

Central Limit Theorem

CSE 312 Spring 21
Lecture 19

Announcements

No responsibilities this summer? Enjoying this class? Consider Taing for it. Kushal will be the instructor! (more info [on Ed](#)).

Why Learn Normals?

When we add together independent normal random variables, you get another normal random variable.

The sum of **any** independent random variables **approaches** a normal distribution.

Central Limit Theorem

Let X_1, X_2, \dots, X_n be i.i.d. random variables, with mean μ and variance σ^2 . Let $Y_n = \frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}}$

As $n \rightarrow \infty$, the CDF of Y_n converges to the CDF of $\mathcal{N}(0, 1)$

Breaking down the theorem

Central Limit Theorem

Let X_1, X_2, \dots, X_n be i.i.d. random variables, with mean μ and variance σ^2 . Let $Y_n = \frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}}$

As $n \rightarrow \infty$, the CDF of Y_n converges to the CDF of $\mathcal{N}(0, 1)$

Proof of the CLT?

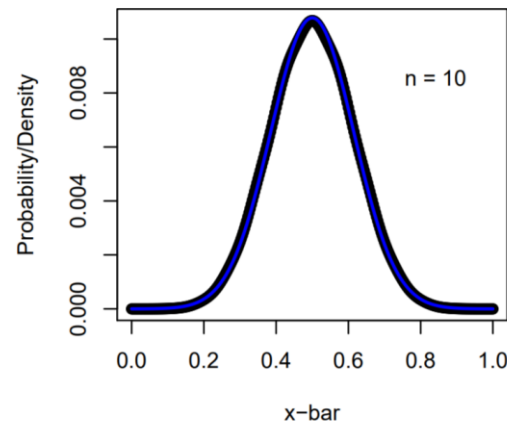
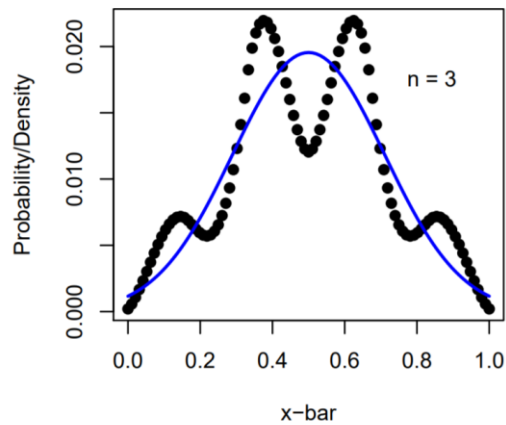
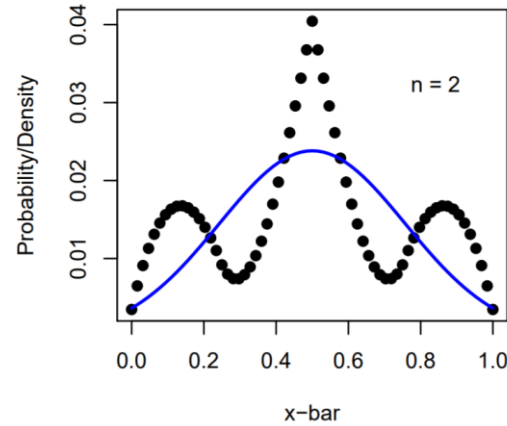
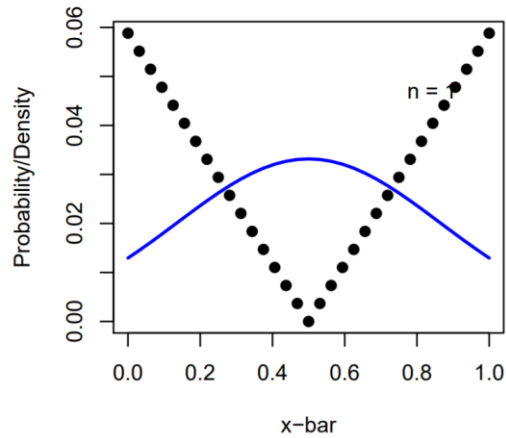
No.

How is the proof done?

