n\}_{n \geq 0}$  is called asymptotic uniformly integrable (u.i.) if:

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} E \left[ |Y_n| \mathbb{I}_{\{Y_n > M\}} \right] = 0$$

The uniform integrability is the missing link between convergence in distribution and convergence of moments.

## Theorem 6.4

Let  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  be measurable and continuous at every point in a set  $C$ ,  $X_n \xrightarrow{d} X$  where  $X$  takes its values in  $C$ . Then  $Ef(X_n) \rightarrow Ef(X)$  if and only if the sequence of r.v.  $f(X_n)$  is asymptotically u.i.

# Moment Approximation

If  $T_n$  has m-th moment exist. Knowing  $\sqrt{n}(T_n - \theta) \xrightarrow{d} T \Rightarrow \sqrt{n}(\phi(T_n) - \phi(\theta)) \xrightarrow{d} \phi'(\theta)T$ , if  $\phi'(\theta) \neq 0$ .

- Can we approximate  $E\phi(T_n)$  by Taylor expansion?

$$\phi(T_n) = \phi(\theta) + \phi'(\theta)(T_n - \theta) + \frac{1}{2}\phi''(\theta)(T_n - \theta)^2 + \dots$$

- So that, do we have the following equations:

$$E\phi(T_n) \approx \phi(\theta) + \phi'(\theta)\text{Bias}(T_n) + \frac{1}{2}\phi''(\theta)\text{MSE}(T_n)$$

$$\text{var}(\phi(T_n)) \approx (\phi'(\theta))^T \text{var}(T_n) (\phi'(\theta))$$

- We need  $\phi(T_n) - \phi(\theta)$  being u.i. If  $T_n - \theta$  is u.i. and  $\phi$  is Lipschitz, then  $\phi(T_n) - \phi(\theta)$  is u.i..
- See also Sargan, J.D. (1976, Econometrica).

## Chapter 6: Moment Estimator (ME)

Let  $X_1, \dots, X_n$  i.i.d.  $F_\theta$ , where  $\theta_0$  is the true parameter,  $f_1, \dots, f_k$  be given known function.

- Moments:

$$E_\theta f_j(X) = \int f_j(x) dF_\theta(x), \quad j = 1, \dots, k$$

- A popular or original choice is:  $f_j(x) = x^j$ . Let  $f = (f_1, \dots, f_k)'$ .

### Definition 7.1 (Moment Estimator (ME))

Match sample moments  $\frac{1}{n} \sum_{i=1}^n f(x_i)$  with its population counterparts:

$$P_n f := \frac{1}{n} \sum_{i=1}^n f(x_i) = e(\theta) = E_\theta f(X) := P_\theta f$$

- If  $e$  is one to one, then the ME is  $\hat{\theta}_n = e^{-1}(P_n f)$ .
- If  $e^{-1}$  is differential and  $E_{\theta_0} f(X) f^T(X) < \infty$  (which implies the AN of  $P_n f$ ), then we can have the AN of  $\hat{\theta}_n$ .
- Note that:

$$(e^{-1}(x_0))' = (e'(\theta_0))^{-1} \Big|_{\theta_0=e^{-1}(x_0)}$$

## Theorem 7.2 (CLT for ME)

If  $e(\theta) = P_\theta f =: E_\theta f(X)$  is 1-1 on an open set  $\Theta \subset \mathbb{R}^k$  and is continuously differentiable at  $\theta_0$  with non-singular  $e'_\theta$  and  $P_{\theta_0} \|f\|^2 < \infty$ , then  $\hat{\theta}_n$  exists with prob approaching to 1 (wpa 1) and:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N\left(0, [(e'(\theta_0))^{-1}](E_{\theta_0} f f^T - E_{\theta_0} f E_{\theta_0} f^T)[(e'(\theta_0))^{-1}]^T\right)$$

We need the inverse function in the proof.

### Lemma 7.3 (The inverse function theorem)

Let  $A$  be open in  $\mathbb{R}^n$ ; let  $g : A \mapsto \mathbb{R}^k$  is continuously differentiable at  $\mathbf{a}$  and differentiable in a neighborhood of  $\mathbf{a} \in A$ . If the Jacobi matrix

$$Dg(x) := \partial g(x) / \partial x^\tau$$

is non-singular at the point  $\mathbf{a}$  and of  $A$ .

