

## 6. The central limit theorem

We will be interested in the distribution of sums  $S_n = X_1 + X_2 + \dots + X_n$ .

For some types of distributions we know the answer but in general this question is difficult to answer.

Moreover, in statistics we need to approximate such distributions even if we do not exactly know the distributions of  $X_1, X_2, \dots$ .

The setup we will look at will

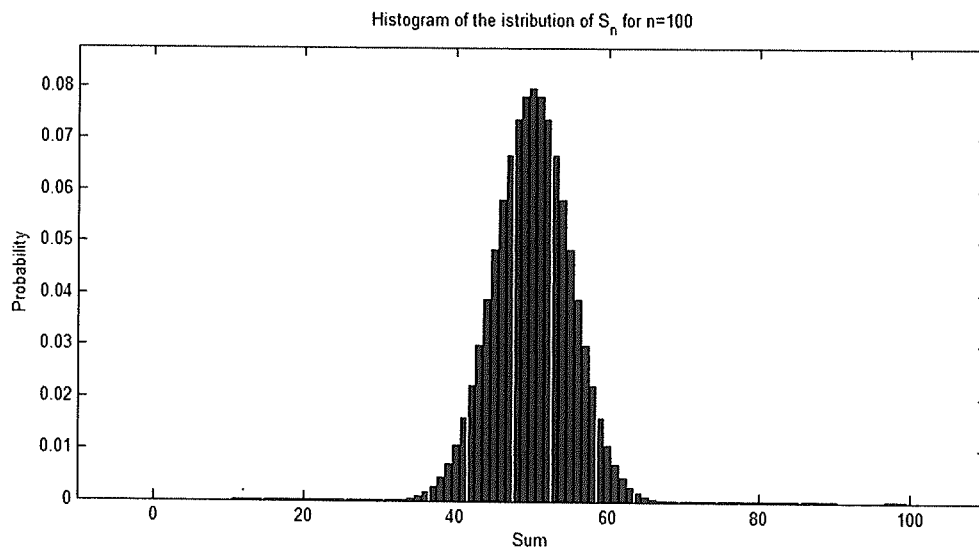
be:  $X_1, X_2, \dots$  are independent, equally distributed random variables.

We denote  $S_n = X_1 + X_2 + \dots + X_n$ .

Let us look at examples of distributions of  $S_n$  for simple distributions.

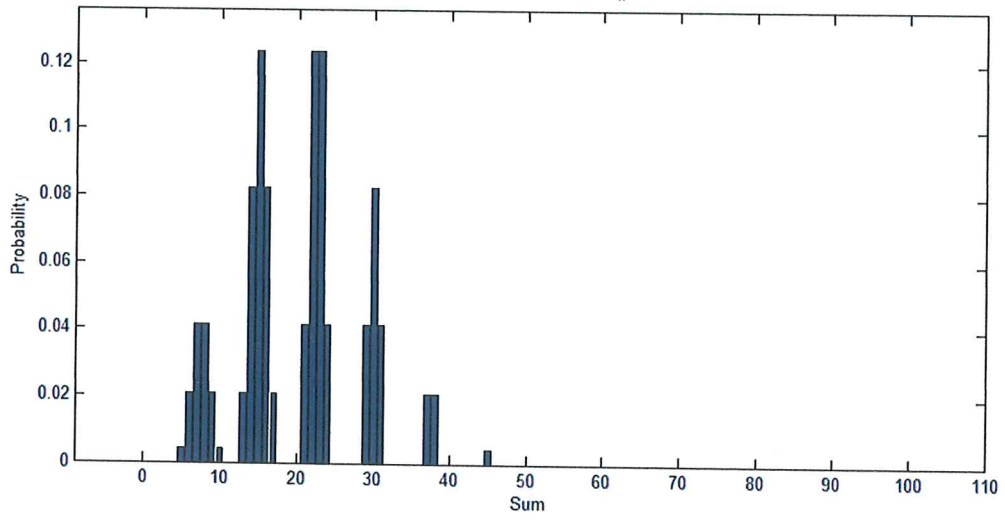
We will look at a few examples of distributions of  $S_n$  for different distributions of  $X_1$  and different  $n$ .

