

ORIE 4580/5580: Simulation Modeling and Analysis

ORIE 5581: Monte Carlo Simulation

Unit 3: Intro to Monte Carlo Simulation

Sid Banerjee

School of ORIE, Cornell University

expectation and variance of sums of rvs

linearity of expectation

for any rvs X and Y , and any constants $a, b \in \mathbb{R}$

$$\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$$

note: no assumptions! (in particular, does not need independence)

- for general X, Y

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab\text{Cov}(XY)$$

'negatively correlated'
can be negative!

- when X and Y are independent

$$\Rightarrow \text{Var}(X+Y) \text{ can be } \leq \text{Var}(X) + \text{Var}(Y)$$

(useful for variance reduction)

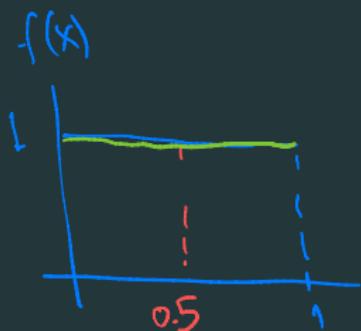
$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y)$$

law of large numbers

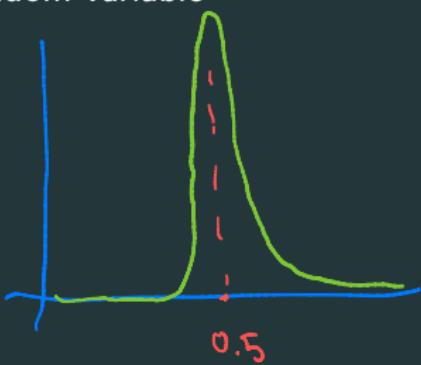
let X_1, X_2, \dots be a sequence of independent rvs with $\mathbb{E}[X_i] = \mu$ for all i
then, “almost” always (almost sure convergence)

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\text{a.s.}} \mu, \quad \text{as } n \rightarrow \infty$$

note: for any finite n , $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ is still a random variable



$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$



central limit theorem

let X_1, X_2, \dots be a sequence of independent rvs with

$$\mathbb{E}[X_i] = \mu, \text{Var}(X_i) = \sigma^2 < \infty \text{ for all } i$$

then,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) = \sqrt{n} \left(\frac{\sum X_i}{n} - \mu \right)$$

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} \sigma \mathcal{N}(0, 1) = \mathcal{N}(0, \sigma^2) \quad , \quad \text{as } n \rightarrow \infty$$

'convergence in distribution'

approximations for large n ,

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{D} \mathcal{N}(\mu, \sigma^2/n)$$

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

QUESTION OF THE DAY

**How can we tell if we are not already in
a computer simulation?**



SingularityNET
community.singularitynet.io

#AGICHAT

clicker question: Staffing a Food Bank

a food bank depends on volunteers for its labor pool

on any given day, the number of workers who show up is $\text{Uniform}(\{1, 2, \dots, 9\})$,

while the number of donations needed to be collected is $\text{Uniform}(\{1, 2, \dots, 29\})$

assuming the work is equally divided among each worker, what is the average load for each worker?

22

- (a) 3

45

- (b) > 3

9

- (c) < 3

24

- (d) can be any of the above on different days

staffing a food bank

a food bank depends on volunteers for its labor pool

on any given day, the number of workers who show up is $\text{Uniform}(\{1, 2, \dots, 9\})$,

while the number of donations needed to be collected is $\text{Uniform}(\{1, 2, \dots, 29\})$

assuming the work is equally divided among each worker, what is the average load for each worker?

$$\mathbb{E}[X] = \frac{1}{9} \sum_{i=1}^9 i = \frac{9 \cdot 10}{2 \cdot 9} = 5, \quad \mathbb{E}[Y] = \frac{1}{29} \cdot \left(\frac{29 \cdot 30}{2} \right) = 15$$

- let X = number of workers, Y = number of donations
- we have $\frac{\mathbb{E}[Y]}{\mathbb{E}[X]} = \frac{30/2}{10/2} = 3$
- on the other hand, $\mathbb{E}[\text{Load}] = \mathbb{E}\left[\frac{Y}{X}\right]$; is this also 3?
- let us simulate and check!

