

Lecture 15: The central limit Thm

The central limit theorem is most likely one of the most fundamental results in probability and statistics.

It is traditionally taught with Fourier methods (i.e. characteristic functions) but, in a way, this method doesn't really extend well for things like martingales, etc.

The method of proof we will give is based on Lindeberg's proof & has proven to be much more powerful and conceptual.

Taylor with remainder for C_c^∞ functions

For $x \in \mathbb{R}$ and $f \in C_c^\infty(\mathbb{R})$ Taylor's Theorem gives

$$f(x+\Delta x) = f(x) + f'(x)\Delta x + \frac{f''(x)}{2}\Delta x^2 + R_f(x, \Delta x)$$

$$\text{where } R_f(x, \Delta x) = \frac{f'''(x_*)}{3!} \Delta x^3$$

$$\text{since } |R_f(x, \Delta x)| = \underbrace{\left| \frac{f'''(x_*)}{3!} \Delta x \right|}_{\Delta x^2}$$

this term $\rightarrow 0$ as
 $\Delta x \rightarrow 0$ uniformly in x
 since $\|f'''\|_\infty < \infty$

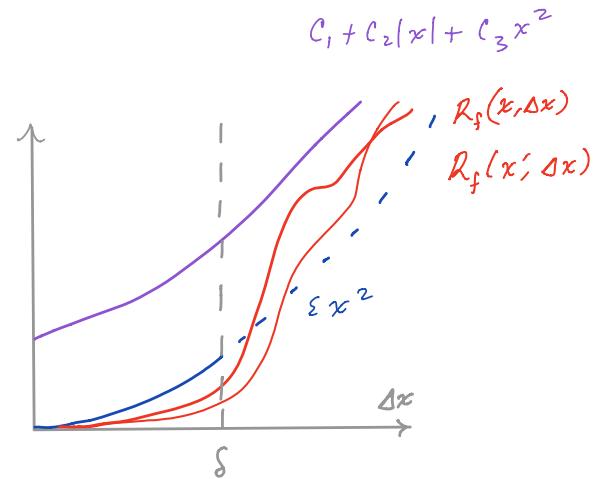
$$\therefore \forall \varepsilon > 0 \exists \delta > 0 \text{ s.t.}$$

$$|\Delta x| \leq \delta \Rightarrow |R_f(x, \Delta x)| \leq \varepsilon (\Delta x)^2$$

(2)
 But we also have $\max(\|f\|_\infty, \|f'\|_\infty, \|f''\|_\infty) < \infty$

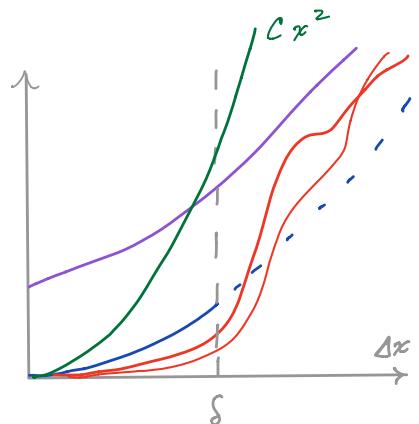
$$\therefore |\Delta x| \geq \delta \Rightarrow |R_f(x, \Delta x)| \leq C_1 + C_2 |\Delta x| + C_3 \Delta x^2$$

Here is the picture



Notice, however, we can find a C large enough so that

$$x \geq \delta \Rightarrow C_1 + C_2 x + C_3 x^2 \leq C x^2$$



$$\therefore \forall \varepsilon > 0 \exists \delta, C > 0 \text{ s.t.}$$

$$f(x+\Delta x) = f(x) + f'(x)\Delta x + \frac{f''(x)}{2}\Delta x^2 + R_f(x, \Delta x)$$

$$\sup_{x \in \mathbb{R}} |R_f(x, \Delta x)| \leq \varepsilon (\Delta x)^2 + C (\Delta x)^2 \int_{\{|x| \geq \delta\}} 1$$

Lindeberg's Method for CLT

(3)

As a warm up suppose X_1, X_2, \dots and Z_1, Z_2, \dots are all independent r.v.s

$$\text{s.t. } E X_i = E Z_i \text{ &}$$

$$E X_i^2 = E Z_i^2 < \infty$$

let $g: \mathbb{R}^n \rightarrow \mathbb{R}$ be a smooth function and notice

$$|E g(X_1, \dots, X_n) - E g(Z_1, \dots, Z_n)|$$

$$\pm E g(Z_1, X_2, \dots, X_n)$$

$$\pm E g(Z_1, Z_2, X_3, \dots, X_n)$$

\vdots

$$\pm E g(Z_1, Z_2, \dots, Z_{n-1}, X_n)$$

$$\leq \sum_{i=1}^n |E g(\dots, Z_{i-1}, X_i, X_{i+1}, \dots) - E g(\dots, Z_{i-1}, Z_i, X_{i+1}, \dots)|$$