Then,

- there is a neighborhood  $U$  of the point  $\mathbf{a}$ , such that

$$g : U \mapsto V \text{ is one-to-one}$$

for an open set  $V$  of  $\mathbb{R}^k$ ;

- and there exists an inverse function  $g^{-1} : V \mapsto U$  which is continuously differentiable with

$$Dg^{-1}(y) := \partial g^{-1}(y) / \partial y^\tau = (Dg(x))^{-1}.$$

## Proof of Theorem 7.2

- Continuous differentiability at  $\theta_0$  presumes differentiability in a neighborhood and the continuity of  $\theta \mapsto e'_\theta$ ;  
the nonsingularity of  $e'_\theta$  implies nonsingularity in a neighborhood.
- Therefore, by the **inverse function theorem** there exist open neighborhoods  $U$  of  $\theta_0$ , and  $V$  of  $P_{\theta_0}f$  such that

$e : U \mapsto V$  is a differentiable bijection (one-to-one)  
with a differentiable inverse  $e^{-1} : V \mapsto U$ .

- By the LLN,  $\mathbb{P}_n f \equiv \frac{1}{n} \sum_{i=1}^n f(X_i) \xrightarrow{a.s.} e(\theta_0) \in V$ . Since  $\mathbb{P}_n f \in V$ ,

$$\hat{\theta}_n := e^{-1}(\mathbb{P}_n f) \xrightarrow{a.s.} e^{-1}(e(\theta_0)) = \theta_0$$

exist with probability tending to 1 by continuous mapping theorem.

- The CLT guarantees asymptotic normality of the sequence  $\sqrt{n}(\mathbb{P}_n f - P_{\theta_0} f)$ . The proof is finished by Delta Method.

## Example : ME for the Beta distribution

Let  $X_1, X_2, \dots, X_n$  be a random sample from the Beta distribution which has a density function  $f(x; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$  for  $x \in (0, 1)$  where  $\alpha > 0$  and  $\beta > 0$  are two unknown parameters, and  $B(\alpha, \beta)$  is the Beta function. Note that

$$\begin{aligned} \mathbb{E}X^k &= \frac{1}{B(\alpha, \beta)} \int x^{\alpha+k-1} (1-x)^{\beta-1} dx = \frac{B(\alpha+k, \beta)}{B(\alpha, \beta)} \\ &= \frac{\Gamma(\alpha+k)}{\Gamma(\alpha)} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha+\beta+k)} = \prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r}, \quad k=1, 2, \dots \end{aligned}$$

then the moment estimator can be solved by the following equations:

$$\bar{X} = \frac{\alpha}{\alpha+\beta}, \quad \overline{X^2} = \frac{\alpha(\alpha+1)}{(\alpha+\beta)(\alpha+\beta+1)}$$

namely,  $\hat{\alpha} = (1 - \bar{X}) \left[ \frac{\bar{X}(\bar{X}-1)}{\bar{X}^2 - \bar{X}^2} - 1 \right]$ ,  $\hat{\beta} = \bar{X} \left[ \frac{\bar{X}(\bar{X}-1)}{\bar{X}^2 - \bar{X}^2} - 1 \right]$  is the solutions.



## Example of Beta distribution: Con.

Let  $\theta = (\alpha, \beta)$  be the set of true parameters vector. The moment function is  $f(x) = (x, x^2)^T$  and the estimated function is:

$$e(\theta) = \mathbb{E}_\theta f(x) = \left( \frac{\alpha}{\alpha + \beta}, \frac{\alpha(\alpha + 1)}{(\alpha + \beta)(\alpha + \beta + 1)} \right)^T$$

It is very easy to verify that  $e^{-1} \in C^\infty(\mathbb{R}_+^2)$ :

$$\hat{\theta} = e^{-1}(\mathbb{P}_n f) \xrightarrow{P} e^{-1}(P_n f) = \theta.$$

