

18: Central Limit Theorem

Lisa Yan and Jerry Cain
October 23, 2020

Quick slide reference

3	i.i.d. random variables	18a_iid
9	Central Limit Theorem	18b_clt
19	CLT example	18c_clt_example
24	Sum/average/max of i.i.d. rvs	18d_sums
30	Exercises	LIVE
43	Extra: History of the CLT	18f_clt_history

i.i.d. random variables

Another big day

Up until this point, we've mostly covered traditional probability topics:

- Equally likely outcomes
- Conditional probability, independence, *random variables*
- Joint probability distributions, conditional expectation

We have done some awesome applications:

- Federalist Papers: Authorship identification
- WebMD: General Inference

Today

- Our last big topic in **traditional probability** before we move onto modern-day statistical analysis!



Independence of multiple random variables

We have independence of n **discrete random variables** X_1, X_2, \dots, X_n if for all x_1, x_2, \dots, x_n :

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = \prod_{i=1}^n P(X_i = x_i)$$

$$p_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p_{X_i}(x_i)$$

We have independence of n **continuous random variables** X_1, X_2, \dots, X_n if for all x_1, x_2, \dots, x_n :

$$P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n) = \prod_{i=1}^n P(X_i \leq x_i)$$

$$f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i)$$

i.i.d. random variables

Consider n variables X_1, X_2, \dots, X_n .

X_1, X_2, \dots, X_n are **independent and identically distributed** if

- X_1, X_2, \dots, X_n are independent, and
- All have the same PMF (if discrete) or PDF (if continuous).
 - $\Rightarrow E[X_i] = \mu$ for $i = 1, \dots, n$
 - $\Rightarrow \text{Var}(X_i) = \sigma^2$ for $i = 1, \dots, n$

Same thing:

i.i.d.

iid

IID

Quick check





Are X_1, X_2, \dots, X_n i.i.d. with the following distributions?

1. $X_i \sim \text{Exp}(\lambda)$, X_i independent
2. $X_i \sim \text{Exp}(\lambda_i)$, X_i independent
3. $X_i \sim \text{Exp}(\lambda)$, $X_1 = X_2 = \dots = X_n$
4. $X_i \sim \text{Bin}(n_i, p)$, X_i independent



Quick check

Are X_1, X_2, \dots, X_n i.i.d. with the following distributions?

1. $X_i \sim \text{Exp}(\lambda)$, X_i independent 
2. $X_i \sim \text{Exp}(\lambda_i)$, X_i independent  (unless λ_i equal)
3. $X_i \sim \text{Exp}(\lambda)$, $X_1 = X_2 = \dots = X_n$  dependent: $X_1 = X_2 = \dots = X_n$
4. $X_i \sim \text{Bin}(n_i, p)$, X_i independent  (unless n_i equal)
Note underlying Bernoulli RVs are i.i.d.!
 $Y_j \sim \text{Ber}(p) \quad j=1, \dots, \sum_{i=1}^n n_i$

Central Limit Theorem



(silent drumroll)

Central Limit Theorem

Consider n **independent and identically distributed (i.i.d.)** variables X_1, X_2, \dots, X_n with $E[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$.

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad \text{As } n \rightarrow \infty$$

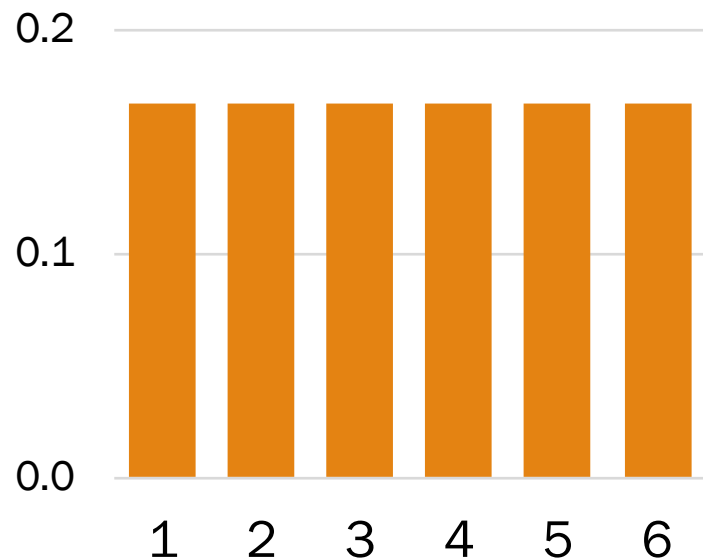
The sum of n **i.i.d.** random variables is normally distributed with mean $n\mu$ and variance $n\sigma^2$.

True happiness



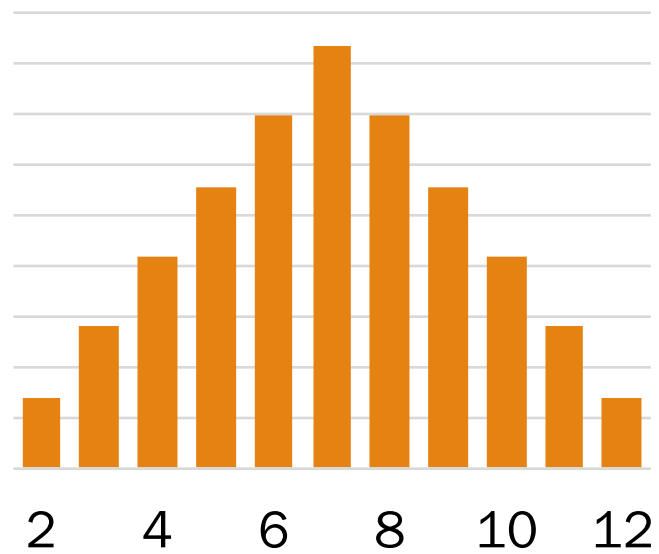
Sum of dice rolls

Roll n independent dice. Let X_i be the outcome of roll i . X_i are i.i.d.