1. Let  $P(X_1 = 0) = P(X_1 = 1) = \frac{1}{2}$ . Take  $n = 100$ . Let  $S_n = X_1 + \dots + X_n$ . The histogram of the distribution of  $S_n$  is:

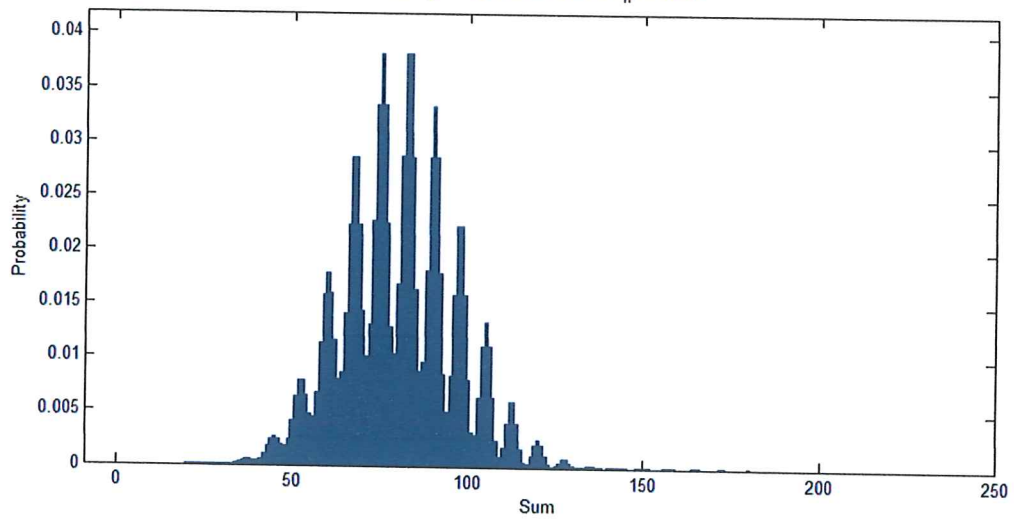


2. Let  $P(X_1 = 1) = P(X_1 = 2) = P(X_1 = 9) = \frac{1}{3}$ . Let  $n = 5, 20, 50, 200$ .

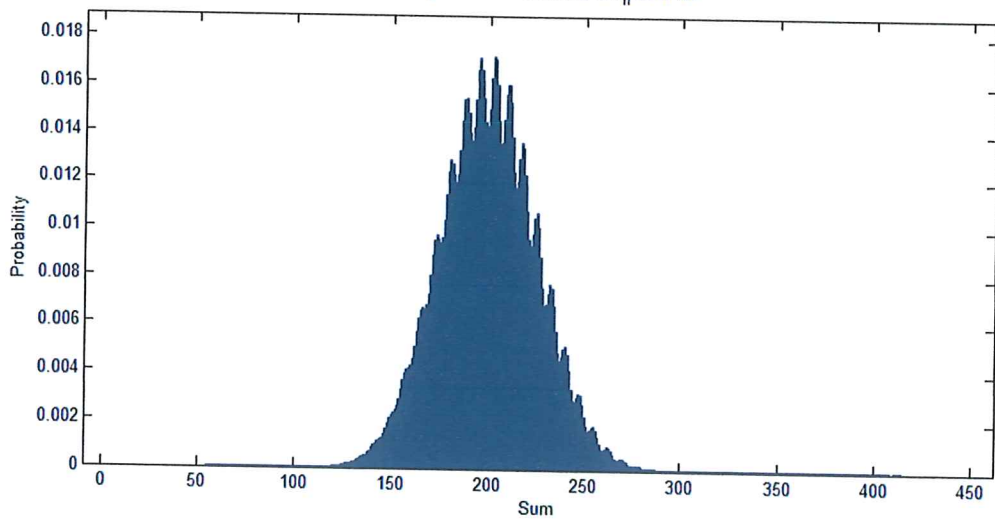
Histogram of the istribution of  $S_n$  for  $n=5$