$$= \sum_{i=1}^n |E g_i(X_i) - E g_i(Z_i)|$$

where g_i depends on \dots, Z_{i-1} & X_{i+1}, \dots

Now let $f \in C_c^\infty(\mathbb{R})$ and look what happens when we set

$$g(x_1, \dots, x_n) = f(x_1 + \dots + x_n)$$

$$\therefore g_i(x) = f(\dots + Z_{i-1} + x + X_{i+1} + \dots)$$

(4)

$$\therefore g_i(X_i) - g_i(0)$$

$$= f(\dots + Z_{i-1} + X_i + X_{i+1} + \dots)$$

$$- f(\dots + Z_{i-1} + 0 + X_{i+1} + \dots)$$

$$= f'(Z_{i-1}) X_i + \underbrace{f''(Z_{i-1}) X_i^2}_{2} + R_f(Z_{i-1}, X_i)$$

$$\text{where } Z_{i-1} := (\dots + Z_{i-1} + X_{i+1} + \dots)$$

We also have

$$g_i(Z_i) - g_i(0)$$

$$= f'(Z_{i-1}) Z_i + \underbrace{f''(Z_{i-1}) Z_i^2}_{2} + R_f(Z_{i-1}, Z_i)$$

Therefore

$$E g_i(X_i) - E g_i(Z_i) \leftarrow \pm E g_i(0)$$

$$= E [f'(Z_{i-1})(X_i - Z_i)]$$

$$+ E \left[\underbrace{f''(Z_{i-1})}_{2} (X_i^2 - Z_i^2) \right]$$

$$+ E R_f(Z_{i-1}, X_i) - E R_f(Z_{i-1}, Z_i)$$

But (X_i, Z_i) is indep of Z_{i-1} so

$$\underbrace{E [f'(Z_{i-1})(X_i - Z_i)]}_{\circ} = 0$$

$$\underbrace{E \left[\underbrace{f''(Z_{i-1})}_{2} (X_i^2 - Z_i^2) \right]}_{\circ} = 0$$

$$\text{since } E(X_i - Z_i) = E(X_i^2 - Z_i^2) = 0$$

$\therefore \forall \varepsilon > 0 \exists \delta, c > 0$ s.t.

(5)

$$\begin{aligned} & |Eg_i(X_i) - Eg_i(z_i)| \\ & \leq E|R_f(z_{X_{(i)}, X_i})| + E|R_f(z_{X_{(i)}, Z_i})| \\ & \leq \varepsilon EX_i^2 + cE(X_i^2 I_{|X_i| \geq \delta}) \\ & \quad + \varepsilon EZ_i^2 + cE(Z_i^2 I_{|Z_i| \geq \delta}) \end{aligned}$$

Let's put these results together into a lemma.

Lemma 1:

If $X_1, X_2, \dots, Z_1, Z_2, \dots$ are indep r.v.s s.t. $EX_i = EZ_i$ & $EX_i^2 = EZ_i^2$ and $f \in C_c^\infty(\mathbb{R})$ then $\forall \varepsilon > 0$ $\exists \delta, c > 0$ s.t.

$$\begin{aligned} & |Ef(X_1 + \dots + X_n) - Ef(Z_1 + \dots + Z_n)| \\ & \leq 23 \sum_{i=1}^n EX_i^2 + c \sum_{i=1}^n E X_i^2 I_{|X_i| \geq \delta} \\ & \quad + c \sum_{i=1}^n EZ_i^2 I_{|Z_i| \geq \delta} \end{aligned}$$

Now the key "reason" the limit in the CLT is Gaussian is the following lemma.

Lemma 2:

If $Z, Z_1, \dots, Z_n \stackrel{iid}{\sim} N(0, \sigma^2)$ then

$$Z \stackrel{D}{=} \frac{Z_1}{\sqrt{n}} + \dots + \frac{Z_n}{\sqrt{n}}.$$

Theorem: (The central limit theorem)

(6)

If X_1, X_2, \dots are iid r.v.s with $EX_i = 0$ and $EX_i^2 = \sigma^2 < \infty$, then

$$\frac{S_n}{\sqrt{n}} \xrightarrow{D} Z \sim N(0, \sigma^2)$$

where $S_n := X_1 + \dots + X_n$.

Proof:

Since $\text{var}\left(\frac{S_n}{\sqrt{n}}\right) = \sigma^2$ Chebyshev's Thm implies

$$\left|\frac{S_n}{\sqrt{n}}\right| = O_p(1) \Rightarrow \left\{\frac{S_n}{\sqrt{n}}\right\}_{n \geq 1} \text{ is tight}$$

Now \mathbb{R} is locally compact and therefore Portmanteau II implies it will be sufficient to show that $\forall f \in C_c^\infty(\mathbb{R})$

$$\left|Ef\left(\frac{S_n}{\sqrt{n}}\right) - Ef(z)\right| \rightarrow 0$$

II Lemma 2

$$\left|Ef\left(\frac{X_1}{\sqrt{n}} + \dots + \frac{X_n}{\sqrt{n}}\right) - Ef\left(\frac{Z_1}{\sqrt{n}} + \dots + \frac{Z_n}{\sqrt{n}}\right)\right| \xrightarrow{iid N(0, \sigma^2)} (x)$$

Now let $c > 0$.

