

AP Statistics in 75 minutes

Data: $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu_0, \sigma_0^2)$
true but
unknown population
parameters

Estimators: sample average $\hat{\mu}_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

sample variance $\hat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu}_n)^2$

Exact
sampling
distributions:

$$\hat{\mu}_n \sim N(\mu_0, \sigma_0^2/n)$$

$$\hat{\sigma}_n^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n-1}{2} \frac{1}{\sigma_0^2}\right)$$

} These random
variables are
independent.

Confidence
interval:

Exact $100 \times (1-\alpha)\%$ CI for the mean

$$\underbrace{\hat{\mu}_n}_{\text{point estimate}} \pm \underbrace{t_{n-1}^* \left(1 - \frac{\alpha}{2}\right)}_{\text{"t score"} \atop \text{margin of error}} \underbrace{\sqrt{\frac{\hat{\sigma}_n^2}{n}}}_{\text{standard error}}$$

Hypothesis
test:

$H_0: \mu_0 = h$ \leftarrow point null

$H_A: \mu_0 \neq h$ \leftarrow two-sided alternative

(test statistic) $T_n = \frac{\hat{\mu}_n - h}{\sqrt{\frac{\hat{\sigma}_n^2}{n}}} \sim t_{n-1}$

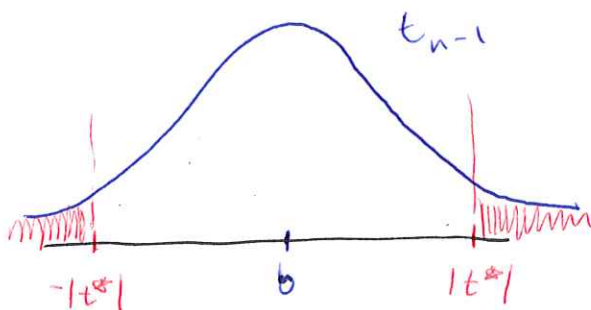
if H_0 were true,
this is the distribution
of the test statistic.
The null distribution.

(observed statistic) t^*

statistic the data gave you.
if it's far in the tails of

(p-value) $p\text{-val} = P(|T_n| \geq t^* | H_0 \text{ true})$

the null distribution,
that's evidence against
the null hypothesis.



Stuff we sweep under the rug in AP Stat or 101

- derivation of sampling distributions
- why is it $n-1$ everywhere and not n ?
- whence the t distribution?
- why are $\hat{\mu}_n$ and $\hat{\sigma}_n^2$ independent?

We will tackle the first 3 points. we defer the fourth to STA 332. You need Calc III and linear algebra.

The tool kit

$$X \sim N(\mu, \sigma^2) \Rightarrow \frac{X - \mu}{\sigma} \sim N(0, 1)$$

$$\Rightarrow cX \sim N(c\mu, c^2\sigma^2)$$

$$Z \sim N(0, 1) \Rightarrow Z^2 \sim \chi_1^2 = \text{Gamma}(\frac{1}{2}, \frac{1}{2})$$

$$X \sim \text{Gamma}(\alpha, \beta) \Rightarrow E(X) = \alpha/\beta$$

$$\text{var}(X) = \alpha/\beta^2$$

$$M(t) = E[e^{tX}] = \left(\frac{\beta}{\beta - t}\right)^\alpha \quad t < \beta$$

$$cX \sim \text{Gamma}(\alpha, \beta/c)$$

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Gamma}(\alpha, \beta) \Rightarrow \sum_{i=1}^n X_i \sim \text{Gamma}(n\alpha, \beta)$$

$$\hat{\mu}_n \sim N(\mu_0, \frac{\sigma_0^2}{n}) \Rightarrow \frac{\hat{\mu}_n - \mu_0}{\frac{\sigma_0^2}{n}} \sim N(0, 1)$$