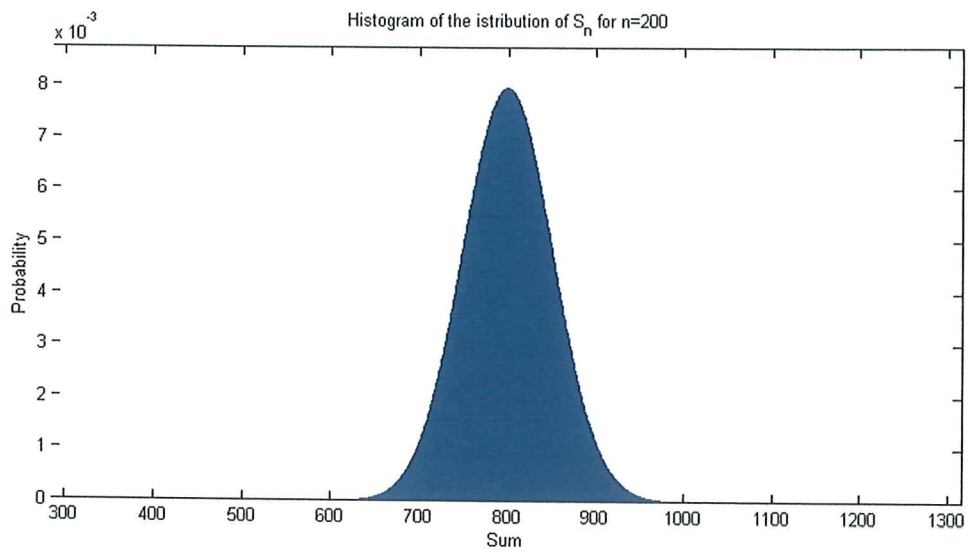


Histogram of the istribution of  $S_n$  for  $n=20$

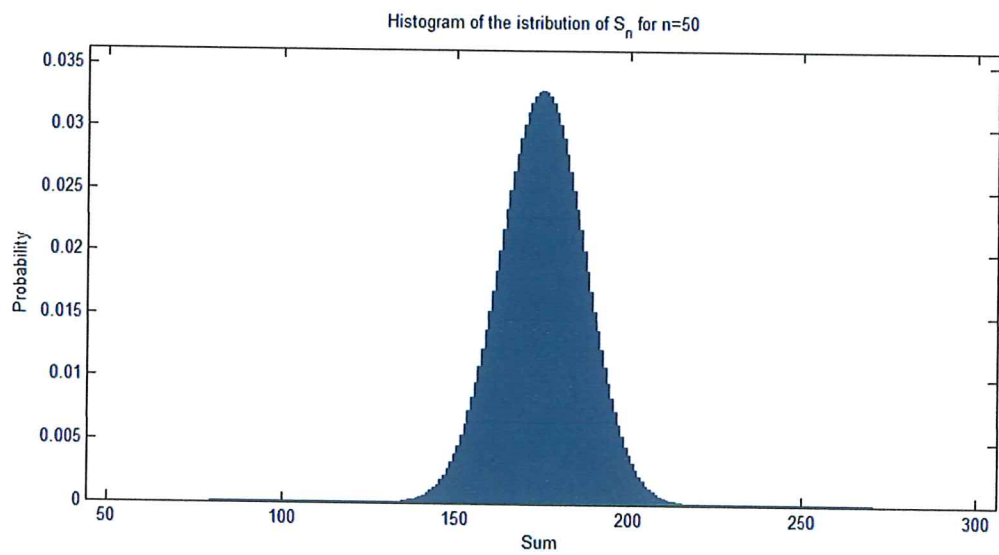


Histogram of the istribution of  $S_n$  for  $n=50$

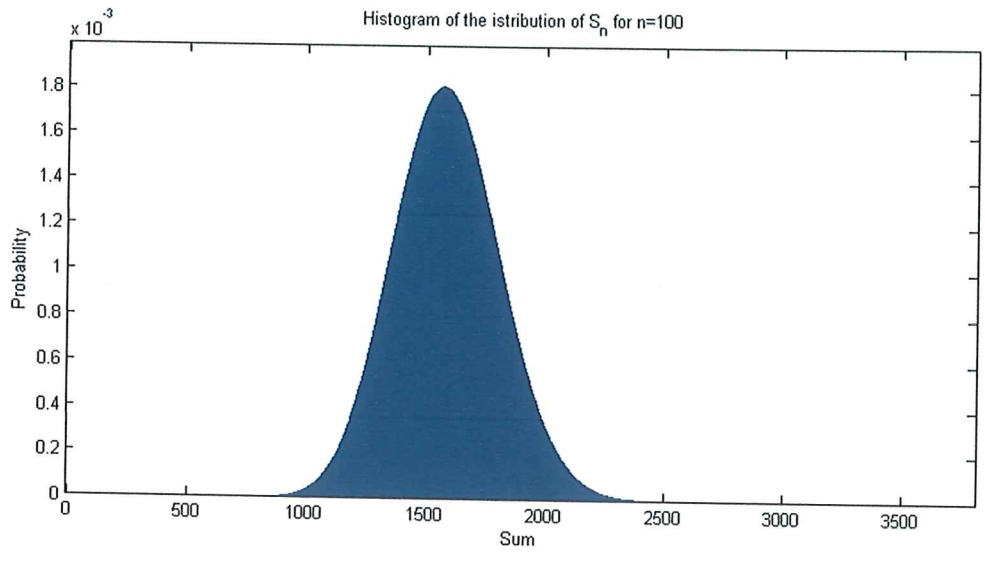




3. Take  $P(X_1 = k) = 1/6$  for  $k = 1, 2, \dots, 6$ . Take  $n = 50$ .



4. Take  $P(X_1 = 2^k) = 1/7$  for  $k = 0, 1, 2, 3, 4, 5, 6$ . Take  $n = 100$ .



From the examples we infer that the distribution of  $S_n$  is similar to the normal distribution. This is not clear in itself. But if we accept this observation we need to find the normal distribution that fits the histogram of  $S_n$  well.