'Assuming X, Y independent'

$$\underbrace{\mathbb{E}[Y]}_{15} \cdot \underbrace{\mathbb{E}\left[\frac{1}{X}\right]}_{?}$$

understanding what happened

for random variables X and Y , and function $g(\cdot, \cdot)$,

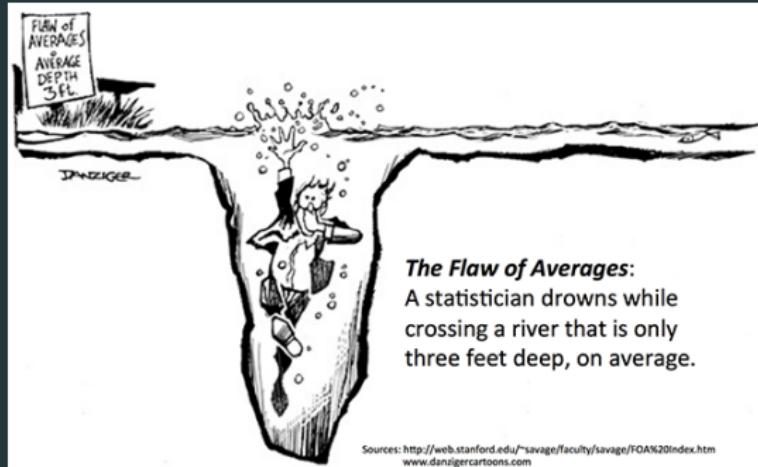
$$\mathbb{E}[g(X, Y)] \neq g(\mathbb{E}[X], \mathbb{E}[Y])!$$

- Equality only if g is linear (LOI)

flaw of averages

moral: in most settings, average inputs don't give average outputs!

$$\mathbb{E}[g(x)] \neq g(\mathbb{E}[x])$$



what can go wrong?

- non-linearities
- correlations between rvs
- 'inspection paradox' (buses take longer to arrive than they should!)
- ...

simulation allows us to avoid such problems!

- Ways of quantifying error

'frequentist' **confidence intervals**

(alternate - Bayesian 'credible' interval)

how many replications?

simulating food bank for many days (**samples/replications**) = distribution of loads
as number of replications increases, **sample average** → **average load** (by LLN)

question: every time we run the simulation model with some fixed number of replications, our estimate of $\mathbb{E}[\text{load}]$ changes.

how 'confident' can we be in our estimate?

- *answer:* use CLT to build a **confidence interval!**

want - true avg \in sim avg $\pm \delta$ with 95% prob

confidence intervals

let X_1, X_2, \dots be a sequence of i.i.d. random variables with mean $\mathbb{E}[X_1] = \mu$ and variance $\text{Var}(X_1) = \sigma^2 < \infty$

want to measure μ from simulations

confidence interval: attempt 1...

fixed

an interval $[a, b]$ is called a 95% confidence interval for $\mathbb{E}[X_1]$ if

$$\mathbb{P}\left[\underbrace{5}_{\text{number}} \leq \underbrace{\mathbb{E}[X_1]}_{\text{number}} \leq \underbrace{10}_{\text{number}}\right] \geq 0.95$$

Q: what is wrong with this?