Note that  $\frac{\partial e(\theta)}{\partial \theta^T} = \begin{bmatrix} \frac{\beta}{(\alpha + \beta)^2} & -\frac{\alpha}{(\alpha + \beta)^2} \\ -\frac{\beta(\alpha + 1)(\alpha + \beta + 1) + \alpha\beta(\alpha + \beta)}{(\alpha + \beta)^2(\alpha + \beta + 1)^2} & -\frac{\alpha(\alpha + 1)(2\alpha + 2\beta + 1)}{(\alpha + \beta)^2(\alpha + \beta + 1)^2} \end{bmatrix}$ , so

$$\mathbb{E}_\theta f f^T - \mathbb{E}_\theta f \mathbb{E}_\theta f^T = \begin{bmatrix} \mathbb{E}X^2 & \mathbb{E}X^3 \\ \mathbb{E}X^3 & \mathbb{E}X^4 \end{bmatrix} - \begin{bmatrix} (\mathbb{E}X)^2 & \mathbb{E}X\mathbb{E}X^2 \\ \mathbb{E}X\mathbb{E}X^2 & (\mathbb{E}X^2)^2 \end{bmatrix}.$$

as a result, we have:

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N\left(0, \left(\frac{\partial e(\theta)}{\partial \theta^T}\right)^{-1} [\mathbb{E}_\theta f f^T - \mathbb{E}_\theta f \mathbb{E}_\theta f^T] \left(\frac{\partial e(\theta)}{\partial \theta^T}\right)^{-1}\right) =: N(0, \Sigma)$$

# Exponential Family

Suppose  $X_1, \dots, X_n$  i.i.d.  $F_\theta$  with density:

$$p_\theta(x) = c(\theta)h(x) \exp\{\theta^T t(x)\}$$

- Likelihood:

$$l_\theta(x) = \log p_\theta(x) = \log c(\theta) + \log h(x) + \theta^T t(x)$$

- Likelihood score:

$$\dot{l}_\theta(x) = \frac{\dot{c}(\theta)}{c(\theta)} + t(x) = t(x) - E_\theta t(X)$$

- $E \dot{l}_\theta(X) = 0$  implies  $\dot{c}(\theta)/c(\theta) = -E_\theta t(X)$ .

Hence, MLE are MEs:

$$\frac{1}{n} \sum_{i=1}^n t(x_i) = E_\theta t(X) = e(\theta)$$

# Exponential Family

Furthermore, if  $e(\theta)$  is continuous and differentiable and a moment condition  $E_{\theta} \|t(X)\|^2 < \infty$  is satisfied, then:

- $\hat{\theta}_{mle} = \hat{\theta}_{me}$  exists wpa 1;
- and

$$\begin{aligned}\sqrt{n}(\hat{\theta}_{me} - \theta) &\xrightarrow{d} N\left(0, [(e'(\theta))^{-1}] \text{var}_{\theta}(t(x)) [(e'(\theta))^{-1}]^T\right) \\ &= N(0, I_{\theta}^{-1})\end{aligned}$$

- $I_{\theta}$  is Fisher Information matrix:

$$I_{\theta} = \text{var}(i_{\theta}(X)) = E\left(i_{\theta}(X)i_{\theta}^T(X)\right) = -E\ddot{l}_{\theta}(X)$$

We can check the equal-sign in the preceding asymptotic distribution.

# Generalized Method of Moments (GMM)

MLE requires parameter model specifications on:

## MLE

- (i) The full density of  $w = (x, y)$ ,  $f(w; \eta)$ ; or
- (ii) The conditional density of  $y$  given  $x$ ,  $f(y|x; \theta)$ ; or
- (iii) The partial density of  $y_t$  given  $x_t$ ,  $f_t(y_t|x_t; \theta)$  in the context of panel data.

GMM was developed by Lars Peter Hansen in 1982 as a generalization of the method of moments. GMM requires less model specifications:

## GMM

- (i) Instead of requiring densities, it asks for only moment restrictions;
- (ii) GMM is largely “semi-parameter”, where the parameter of interest is finite-dimensional [the full data’s distribution function may not be known (infinite-dimensional)], and therefore MLE is not applicable.

## Examples: Parametric regressions

The regression data  $\{w_i = (Y_i, x_i)\}_{i=1}^n$  ( $Y_i \in \mathbb{R}$ : response,  $x_i$ : covariate)

$$y_i = m(x_i, \theta) + \varepsilon_i, E(\varepsilon_i | x_i) = 0, \text{Var}(\varepsilon_i | x) = \sigma^2(x_i) < \infty$$

where  $\{\varepsilon_i\}_{i=1}^n$  are the independent error variable.

By least square method, the score functions is

$$g(w_i, \theta) = \frac{\partial m(x_i, \theta)}{\partial \theta} (y_i - m(x_i, \theta)).$$

and the weighted least square method leads to

$$g(w_i, \theta) = \frac{\partial m(x_i, \theta)}{\partial \theta} \cdot \frac{y_i - m(x_i, \theta)}{\sigma^2(x_i, r_0)}.$$

if  $\sigma^2(x_i, \gamma_0)$  known.

