

## class poll: (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 (95%). CI ← CI for my mean estimate  
therefore, the rainfall next year will be between 25 feet and 27 feet with 95% probability.

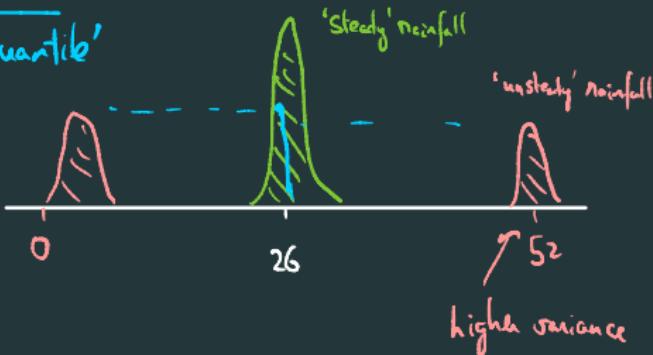
what you want is a 'Quantile'

(a) true

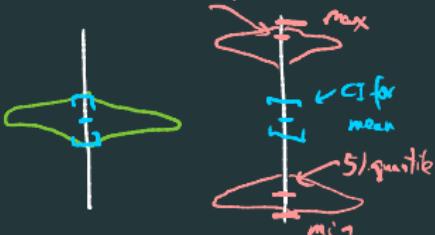
✓ (b) false

(c) where is New Zealand?

Better way of representing data - Violin plot



both these examples have the same mean!



Last class - 95% CI from data  $(X_1, X_2, \dots, X_n)$

$$\approx \left[ \bar{X}_n - 1.96 \sqrt{\text{Var}(\bar{X}_n)}, \bar{X}_n + \underbrace{1.96}_{\approx 2} \sqrt{\text{Var}(\bar{X}_n)} \right]$$

from the Gaussian CLT

$$\frac{\text{std-dev}(X_i)}{\sqrt{n}} \approx \frac{1}{\sqrt{n}} \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2}$$

- Use 'pilot runs' to estimate  $\text{Var}(X_i)$ , decide # of replicates n empirical std-dev

## three advanced probability tools

---

Today

- 3 advanced tools - 'better' understand CIs
- Unit 5 - How to generate complex random variables

1) How can we 'guess'  $\underbrace{\mathbb{E}[f(x)]}_{\text{may be unknown}}$  using  $\underbrace{\mathbb{E}[x]}_{\text{known}}$ ?

- Jensen's Inequality (when is  $\mathbb{E}[f(x)]$  'related to'  $f(\mathbb{E}[x])$ )

2) How we deal with multiple CI (how to combine them)?

Eg- Testing multiple hypothesis  $\equiv$  Union bound

3) Can we get valid CIs if we do not believe in the CLT?

$$\text{More generally } \sigma^2 = \text{Va}(\bar{x}_n)$$

$$P[|\bar{x}_n - \mu| > c\sigma] \leq \frac{1}{c^2}$$

$$CI \equiv \bar{x}_n \pm \underline{\frac{1.96}{\sqrt{\text{Va}(\bar{x}_n)}}}$$

Gaussian  
approximation

$\Rightarrow$  Can always use CIs of the form  $\bar{x}_n \pm \frac{c}{\sqrt{\text{Va}(\bar{x}_n)}}$   
 $\approx 4 \dots$  for 95%.