$$\underbrace{\frac{1}{n} \sum_{i=1}^n X_i}_{\text{Sample mean}} = \bar{X}_n$$

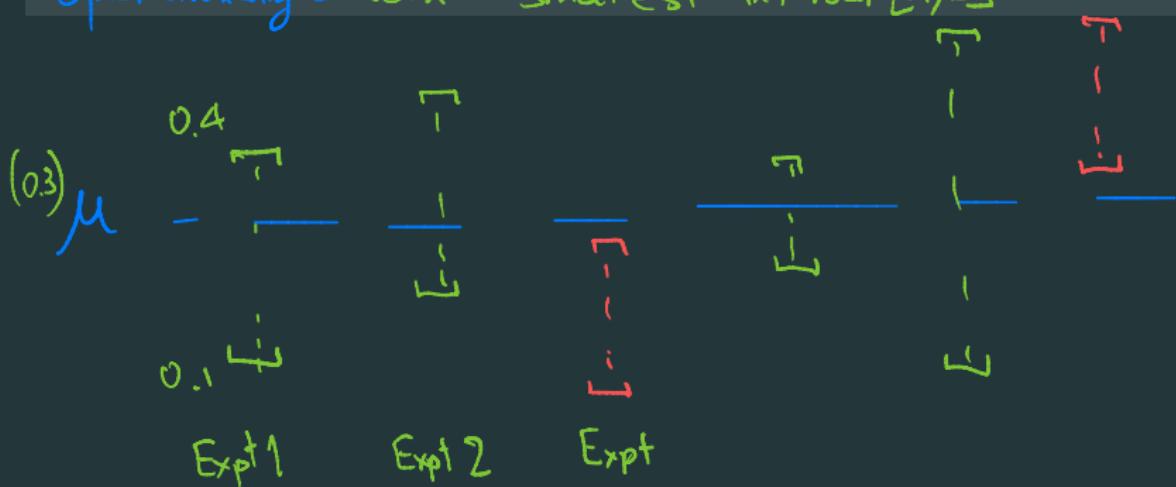
confidence intervals: definition

function of the simulation

an random interval $[A, B]$ (computed from data/experiments) is called a **95%** confidence interval for some (deterministic) quantity μ if

$$\mathbb{P}[A \leq \mu \leq B] \geq 0.95$$

operationally - want smallest interval $[A, B]$



confidence intervals for population mean

X_1, X_2, \dots are i.i.d. rvs with $\mathbb{E}[X_1] = \mu$ and $\text{Var}(X_1) = \sigma^2 < \infty$; $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

Pop mean *Want to measure*

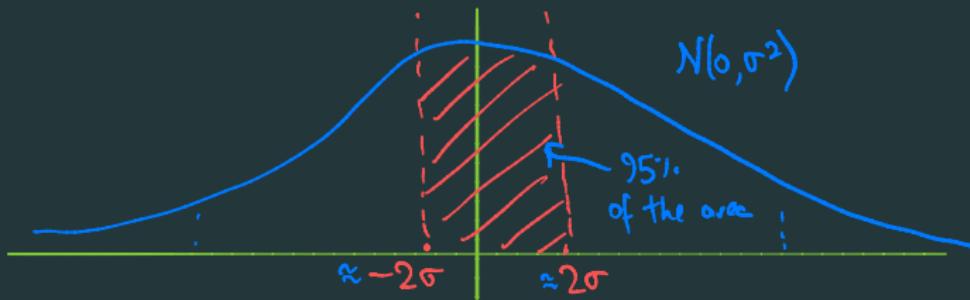
↑
unbiased statistic

- from the central limit theorem:

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} \sigma \mathcal{N}(0, 1) = \mathcal{N}(0, \sigma^2)$$

- from the inverse cdf of $\mathcal{N}(0, 1)$, we can compute

$$\mathbb{P}[-1.96 \leq \mathcal{N}(0, 1) \leq 1.96] \geq 0.95.$$



confidence intervals

want to measure $\mu = \mathbb{E}[X_1]$ from simulations

- from the central limit theorem: $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} \sigma \mathcal{N}(0, 1)$ ↪
- from the cdf of $\mathcal{N}(0, 1)$, we have $\mathbb{P}[-1.96 \leq \mathcal{N}(0, 1) \leq 1.96] \geq 0.95$

putting these together, we have:

$$\boxed{\mathbb{P}\left[\bar{X}_n - \frac{1.96\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + \frac{1.96\sigma}{\sqrt{n}}\right] \geq 0.95}$$

constant (unknown)
random random

confidence intervals: problems

- the confidence interval is approximate because

$\sqrt{n}(\bar{X}_n - \mu)$ is only approximately Gaussian

- the confidence interval is 'exact' when

$\sqrt{n}(\bar{X}_n - \mu)$ is 'exactly' Gaussian
(ie., $X_i \sim N(\mu, \sigma^2)$)

- the confidence interval above requires knowledge of σ^2

Soln 1) use upper bound for σ^2 ($Eg -$ if $X \sim Ber(p)$
 $V_n(x) = p(1-p)$
 $\leq \frac{1}{4}$
 (≤ 1))

confidence intervals: problems

- the confidence interval is approximate because
- the confidence interval is ‘exact’ when
- the confidence interval above requires knowledge of σ^2

can replace σ^2 with its sample estimator

$$s_n^2 = \underbrace{\frac{1}{n-1}}_{\text{correction}} \sum_{i=1}^n (X_i - \underbrace{\bar{X}_n}_{\text{sample mean}})^2.$$

