

Last time:

Sec. 6.4 Chebychev's inequality

Sec. 7.1, 7.2 Variance of sums (Binomial, Poisson, NegBinomial, Hypergeometric)

Sec. 7.3 Law of averages / Law of large number

Today: the Central Limit Theorem "CLT"

Sec. 8.1 The dist. of an i.i.d. sum.

$X_1, X_2, X_3, \dots, X_n$ i.i.d. with mean μ , sd σ

Let $S_n := \sum_{i=1}^n X_i$

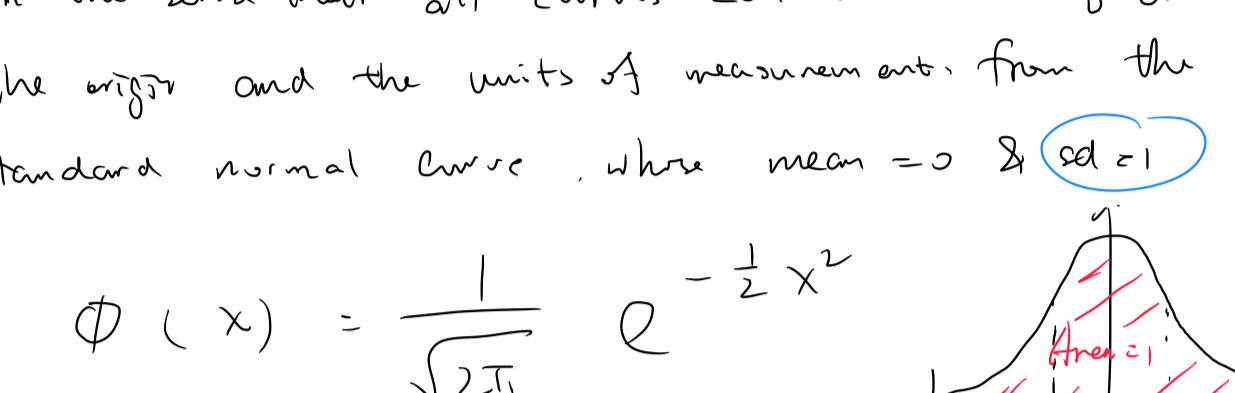
We want to see the shape of the dist. of S_n .

$$\mathbb{E} S_n = n\mu \quad \text{Var}(S_n) = n\text{Var}(X_i) = n\sigma^2$$

$$\text{SD}(S_n) = \sqrt{n}\sigma$$

Roughly speaking, CLT tells us that S_n will eventually

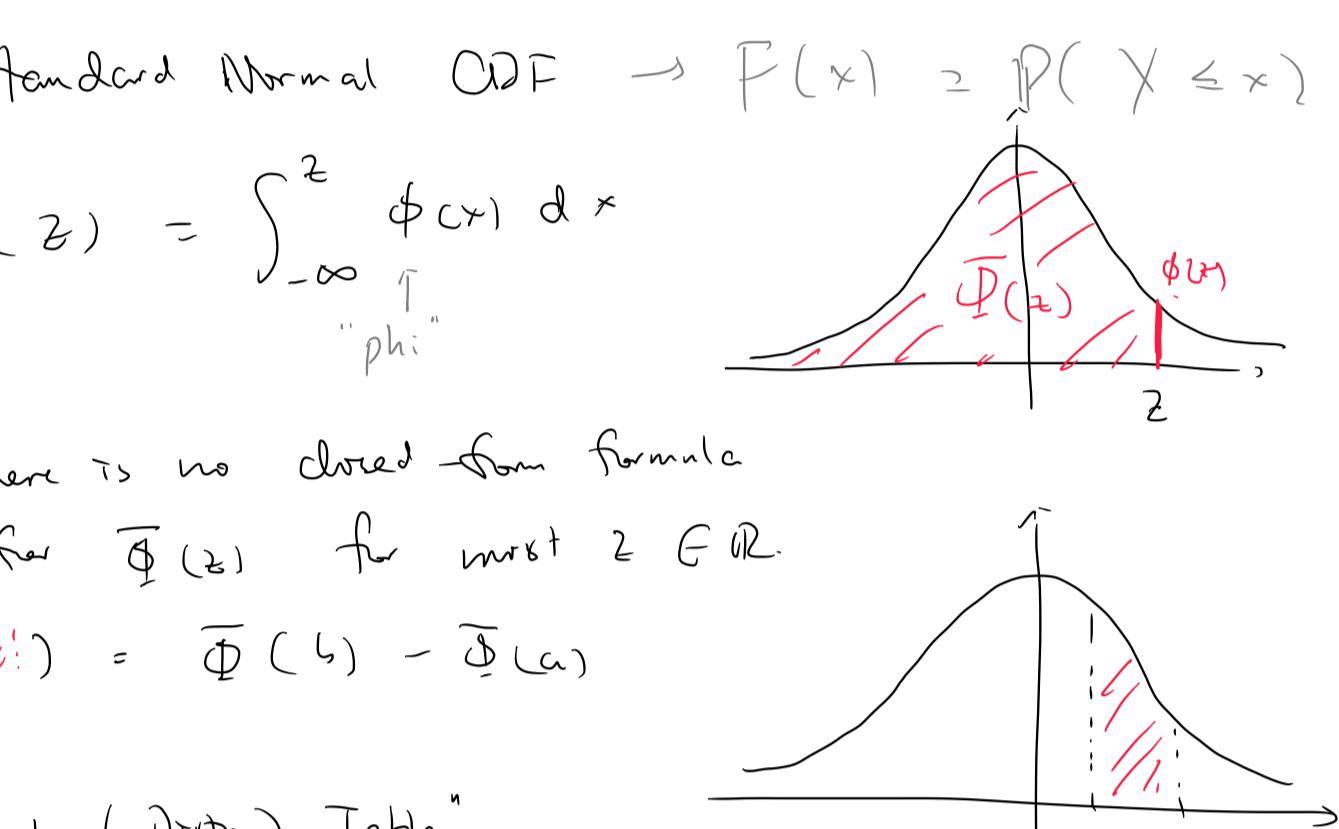
become bell-shaped when n is large, regardless of the dist. of X_i .