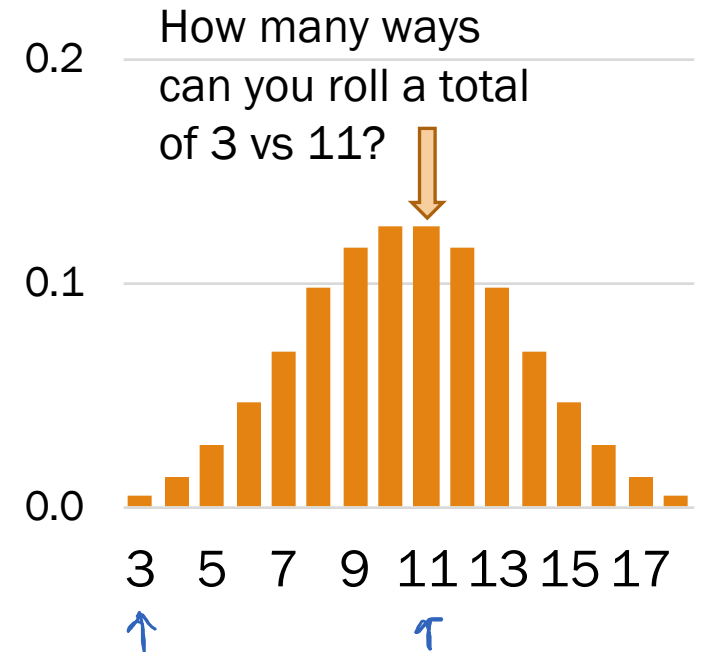
$$\sum_{i=1}^1 X_i$$

Sum of 1 die roll



$$\sum_{i=1}^2 X_i$$

Sum of 2 dice rolls



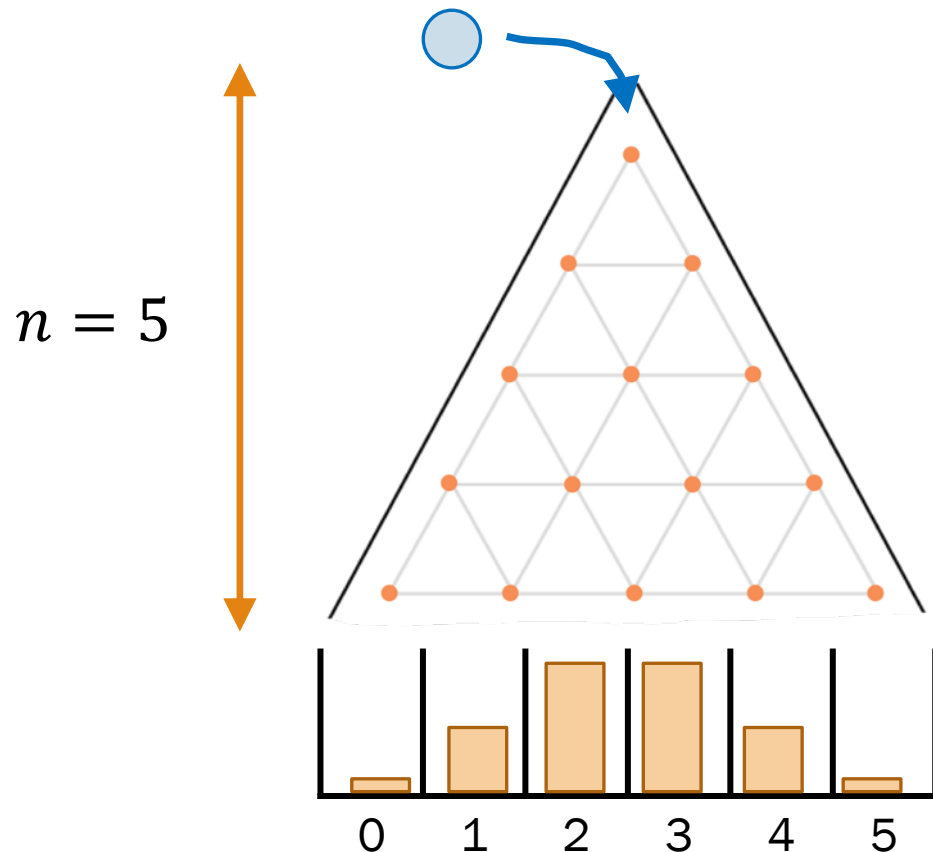
$$\sum_{i=1}^3 X_i$$

Sum of 3 dice rolls

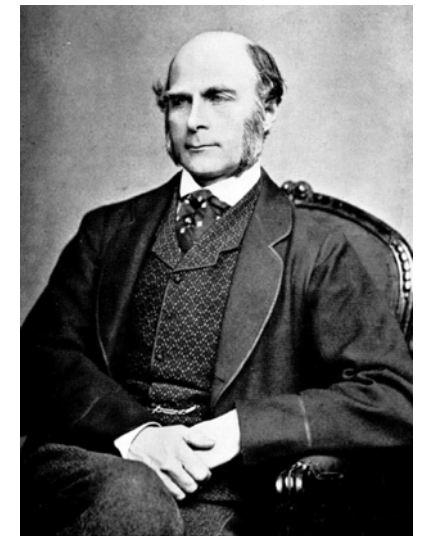
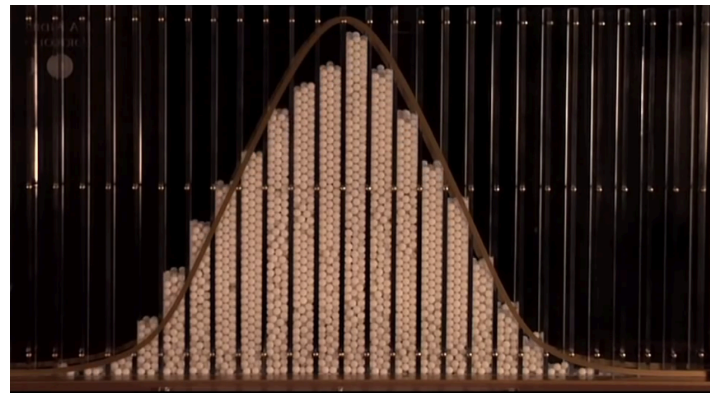
CLT explains a lot

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad \text{As } n \rightarrow \infty$$

The sum of n **i.i.d.** random variables is normally distributed with mean $n\mu$ and variance $n\sigma^2$.



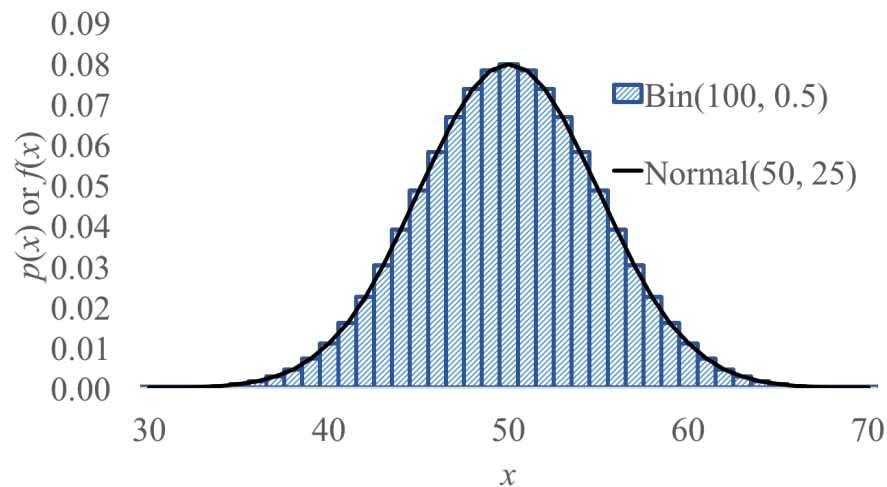
Galton Board, by Sir Francis Galton (1822-1911)



CLT explains a lot

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad \text{As } n \rightarrow \infty$$

The sum of n **i.i.d.** random variables is normally distributed with mean $n\mu$ and variance $n\sigma^2$.



