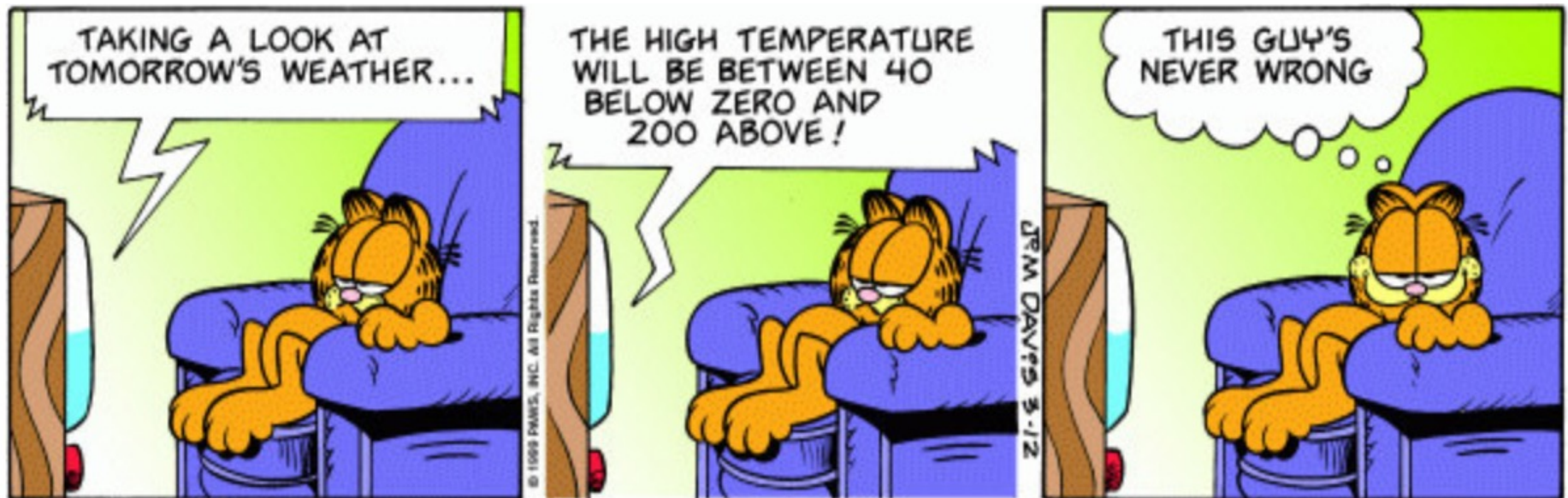# Stat 88: Prob. & Mathematical Statistics in Data Science



Lecture 31: 4/10/2024

Section 9.1, 9.2

Confidence intervals

Goal : Estimate a population parameter.

Using $\bar{X}$ to estimate $\mu$     $X_1, X_2, \; -- \; X_n \; iid : \mu, \sigma^2$

- $\bar{X}$ is an unbiased estimator of $\mu$ (what does that mean?)   $E(X_k) = \mu$
- If we also know that each of the $X_k$ had SD $\sigma$, what can we say about $SD(\bar{X})$?

$$SD(\bar{X}) = \sigma/\sqrt{n}$$

- What does the Central Limit theorem say about the sample mean?

For n large enough $\bar{X} \approx$

- We will use the CLT and the sample mean to define a random interval (why is it random?) that will *cover* the true mean with a specified probability, say 95%

- Based on data from a *random sample*, we will construct an interval of estimates for some unknown (but fixed) population parameter.

$\approx$ "approximately distributed as"

$$X \sim Bin(n = 100, p = 0.5)$$     $\sigma = \sqrt{npq} = \sqrt{100 \times \frac{1}{2} \times \frac{1}{2}}$

$$X \approx N(\mu = 50, \sigma = 5) \quad \leftarrow \quad example$$     $= 5$

Given an iid sample $X_1, X_2, \dots, X_n$

By CLT, $\bar{X} = \dfrac{X_1 + X_2 + \dots + X_n}{n}$ is

approx normal with $E(\bar{X}) = \mu$, $Var(\bar{X}) = \dfrac{\sigma^2}{n}$

$Z = \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}} \approx N(0,1)$ (by CLT)

$$P\left(-2 \leq Z \leq 2\right) \approx 0.95$$

$$P\left(|Z| \leq 2\right) = 0.95$$

$\uparrow$ dist b/w $Z$ & $0 \leq 2$



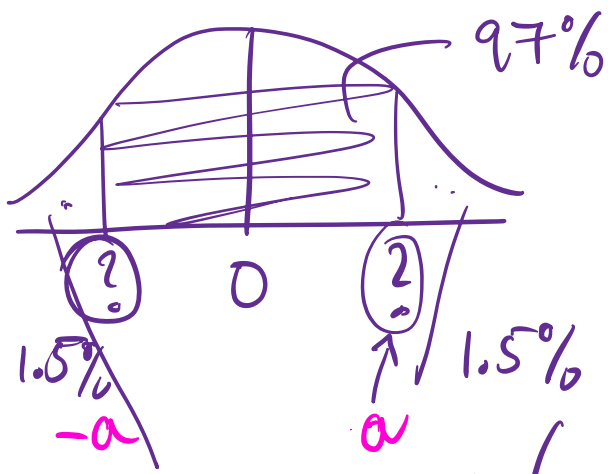$-2 \quad 0 \quad Z \quad 2$

$$Z = \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

$$P\left(-2 \leq Z \leq 2\right) \approx 0.95$$

$$P\left(-2 \leq \dfrac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq 2\right) \approx 0.95$$

$$P\left(-2 * \dfrac{\sigma}{\sqrt{n}} - \bar{X} \leq -\mu \leq 2 * \dfrac{\sigma}{\sqrt{n}} - \bar{X}\right) \approx 0.95$$

$$\Rightarrow P\left(\bar{X} - 2\dfrac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + 2\dfrac{\sigma}{\sqrt{n}}\right) = 0.95$$