In both methods,  $\gamma_0 = p$ . And  $\hat{\theta}$  is directly solved by estimating equation  $\frac{1}{n} \sum_{i=1}^n g(w_i, \theta) = 0$ .

## Example: Poisson regressions

Consider the equal-dispersion assumption for  $Y_i$  being count data

$$E(Y_i|x_i) = \mu(x_i, \theta), \quad \sigma^2(Y_i|x_i) = \mu(x_i, \theta).$$

The  $\mu(x_i, \theta)$  is a known function, for example:  $\mu(x_i, \theta) = e^{x_i^\top \theta}$ . Define

$$g(w_i, \theta) = \begin{pmatrix} g_1(w_i, \theta) \\ g_2(w_i, \theta) \end{pmatrix} =: \begin{pmatrix} \frac{\partial \mu(x_i, \theta)}{\partial \theta} (y_i - \mu(x_i, \theta)) \\ a(x_i, \theta) \{ [y_i - \mu(x_i, \theta)]^2 - \mu(x_i, \theta) \} \end{pmatrix}$$

We can choose  $a(x_i, \theta)$  almost freely to satisfy, but we need choose one that  $\hat{\theta}$  is most efficient.

Let  $\{w_i\}_{i=1}^n$  be IID r.vs in  $\mathbb{R}^m$ ,  $g(w_i; \theta) \in \mathbb{R}^r$  be  $r$ -dimensional known function of  $w_i$  and  $\theta \in \Theta \subset \mathbb{R}^p$ . So that

$$\exists \theta_0 \in \Theta, \quad E\{g(w_i; \theta_0)\} = 0$$

- If  $r = p$ , we call it “just-identified”;
- If  $r = p$ , the  $\hat{\theta}$  can be made by solving directly:

$$\frac{1}{n} \sum_{i=1}^n g(w_i; \theta) = 0$$

- When  $r > p$ , we call it “over-identified”.

## Definition 7.4

Given  $g(w_i; \theta)$  s.t.  $Eg(w_i; \theta_0) = 0$  for some  $\theta_0 \in \Theta$ . The GMM estimator  $\hat{\theta}_{GMM}$  of  $\theta$  is

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \left( \frac{1}{n} \sum_{i=1}^n g(w_i; \theta) \right)^T \widehat{W}_n \left( \frac{1}{n} \sum_{i=1}^n g(w_i; \theta) \right)$$

for contain  $r \times r$  non-negative definite matrices  $\widehat{W}_n$ , which satisfied that  $\widehat{W}_n \xrightarrow{P} W_0 > 0$ ,  $W_0$  is deterministic and may depend on  $\theta_0$ .

The GMM estimator above is asymptotically equivalent to

$$\hat{\theta}_{n0} = \operatorname{argmin}_{\theta \in \Theta} \left( \frac{1}{n} \sum_{i=1}^n g^T(w_i; \theta) W_0 \frac{1}{n} \sum_{i=1}^n g(w_i; \theta) \right)$$

which is a M-estimator.



- To ensure identification of  $\theta_0$ , we assume  $\theta_0$  is the “unique”

$$\theta \in \Theta, \quad s.t. \quad E g(w_i; \theta) = 0$$

- As  $W_0 > 0$ ,  $\theta_0$  is also the unique  $\theta$  which minimizes

$$E\{g^T(w_i; \theta)\} W_0 E\{g(w_i; \theta)\}$$

- Under certain conditions, we have  $\hat{\theta}_{GMM} \xrightarrow{P} \theta_0$ .

# Asymptotic Normality

Suppose:

- (i)  $g(w, \cdot)$  is a continuous differentiable function on  $\theta \in \text{Int}(\Theta)$ ;
- (ii)  $G_0 = E \left( \frac{\partial g(w, \theta_0)}{\partial \theta} \right)_{r \times p}$  exists and its has full rank  $p$ .