Proof:

Let $X_i \sim \text{Ber}(p)$ for $i = 1, \dots, n$, where X_i are i.i.d.
 $E[X_i] = p$, $\text{Var}(X_i) = p(1 - p)$

$$X = \sum_{i=1}^n X_i \quad (X \sim \text{Bin}(n, p))$$

$$X \sim \mathcal{N}(n\mu, n\sigma^2) \quad (\text{CLT, as } n \rightarrow \infty)$$

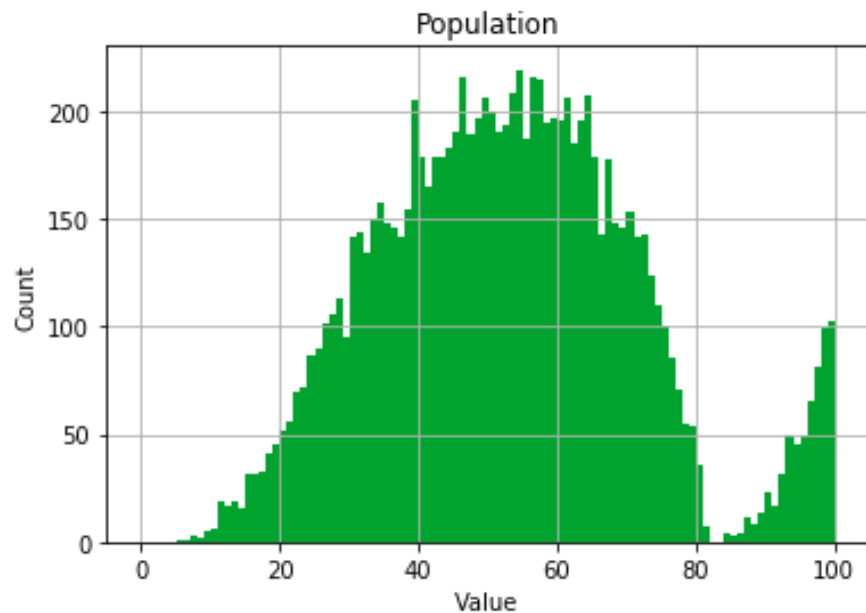
$$X \sim \mathcal{N}(np, np(1 - p)) \quad (\text{substitute mean, variance of Bernoulli})$$

Normal approximation of Binomial
Sum of i.i.d. Bernoulli RVs \approx Normal

CLT explains a lot

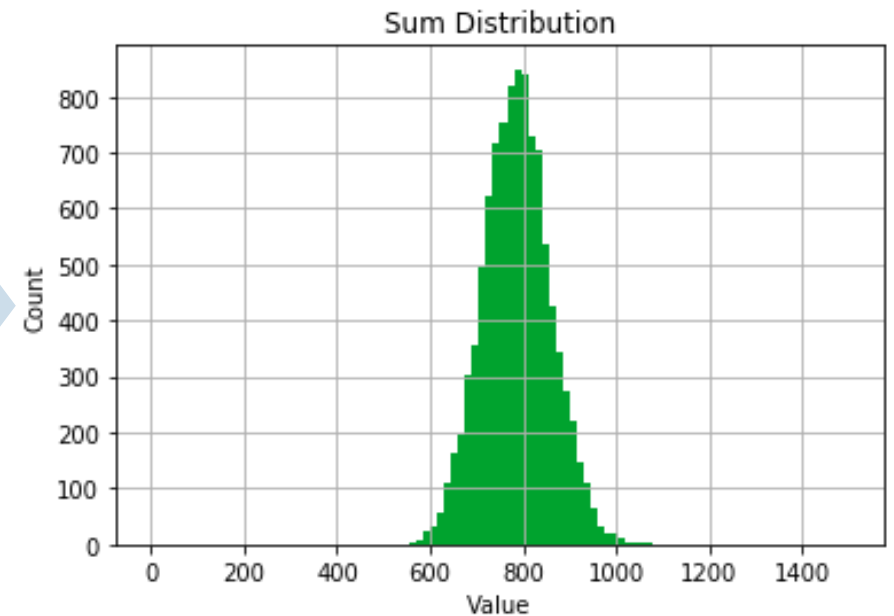
$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad \text{As } n \rightarrow \infty$$

The sum of n **i.i.d.** random variables is normally distributed with mean $n\mu$ and variance $n\sigma^2$.



Distribution of X_i

Sample of
size 15,
sum values

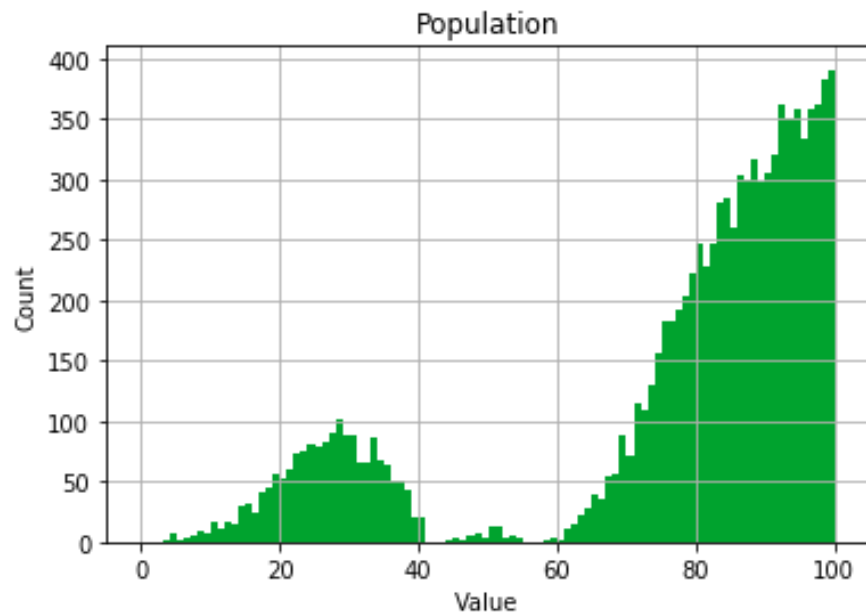


Distribution of $\sum_{i=1}^{15} X_i$

CLT explains a lot

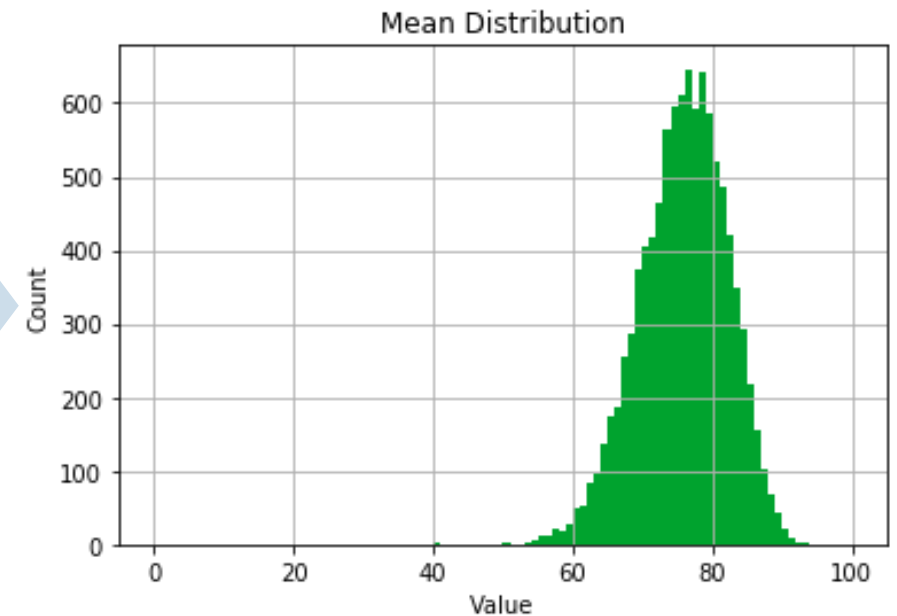
$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad \text{As } n \rightarrow \infty$$

The sum of n **i.i.d.** random variables is normally distributed with mean $n\mu$ and variance $n\sigma^2$.



