

# SDS 387 Linear Models

Fall 2024

Lecture 9 - Tue, Sep 24, 2024

Instructor: Prof. Ale Rinaldo

CLT: Lindeberg-Feller (LF) r.v.'s on the same row, are indep.

Let  $\{X_{j;n}\}_{1 \leq j \leq n}$ ,  $n=1, 2, \dots$  be a triangular array (I wrote  $X_{n,j}$ )

$\hookrightarrow$  different notation than last time

$n$ : index for  $n^{\text{th}}$  row of the triangular array

s.t.  $\mathbb{E}[X_{j;n}] = 0$   $\forall j, n$  and let

$$B_n^2 = \sum_{j=1}^n \sigma_{n,j}^2 \quad \sigma_{n,j}^2 = \text{Var}[X_{n,j}]$$

Then 
$$\frac{\sum_{j=1}^n X_{j;n}}{B_n} \xrightarrow{d} N(0,1) \quad \text{if}$$

$$(LF) \quad \frac{1}{B_n^2} \sum_{j=1}^n \mathbb{E} \left[ X_{j;n}^2 \mathbb{1}_{\{|X_{j;n}| > \varepsilon B_n\}} \right] \rightarrow 0$$

as  $n \rightarrow \infty$ ,  $\forall \varepsilon > 0$ . (1)

• Conversely, if  $\frac{\sum_{j=1}^n X_{j,n}}{B_n} \rightarrow N(0,1)$  and

$\max_{j=1, \dots, n} \frac{\sigma_{j,n}^2}{B_n^2} \rightarrow 0$  as  $n \rightarrow \infty$ , then (LF) holds.

• The proof is based on use of ch.f. See, e.g., Petrov.

• A stronger condition than (LF) is Lyapunov:

$\exists \delta > 0$  s.t.

$$\frac{1}{B_n^{2+\delta}} \sum_{j=1}^n \mathbb{E} [ |X_{j,n}|^{2+\delta} ] \rightarrow 0 \text{ as } n \rightarrow \infty$$

Example  $Y_j = \text{Bernoulli}(p_j)$  independent

Let  $X_j = Y_j - p_j$  ( $\mathbb{E}[X_j] = 0$ ).

Under what conditions on the sequence  $\{p_j\}$  do we have a CLT?

Use Lyapunov with  $\delta = 1$  (we control 3<sup>rd</sup> moment)

Let  $\sigma_j^2 = \text{Var}[X_j] = p_j(1-p_j)$ . Then

$$\mathbb{E}[|X_j|^3] \leq \sigma_j^2. \quad \square$$

$$\frac{\sum_{j=1}^n \mathbb{E} [ |X_{j,1}|^3 ]^{\overrightarrow{2+\delta}}}{B_n^3} \leq \frac{B_n^2}{B_n^3} = \frac{1}{B_n}$$

$$B_n^2 = \sum_{j=1}^n \sigma_j^2 = \sum_{j=1}^n p_j (1-p_j)$$

This quantity will vanish as  $n \rightarrow \infty$  if

$$\sum_{j=1}^n p_j (1-p_j) \rightarrow \infty$$

• The multivariate case:  $d > 1$  but fixed (not increasing with  $n$ !)

Consider a triangular array of centered  $d$ -dimensional random vectors  $\{X_{j,n}\}_{j \leq n}$  s.t.

$\text{Var}[X_{j,n}]$  exists. Let

$$Y_{j,n} = \left( \sum_{i=1}^n \text{Var}[X_{i,n}] \right)^{-1/2} X_{j,n}$$

(LF) If  $\lim_{n \rightarrow \infty} \sum_{j=1}^n \mathbb{E} \left[ \|Y_{j,n}\|^2 \mathbb{1}_{\{\|Y_{j,n}\| > \varepsilon\}} \right] = 0$   
 $\forall \varepsilon > 0$

Then  $\sum_{j=1}^n Y_{j,n} \xrightarrow{d} N_d(0, I)$

PA/

Use Cramer-Wold device. Let's first consider points  $t \in \mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\| = 1\}$  (3)

So we need to show that,  $\forall t \in \mathbb{S}^{d-1}$

$$t^T \sum_{j=1}^n Y_{j,n} \xrightarrow{d} N(0, 1)$$

(if  $\|t\| \neq 1$ , the limiting distribution will be  $N(0, \|t\|^2)$ )

First, notice that for  $t \in \mathbb{S}^{d-1}$ ,

$$\sum_{j=1}^n \text{Var}[t^T Y_{j,n}] = 1$$

So, using the CLT for univariate CLT we only need to look at the sequence

$$\sum_{j=1}^n \mathbb{E} \left[ \underbrace{(t^T Y_{j,n})^2}_{\leq \underbrace{\|t\|^2}_{=1} \|Y_{j,n}\|^2} \mathbb{1}_{\{|t^T Y_{j,n}| > \varepsilon\}} \right]$$

$\downarrow$   
arbitrary  $> 0$

and show that it is vanishing as  $n \rightarrow \infty$ .