Idea: We match the top of the two distributions which means that the mean of the normal distribution will be  $E(S_n)$ . We match the dispersion by choosing the second parameter to be  $\text{var}(S_n)$ .

Denote :

$$E(X_1) = \mu \quad \text{var}(X_1) = \sigma^2$$
$$E(S_n) = n E(X_1) = \nu$$
$$\text{var}(S_n) = n \text{var}(X_1) = \tau^2$$

To approximate the distribution we can say

$$P(a \leq S_n \leq b) \approx \frac{1}{\sqrt{2\pi} \cdot \tau} \int_a^b e^{-\frac{(x-\nu)^2}{2\tau^2}} dx$$

The area of columns in the histogram between  $a$  and  $b$  ~~are~~ is exactly  $P(a \leq S_n \leq b)$ . We superimpose a curve closely following the histogram and replace the area of columns with the integral under the curve.

To turn the above into a mathematical theorem we will reformulate.

Take  $a = \nu + \alpha \cdot \tau$  and  $b = \nu + \beta \cdot \tau$ .

We compute

$$P(a \leq S_u \leq b)$$

$$= P(\nu + \alpha \cdot \tau \leq S_u \leq \nu + \beta \cdot \tau)$$

$$\approx \frac{1}{\sqrt{2\pi} \tau} \int_{\nu + \alpha \cdot \tau}^{\nu + \beta \cdot \tau} e^{-\frac{(x - \nu)^2}{2\tau^2}} dx$$

○ New variable:  $\frac{x - \nu}{\tau} = u$

$$= \frac{1}{\sqrt{2\pi}} \int_{\alpha}^{\beta} e^{-u^2/2} du.$$

On the other hand we have

○  $P(\nu + \alpha \cdot \tau \leq S_u \leq \nu + \beta \cdot \tau)$

$$= P\left(\alpha \leq \frac{S_u - \nu}{\tau} \leq \beta\right)$$

$$= P\left(\alpha \leq \frac{S_u - E(S_u)}{\sqrt{\text{var}(S_u)}} \leq \beta\right)$$



Definition: The expression

$$\tilde{S}_n = \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}}$$

is called the standardized sum.

We noticed that the approximation

is "better" if  $n$  is "large". We

expect the mathematical form

to include limits.

Theorem 6.1 (central limit theorem)

Let  $X_1, X_2, \dots$  be independent

equally distributed random

variables with  $E(X_i) = \mu$  and

$\text{var}(X_i) = \sigma^2 < \infty$ . Let  $S_n = X_1 + \dots + X_n$ .

For any  $\alpha < \beta$  we have

$$\lim_{n \rightarrow \infty} P\left(\alpha \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \beta\right)$$

$$= \Phi(\beta) - \Phi(\alpha)$$

where  $\Phi$  is the distribution function of the standard normal distribution.

### Comments:

(i) we will prove the theorem in several steps. It is true as it is formulated but we will impose the additional assumption  $E(|X_1|^3) < \infty$ .

(ii) It is enough to prove

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \beta\right) = \Phi(\beta)$$

for  $\beta \in \mathbb{R}$ .

(iii) In the limit we get equality.

For finite  $n$  we use the limit as an approximation.

To prove the central limit theorem we need the following result.

Theorem 6.2 (Lindeberg-Bergman)

Let  $X_1, X_2, \dots, X_n$  be independent and such that  $\text{var}(X_1 + X_2 + \dots + X_n) = 1$   
and  $E(X_1 + \dots + X_n) = 0$ . Assume that  
 $E(|X_k|^3) < \infty$  for all  $k = 1, 2, \dots, n$ .

Let  $f$  be a three times continuously differentiable function such that  $|f(x)|, |f'(x)|, |f''(x)|, |f'''(x)| \leq M$   
for some  $M < \infty$  and all  $x \in \mathbb{R}$ .