Distribution of X_i

Sample of
size 15,
average values



Distribution of $\frac{1}{15} \sum_{i=1}^{15} X_i$

Proof of CLT

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad \text{As } n \rightarrow \infty$$

The sum of n **i.i.d.** random variables is normally distributed with mean $n\mu$ and variance $n\sigma^2$.

Proof:

- The Fourier Transform of a PDF is called a **characteristic function**.
- Take the characteristic function of the probability mass of the sample distance from the mean, divided by standard deviation $\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{X_i - \mu}{\sigma}$
- Show that this approaches an exponential function in the limit as $n \rightarrow \infty$: $f(x) = e^{-\frac{x^2}{2}}$
- This function is in turn the characteristic function of the Standard Normal, $Z \sim \mathcal{N}(0,1)$.

(this proof is beyond the scope of CS109)

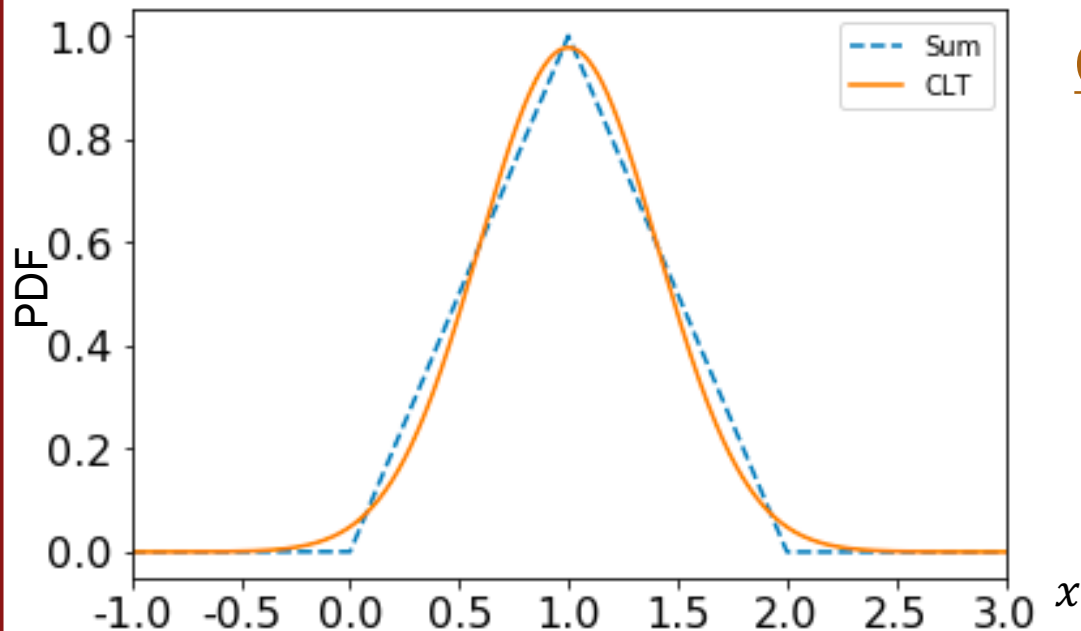
CLT example

Sum of n independent Uniform RVs

Let $X = \sum_{i=1}^n X_i$ be sum of i.i.d. RVs, where $X_i \sim \text{Uni}(0,1)$. $\mu = E[X_i] = 1/2$
 $\sigma^2 = \text{Var}(X_i) = 1/12$

For different n , how close is the CLT approximation of $P(X \leq n/3)$?

$n = 2$:



Exact

$$P(X \leq 2/3) \approx 0.2222$$

CLT approximation

$$X \approx Y \sim \mathcal{N}(n\mu, n\sigma^2) \implies Y \sim \mathcal{N}(1, 1/6)$$

$$P(X \leq 2/3) \approx P(Y \leq 2/3)$$

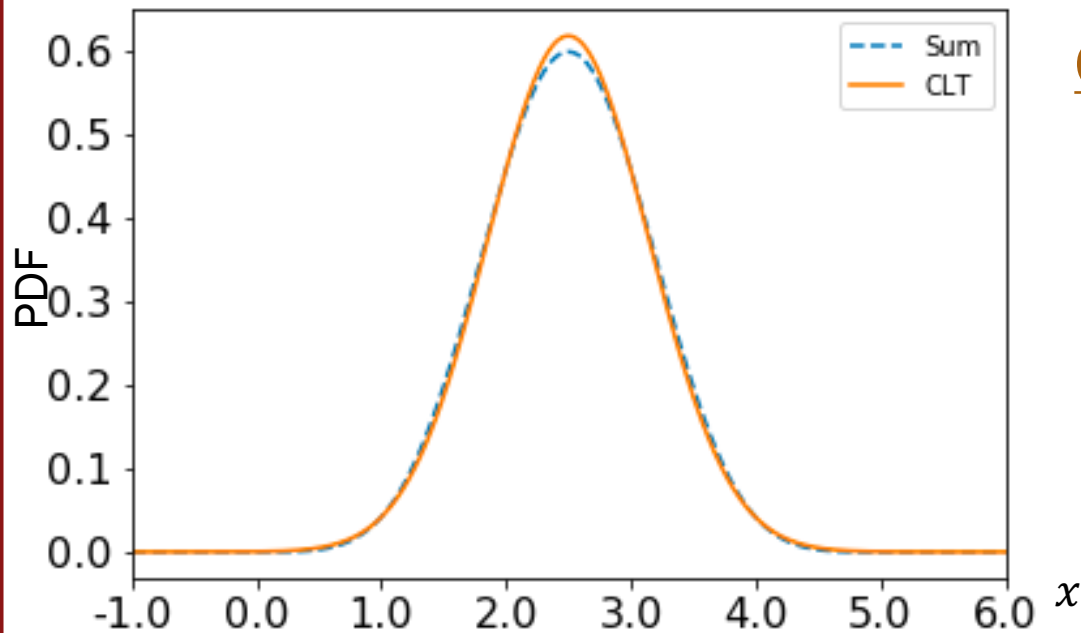
$$= \Phi\left(\frac{2/3 - 1}{\sqrt{1/6}}\right) \approx 0.2071$$

Sum of n independent Uniform RVs

Let $X = \sum_{i=1}^n X_i$ be sum of i.i.d. RVs, where $X_i \sim \text{Uni}(0,1)$. $\mu = E[X_i] = 1/2$
 $\sigma^2 = \text{Var}(X_i) = 1/12$

For different n , how close is the CLT approximation of $P(X \leq n/3)$?