In practice
 $n \approx n-1$

fixed sample-size: recipe for CI

pilot approximate $100(1 - \alpha)\%$ Gaussian CI for $\mathbb{E}X$

- 1. select a sample size N
2. generate N i.i.d. samples X_1, X_2, \dots, X_N of X (replicates)
3. compute the estimators \bar{X}_N, s_N^2

$$\bar{X}_N = \frac{1}{N} \sum_{n=1}^N X_n, \quad s_N^2 = \frac{1}{N-1} \sum_{n=1}^N (X_n - \bar{X}_N)^2$$

4. look up the value of $z_{\alpha/2}$ such that

$$\mathbb{P}[-z_{\alpha/2} \leq N(0, 1) \leq z_{\alpha/2}] = 1 - \alpha$$

For this course
 $\alpha = 0.05$
 $z_{\alpha/2} = 1.96 \approx 2$

5. the approximate $100(1 - \alpha)\%$ CI for $\mathbb{E}X$ is given by

$$\bar{X}_N \mp z_{\alpha/2} \frac{s_N}{\sqrt{N}}$$

selecting the sample size

how large should N be so that the resulting $100(1 - \alpha)\%$ confidence interval will have a pre-specified width?

- CI $\Rightarrow \bar{X}_N \mp z_{\alpha/2} \left(\sigma / \sqrt{N} \right)$
- half-width $\Rightarrow \frac{z_{\alpha/2}}{2} \frac{\sigma}{\sqrt{N}} = \ell$
- ℓ be the desired half-width
- Set $N = \frac{4 \sigma^2}{\ell^2} \left(\frac{(z_{\alpha/2})^2 \sigma^2}{\ell^2} \right)$

selecting the sample size

- problem: σ^2 is unknown!
- estimating σ^2 through s_N^2 requires simulation!
- solution: ‘pilot runs’

perform k simulation runs to get $[X'_n : n = 1, \dots, k]$ as outcomes
compute

$$\tilde{X}_k = \frac{1}{k} \sum_{n=1}^k X'_n, \quad \tilde{s}_k^2 = \underbrace{\frac{1}{k-1}}_{\text{correction}} \sum_{n=1}^k \underbrace{\left(X'_n - \tilde{X}_k \right)^2}_{\text{empirical variance}}$$

use \tilde{s}_k^2 to estimate σ^2

for confidence level α , half-width ℓ , set

$$N = \left\lceil \frac{z_{\alpha/2}^2 \tilde{s}_k^2}{\ell^2} \right\rceil$$

clicker question

we want to estimate $\mathbb{E}[X]$ with ∓ 0.01 accuracy with 95% confidence from 10 replications, $\tilde{s}_{10}^2 = 0.5$

how many replications should we use?

- A** 200
- B** 2,000
- C** 20,000
- D** 40,000
- E** 200,000

clicker question

we want to estimate $\mathbb{E}[X]$ with ∓ 0.01 accuracy with 95% confidence from 10 replications, $\tilde{s}_{10}^2 = 0.5$

thus we have

$$\left\lceil \frac{z_{\alpha/2}^2 \tilde{s}_k^2}{\ell^2} \right\rceil = \frac{2^2 \times 0.5}{0.01^2} = 20,000$$

how many replications should we use?

- A 200
- B 2,000
- C 20,000
- D 40,000
- E 200,000

basic simulation workflow

X_i

$E[X_i] = \mu$ ← what we want to compute

Typical Sim : Choose 'statistic' which is an 'unbiased estimator'

- perform **pilot run** of k simulations (sufficient but not large k) — compute \tilde{s}_k^2
- compute required sample-size N for desired confidence interval
- run N additional simulations \Rightarrow production runs
- form fixed-sample confidence intervals from these N samples
- **note:** final CI may be different than desired, because it is constructed by using s_N^2 (may be larger/smaller than \tilde{s}_k^2)
- for the final confidence interval, discard the information from the trial runs not a problem, since $N \gg k$ usually

clicker question: abusing CI (?)

in a homework assignment in Simulation, students were asked to do a Monte Carlo simulation to compute π up to 2 decimal places

there were 100 homework submissions, and *all* submissions reported 95% CIs which included 3.14

the probability that this happened 'by chance' is approximately

5

(a) 0.4

$$P[3.14 \in \text{CI}] \approx 0.95$$

44

(b) 0.05

$$\Rightarrow P[\text{All students have } 3.14 \in \text{CI}] \approx (0.95)^{100}$$

20

(c) 0.006

Reasons -

- Everyone chooses same 'seed'
- Everyone has 'prior knowledge'

8

(d) 0.0007

23

(e) 0

confidence intervals as a social contract

a random interval $[A, B]$ (computed from data/experiments) is a 95% confidence interval for some unknown μ if **before the experiment is done**

$$\mathbb{P}[A \leq \mu \leq B] \geq 0.95$$

Good practices →

- Announce questions before hand
- Design exp't, let someone else perform it

clicker question: (mis)interpreting CI

a 95% confidence interval for the mean annual rainfall in Milford Sound, New Zealand is 26 feet plus or minus 1 foot
therefore, the rainfall next year will be between 25 feet and 27 feet with 95% probability.