Step 1: Prove that for all positive integers k , $\mathbb{E}[(Y_n)^k] \rightarrow \mathbb{E}[Z^k]$

Step 2: Prove that if $\mathbb{E}[(Y_n)^k] = \mathbb{E}[Z^k]$ for all k then $F_{Y_n}(z) = F_Z(z)$

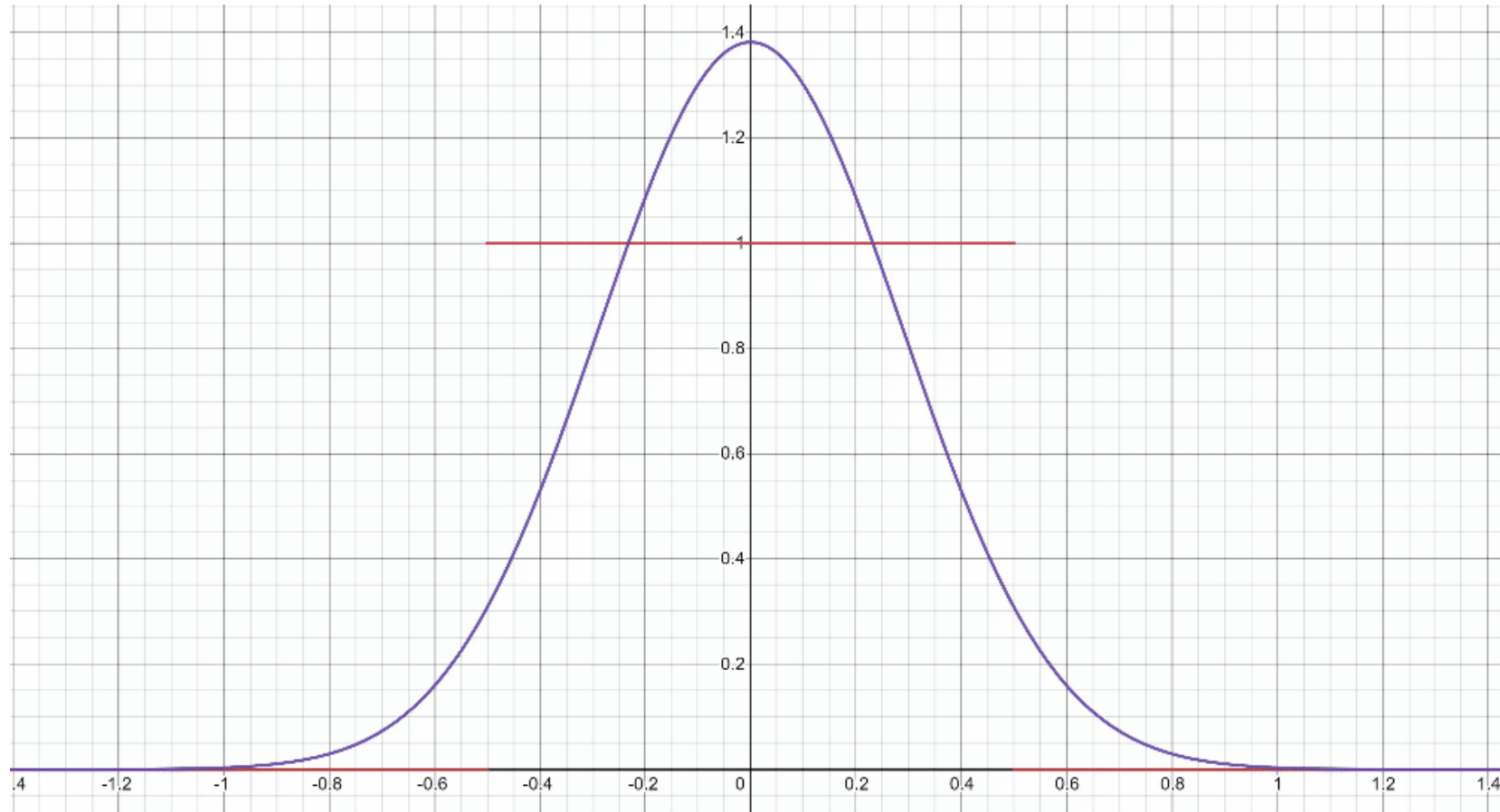
“Proof by example”



The dotted lines show an “empirical pmf” – a pmf estimated by running the experiment a large number of times. The blue line is the normal rv that the CLT predicts.