$n = 5$:



Exact

$$P(X \leq 5/3) \approx 0.1017$$

CLT approximation

$$X \approx Y \sim \mathcal{N}(n\mu, n\sigma^2) \implies Y \sim \mathcal{N}(5/2, 5/12)$$

$$P(X \leq 5/3) \approx P(Y \leq 5/3)$$

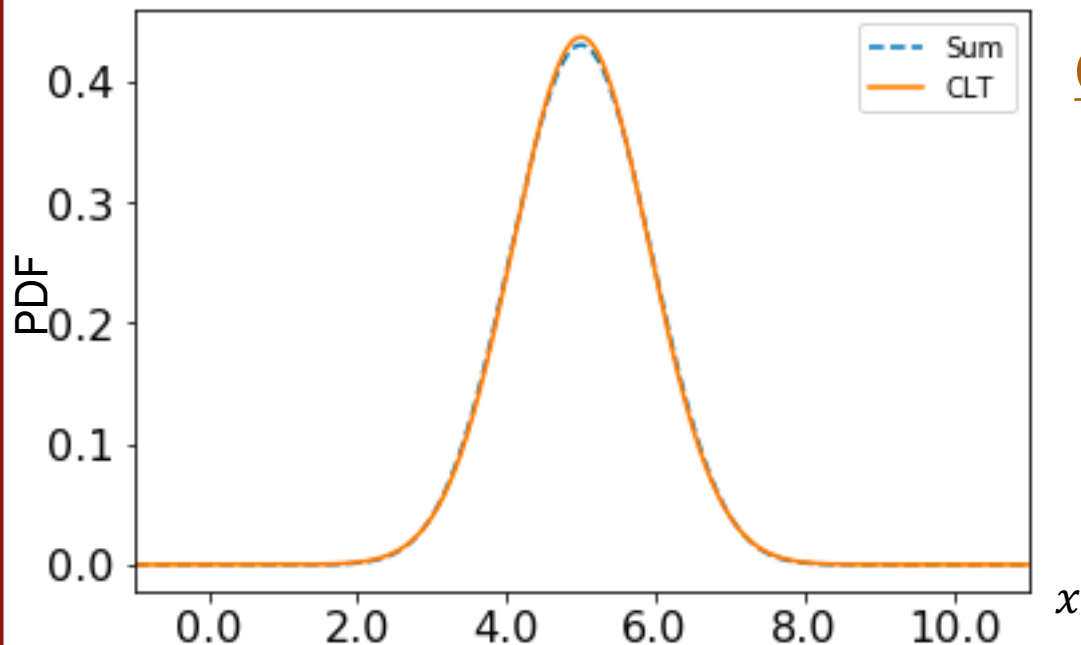
$$= \Phi\left(\frac{5/3 - 5/2}{\sqrt{5/12}}\right) \approx 0.0984$$

Sum of n independent Uniform RVs

Let $X = \sum_{i=1}^n X_i$ be sum of i.i.d. RVs, where $X_i \sim \text{Uni}(0,1)$. $\mu = E[X_i] = 1/2$
 $\sigma^2 = \text{Var}(X_i) = 1/12$

For different n , how close is the CLT approximation of $P(X \leq n/3)$?

$n = 10$:



Exact

$$P(X \leq 10/3) \approx 0.0337$$

CLT approximation

$$X \approx Y \sim \mathcal{N}(n\mu, n\sigma^2) \implies Y \sim \mathcal{N}(5, 5/6)$$

$$P(X \leq 10/3) \approx P(Y \leq 10/3)$$

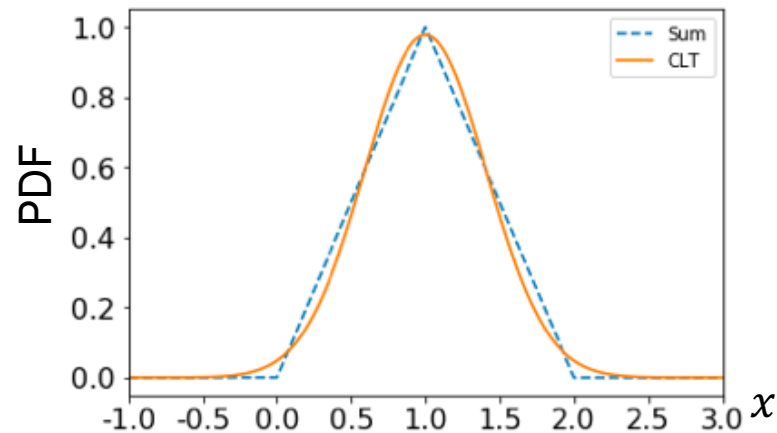
$$= \Phi\left(\frac{10/3 - 5}{\sqrt{5/6}}\right) \approx 0.0339$$

Sum of n independent Uniform RVs

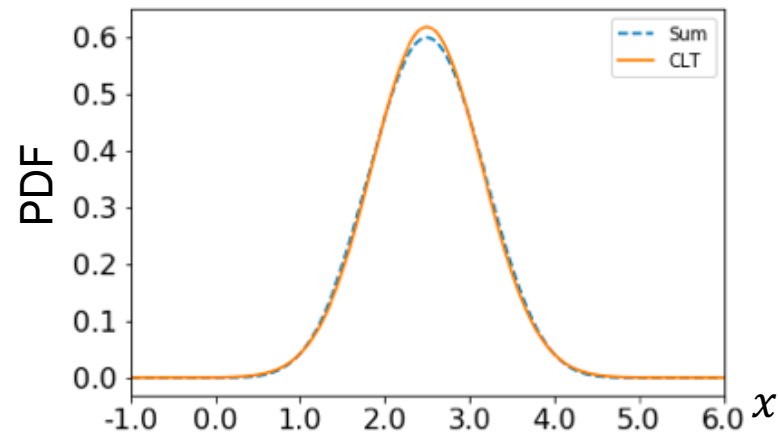
Let $X = \sum_{i=1}^n X_i$ be sum of i.i.d. RVs, where $X_i \sim \text{Uni}(0,1)$. $\mu = E[X_i] = 1/2$
 $\sigma^2 = \text{Var}(X_i) = 1/12$

For different n , how close is the CLT approximation of $P(X \leq n/3)$?

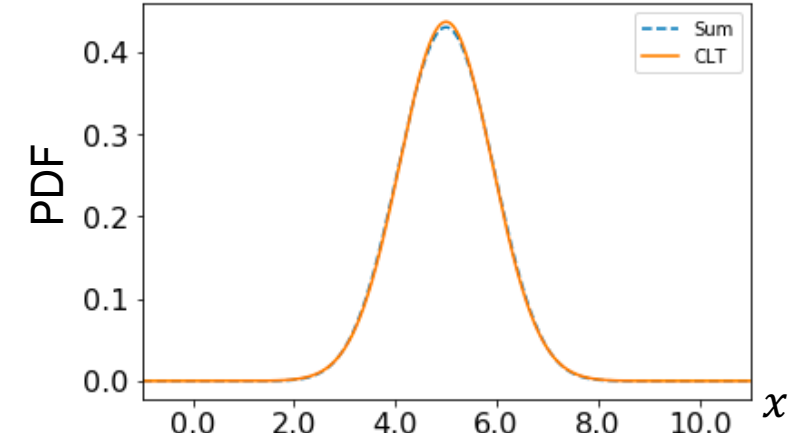
$n = 2$:



$n = 5$:



$n = 10$:



Most books will tell you that CLT holds if $n \geq 30$, but it can hold for smaller n depending on the distribution of your i.i.d. X_i 's.

Sum/average/
max of i.i.d.
random
variables

What about other functions?

Let X_1, X_2, \dots, X_n be i.i.d., where $E[X_i] = \mu$, $\text{Var}(X_i) = \sigma^2$. As $n \rightarrow \infty$:

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

Sum of i.i.d. RVs

?

Average of i.i.d. RVs
(sample mean)

?

Max of i.i.d. RVs

What about other functions?

Let X_1, X_2, \dots, X_n be i.i.d., where $E[X_i] = \mu, \text{Var}(X_i) = \sigma^2$. As $n \rightarrow \infty$:

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

Sum of i.i.d. RVs

?