Then under the assumption  $\hat{\theta}_n \xrightarrow{P} \theta_0$ , we have

- AN:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N \left( 0, \left( G_0^T W_0 G_0 \right)^{-1} G_0^T W_0 \Lambda_0 W_0 G_0 \left( G_0^T W_0 G_0 \right)^{-1} \right)$$

where

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n g(w_i; \theta_0) \xrightarrow{d} N(0, \Lambda_0), \quad \Lambda_0 = \text{var} \{g(w_i; \theta_0)\}$$

From the definition of  $\hat{\theta}_n$  (Definition 6.4),

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial g^\top(w_i, \hat{\theta}_n)}{\partial \theta} \widehat{W}_n \left( \frac{1}{n} \sum_{i=1}^n g(w_i, \hat{\theta}_n) \right) = 0 \quad (34)$$

Note that:

$$\frac{\partial g^\top W g}{\partial \theta_i} = \frac{\partial g^\top}{\partial \theta_i} \frac{\partial g^\top W g}{\partial g} = \frac{\partial g^\top}{\partial \theta_i} 2Wg$$

So,

$$\frac{\partial g^\top W g}{\partial \theta} = \frac{\partial g^\top}{\partial \theta} \frac{\partial g^\top W g}{\partial g} = 2 \frac{\partial g^\top}{\partial \theta} Wg$$

By using Taylor formula on (34) around  $\theta_0$  (assuming  $\hat{\theta}_n \xrightarrow{P} \theta_0$ ):

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial g^\top(w_i, \hat{\theta}_n)}{\partial \theta} \widehat{W}_n \left\{ \frac{1}{n} \sum_{i=1}^n g(w_i, \theta_0) + \frac{1}{n} \sum_{i=1}^n \frac{\partial g(w_i, \hat{\theta}_n^*)}{\partial \theta} (\hat{\theta}_n - \theta_0) \right\} = 0 \quad (35)$$

where  $\hat{\theta}_n^*$  is between  $\theta_0$  and  $\hat{\theta}_n$ .

# Proof

Since  $\hat{\theta}_n \xrightarrow{P} \theta_0$ , it is easy to verify that:

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial g^\top(w_i, \hat{\theta}_n)}{\partial \theta}, \frac{1}{n} \sum_{i=1}^n \frac{\partial g^\top(w_i, \hat{\theta}_n^*)}{\partial \theta} \xrightarrow{P} G_0^\top = E \left[ \frac{\partial g(w, \theta_0)}{\partial \theta} \right]$$

Note that  $\widehat{W}_n \xrightarrow{P} W_0$ , from (35),

$$A_0(\hat{\theta}_n - \theta_0) := (G_0^\top W_0 G_0)(\hat{\theta}_n - \theta_0) = -G_0 W_0 \frac{1}{n} \sum_{i=1}^n g(w_i, \theta_0) \{1 + o_p(1)\}$$

and,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = -A_0^{-1} G_0^\top W_0 \frac{1}{\sqrt{n}} \sum_{i=1}^n g(w_i, \theta_0) + o_p(1)$$

As  $\frac{1}{\sqrt{n}} \sum_{i=1}^n g(w_i, \theta_0) \xrightarrow{d} N(0, \Lambda_0)$  where  $\Lambda_0 = \text{var}(g(w_i, \theta_0))$ , we obtain:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, A_0^{-1} G_0^\top W_0 \Lambda_0 W_0 G_0 A_0^{-1})$$

Thus establish the asymptotic normality of  $\hat{\theta}_n$ .

If we use  $W_0$  instead of  $W_n$ , then

- We can attain the same asymptotic normal distribution.
- Denote  $A_0 = (G_0^T W_0 G_0)$  and  $B_0 = G_0^T W_0 \Lambda_0 W_0 G_0$ , Then  $\text{Avar}(\hat{\theta}_{GMM}) = A_0^{-1} B_0 A_0^{-1} / n$ .
- It can be estimated by  $\hat{A}_0^{-1} \hat{B}_0 \hat{A}_0^{-1} / n$ , where

$$\hat{A}_0 = \hat{G}^T \hat{W}_n \hat{G}, \quad \hat{B}_0 = \hat{G}^T \hat{W}_n \hat{\Lambda} \hat{G} \hat{W}_n$$

$$\hat{G} = \frac{1}{n} \sum_{i=1}^n \frac{\partial g(w_i; \hat{\theta})}{\partial \theta}, \quad \hat{\Lambda} = \frac{1}{n} \sum_{i=1}^n g(w_i; \hat{\theta}) g^T(w_i; \hat{\theta})$$

# Optional Weighting Matrix

The GMM estimator depends on the choice of  $\widehat{W}_n$ , the weighting matrix.

