- Last time CLT in $\mathbb{R}^d$:

Consider a triangular array of random vectors in $\mathbb{R}^d$

$\{ X_{n,j}, \; j=1,\ldots,n \}_{n=1,2,\ldots}$    s.t.   $\mathbb{E}[X_{n,j}] = 0 \in \mathbb{R}^d$

and that   $\text{Var}[X_{n,j}]$ exists     $\forall n$ and $j$

       $d \times d$ matrix      and is invertible

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

If

$(*)$    $\lim_{n \to \infty} \sum_{i=1}^{n} \mathbb{E}\left[ \| Y_{n,j} \|^2 \, \mathbb{1}\{\| Y_{n,j} \| \geq \varepsilon \} \right] = 0$

                                      $\forall \varepsilon > 0$

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

                                $\hookrightarrow$ identity matrix

                as $n \to$

PF/ Use the Cramer Wold device. We need to
show that $\forall t \in \mathbb{R}^d$

$$t^T \sum_{j=1}^{n} Y_{n,j} \xrightarrow{d} t^T Z$$

$$\hookrightarrow \ Z \sim N(0, I_d)$$

It is easy to see that

<span style="color:red">Exercise or HW</span> $\longleftarrow$ $\quad t^T \sum_{j=1}^{n} Y_{n,j} \sim 0, \|t\|^2$

$\underbrace{\qquad}$ mean is 0 and variance is $\|t\|^2$

To show that $\dfrac{\left( t^T \sum_{j=1}^{n} Y_{n,j} \right)}{\|t\|} \xrightarrow{d} N(0,1)$ we will

check that the LF condition holds. So $\forall \varepsilon > 0$

$$\frac{1}{\|t\|^2} \sum_{j=1}^{n} \mathbb{E}\left[ (t^T Y_{n,j})^2 \, \mathbb{1}\left\{ |t^T Y_{n,j}| > \varepsilon \|t\| \right\} \right] \qquad \begin{array}{l} \text{\color{blue}By Cauchy} \\ \text{\color{blue}Schwartz} \\ \text{\color{blue}$|t^T Y_{n,j}| \leq \|t\|$} \\ \text{\color{blue}$\|Y_{n,j}\|$} \end{array}$$

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

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

as $n \to \infty$
by assumption. (*)

∎

■ BERRY - ESSEEN BOUNDS $\qquad$ <span style="color:red">See Petrov chapter 5</span>

Let $X_1, X_2, \ldots, X_n$ be independent <u>univariate</u> r.v.'s
s.t. $\mathbb{E}[X_n] = 0$, $\sigma_n^2 = \text{Var}[X_n]$ and $\mathbb{E}|X_n|^3 = \mu_{n,3} < \infty$

Then

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P}\left( \frac{\sum_{n=1}^{n} X_n}{B_n} \leq x \right) - \Phi(x) \right| \leq C \frac{\sum_{n=1}^{n} \mu_{n,3}}{B_n^3}$$

②

where $\quad B_n^2 = \sum_{i=1}^{n} \sigma_i^2 \quad$ and $\quad \Phi \quad$ is the cdf of
$$N(0,1)$$

$C$ is an universal constant $< \frac{1}{2}$

- Assume $\sigma_i^2 = \sigma^2$ and $\mathbb{E}\left[|X_n|^3\right] = \mu_3$ all $i$. Then

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P}\left( \frac{\sqrt{n} \, \bar{X}_n}{\sigma} \leq x \right) - \Phi(x) \right| \leq c \, \frac{\sqrt{n} \, \mu_3}{\sigma^3 \, n^{3/2}}$$

$$\mathbb{P}\left( \sqrt{n} \, \frac{\bar{Z}_n}{\sigma} \leq x \right) = c \, \frac{\mu_3}{\sigma^3} \, \frac{1}{\sqrt{n}}$$

$Z_1, \ldots, Z_n \overset{iid}{\sim} N(0, \sigma^2)$

- This requires a $3^{rd}$ moment!

<u>Example</u> $\qquad X_1, X_2, \ldots, X_n$ independent with $X_i \sim$ Bernoulli $(p_i)$

$p_i \in (0,1)$

Then the Berry-Esseen bound is as $\qquad$ all $i$

follows:

$$\mathbb{E}\left[|X_i - p_i|^3\right] = p_i (1-p_i) \underbrace{\left[(1-p_i)^2 + p_i^2\right]}_{\leq 1}$$

$$\leq p_i (1-p_i)$$

So, assuming that $\quad p_i \in [\varepsilon, 1-\varepsilon] \quad 0 < \varepsilon < \frac{1}{2}$

all $i$

The RHS of Berry-Esseen bound is

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

$$\boxed{3}$$

Next, for all $i$, $\dfrac{1}{p_i(1-p_i)} \leq \dfrac{1}{\varepsilon(1-\varepsilon)}$ . To see this, e.g.

look at the graph of the function $x \in [\varepsilon, 1-\varepsilon] \longmapsto x(1-x)$

(formally, use
concavity
of the function)

Then,

$$\dfrac{1}{\sqrt{\displaystyle\sum_{i=1}^{n} p_i(1-p_i)}} \leq \dfrac{1}{\sqrt{n \min_{i} p_i(1-p_i)}} \leq \dfrac{1}{\sqrt{n \, \varepsilon(1-\varepsilon)}}$$

If we let $\varepsilon = \varepsilon_n \to 0$ as $n \to \infty$ the we have a CLT as long as $\frac{1}{\sqrt{n}} = o\left(\sqrt{\varepsilon_n(1-\varepsilon_n)}\right)$

Equivalently $\varepsilon_n$ can go to zero but slower than $\frac{1}{n}$

For example, if $\varepsilon_n = n^{-\alpha}$ for $\alpha \in (0,1)$ the Berry-Esseen bound is of order $n^{(\alpha-1)/2}$.

# HIGH-DIM BERRY ESSEEN BOUNDS

Let $X_1, \ldots, X_n$ are independent centered r.v.'s in $\mathbb{R}^d$ s.t. $\text{Cov}[X_n] = \Sigma_i$. Let $z_1, \ldots, z_n$ be independent centered Gaussians s.t. $\text{Var}[z_n] = \underset{\shortparallel}{\Sigma_i}$

$$\text{Var}[x_i]$$

Let $\mathcal{A}$ be a collection of subsets of $\mathbb{R}^d$. Examples:

· set of all convex set

· set of all balls or ellipsoids

· set of all hyper-rectangles

We want to establish the bound.

$$\underset{A \in \mathcal{A}}{\sup} \left| \mathbb{P}\left( \frac{\sum_i X_i}{\sqrt{n}} \in A \right) - \mathbb{P}\left( \frac{\sum_i z_i}{\sqrt{n}} \in A \right) \right|$$

$$\leq C(d, \mathcal{A}) \frac{1}{\sqrt{n}} \qquad \text{"Third moment term"}$$