To draw the graph:

① the center of the bell is $\mathbb{E} S_n = n\mu$

② $\mathbb{E} S_n \pm \text{SD}(S_n)$ are the two points of inflection



Sec. 8.2 the Standard Normal Curve

For different $n\mu$ & $\sqrt{n}\sigma$, you may result in different bell-shaped curves, but they are essentially the same in the sense that all curves can be derived by changing the origin and the units of measurement from the standard normal curve, where mean = 0 & $\text{sd} = 1$.

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$$

Properties:

- ① bell-shaped
- ② symmetric about 0 / even function
- ③ two points of inflection ± 1 — $\Phi(-1) = 0$
- ④ pretty close to 0 outside of $(-3, +3)$
- ⑤ a prob. dist. histogram (prob. density function) from next week.

The standard Normal CDF $\rightarrow F(x) = P(X \leq x)$

$$\Phi(z) = \int_{-\infty}^z \phi(x) dx$$

There is no closed form formula

for $\Phi(z)$ for most $z \in \mathbb{R}$

$$P(a) = \Phi(b) - \Phi(a)$$



"Normal (Distr.) Table"

$$z \mid \Phi(z)$$

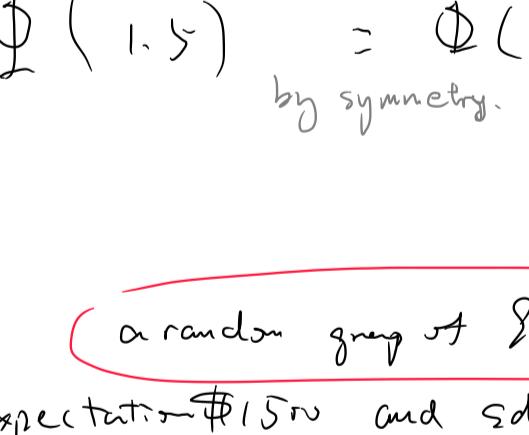
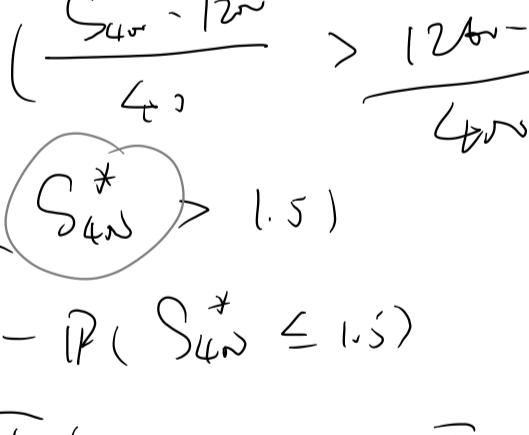
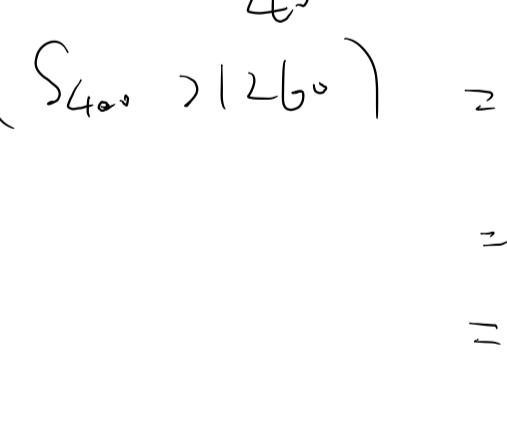
$$1 \mid 0.84$$