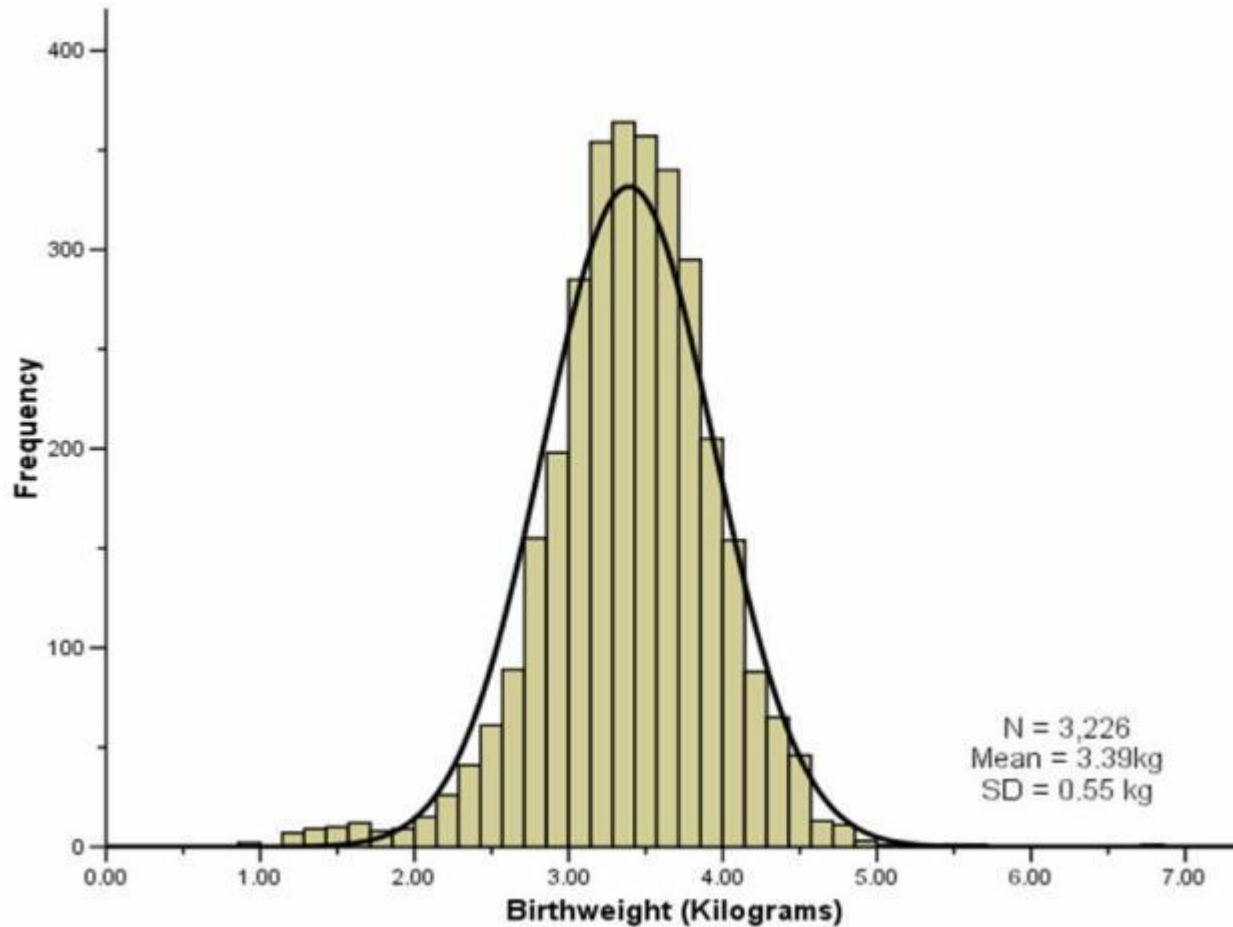
Shown are $n = 1, 2, 3, 10$

"Proof by example" -- uniform



<https://www.desmos.com/calculator/2n2m05a9km>

“Proof by real-world”



birthweight

A lot of real-world bell-curves can be explained as:

1. The random variable comes from a combination of independent factors.
2. The CLT says the distribution will become like a bell curve.

Theory vs. Practice

The formal theorem statement is “in the limit”

You might not get exactly a normal distribution for any finite n (e.g. if you sum indicators, your random variable is always discrete and will be discontinuous for every finite n).

In practice, the approximations get very accurate very quickly (at least with a few tricks we’ll see soon).

They won’t be exact (unless the X_i are normals) but it’s close enough to use even with relatively small n .

Using the Central Limit Theorem

Suppose you are managing a factory, that produces widgets. Each widget produced is defective (independently) with probability 5%.