Let  $S_n = X_1 + X_2 + \dots + X_n$ . Then

$$\begin{aligned} & |E(f(S_n)) - E(f(z))| \\ & \leq \frac{1}{6} M \left(1 + \sqrt{\frac{8}{\pi}}\right) E(|X_1|^3 + \dots + |X_n|^3). \end{aligned}$$

for  $z \sim N(0, 1)$ .

Proof: Without loss of generality we can assume  $E(X_k) = 0$  for all  $k = 1, 2, \dots, n$ . Let  $Z_1, Z_2, \dots, Z_n$  be independent and independent of  $X_1, X_2, \dots, X_n$  and such that  $Z_k \sim N(0, \text{var}(X_k))$ ,  $k = 1, 2, \dots, n$ .

Since  $\text{var}(X_1) + \dots + \text{var}(X_n) = 1$  by assumption we have that

$$Z = Z_1 + Z_2 + \dots + Z_n \sim N(0, 1).$$

Define

$$a_1 = E[f(Z_1 + Z_2 + \dots + Z_n)] - E[f(X_1 + Z_2 + \dots + Z_n)]$$

$$a_2 = E[f(X_1 + Z_2 + \dots + Z_n)] - E[f(X_1 + X_2 + \dots + Z_n)]$$

⋮

$$a_n = E[f(X_1 + \dots + X_{n-1} + Z_n)] - E[f(X_1 + X_2 + \dots + X_n)]$$

By triangle inequality we have

$$|E[f(X_1 + \dots + X_n)] - E[f(Z_1 + \dots + Z_n)]| \leq \sum_{k=1}^n |a_k|$$

By Taylor we have

$$f(x+h) - f(x) = f'(x) \cdot h + \frac{1}{2} f''(x) h^2 + r$$

where  $r = \frac{1}{6} f'''(\xi) h^3$  for some

$\xi$  between  $x$  and  $x+h$ . By our

assumption  $|r| \leq \frac{1}{6} \cdot M \cdot |h|^3$ .

○ Define

$$Y_1 = z_2 + z_3 + \dots + z_n$$

$$Y_2 = x_1 + z_3 + \dots + z_n$$

$$Y_3 = x_1 + x_2 + z_4 + \dots + z_n$$

○ 
$$Y_n = x_1 + x_2 + \dots + x_{n-1}$$

Note that  $Y_k$  is independent of  $(x_k, z_k)$  for all  $k = 1, 2, \dots, n$ .

We use Taylor's expansion around  $Y_k$  to get

$$E \left[ f(x_1 + \dots + x_{k-1} + z_k + \dots + z_n) \right]$$

$$= E \left[ f(y_k) + f'(y_k) z_k + \frac{1}{2} f''(y_k) z_k^2 + R_k \right]$$

and

$$E \left[ f(x_1 + \dots + x_k + z_{k+1} + \dots + z_n) \right]$$

$$= E \left[ f(y_k) + f'(y_k) x_k + \frac{1}{2} f''(y_k) x_k^2 + \tilde{R}_k \right]$$

Subtracting we get

$$\alpha_k = E \left[ f'(y_k)(z_k - x_k) + \frac{1}{2} f''(y_k)(z_k^2 - x_k^2) + R_k - \tilde{R}_k \right]$$

$$= E \left[ f'(y_k)(z_k - x_k) \right]$$

$$+ E \left[ \frac{1}{2} f''(y_k)(z_k^2 - x_k^2) \right]$$

$$+ E \left[ R_k - \tilde{R}_k \right]$$

By independence

$$\begin{aligned} E [ f'(Y_k) (Z_k - X_k) ] \\ &= E [ f'(Y_k) ] \underbrace{E [ Z_k - X_k ]}_{= 0 \text{ by assumption}} \\ &= 0 \end{aligned}$$

and

$$\begin{aligned} E [ f''(Y_k) (Z_k^2 - X_k^2) ] \\ &= E [ f''(Y_k) ] \underbrace{E [ Z_k^2 - X_k^2 ]}_{= 0 \text{ by assumption}} \\ &= 0. \end{aligned}$$