$$2 \mid 0.977$$

Easier to do with Python by

stats.norm.cdf(z)

$$= \Phi(z)$$



Percentiles:

the 84th percentile is ≈ 1

the 98th percentile is ≈ 2

the pth percentile is given by

$$z = \Phi^{-1}(p\%)$$

$$\text{since } \Phi(z) = p\% \Leftrightarrow z = \Phi^{-1}(p\%)$$

In python, use the percent point function (ppf)

$$\text{stats.norm.ppf}(p)$$

Sec. 8.3 Normal Approximation

How could we use CLT and the normal curve

to approx. a prob. dist.?

the restatement of

First, we consider standard unit \Rightarrow Recall Chebychev's inequality

use $\mu = \mathbb{E} X$ as origin

use $\sigma = \text{SD}(X)$ as the unit distance

to create a new R.V. $X^* = X - \mu$ in standard units

$$X^* = \frac{X - \mu}{\sigma} = \frac{X - \mu}{\sigma} \Leftrightarrow X = \mu + \sigma X^*$$

$$\left\{ \begin{array}{l} \mathbb{E} X^* = \frac{1}{\sigma}(\mathbb{E}(X) - \mu) = \frac{1}{\sigma}(\mu - \mu) = 0 \\ \text{SD}(X^*) = \frac{1}{\sigma} \text{SD}(X - \mu) = \frac{1}{\sigma} \text{SD}(X) = \frac{1}{\sigma} \cdot \sigma = 1 \end{array} \right.$$

Example: X random weight in pounds (lbs)

$$\mathbb{E} X = 150$$

$$\text{SD}(X) = 20$$

$$X = 150 + 20 \Leftrightarrow X^* = \frac{150 - 150}{20} = 0$$

X is one SD to the right of mean

$$X^* = -1.5 \Leftrightarrow X = 150 + 20 \cdot (-1.5) = 120$$

negative means smaller than expectation

Percentile:

the 84th percentile is ≈ 1

the 98th percentile is ≈ 2

the pth percentile is given by

$$z = \Phi^{-1}(p\%)$$

$$\text{since } \Phi(z) = p\% \Leftrightarrow z = \Phi^{-1}(p\%)$$

In python, use the percent point function (ppf)

$$\text{stats.norm.ppf}(p)$$

Sec. 8.4 How large is large?

- If the dist. is good $\left\{ \begin{array}{l} \text{Smooth} \\ \text{uni-modal} \\ \text{balanced} \end{array} \right.$

\Rightarrow not very large n needed

(30)

- If the dist. is bad $\left\{ \begin{array}{l} \text{sharp} \\ \text{many curves} \\ \text{skewed} \end{array} \right.$

\Rightarrow need larger n .

(100 or 1000)

In this course: Consider good dist. only

Calculate the mean and sd. and check if

mean $\pm 3\sigma$ are within reasonable / possible region

Example: $X_i \sim \text{Bern}(\frac{1}{100})$ $\mathbb{E} X_i = \frac{1}{100}$ $\text{SD}(X_i) = \sqrt{\frac{1}{100} \cdot \frac{99}{100}} \approx \frac{1}{10}$

$$S_n = X_1 + X_2 + \dots + X_{100}$$

$$\mathbb{E} S_n = 1$$

$$\text{SD}(S_n) \approx \sqrt{100} \cdot \frac{1}{10} = 1$$

$$1 \pm 3 = (-2) \text{ or } 4$$

not reasonable

$$\Rightarrow n = 100 \text{ not large enough.}$$

$$n = 1000$$

$$\mathbb{E} S_n = 10$$

$$\text{SD}(S_n) = \sqrt{1000} \cdot \frac{1}{10} = 10$$

$$10 \pm 3 \cdot 10 = 70 \text{ or } 130$$

$$\Rightarrow n = 1000 \text{ large enough.}$$

$$n = 10000$$

$$\mathbb{E} S_n = 100$$

$$\text{SD}(S_n) = \sqrt{10000} \cdot \frac{1}{10} = 100$$

$$100 \pm 3 \cdot 100 = 700 \text{ or } 1300$$

$$\Rightarrow n = 10000 \text{ large enough.}$$