Question: which  $\widehat{W}_n$  or  $W_0$  is the best?

- If we choose  $W_0 = \Lambda_0^{-1}$ , then  $B_0 = A_0$  and

$$\text{Avar}(\hat{\theta}_{GMM}) = A_0^{-1} = (G_0^T \Lambda_0^{-1} G_0)^{-1}$$

- In Hansen(1982), it can be shown that for any  $W_0 > 0$ :

$$\begin{aligned}\text{Avar}[\hat{\theta}_n(W_0)] &:= (G_0^T W_0^{-1} G_0)^{-1} (G_0^T W_0 \Lambda_0 W_0 G_0) (G_0^T W_0^{-1} G_0)^{-1} \\ &\geq (G_0^T \Lambda_0^{-1} G_0)^{-1} =: \text{Avar}[\hat{\theta}_n(\Lambda_0^{-1})]\end{aligned}$$

- Hence, the choice of  $W_0^* = \Lambda_0^{-1}$  is the optimal.

# Weight Matrix

- The efficient GMM estimator will use  $\widehat{W}_n = \widehat{\Lambda}^{-1}$  as the weight matrix, where

$$\widehat{\Lambda} = \frac{1}{n} \sum_{i=1}^n g(w_i; \widehat{\theta}) g^T(w_i; \widehat{\theta})$$

- The  $\widehat{\theta}$  is an initial estimator, any consistent estimator of  $\theta$  can be used.
- For instance, it can be derived by minimizing

$$\frac{1}{n} \sum_{i=1}^n g^T(w_i; \theta) \frac{1}{n} \sum_{i=1}^n g(w_i; \theta)$$

- So  $\widehat{\theta}$  is a GMM with  $\widehat{W}_n = I_r$ .

# The GMM estimator

**Step 1:** Construct an initial estimator  $\hat{\theta}$ , which is a GMM with any weight  $\widehat{W}_n > 0$ , for example  $\widehat{W}_n = I_r$ .

**Step 2:** Obtain the optimal weight matrix

$$\widehat{W}_n^* = \left( \frac{1}{n} \sum_{i=1}^n g(w_i; \hat{\theta}) g^T(w_i; \hat{\theta}) \right)^{-1}$$

Then the GMM estimator with  $\widehat{W}_n^*$  as the weight matrix satisfies:

$$\sqrt{n} \left( \hat{\theta}_{GMM}^* - \theta_0 \right) \xrightarrow{d} N \left( 0, \left( G_0^T \Lambda_0^T G_0 \right)^{-1} \right)$$



# Sargan-Hansen Test

It can be shown with  $\hat{\theta} = \hat{\theta}_{GMM}$ , the objective function satisfies (Homework):

$$T_n(\hat{\theta}) = n^{-1/2} \sum_{i=1}^n g^T(w_i; \hat{\theta}) \widehat{W}_n^* n^{-1/2} \sum_{i=1}^n g(w_i; \hat{\theta}) \xrightarrow{d} \chi_{r-p}^2$$

From this asymptotic distribution:

- We also need the condition  $r - p \geq 1$ .
- Hypothesis Testing: reject  $H_0 : Eg(w_i; \theta_0) = 0$  if  $T_n(\theta) > \chi_{r-p, 1-\alpha}^2$ .
- $\hat{\theta}_{GMM}$  also is named minimum  $\chi^2$ -estimator.

# Hypothesis Testing

Test for  $H_0 : c(\theta_0) = 0$ , where  $c(\theta) \in \mathbb{R}^Q$  with  $Q \leq q$ .

Wald Test:

$$c^T(\hat{\theta}) \left( \hat{V}(\hat{\theta}) \right)^{-1} c(\hat{\theta}), \quad \hat{V}(\hat{\theta}) = \text{var} \left( c(\hat{\theta}) \right)$$

LM Test:

$$\tilde{T}_n = \frac{1}{n} \left[ \sum_{i=1}^n g^T(w_i; \tilde{\theta}_n) \widehat{W}_n^* \sum_{i=1}^n g(w_i; \tilde{\theta}_n) - \sum_{i=1}^n g^T(w_i; \hat{\theta}_n) \widehat{W}_n^* \sum_{i=1}^n g(w_i; \hat{\theta}_n) \right]$$