Your factory will produce 1000 (possibly defective) widgets. You want to know what the chances are of having a “very bad day” where “very bad” means producing at most 940 non-defective widgets.
(In expectation, you produce 950 non-defective widgets)

What is the probability?

True Answer

Let $X \sim \text{Bin}(1000, .95)$

What is $\mathbb{P}(X \leq 940)$?

The cdf is ugly...and that's a big summation.

$$\sum_{k=0}^{940} \binom{1000}{k} (.95)^k \cdot (.05)^{1000-k} \approx .08673$$

What does the CLT give?

CLT setup

$\text{Bin}(1000, .95)$ is the sum of a bunch of independent random variables (the indicators/Bernoullis we summed to get the binomial)

So, let's use the CLT instead

$$\mathbb{E}[X_i] = p = .95.$$

$$\text{Var}(X_i) = p(1 - p) = .0475$$

$$Y_{1000} = \frac{\sum_{i=1}^{1000} X_i - 1000 \cdot .95}{\sqrt{1000 \cdot .0475}} \text{ is approximately } \mathcal{N}(0,1).$$

With the CLT.

The event we're interested in is $\mathbb{P}(X \leq 940)$

$$\mathbb{P}(X \leq 940)$$

$$= \mathbb{P}\left(\frac{X - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}} \leq \frac{940 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right)$$

$$= \mathbb{P}\left(Y_{1000} \leq \frac{940 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right)$$

$$\approx \mathbb{P}\left(Y \leq \frac{940 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right) \text{ by CLT}$$

$$= \Phi\left(\frac{940 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right)$$

$$\approx \Phi(-1.45) = 1 - \Phi(1.45)$$

$$\approx 1 - .92647 = .07353.$$

It's an approximation!

The true probability is

$$1 - \sum_{k=941}^{1000} \binom{1000}{k} (.95)^k \cdot (.05)^{1000-k} \approx .08673$$

The CLT estimate is off by about 1.3 percentage points.

We can get a better estimate if we fix a subtle issue with this approximation.

A problem

What's the probability that $X = 950$? (exactly)

True value, we can get with binomial:

$$\binom{1000}{950} \cdot (.95)^{950} \cdot (.05)^{50} \approx .05779$$

What does the CLT say?

$$= \mathbb{P}\left(\frac{X - 1000 \cdot .95}{\sqrt{1000 \cdot .0475}} = \frac{950 - 1000 \cdot .95}{\sqrt{1000 \cdot .0475}}\right)$$

$$\approx \mathbb{P}(Y = 0)$$

$$= 0$$

Uh oh.

Continuity Correction

The binomial distribution is discrete, but the normal is continuous.

Let's correct for that (called a "continuity correction")

Before we switch from the binomial to the normal, ask "what values of a continuous random variable would round to this event?"

Applying the continuity correction

$$\begin{aligned}\mathbb{P}(X = 950) \\ = \mathbb{P}(949.5 \leq X < 950.5)\end{aligned}$$

Continuity correction.
This is an "exactly equal to"
The discrete rv X can't equal 950.2.

$$\begin{aligned}&= \mathbb{P}\left(\frac{949.5-950}{\sqrt{1000 \cdot 0.0475}} \leq \frac{X-950}{\sqrt{1000 \cdot 0.0475}} < \frac{950.5-950}{\sqrt{1000 \cdot 0.0475}}\right) \\ &\approx \mathbb{P}\left(\frac{949.5-950}{\sqrt{1000 \cdot 0.0475}} \leq Y < \frac{950.5-950}{\sqrt{1000 \cdot 0.0475}}\right) \text{ By CLT} \\ &= \Phi\left(\frac{950.5-950}{\sqrt{1000 \cdot 0.0475}}\right) - \Phi\left(\frac{949.5-950}{\sqrt{1000 \cdot 0.0475}}\right) \\ &\approx \Phi(0.07) - \Phi(-0.07) = \Phi(0.07) - (1 - \Phi(0.07)) \\ &\approx 0.5279 - (1 - 0.5279) = 0.0558\end{aligned}$$

Still an Approximation

$\binom{1000}{950} \cdot (.95)^{950} \cdot (.05)^{50} \approx .05779$ is the true value

The CLT approximates: 0.0558

Very close! But still not perfect.

Let's fix that other one

Question was "what's the probability of seeing at most 940 non-defective widgets?"

With the CLT.