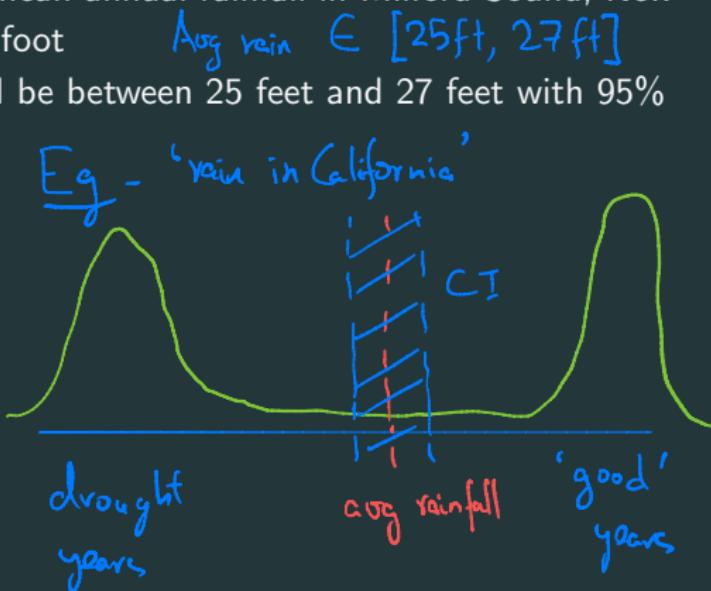
32:1.

(a) true

63:1.

(b) false

(c) where is New Zealand?

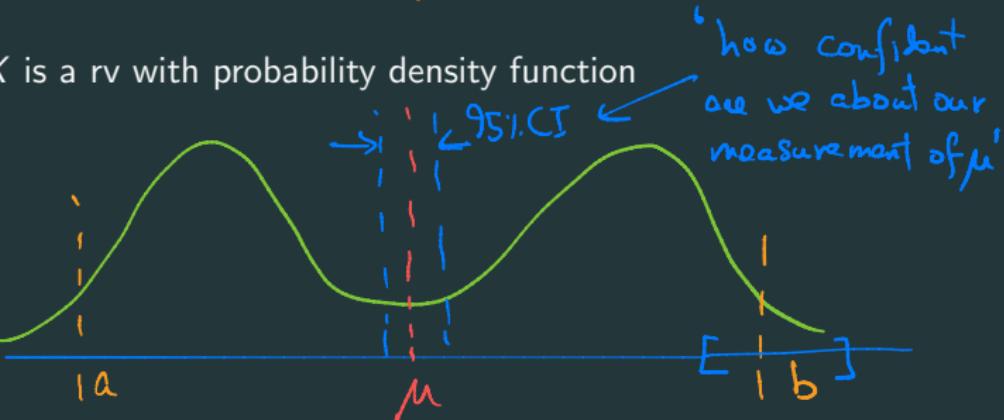


confidence intervals vs quantiles

the CI for the mean is NOT the same as the quantiles of a random variable.

- suppose that X is a rv with probability density function

$[a, b]$ is a
95% quantile
if $X \in [a, b]$ w.p 95%



- we can select q_1 and q_2 so that

$$\mathbb{P}[q_1 \leq X \leq q_2] = 0.95,$$

$b \equiv 95\%$ quantile
 $\in [b-\delta, b+\delta]$

but $[q_1, q_2]$ is not a 95% confidence interval for $\mathbb{E}X$.