We have  $\tilde{T}_n \xrightarrow{d} \chi_Q^2$  under  $H_0$ , where:

$$\tilde{\theta}_n = \underset{\theta \in \Theta, c(\theta)=0}{\text{argmin}} \sum_{i=1}^n g^T(w_i; \theta) \widehat{W}_n^* \sum_{i=1}^n g(w_i; \theta)$$

# An important feature

If we have  $r$  restrictions, would more moment restrictions lead to more efficiency?  
To appreciate this question, let

$$g(w, \theta) = \underbrace{(g_{(r-1)}(w, \theta))}_r, \underbrace{g_r(w, \theta)}_1)^T, E g(w, \theta) = 0.$$

The asymptotic variance of  $\hat{\theta}$  based on  $g(\cdot, \cdot) \in \mathbb{R}^r$  is:

$$V_r^{-1} := \left[ E \left( \frac{\partial g^T}{\partial \theta} \right) E^{-1}(g g^T) E \left( \frac{\partial g}{\partial \theta} \right) \right]^{-1}$$

And asymptotic variance of  $\hat{\theta}_{(r-1)}$  based on  $g_{(r-1)}$  is:

$$V_{r-1}^{-1} := \left[ E \left( \frac{\partial g_{(r-1)}^T}{\partial \theta} \right) E^{-1}(g_{(r-1)} g_{(r-1)}^T) E \left( \frac{\partial g_{(r-1)}}{\partial \theta} \right) \right]^{-1}$$

An important feature in GMM:

- $V_r^{-1} \leq V_{r-1}^{-1}$ .
- Hence, GMM with  $r$  restrictions is at least as efficient as GMM with  $r - 1$  restrictions.

# Proof

Note that at  $\theta = \theta_0$

$$E(gg^T) = \begin{pmatrix} E(g_{(r-1)}g_{(r-1)}^T) & E(g_{(r-1)}g_r^T) \\ E(g_{(r-1)}^Tg_r) & E(g_r^2) \end{pmatrix} := \begin{pmatrix} A & B \\ B^T & C \end{pmatrix}$$

Let  $D = C - B^T A B$ . The inverse matrix of block matrices formula gives:

$$\begin{pmatrix} A & B \\ B^T & C \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} A^{-1}B \\ -I \end{pmatrix} D^{-1} \begin{pmatrix} B^T A, & -I \end{pmatrix}$$

Write  $E(gg^T) := \begin{pmatrix} A & B \\ B^T & C \end{pmatrix}$  and  $A := E(g_{(r-1)}g_{(r-1)}^T)$ . We have

$$\left(E(gg^T)\right)^{-1} = \begin{bmatrix} \left(E(g_{(r-1)}g_{(r-1)}^T)\right)^{-1} & 0 \\ 0 & 0 \end{bmatrix} + H$$

where  $H \geq 0$ .

# Proof

Therefore,

$$\begin{aligned} V_r &= \begin{pmatrix} E\left(\frac{\partial \mathbf{g}_{(r-1)}}{\partial \theta}\right) \\ E\left(\frac{\partial \mathbf{g}_r}{\partial \theta}\right) \end{pmatrix}^\top \left\{ \begin{bmatrix} \left(E\left(\mathbf{g}_{(r-1)}\mathbf{g}_{(r-1)}^\top\right)\right)^{-1} & 0 \\ 0 & 0 \end{bmatrix} + H \right\} \begin{pmatrix} E\left(\frac{\partial \mathbf{g}_{(r-1)}}{\partial \theta}\right) \\ E\left(\frac{\partial \mathbf{g}_r}{\partial \theta}\right) \end{pmatrix} \\ &= E^\top\left(\frac{\partial \mathbf{g}_{(r-1)}}{\partial \theta}\right) \left(E\left(\mathbf{g}_{(r-1)}\mathbf{g}_{(r-1)}^\top\right)\right)^{-1} E\left(\frac{\partial \mathbf{g}_{(r-1)}}{\partial \theta}\right) + \tilde{H} \\ &= V_{r-1} + \tilde{H} \geq V_{r-1}, \quad V_r^{-1} \leq V_{r-1}^{-1} \end{aligned}$$

## Remark 21

From the above inequality, we know that if

$$E\left(\frac{\partial \mathbf{g}^\top}{\partial \theta}\right) H E\left(\frac{\partial \mathbf{g}}{\partial \theta}\right) \neq 0$$

then there will be reduction of the asymptotic variance in using the moment restriction in some directions or combination of the parameter space. This is exactly the attraction of GMM.