The event we're interested in is $\mathbb{P}(X \leq 940)$

$$\mathbb{P}(X \leq 940)$$

$$= \mathbb{P}\left(\frac{X - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}} \leq \frac{940 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right)$$

$$\approx \mathbb{P}\left(Y \leq \frac{940 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right) \text{ By CLT}$$

$$= \Phi\left(\frac{940 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right)$$

$$\approx \Phi(-1.45) = 1 - \Phi(1.45)$$

$$\approx 1 - .92647 = .07353.$$

$$\mathbb{P}(X \leq 940.5)$$

$$= \mathbb{P}\left(\frac{X - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}} \leq \frac{940.5 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right)$$

$$\approx \mathbb{P}\left(Y \leq \frac{940.5 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right) \text{ By CLT}$$

$$= \Phi\left(\frac{940.5 - 1000 \cdot 0.95}{\sqrt{1000 \cdot 0.0475}}\right)$$

$$\approx \Phi(-1.38) = 1 - \Phi(1.38)$$

$$\approx 1 - .91621 = .08379.$$

True answer: .08673

Approximating a continuous distribution

You buy lightbulbs that burn out according to an exponential distribution with parameter of $\lambda = 1.8$ lightbulbs per year.

You buy a 10 pack of (independent) light bulbs. What is the probability that your 10-pack lasts at least 5 years?

Let X_i be the time it takes for lightbulb i to burn out.

Let X be the total time. Estimate $\mathbb{P}(X \geq 5)$.

Where's the continuity correction?

There's no correction to make – it was already continuous!!

$$\mathbb{P}(X \geq 5)$$

$$= \mathbb{P}\left(\frac{X-10/1.8}{\sqrt{10/1.8^2}} \geq \frac{5-10/1.8}{\sqrt{10/1.8^2}}\right)$$

$$\approx \mathbb{P}\left(Y \geq \frac{5-10/1.8}{\sqrt{10/1.8^2}}\right) \text{ By CLT}$$

$$\approx \mathbb{P}(Y \geq -0.32)$$

$$= 1 - \Phi(-0.32) = \Phi(0.32)$$

$$\approx .62552$$

True value (needs a distribution not in our zoo) is ≈ 0.58741

Outline of CLT steps

1. Write event you are interested in, in terms of sum of random variables.
2. Apply continuity correction if RVs are discrete.
3. Normalize RV to have mean 0 and standard deviation 1.
4. Replace RV with $\mathcal{N}(0,1)$.
5. Write event in terms of Φ
6. Look up in table.

Polling

Suppose you know that 60% of CSE students support you in your run for SAC. If you draw a sample of 30 students, what is the probability that you don't get a majority of their votes.

How are you sampling?

Method 1: Get a uniformly random subset of size 30.

Method 2: Independently draw 30 people with replacement.

Which do we use?

Polling

Method 1 is what's accurate to what is actually done...

...but we're going to use the math from Method 2.

Why?

Hypergeometric variable formulas are rough, and for increasing population size they're very close to binomial.

And we're going to approximate with the CLT anyway, so...the added inaccuracy isn't a dealbreaker.

If we need other calculations, independence will make any of them easier.

Polling

Let X_i be the indicator for “person i in the sample supports you.”

$\bar{X} = \frac{\sum_{i=1}^n X_i}{30}$ is the fraction who support you.

We’re interested in the event $\mathbb{P}(\bar{X} \leq .5)$.

What is $\mathbb{E}[\bar{X}]$? What is $\text{Var}(\bar{X})$?

$$\mathbb{E}[\bar{X}] = \frac{1}{30} \mathbb{E}[\sum X_i] = \frac{.6 \cdot 30}{30} = \frac{3}{5}.$$

$$\text{Var}(\bar{X}) = \frac{1}{30^2} \text{Var}(\sum X_i) = \frac{1}{30} \cdot .6 \cdot .4 = \frac{1}{125}.$$

Using the CLT

$$\mathbb{P}(\bar{X} \leq .5)$$

$$= \mathbb{P}\left(\frac{\bar{X} - .6}{1/\sqrt{125}} \leq \frac{.5 - .6}{1/\sqrt{125}}\right)$$

$$\approx \mathbb{P}\left(Y \leq \frac{.5 - .6}{1/\sqrt{125}}\right) \text{ where } Y \sim \mathcal{N}(0,1)$$

$$\approx \mathbb{P}(Y \leq -1.12)$$

$$= \Phi(-1.12) = 1 - \Phi(1.12) \approx 1 - 0.86864 = 0.13136$$