We are left with

$$a_k = E [ R_k - \tilde{R}_k ]$$

Since  $-|x| \leq x \leq |x|$  we have  $|E(x)| \leq E(|x|)$

so

$$|a_k| \leq E [ |R_k - \tilde{R}_k| ]$$

But

$$|R_k| \leq \frac{1}{6} M \cdot |X_k|^3$$

$$|\tilde{R}_k| \leq \frac{1}{6} M |Z_k|^3$$

so

$$|R_k - \tilde{R}_k| \leq \frac{1}{6} M (|X_k|^3 + |Z_k|^3).$$

It follows

$$E[|R_k - \tilde{R}_k|] \leq \frac{1}{6} \cdot M (E(|X_k|^3) + E(|Z_k|^3)).$$

A standard calculation gives that for  $Z \sim N(0, \sigma^2)$  we have

$$E(|Z|^3) = \sqrt{\frac{8}{\pi}} \cdot \sigma^3$$

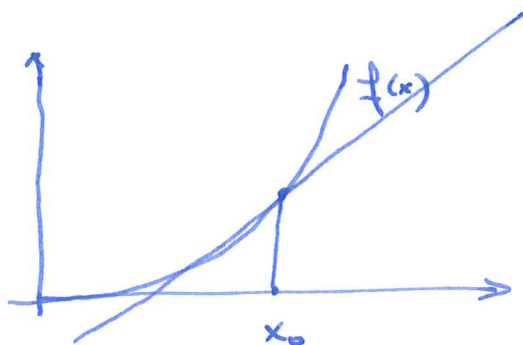
We then have

$$\begin{aligned} E(|Z_k|^3) &= \sqrt{\frac{8}{\pi}} \text{var}(X_k)^{3/2} \\ &= \sqrt{\frac{8}{\pi}} \cdot E(X_k^2)^{3/2} \\ &= \sqrt{\frac{8}{\pi}} \cdot E(|X_k|^2)^{3/2} \end{aligned}$$



Since  $f(x) = x^{3/2}$  is convex,  
the function is above its tangent.

Figure :



We have for  $x_0 > 0$

$$f(x) = x^{3/2} \geq \underbrace{f'(x_0)(x - x_0) + f(x_0)}_{\text{tangent}}$$

Take  $x_0 = E(|X_k|^2)$ . We have

$$(|X_k|^2)^{3/2} \geq f'(x_0)(|X_k|^2 - x_0) + x_0^{3/2}$$

Taking expectations we get

$$\begin{aligned} E(|X_k|^3) &\geq E(|X_k|^2)^{3/2} \\ &= (\text{var}(X_k))^{3/2} \end{aligned}$$

Taking all inequalities we get

$$E(|x_k|^3) + E(|z_k|^3)$$

$$\leq \left(1 + \sqrt{\frac{8}{\pi}}\right) E(|x_k|^3)$$

Finally,

$$\sum_{k=1}^n |a_k| \leq \sum_{k=1}^n \frac{1}{6} \cdot M \left(1 + \sqrt{\frac{8}{\pi}}\right) E(|x_k|^3) \quad \square$$

The inequality is valid for arbitrary  $x_1, x_2, \dots, x_n$  provided they are independent. If

$x_1, x_2, \dots, x_n$  are independent and

equally distributed define

$$x_k' = \frac{x_k - E(x_k)}{\sqrt{\text{var}(S_n)}}$$

for  $S_n = x_1 + x_2 + \dots + x_n$ . Note that

$$x_1' + x_2' + \dots + x_n' = \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} = \tilde{S}_n$$

We have that  $X_1', X_2', \dots, X_n'$

are independent,  $E(X_k') = 0$

and  $\text{var}(X_1' + \dots + X_n') = 1$ .

Theorem 6.2 implies

$$|E[f(\hat{S}_n)] - E[f(z)]|$$

$$\leq \frac{1}{6} M \left(1 + \sqrt{\frac{8}{\pi}}\right) \cdot n \cdot E(|X_1'|^3)$$

But

$$E(|X_1'|^3) = E\left(\frac{|X_1 - E(X_1)|^3}{\sqrt{n} \sqrt{\text{var}(X_1)}}\right)$$

$$= \frac{1}{n^{3/2} \text{var}(X_1)^{3/2}} E(|X_1 - E(X_1)|^3)$$

If  $\gamma = E(|X_1 - E(X_1)|^3) < \infty$  we have

$$|E[f(\hat{S}_n)] - E[f(z)]|$$

$$\leq \frac{1}{6} \cdot M \left(1 + \sqrt{\frac{8}{\pi}}\right) \cdot \frac{1}{\sqrt{n}} \cdot \frac{\gamma}{\text{var}(X_1)^{3/2}}$$

$\rightarrow 0$ , as  $n \rightarrow \infty$ .

We will prove the central limit theorem under the additional assumption that  $\gamma = E(|X_1 - E(X_1)|^3) < \infty$ .