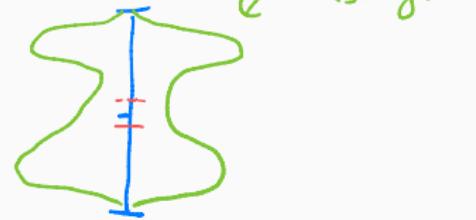


Useful visualizations

SMORE plot



Violin plot



advanced probability tools

tool 1: Jensen's inequality

$$\text{Recall - LoI} - \mathbb{E}[g(x)] = g(\mathbb{E}[x])$$

if g is linear

Q: can we say something about $\mathbb{E}[f(X)]$ vs $f(\mathbb{E}[X])$ (in particular, \geq or \leq) without simulating?

Jensen's inequality

if X is a random variable and f is a convex function, then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$



Example – in our food-bank staffing problem, since $f(x) = \frac{1}{x}$, $x > 0$ is convex:

$$\mathbb{E}[Y]\mathbb{E}\left[\frac{1}{X}\right] \geq \frac{\mathbb{E}[Y]}{\mathbb{E}[X]}$$

$\Rightarrow \mathbb{E}\left[\frac{1}{X}\right] \geq \frac{1}{\mathbb{E}[X]}$ if $X > 0$

$$\mathbb{E}[y] - \mathbb{E}[x^2] \geq (\mathbb{E}[x])^2$$

- If g is concave (\cap), then $\mathbb{E}[g(x)] \leq g(\mathbb{E}[x])$

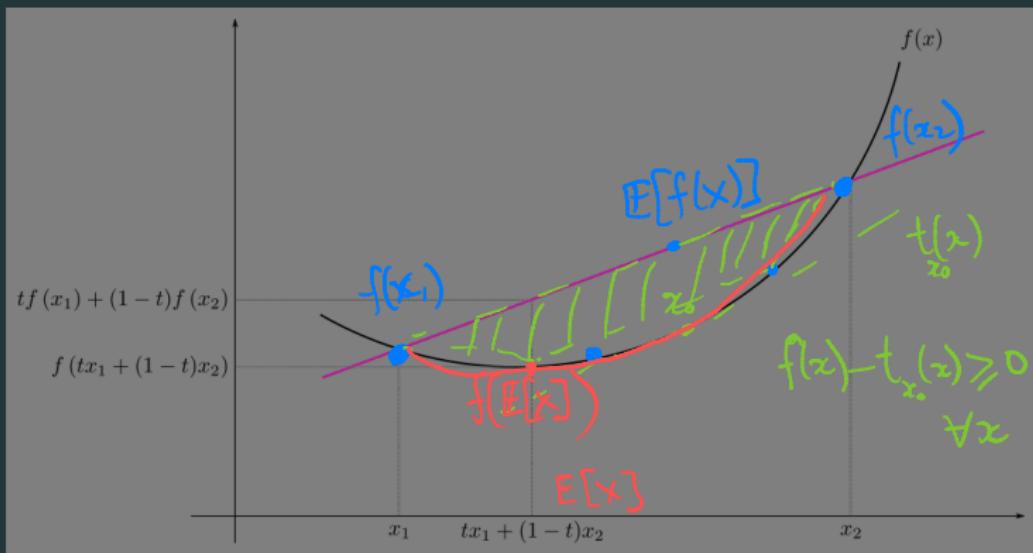
tool 1: Jensen's inequality

Jensen's inequality

if X is a random variable and f is a **convex function**, then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

proof sketch (plus way to remember)



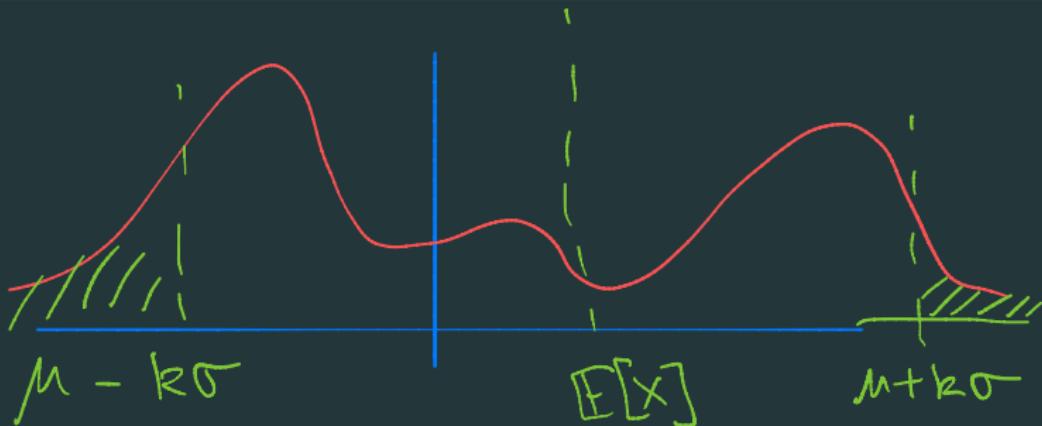
tool 2: Chebyshev's Inequality

Q: since the CLT convergence is faster/slower for different rvs, can we be sure that CIs based on variance always make sense?