## staffing a food bank (example from 2 classes back)

a food bank depends on volunteers for its labor pool

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

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

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

i.e., Number of donations they need to pick up

- let  $X = \text{number of workers}$ ,  $Y = \text{number of donations}$

- 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!

$$\mathbb{E}[X] = \frac{1}{9} \sum_{i=1}^9 i = \frac{9.10}{9.2}$$

$$\mathbb{E}[Y] = \frac{1}{29} \sum_{i=1}^{29} i = \frac{30}{2}$$

found  $\mathbb{E}\left[\frac{Y}{X}\right] > 3 = \frac{\mathbb{E}[Y]}{\mathbb{E}[X]}$

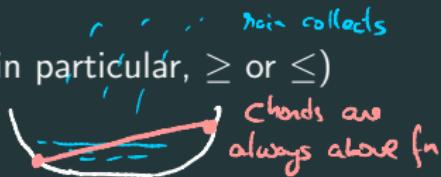
Load

# tool 1: Jensen's inequality

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

•  $f(x)$  is convex in  $x$  if

$$(\frac{d^2}{dx^2} f) \geq 0$$



## Jensen's inequality

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

$g$  is a concave fn



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

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

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

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

convex fn

$$f(0)$$

$$f(1)$$

$$X \sim \text{Beta}(0.5)$$

$$\mathbb{E}[X]$$

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

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

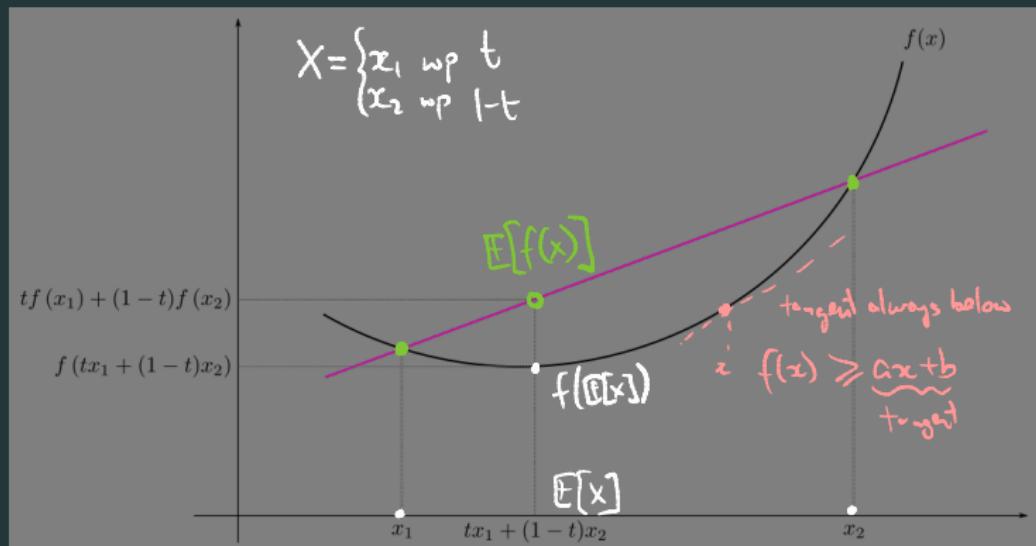
# tool 1: Jensen's inequality

## Jensen's inequality

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

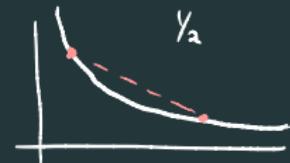
$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

proof sketch (plus way to remember)



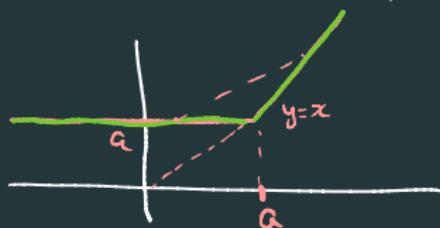
Convex fns

-  $f(x) = x^2$



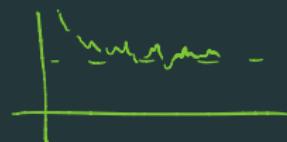
$f(x) = \frac{1}{x}, x > 0$

$f(x) = \max\{x, a\}$



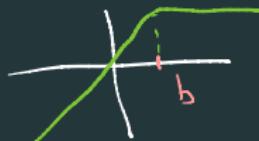
Eg -  $\mathbb{E}[X^2] \geq (\mathbb{E}[X])^2 \quad (\Rightarrow \text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \geq 0)$