Average of i.i.d. RVs
(sample mean)

?

Max of i.i.d. RVs

Distribution of sample mean

Let X_1, X_2, \dots, X_n be i.i.d., where $E[X_i] = \mu$, $\text{Var}(X_i) = \sigma^2$. As $n \rightarrow \infty$:

Define: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ (sample mean) $Y = \sum_{i=1}^n X_i$ (sum)

$Y \sim \mathcal{N}(n\mu, n\sigma^2)$ (CLT, as $n \rightarrow \infty$)

$$\bar{X} = \frac{1}{n} Y$$

$\bar{X} \sim \mathcal{N}(\text{?}, \text{?})$ (Linear transform of a Normal)

$$\begin{aligned} \mathbb{E}[\bar{X}] &= \frac{1}{n} \mathbb{E}[Y] = \mu \\ \text{Var}(\bar{X}) &= \left(\frac{1}{n}\right)^2 \text{Var}(Y) = \left(\frac{1}{n}\right)^2 n \sigma^2 = \frac{\sigma^2}{n} \\ \bar{X} &\sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) \end{aligned}$$

Distribution of sample mean

Let X_1, X_2, \dots, X_n be i.i.d., where $E[X_i] = \mu$, $\text{Var}(X_i) = \sigma^2$. As $n \rightarrow \infty$:

Define: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ (sample mean) $Y = \sum_{i=1}^n X_i$ (sum)

$Y \sim \mathcal{N}(n\mu, n\sigma^2)$ (CLT, as $n \rightarrow \infty$)

$$\bar{X} = \frac{1}{n} Y$$

$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ (Linear transform of a Normal)

$$\frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

The average of i.i.d. random variables (i.e., **sample mean**) is normally distributed with mean μ and variance σ^2/n .

Demo: http://onlinestatbook.com/stat_sim/sampling_dist/

What about other functions?

Let X_1, X_2, \dots, X_n be i.i.d., where $E[X_i] = \mu$, $\text{Var}(X_i) = \sigma^2$. As $n \rightarrow \infty$:

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

Sum of i.i.d. RVs

$$\frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Average of i.i.d. RVs
(sample mean)

Gumbel

Max of i.i.d. RVs

(see Fisher-Tippett Gnedenko Theorem)

(live)

18: Central Limit Theorem

Lisa Yan and Jerry Cain
October 23, 2020

Think

Slide 36 has a question to go over by yourself.

Post any clarifications here!

<https://us.edstem.org/courses/2678/discussion/153773>

Think by yourself: 2 min



Quick check

What dimensions are the following RVs?
(Let X_i be i.i.d. with mean μ)

1. X_1

2. (X_1, X_2, \dots, X_n)

3. $\sum_{i=1}^n X_i$

4. $\frac{1}{n} \sum_{i=1}^n X_i$

5. $\frac{1}{n} \sum_{i=1}^n \mu$

- A. 1-D random variable
- B. n -D random variable (a vector)
- C. not a random variable



Quick check

What dimensions are the following RVs?
(Let X_i be i.i.d. with mean μ & variance σ^2)

- A. 1-D random variable
- B. n -D random variable (a vector)
- C. not a random variable

1. X_1 A

2. (X_1, X_2, \dots, X_n) (aka a **sample**) B

3. $\sum_{i=1}^n X_i$ A

$n \rightarrow \infty \rightarrow N(n\mu, n\sigma^2)$ of $(X_1, X_2, \dots, X_n) \rightarrow$ one random value

4. $\frac{1}{n} \sum_{i=1}^n X_i$ A = \bar{X}

(aka the **sample mean**)

$\frac{1}{n}(Y) \rightarrow \mathbb{E}[\frac{1}{n}Y] = \frac{1}{n}\mathbb{E}[Y]$
 $\text{var}(\frac{1}{n}Y) = (\frac{1}{n})^2 \text{var}(Y)$
 $n \rightarrow \infty \rightarrow N(\mu, \frac{\sigma^2}{n})$

5. $\frac{1}{n} \sum_{i=1}^n \mu$ C

Dice game

$$\text{As } n \rightarrow \infty: \sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

You will roll 10 6-sided dice (X_1, X_2, \dots, X_{10}).

- Let $X = X_1 + X_2 + \dots + X_{10}$, the total value of all 10 rolls.
- You win if $X \leq 25$ or $X \geq 45$.



[To the demo!](#)



Breakout Rooms

Check out the question on the next slide (Slide 36). Post any clarifications here!

<https://us.edstem.org/courses/2678/discussion/153773>

Breakout rooms: 3 min



Dice game

$$\text{As } n \rightarrow \infty: \sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

You will roll 10 6-sided dice (X_1, X_2, \dots, X_{10}).

- Let $X = X_1 + X_2 + \dots + X_{10}$, the total value of all 10 rolls.
- You win if $X \leq 25$ or $X \geq 45$.



And now the truth (according to the CLT)...

1. Define RVs and state goal.

$$E[X_i] = 3.5,$$
$$\text{Var}(X_i) = 35/12$$

Want: $X \leq 25 \cup X \geq 45$
 $P(X \leq 25 \text{ or } X \geq 45)$

Approximate:

?

2. Solve.

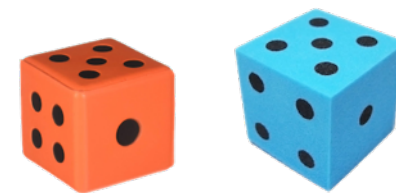


Dice game

$$\text{As } n \rightarrow \infty: \sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

You will roll 10 6-sided dice (X_1, X_2, \dots, X_{10}).

- Let $X = X_1 + X_2 + \dots + X_{10}$, the total value of all 10 rolls.
- You win if $X \leq 25$ or $X \geq 45$.



And now the truth (according to the CLT)...

$$X = \sum_{i=1}^{10} X_i$$

1. Define RVs and state goal.

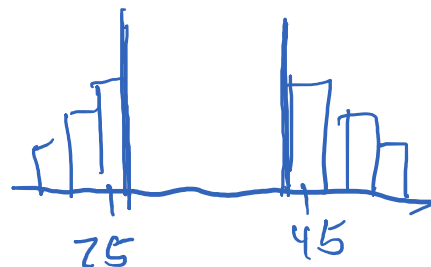
$$E[X_i] = 3.5, \\ \text{Var}(X_i) = 35/12$$

Want: $P(X \leq 25 \text{ or } X \geq 45)$

Approximate:

$$\underbrace{X}_{\text{discrete}} \approx \underbrace{Y}_{\text{continuous}} \sim \mathcal{N}(10(3.5), 10(35/12))$$

2. Solve.



$$P(Y \leq 25.5) + P(Y \geq 44.5) \rightarrow \text{or}$$

$$1 - P(25.5 \leq Y \leq 44.5)$$



continuity
correction

Dice game

$$\text{As } n \rightarrow \infty: \sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

You will roll 10 6-sided dice (X_1, X_2, \dots, X_{10}).

- Let $X = X_1 + X_2 + \dots + X_{10}$, the total value of all 10 rolls.
- You win if $X \leq 25$ or $X \geq 45$.



And now the truth (according to the CLT)...