97%

$$\Phi^{-1}(0.015) = -a$$

stats.norm.ppf(0.015)

1.5%   1.5%

−a    a

$$P\left(\overline{X} - a * \frac{\sigma}{\sqrt{n}} \leq \mu \leq \overline{X} + \frac{a\sigma}{\sqrt{n}}\right)$$

$$\approx 0.97 = 97\%$$

$$\left(\overline{X} - 2\frac{\sigma}{\sqrt{n}}, \overline{X} + 2\frac{\sigma}{\sqrt{n}}\right) \text{ is called}$$

95% CI

Once I plug in my observed value for $\overline{X}$. (Call this $\overline{x}$)

$$\left(\overline{x} - 2\frac{\sigma}{\sqrt{n}}, \overline{x} + 2\frac{\sigma}{\sqrt{n}}\right) \quad 95\% \text{ CI for } \mu.$$

Example Say the ht of randmly selected Data88 Students is 67", $\sigma = 5$"
Take a sample & construct a 95% CI to get (60", 70")

# Confidence intervals

The graph at top shows a normal distribution with shaded tails labeled $\alpha/2$ on each side, $1-\alpha$ in the middle, and x-axis markers $-z_{\alpha/2}$, $0$, $z_{\alpha/2}$.

- In the previous slide, we derived an ***approximate 95% Confidence Interval for the population mean μ***

$$100(1-\alpha)\% \text{ Conf. Int for } \mu$$

$$\left( \bar{X} - z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}, \ \bar{X} + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \right)$$

- Why is the interval random?

- ***A confidence interval is an interval on the real line, that is, a collection of values, that are plausible estimates for the true mean μ.***

- Using the CLT, we can estimate the chance that this interval contains the true mean. If we want the chance to be higher, we make the interval bigger. The interval is like a net. We are trying to catch the true mean in our net.

$$\text{width of interval} = 2 \times \text{margin of error}$$

- The CLT takes the form: $\bar{X} \pm$ *margin of error,* where the margin of error tells us how big our interval is, and depends on the SD of the sample mean.

- The margin of error $= z_{\alpha/2} \times SD(\bar{X})$, where $z_{\alpha/2}$ is the quantile we need to have an area of $1 - \alpha$ in the middle, that is, a ***coverage probability*** of $1 - \alpha$

# Example

$\sigma = 20 \quad \mu = ?$

- A population distribution is known to have an SD of 20. The average of an iid sample of 64 observations is 55. What is your 95% confidence interval for the population mean?

$n = 64 \quad \bar{x} = 55$

$\boxed{1.96} = Z\text{-score for } 95\%$

$$I = \left( 55 - 2 * \frac{20}{8}, \quad 55 + 2 * \frac{20}{8} \right) = (55-5, 55+5)$$

CI

$= (50, 60)$

What is the prob that this interval $I$ contains $\mu$. $(0 \text{ or } 1)$

$\mu$

# Confidence levels

- The probability with which our *random* interval will cover the mean is called the confidence level.

- In reality (vs theory), we will have just one *realization* (observed value) of the sample mean (from our data sample), and we use that value to write down the **realization** of our random interval.

- What would we do differently if we wanted a 68% CI? 99.7% CI?

- What about an 80% CI? 99% CI?

*Exercise*

# Dealing with proportions

- A sample proportion is just the sample mean of a special population of 0's and 1's.

- This kind of population is so common since many of our problems deal with *classifying* and *counting*.

- We have a population of 1 million in a town. We take a SRS of size 400 and find that 22% of the sample is unemployed. Estimate the percentage of unemployed people in the town.

When the $X_k$'s are 0 or 1, $\quad X_1, X_2 \ldots X_{400}$ pretend Bin.

$\overline{X}$ is a proportion & $\overline{X}$ is just

called $\hat{p}$

$X \sim Bin(400, \hat{p})$?

$E(\overline{X}) \quad Var(\overline{X}) \quad ??$

$$\hat{p} = 0.22 \qquad E(X) = np \qquad Var(X) = npq$$

$$X = X_1 + X_2 + \cdots + X_{400}$$

$$E(X_k) = p = ? \qquad Var(X_k) = p(1-p)$$

$$\overset{\approx \hat{p}(1-\hat{p})}{SD(\overline{X}) = SD(\hat{p}) \approx \frac{\sqrt{\hat{p}(1-\hat{p})}}{\sqrt{n}} = \frac{\sqrt{\hat{p}(1-\hat{p})}}{\sqrt{400}}}$$

Exercise  Complete-this

$$\hat{p} \pm 2 * \frac{\sigma}{\sqrt{n}} \xleftarrow{\sqrt{\hat{p}(1-\hat{p})}}$$

# Example

- In a simple random sample of 400 voters in a state, 23% are undecided about which way they will vote. Find a 95% CI for the proportion of undecided voters in the state.

- In the above problem, find 99.7% confidence interval.