$\mathbb{E}[\frac{1}{X}] \geq \frac{1}{\mathbb{E}[X]} \leftarrow \text{showed up for Buffon's needle}$



$\mathbb{E}[\max(X, a)] \geq \max(\mathbb{E}[X], a)$

$\mathbb{E}[\min(X, b)] \leq \min(\mathbb{E}[X], b)$



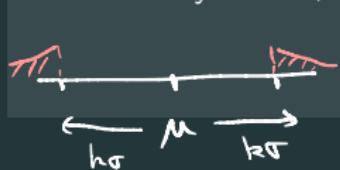
## tool 2: Chebyshev's Inequality (similar to linearity of indep in that NO FEW ASSUMPTIONS)

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

What we need for CLT

let  $X$  be any rv with finite mean  $\mu$  and finite variance  $\sigma^2 > 0$  ie  $|E[X] - \mu| < \epsilon$   
then for any  $k > 0$ ,  $V_n(X) = \sigma^2 < \epsilon$



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

absolute deviation from mean

- Gaussian CI  $\equiv \bar{X}_n \pm 1.95 \sqrt{V_n(\bar{X}_n)}$  (for 95%)

Always valid CI  $\equiv \bar{X}_n \pm 5 \sqrt{V_n(\bar{X}_n)}$  for 95% CI

- Is Chebyshev tight? YES if  $X$  is Bernoulli (extreme values)

Why is this useful

$$- (1-\alpha)\text{-CI} \equiv P\left[\mu \in [A, B] \right] \geq 1-\alpha$$

random interval

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

$$\Rightarrow P\left[|\bar{X}_n - \overset{\text{true mean}}{\hat{\mu}}| \geq k\sqrt{V_n(\bar{X}_n)}\right] \quad \text{ALWAYS TRUE}$$

$$= P\left[\mu \notin \left[\bar{X}_n - k\sqrt{V_n(\bar{X}_n)}, \bar{X}_n + k\sqrt{V_n(\bar{X}_n)}\right]\right] \leq \frac{1}{k^2}$$

at most  $\alpha$

For  $(1-\alpha)\text{-CI}$ , can choose  $k \geq \frac{1}{\sqrt{\alpha}}$

$$\text{Eg- } 75\% \text{-CI} \Rightarrow P[\dots] \geq 0.75 \Rightarrow k = \frac{1}{\sqrt{0.25}} = 2$$