1. Define RVs and state goal.

$$E[X_i] = 3.5, \\ \text{Var}(X_i) = 35/12$$

Want: $P(X \leq 25 \text{ or } X \geq 45)$

Approximate:

$$X \approx Y \sim \mathcal{N}(10(3.5), 10(35/12))$$

2. Solve.

$$P(Y \leq 25.5) + P(Y \geq 44.5) = \Phi\left(\frac{25.5 - 35}{\sqrt{10(35/12)}}\right) + \left(1 - \Phi\left(\frac{44.5 - 35}{\sqrt{10(35/12)}}\right)\right)$$

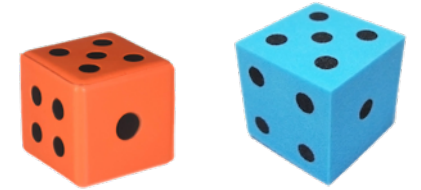
$$\approx \Phi(-1.76) + (1 - \Phi(1.76)) \approx (1 - 0.9608) + (1 - 0.9608) = \mathbf{0.0784}$$

Dice game

$$\text{As } n \rightarrow \infty: \sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

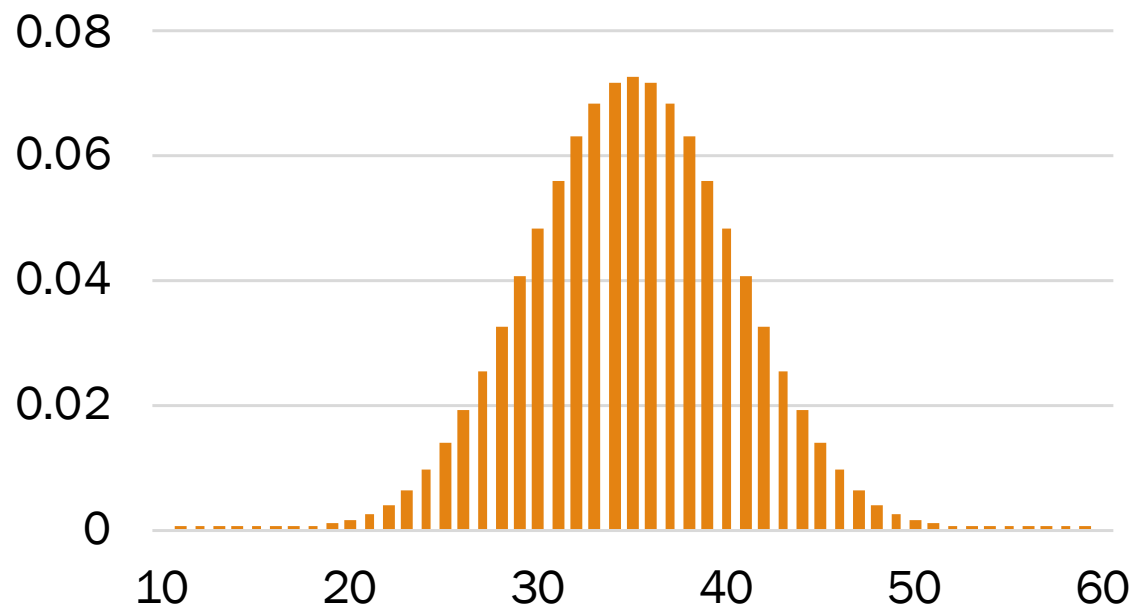
You will roll 10 6-sided dice (X_1, X_2, \dots, X_{10}).

- Let $X = X_1 + X_2 + \dots + X_{10}$, the total value of all 10 rolls.
- You win if $X \leq 25$ or $X \geq 45$.



And now the truth (according to the CLT)...

Check out the [code!](#)



(by CLT)

$$\approx P(Y \leq 25.5) + P(Y \geq 44.5) \approx 0.0786$$

(exact, by computer)

$$P(X \leq 25 \text{ or } X \geq 45) = 0.0780$$

simulated
~~(exact, by computer)~~

$$P(X \leq 25 \text{ or } X \geq 45) \approx 0.0776$$

Summary: Working with the CLT

Let X_1, X_2, \dots, X_n i.i.d., where $E[X_i] = \mu, \text{Var}(X_i) = \sigma^2$. As $n \rightarrow \infty$:

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2)$$

Sum of i.i.d. RVs

$$\frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Average of i.i.d. RVs
(sample mean)



If X_i is discrete:
Use the **continuity correction** on Y !

microwaves



what washes up on
tiny beaches?

Interlude for jokes/announcements

Announcements

Quiz #2 Review session

When:

Monday 10/26 7pm-9pm PT

Recorded:

yes

[Zoom link](#)

Quiz #2 Info and practice:

[Exam page link](#)

Covers PS3, PS4 (i.e., up to and including Lecture 15)

Think

Slide 43 has a question to go over by yourself.

Post any clarifications here!

<https://us.edstem.org/courses/2678/discussion/153773>

Think by yourself: 2 min



Crashing website

- Let X = number of visitors to a website, where $X \sim \text{Poi}(100)$.
- The server crashes if there are ≥ 120 requests/minute.

What is $P(\text{server crashes in next minute})$?

Strategy:

Poisson (exact)

$$P(X \geq 120) = \sum_{k=120}^{\infty} \frac{(100)^k e^{-100}}{k!} \approx 0.0282$$

Strategy:

CLT

(approx.)

How would we involve CLT here?

(Hint: Is there a way to represent X as a sum of i.i.d. RVs?)



Crashing website

- Let X = number of visitors to a website, where $X \sim \text{Poi}(100)$.
- The server crashes if there are ≥ 120 requests/minute.

What is $P(\text{server crashes in next minute})$?

Strategy:

Poisson (exact)

$$P(X \geq 120) = \sum_{k=120}^{\infty} \frac{(100)^k e^{-100}}{k!} \approx 0.0282$$

Strategy:

CLT
(approx.)

State
approx.
goal

sum of indep Poi is Poi

$$\text{Poi}(100) \sim \sum_{i=1}^n \text{Poi}(100/n)$$

X_i

$$X \approx Y \sim \mathcal{N}(n\mu, n\sigma^2)$$

$Y \sim \mathcal{N}(\underbrace{100}, \underbrace{100})$

$n \cdot \frac{100}{n}$ $n \cdot \frac{100}{n}$

$$P(X \geq 120) \approx P(Y \geq 119.5)$$

Check out
the [code!](#)

Solve

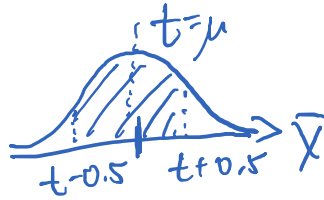
$$P(Y \geq 119.5) = 1 - \Phi\left(\frac{119.5 - 100}{\sqrt{100}}\right) = 1 - \Phi(1.95) \approx 0.0256$$

Clock running time

$$\text{As } n \rightarrow \infty: \frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Want to find the mean (clock) runtime of an algorithm, $\mu = t$ sec.

- Suppose variance of runtime is $\sigma^2 = 4 \text{ sec}^2$.



Run algorithm repeatedly (i.i.d. trials):

- $X_i =$ runtime of i -th run (for $1 \leq i \leq n$)
- Estimate runtime to be **average** of n trials, $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

How many trials do we need s.t. estimated time = $t \pm 0.5$ with **95% certainty**?