By Lemma 1, $\exists \delta, c > 0$ s.t.

$$\begin{aligned} (x) & \leq 23 \sum_{i=1}^n EX_i^2 + c \sum_{i=1}^n E \frac{X_i^2}{n} I_{|X_i| \geq \sqrt{n}\delta} \\ & \quad + c \sum_{i=1}^n E \frac{Z_i^2}{n} I_{|Z_i| \geq \sqrt{n}\delta} \\ & = 2c\sigma^2 + cEX_1^2 I_{|X_1| \geq \sqrt{n}\delta} \\ & \quad + cEZ_1^2 I_{|Z_1| \geq \sqrt{n}\delta} \end{aligned}$$

since $X_i \sim X_1$, $Z_i \sim Z_1$ & $EX_1^2 = EZ_1^2 = \sigma^2$

Now by the DCT we get

(7)

$$E X_1^2 I_{|X_1| \geq \sqrt{n}\delta} \xrightarrow{n \rightarrow \infty} 0$$

$$E Z_1^2 I_{|Z_1| \geq \sqrt{n}\delta} \xrightarrow{n \rightarrow \infty} 0.$$

Therefore

$$\limsup_{n \rightarrow \infty} \left| E f\left(\frac{S_n}{\sqrt{n}}\right) - E f(z) \right| \leq 2\epsilon\sigma^2$$

Since $\epsilon > 0$ was arbitrary

$$\left| E f\left(\frac{S_n}{\sqrt{n}}\right) - E f(z) \right| \xrightarrow{n \rightarrow \infty} 0$$

as was to be shown.

QED

Corollary:

If Z is a r.v. that satisfies

$$(i) \text{ var}(z) := \sigma^2 < \infty$$

$$(ii) Z = \frac{z_1}{\sqrt{n}} + \dots + \frac{z_n}{\sqrt{n}}$$

where z_1, \dots, z_n are indep. copies of z

then Z must be Gaussian:

$$Z \sim N(0, \sigma^2)$$

Proof:

Properties (i) & (ii) were the only properties of Z we use in the proof of the CLT.

By uniqueness of distributional limits it must be the only distribution with this property. QED

Theorem: (CLT for random vectors)

(8)

If X_1, X_2, \dots are iid d-dimensional r.v.s s.t.

$$E(X_i) = 0$$

$$E(X_i X_i^T) = \Sigma \in \mathbb{R}^{d \times d}$$

Then

$$\frac{S_n}{\sqrt{n}} \xrightarrow{\mathcal{D}} Z \sim N_d(0, \Sigma)$$

where $S_n := X_1 + X_2 + \dots + X_n$.

Proof:

One can use the same method of proof as in the CLT... but it requires d-dimensional Taylor approximations. Instead let's use the Cramér-Wold device & show

$$\text{If } k \in \mathbb{R}^d \quad \langle k, \frac{S_n}{\sqrt{n}} \rangle \xrightarrow{\mathcal{D}} \langle k, Z \rangle.$$

Now since

$$\langle k, \frac{S_n}{\sqrt{n}} \rangle = \underbrace{\langle k, X_1 \rangle}_{\text{iid}} + \dots + \underbrace{\langle k, X_n \rangle}_{\text{iid}}$$

where $E \langle k, X_i \rangle = \langle k, E X_i \rangle = 0$

$$E \langle k, X_i \rangle^2 = k^T \Sigma k < \infty$$

$$\langle k, Z \rangle \sim N(0, k^T \Sigma k).$$

the univariate CLT applies.

$$\therefore \langle k, \frac{S_n}{\sqrt{n}} \rangle \xrightarrow{\mathcal{D}} \langle k, Z \rangle \text{ as } n \rightarrow \infty$$

as was to be shown.

QED

Theorem: (Lindeberg CLT generalization)

(9)

Let $X_{i,n}$ be a triangular array
of independent r.v.s.

$$X_{11}$$

$$X_{12}, X_{22}$$

$$X_{13}, X_{23}, X_{33}$$

$$\vdots \quad \ddots$$

which satisfy

(a) $E X_{i,n} = 0, \forall n \forall i \leq n$

(b) $\sum_{i=1}^n E X_{i,n}^2 = \sigma^2, \forall n$

(c) $\sum_{i=1}^n E X_{i,n}^2 I_{|X_{i,n}| \geq s} \xrightarrow{n \rightarrow \infty} 0, \forall s > 0$

Then

$$X_{1n} + \dots + X_{nn} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} N(0, \sigma^2).$$

Proof:

This proof is nearly identical to the
proof we gave for the univariate CLT.
I'll leave the details for a homework.

QED

Kolmogorov's 3 series theorem

(10)

Stem's Method