vs 95%  
CI for Gaussian CI

## 'always-valid' confidence intervals

### Chebyshev's inequality

let  $X$  be any rv with finite mean  $\mu$  and finite variance  $\sigma^2 > 0$

then for any  $k > 0$ ,

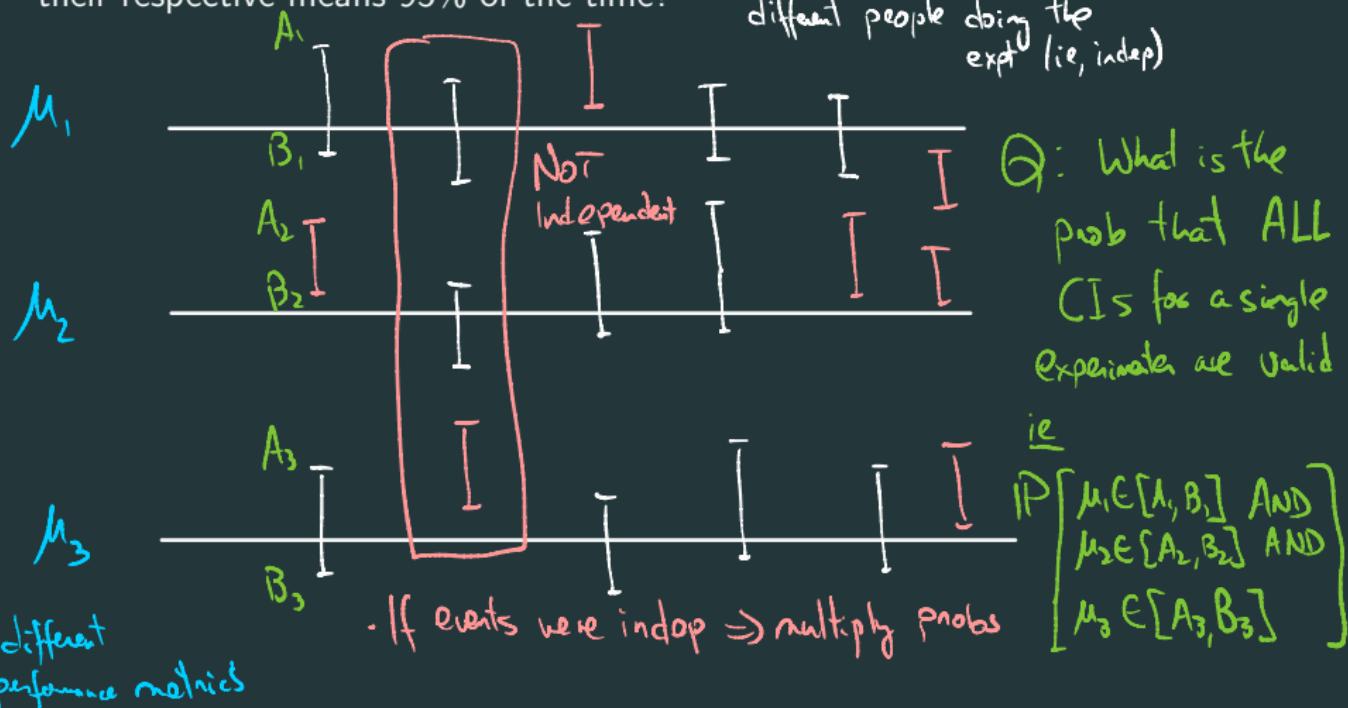
$$\mathbb{P}[|X - \mu| \geq k\sigma] \leq \frac{1}{k^2}$$

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

### tool 3: the union bound

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

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



### tool 3: the union bound

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

the union bound  $\leq \Pr[A_1 \text{ OR } A_2 \text{ OR } A_3] \leq \Pr[A_1] + \Pr[A_2] + \Pr[A_3]$

let  $A_1, A_2, \dots, A_k$  be events; then  $\Pr[A_1 \cup A_2 \cup \dots \cup A_k] \leftarrow \text{de Morgan's Law}$

$$\Pr[A_1 \cap A_2 \cap \dots \cap A_k] = 1 - \Pr[A_1^c \cup A_2^c \cup \dots \cup A_k^c] \quad (A \cup B)^c = A^c \cap B^c$$

AND  $\geq 1 - (\Pr[A_1^c] + \Pr[A_2^c] + \dots + \Pr[A_k^c])$

$\Pr[\text{CI is BAD}]$

let  $A_i = \text{event that the } i\text{th CI contains its true mean...}$   $\left(\begin{array}{l} \text{E.g. 95% CI for} \\ \mu_1, \mu_2 \text{ and } \mu_3 \end{array}\right)$

$$\Pr \left[ \begin{array}{l} \mu_1 \in [A_1, B_1] \text{ AND} \\ \mu_2 \in [A_2, B_2] \text{ AND} \\ \mu_3 \in [A_3, B_3] \end{array} \right] \geq 1 - (\alpha + \alpha + \alpha) = 1 - 3\alpha$$

'bad probabilities add up'

Best case scenario  $\equiv$  If all were independent,  $(1-\alpha)^3 \geq 1 - 3\alpha$

## 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
- we can select  $q_1$  and  $q_2$  so that

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

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