1. Define RVs and state goal.

2. Solve.

$$\text{(CLT)} \quad \bar{X} \sim \mathcal{N}\left(t, \frac{4}{n}\right)$$

$$\text{Want: } P(t - 0.5 \leq \bar{X} \leq t + 0.5) = 0.95$$



(linear transform of a normal)

$$\bar{X} - t \sim \mathcal{N}\left(0, \frac{4}{n}\right)$$

$$P(-0.5 \leq \bar{X} - t \leq 0.5) = 0.95$$

Clock running time

$$\text{As } n \rightarrow \infty: \frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Want to find the mean (clock) runtime of an algorithm, $\mu = t$ sec.

- Suppose variance of runtime is $\sigma^2 = 4 \text{ sec}^2$.

Run algorithm repeatedly (i.i.d. trials):

- $X_i =$ runtime of i -th run (for $1 \leq i \leq n$)
- Estimate runtime to be **average** of n trials, \bar{X}

How many trials do we need s.t. estimated time = $t \pm 0.5$ with **95% certainty**?

1. Define RVs and state goal.

$$\bar{X} - t \sim \mathcal{N}\left(0, \frac{4}{n}\right)$$

$$0.95 = P(-0.5 \leq \bar{X} - t \leq 0.5)$$

2. Solve.

$$\begin{aligned} 0.95 &= F_{\bar{X}-t}(0.5) - F_{\bar{X}-t}(-0.5) \\ &= \Phi\left(\frac{0.5 - 0}{\sqrt{4/n}}\right) - \Phi\left(\frac{-0.5 - 0}{\sqrt{4/n}}\right) = 2\Phi\left(\frac{\sqrt{n}}{4}\right) - 1 \end{aligned}$$

Handwritten notes:
Under the second term: $1 - \Phi\left(\frac{0.5}{\sqrt{4/n}}\right)$
To the right: $\frac{0.5}{\sqrt{4/n}} = \frac{1}{2} \cdot \frac{\sqrt{n}}{2}$

Clock running time

$$\text{As } n \rightarrow \infty: \frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Want to find the mean (clock) runtime of an algorithm, $\mu = t$ sec.

- Suppose variance of runtime is $\sigma^2 = 4 \text{ sec}^2$.

Run algorithm repeatedly (i.i.d. trials):

- $X_i =$ runtime of i -th run (for $1 \leq i \leq n$)
- Estimate runtime to be **average** of n trials, \bar{X}

How many trials do we need s.t. estimated time = $t \pm 0.5$ with **95% certainty**?

1. Define RVs and state goal.

$$\bar{X} - t \sim \mathcal{N}\left(0, \frac{4}{n}\right)$$

$$0.95 =$$

$$P(-0.5 \leq \bar{X} - t \leq 0.5)$$

2. Solve.

$$0.95 = F_{\bar{X}-t}(0.5) - F_{\bar{X}-t}(-0.5)$$

$$= \Phi\left(\frac{0.5 - 0}{\sqrt{4/n}}\right) - \Phi\left(\frac{-0.5 - 0}{\sqrt{4/n}}\right) = 2\Phi\left(\frac{\sqrt{n}}{4}\right) - 1$$

$$0.975 = \Phi(\sqrt{n}/4)$$

$$\sqrt{n}/4 = \Phi^{-1}(0.975) \approx 1.96 \quad \Rightarrow \quad n \approx 62$$

Clock running time

$$\text{As } n \rightarrow \infty: \frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

Want to find the mean (clock) runtime of an algorithm, $\mu = t$ sec.

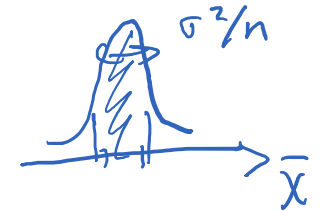
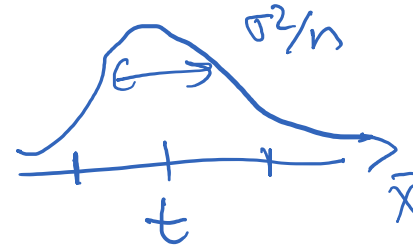
- Suppose variance of runtime is $\sigma^2 = 4 \text{ sec}^2$.

Run algorithm repeatedly (i.i.d. trials):

- $X_i =$ runtime of i -th run (for $1 \leq i \leq n$)
- Estimate runtime to be **average** of n trials, \bar{X}

How many trials do we need s.t. estimated time = $t \pm 0.5$ with **95% certainty**?

$$n \approx 62$$



Interpret: As we increase n (the size of our sample):

- The variance of our sample mean, σ^2/n decreases
- The probability that our sample mean \bar{X} is *close* to the true mean μ increases

Wonderful form of cosmic order

I know of scarcely anything so apt to impress the imagination as the wonderful form of cosmic order expressed by the "[Central limit theorem]". The law would have been personified by the Greeks and deified, if they had known of it. It reigns with serenity and in complete self-effacement, amidst the wildest confusion. The huger the mob, and the greater the apparent anarchy, the more perfect is its sway. It is the supreme law of Unreason.

Whenever a large sample of chaotic elements are taken in hand and marshalled in the order of their magnitude, an unsuspected and most beautiful form of regularity proves to have been latent all along.

– Sir Francis Galton
(of the Galton Board)

Next time

Central Limit Theorem:

- Sample mean $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$
- If we know μ and σ^2 , we can compute probabilities on sample mean \bar{X} of a given sample size n

In real life:

- Yes, the CLT still holds....
- But we **often don't know** μ or σ^2 of our original distribution
- However, we can collect data (a sample of size n)!
- How can we **estimate** the values μ and σ^2 from our sample?

...until next time!

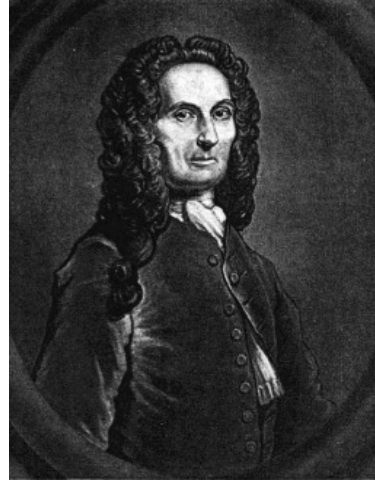
Extra: History of the CLT

Once upon a time...

THE
DOCTRINE
OF
CHANCES:
OR,
A Method of Calculating the Probability
of Events in Play.



By *A. De Moivre*. F. R. S.
L O N D O N:
Printed by *W. Pearson*, for the Author. MDCCLXVIII.



Abraham de Moivre
CLT for $X \sim \text{Ber}(1/2)$
1733



Aubrey Drake Graham
(Drake)

A short history of the CLT

1700



1733: CLT for $X \sim \text{Ber}(1/2)$
postulated by Abraham de Moivre

1800



1823: Pierre-Simon Laplace extends de Moivre's
work to approximating $\text{Bin}(n, p)$ with Normal

1900



1901: Alexandr Lyapunov provides precise
definition and rigorous proof of CLT

2000

1823: Pierre-Simon Laplace extends de Moivre's
work to approximating $\text{Bin}(n, p)$ with Normal

2018: Drake releases *Scorpion*

- It was his 5th studio album, bringing his total # of songs to 190
- Mean quality of subsamples of songs is normally distributed (thanks to the Central Limit Theorem)