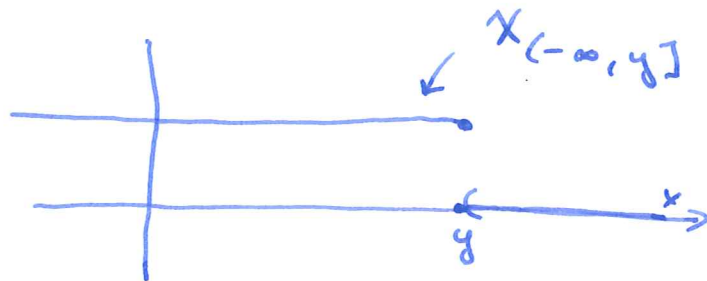
Comment: If  $E(|X_1|^3) < \infty$  then  $\gamma < \infty$ .

Proof: Let  $\varepsilon > 0$ . The distribution

function  $\Phi$  of  $Z \sim N(0,1)$  is continuous so for a fixed  $\beta \in \mathbb{R}$  there is a  $\delta > 0$  such that  $|\Phi(x) - \Phi(\beta)| < \varepsilon$  for  $|x - \beta| < \delta$ .

Denote by  $X_{(-\infty, y]}$  the indicator function of the interval  $(-\infty, y]$ .

Figure:



Analysis 1 gives: there are functions  $f^{-\varepsilon}$  and  $f^{\varepsilon}$  with values on  $[0,1]$  such that:

(i)  $f^{-\varepsilon}, f^{\varepsilon}$  are three times continuously differentiable.

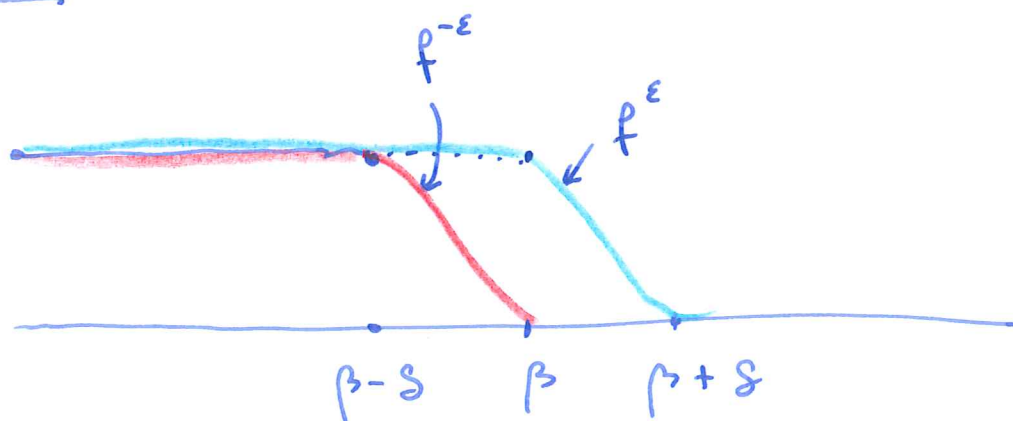
(ii) derivatives up to the third are bounded by  $M < \infty$

(iii)

$$\chi_{(-\infty, \beta - \delta]} \leq f^{-\varepsilon} \leq \chi_{(-\infty, \beta]}$$

$$\leq f^{\varepsilon} \leq \chi_{(-\infty, \beta + \delta]}$$

Figure:



We have

$$\begin{aligned} E(\chi_{(-\infty, y]}(\tilde{S}_n)) \\ = P(\tilde{S}_n \leq y) \end{aligned}$$

and similarly

$$E(\chi_{(-\infty, y]}(z)) = P(z \leq y) = \Phi(y).$$

Inequalities in (iii) imply

$$\begin{aligned} \Phi(\beta - \delta) &\leq E(\varphi^{-\varepsilon}(z)) \\ &\leq \Phi(\beta) \leq E(\varphi^{\varepsilon}(z)) \\ &\leq \Phi(\beta + \delta) \end{aligned}$$

$$\text{But } E(\varphi^{-\varepsilon}(\tilde{S}_n)) \rightarrow E(\varphi^{-\varepsilon}(z))$$

$$E(\varphi^{\varepsilon}(\tilde{S}_n)) \rightarrow E(\varphi^{\varepsilon}(z))$$

as  $n \rightarrow \infty$ .