Chebyshev's inequality

let X be any rv with finite mean μ and finite variance $\sigma^2 > 0$
then for any $k > 0$,

$$\mathbb{P}[|X - \mu| \geq k\sigma] \leq \frac{1}{k^2}$$



'always-valid' confidence intervals

Chebyshev's inequality

let X be any rv with finite mean μ and finite variance $\sigma^2 > 0$
then for any $k > 0$,

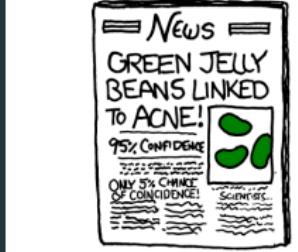
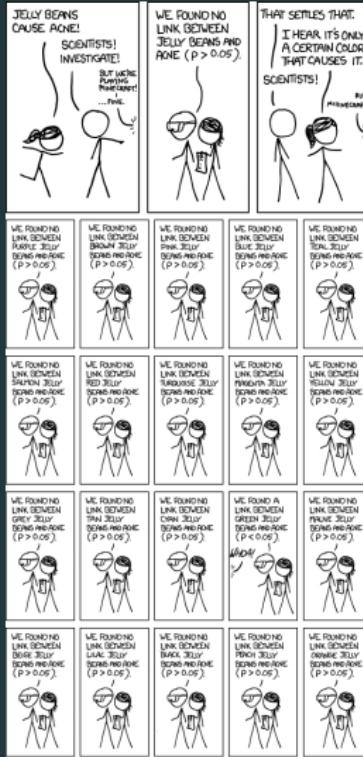
$$\mathbb{P}[|X - \mu| \geq k\sigma] \leq \frac{1}{k^2} = \frac{1}{20}$$

worst-case CI: if we choose $k = 2$, then we **always** have 75% confidence intervals

For 95%, choose k s.t. $\frac{1}{k^2} = 0.05 = \frac{1}{20}$

$$\Rightarrow k = \sqrt{20} \leq 5$$

(i.e., $\mu \in [\bar{x}_n - \frac{5\sigma}{\sqrt{n}}, \bar{x}_n + \frac{5\sigma}{\sqrt{n}}]$ w.p. > 95%)
always!



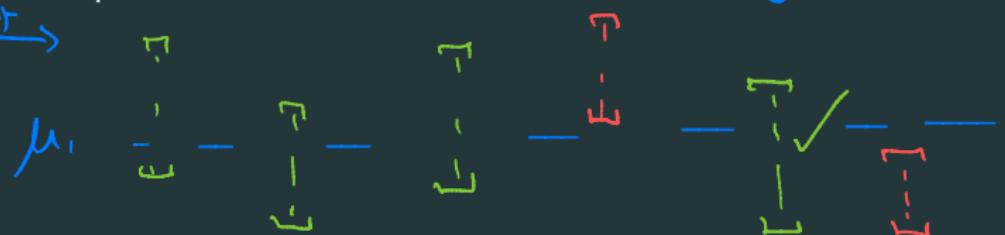
tool 3: the union bound

Q: how can we get simultaneous confidence intervals for multiple hypothesis?

- Eg. I give you five 95% confidence intervals; do they simultaneously contain their respective means 95% of the time?



exp →



tool 3: the union bound

the union bound

let A_1, A_2, \dots, A_k be events; then

$$P(A_1 \cap A_2 \cap \dots \cap A_k) = P(\overline{A_1 \cup A_2 \cup \dots \cup A_k}) \geq 1 - (P(A_1^c) + P(A_2^c) + \dots + P(A_k^c))$$

All A_i are true

let A_i = event that the i th CI contains its true mean... α CI for each

$$\begin{aligned} & P[\text{All } A_i \text{ hold simultaneously}] \geq 1 - k\alpha \\ \Rightarrow & (1-\alpha) \text{ CI for each } A_i \rightarrow (1-k\alpha) \text{ CI for all simultaneously} \end{aligned}$$