But this easily follows from bounding the last expression by

$$\sum_{j=1}^n \mathbb{E} \left[ \|Y_{j,n}\|^2 \mathbb{1}_{\{\|Y_{j,n}\| > \varepsilon\}} \right] \rightarrow 0$$

by assumption.

check the calculations for the case  $\|t\| \neq 1$ .

## BERRY-ESSEEN BOUNDS

→ See Petrov's book

This result indicates how fast the CLT approximation works.  
finite sample

Let  $X_1, X_2, \dots$  be independent r.v.'s s.t.  $\mathbb{E}[X_i] = 0$   
 $\text{Var}[X_i] < \sigma_i^2 < \infty$  and  $\mathbb{E}[|X_i|^3] < \infty$ , all  $i$ .

↳ stronger condition than (LF)

Then

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P}\left(\frac{\sum_{i=1}^n X_i}{B_n} \leq x\right) - \Phi(x) \right| \leq C \frac{\sum_{i=1}^n \mathbb{E}[|X_i|^3]}{B_n^{3/2}}$$

$\downarrow$  cdf of  $N(0,1)$        $\downarrow$  universal constant  $< \frac{1}{2}$

$$\sqrt{\sum_{i=1}^n \sigma_i^2}$$

- When  $\sigma_i^2 = \sigma^2$  all  $i$  and  $\mathbb{E}[|X_i|^3] \leq \mu_3$  all  $i$   
then the RHS of this bound is upper bounded by

$$C \frac{n \mu_3}{\sqrt{n} \sigma^{3/2}} = \frac{C}{\sqrt{n}} \frac{\mu_3}{\sigma^3}$$

- Example:  $X_1, X_2, \dots$  independent with  $X_i \sim \text{Bernoulli}(p_i)$   
where  $p_i \in [\varepsilon, 1-\varepsilon]$  all  $i$ .  
Then, we saw that,  
 $\mathbb{E}[|X_i - p_i|^3] \leq p_i(1-p_i)$  so the  
Berry-Esseen bound is

$$\frac{\sum_{i=1}^n p_i (1-p_i)}{\left( \sum_{i=1}^n p_i (1-p_i) \right)^{3/2}} \leq \frac{1}{\sqrt{n} \varepsilon (1-\varepsilon)}$$

We can allow  $\varepsilon = \varepsilon_n$  to also depend on  $n$ :  $\varepsilon_n \rightarrow 0$

CLT works as long as  $n \varepsilon_n (1-\varepsilon_n) \rightarrow \infty$

So  $\varepsilon_n$  can  $\downarrow 0$  but slower than  $1/n$ .

- Weakest version of Berry-Esseen is (see Petrov).

Assume only that  $\mathbb{E}[|X_i|^{2+\delta}] < \infty$  all  $i$ ,  $\delta \in (0, 1]$

Then

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P}\left( \frac{\sum_{i=1}^n X_i}{B_n} \leq x \right) - \Phi(x) \right| \leq C \frac{\sum_{i=1}^n \mathbb{E}[|X_i|^{2+\delta}]}{B_n^{1+\delta/2}}$$

- General way of thinking about CLT: replace each  $X_i$  by a  $Z_i \sim N(\mathbb{E}[X_i], \text{Var}[X_i])$   
Then "well-behaved" functions of the  $X_i$ 's have a distribution that is close to that of the same function of the  $Z_i$ 's.
- In general we can frame this task (i.e. the task of demonstrating a Gaussian approximation) as follows.

Let  $F$  be a class of functions. Then we want to bound:

$$\sup_{f \in F} \left| \mathbb{E} [f(x_1, \dots, x_n)] - \mathbb{E} [f(z_1, \dots, z_n)] \right|$$

where  $z_i \sim N(\mathbb{E}[x_i], \text{Var}[x_i])$

For example, take  $F = \{1_{(-\infty, z]}, z \in \mathbb{R}\}$  and bound

$$\sup_{f \in F} \left| \mathbb{E} \left[ f \left( \frac{\sum x_i}{\sqrt{n}} \right) \right] - \mathbb{E} [f(z)] \right|$$

Berry - Esseen bound

- We can take  $F$  to be any class of functions we like.