For sufficiently large  $n$   
we will have

$$\Phi(\beta - \delta) - \varepsilon \leq \mathbb{P}(\tilde{S}_n \leq \beta) \leq \Phi(\beta + \delta) + \varepsilon.$$

or

$$\Phi(\beta) - 2\varepsilon \leq \mathbb{P}(\tilde{S}_n \leq \beta) \leq \Phi(\beta) + 2\varepsilon.$$

○ This proves that

$$\lim_{n \rightarrow \infty} \mathbb{P}(\tilde{S}_n \leq \beta) = \Phi(\beta).$$

Examples: Typically we want to

○ estimate probabilities of the form

$$\mathbb{P}(a \leq S_n \leq b).$$

we compute

$$\mathbb{P}(a \leq S_n \leq b)$$

$$= \mathbb{P}(a - E(S_n) \leq S_n - E(S_n) \leq b - E(S_n))$$

$$= P \left( \underbrace{\frac{a - E(S_n)}{\sqrt{\text{var}(S_n)}}}_{\alpha} \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \underbrace{\frac{b - E(S_n)}{\sqrt{\text{var}(S_n)}}}_{\beta} \right)$$

$$= P \left( \alpha \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \beta \right)$$

CLT  
x

$$\Phi(\beta) - \Phi(\alpha)$$

(i) Let  $X_1, X_2, \dots$  be independent

and  $X_k \sim \text{Bernoulli}(p)$ . We know

that  $S_n = X_1 + X_2 + \dots + X_n \sim \text{Bin}(n, p)$

so  $E(S_n) = n \cdot p$  and  $\text{var}(S_n) = npq$ .

Assume  $n = 10,000$  and  $p = \frac{1}{2}$

and  $a = 4950$  and  $b = 5050$ . We

have

$$P(4950 \leq S_n \leq 5050)$$

$$= P \left( \frac{4950 - 5000}{50} \leq \frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq \frac{5050 - 5000}{50} \right)$$

$$\approx P(-1 \leq Z \leq 1)$$



Statistical software gives

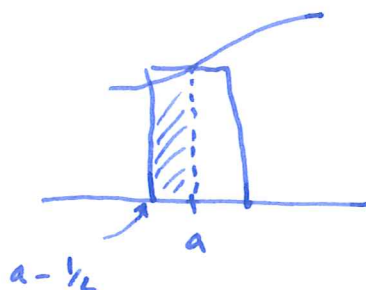
$$P(-1 \leq z \leq 1) = \Phi(1) - \Phi(-1)$$

$$= 0.6827$$

The exact probability is 0.6825.

If  $x_1, x_2, \dots$  are integer valued we can improve the approximation by changing  $a$  to  $a - 1/2$  and  $b$  to  $b + 1/2$ .

Figure:



Changing  $a$  to  $a - 1/2$  adds the "half" of the column over  $a$ . This correction is called correction for continuity.

Using this correction we get

$$P(4950 \leq S_n \leq 5500)$$

$$\approx \Phi(1.0100) - \Phi(-1.0100)$$

$$\doteq 0.6875.$$

This is accurate to 4 decimals!

Q ii) Let  $X_1, X_2, \dots$  be independent  
and  $P(X_i = 1) = P(X_i = 2) = P(X_i = 3) = 1/3$

Let  $n = 300$ . We find

$$E(S_{300}) = 300 \cdot 4 = 1200$$

$$\text{var}(S_{300}) = 300 \cdot \frac{38}{3} = 3800$$

We approximate

$$P(1,100 \leq S_{300} \leq 1,300)$$

$$= P\left(\frac{1100 - 1200}{\sqrt{3800}} \leq \frac{S_{300} - E(S_{300})}{\sqrt{3800}} \leq \frac{1300 - 1200}{\sqrt{3800}}\right)$$

$$\approx \Phi(1.622) - \Phi(-1.622)$$

$$\doteq 0.8952$$

The exact probability using the fast Fourier transform turns out to be 0.8970. If we include the continuity correction we get 0.8970!

Can we say anything about the accuracy approximation? The answer is yes but the proof is demanding.

Theorem 6.3 (Berry-Esséen) Let

$$g = E(|X_1 - E(X_1)|^3) \text{ and keep all}$$

the assumptions of Theorem 6.1.

Then

$$\sup_{x \in \mathbb{R}} \left| P\left(\frac{S_n - E(S_n)}{\sqrt{\text{var}(S_n)}} \leq x\right) - \Phi(x) \right|$$

$$\leq \frac{C \cdot g}{\sqrt{n} \text{var}(X_1)^{3/2}}$$

where  $C < 0.4748$ .

For a proof see

Shvertsova, I., On the accuracy  
of the normal approximation  
for sums of independent symmetric  
random variables, Dokl. Akad. Nauk  
443 (2012), no. 6, 671-676.

Ⓜ THE END Ⓜ