Recall the Student's t derivation

$$\begin{array}{l} \text{independent } \left\{ \begin{array}{l} V \sim \text{Gamma}(\frac{df}{2}, \frac{df}{2}) \\ Z \sim N(0, 1) \\ T = Z/\sqrt{V} \end{array} \right\} \Rightarrow \\ \quad \parallel \\ \quad \checkmark \\ T|V=v \sim N(0, 1/v) \end{array}$$

$$\begin{array}{l} V \sim \text{Gamma}(\frac{df}{2}, \frac{df}{2}) \\ T|V=v \sim N(0, 1/v) \end{array}$$

$$T \sim t_{df}$$

Marginal
distribution
of T

joint dist
of (V, T)
written
hierarchically

Let's show: $\hat{\sigma}_n^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n-1}{2} \frac{1}{\sigma_0^2}\right)$

BTW: $E\left(\hat{\sigma}_n^2\right) = \frac{\cancel{\frac{n-1}{2}} \frac{1}{\cancel{\frac{n-1}{2}} \frac{1}{\sigma_0^2}}}{\cancel{\frac{n-1}{2}} \frac{1}{\cancel{\frac{n-1}{2}} \frac{1}{\sigma_0^2}}} = \sigma_0^2$ *unbiased!*

$\text{var}\left(\hat{\sigma}_n^2\right) = \frac{\frac{n-1}{2}}{\left(\frac{n-1}{2} \frac{1}{\sigma_0^2}\right)^2} = \frac{4 \sigma_0^2}{n-1} \rightarrow 0$ *consistent!*

Itchy calculation:

$$\begin{aligned} \sum_{i=1}^n (X_i - \mu_0)^2 &= \sum_{i=1}^n (X_i - \mu_0 + \hat{\mu}_n - \hat{\mu}_n)^2 \\ &= \sum_{i=1}^n \left[(X_i - \hat{\mu}_n) + (\hat{\mu}_n - \mu_0) \right]^2 \\ &= \sum_{i=1}^n \left[(X_i - \hat{\mu}_n)^2 + 2(\hat{\mu}_n - \mu_0)(X_i - \hat{\mu}_n) + (\hat{\mu}_n - \mu_0)^2 \right] \\ &= \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 + 2(\hat{\mu}_n - \mu_0) \sum_{i=1}^n (X_i - \hat{\mu}_n) + \sum_{i=1}^n (\hat{\mu}_n - \mu_0)^2 \\ &= \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 + 2(\hat{\mu}_n - \mu_0) \left[\sum_{i=1}^n X_i - \sum_{i=1}^n \hat{\mu}_n \right] + \sum_{i=1}^n (\hat{\mu}_n - \mu_0)^2 \\ &= \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 + 2(\hat{\mu}_n - \mu_0) \left[\sum_{i=1}^n X_i - n \hat{\mu}_n \right] + \sum_{i=1}^n (\hat{\mu}_n - \mu_0)^2 \\ &= \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 + 2(\hat{\mu}_n - \mu_0) \left[\sum_{i=1}^n X_i - n \frac{1}{n} \sum_{i=1}^n X_i \right] + \sum_{i=1}^n (\hat{\mu}_n - \mu_0)^2 \\ &= \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 + 2(\hat{\mu}_n - \mu_0) \left[\sum_{i=1}^n X_i - \sum_{i=1}^n X_i \right] + \sum_{i=1}^n (\hat{\mu}_n - \mu_0)^2 \end{aligned}$$

$$\sum_{i=1}^n (X_i - \mu_0)^2 = (n-1) \hat{\sigma}_n^2 + n (\hat{\mu}_n - \mu_0)^2$$

$$\frac{1}{\sigma_0^2} \sum_{i=1}^n (X_i - \mu_0)^2 = \frac{(n-1)}{\sigma_0^2} \hat{\sigma}_n^2 + n \frac{(\hat{\mu}_n - \mu_0)^2}{\sigma_0^2}$$

$$\sum_{i=1}^n \frac{(X_i - \mu_0)^2}{\sigma_0^2} = \frac{(n-1)}{\sigma_0^2} \hat{\sigma}_n^2 + \frac{(\hat{\mu}_n - \mu_0)^2}{\sigma_0^2/n}$$

$$\sum_{i=1}^n \left(\frac{X_i - \mu_0}{\sigma_0} \right)^2 = \frac{n-1}{\sigma_0^2} \hat{\sigma}_n^2 + \left(\frac{\hat{\mu}_n - \mu_0}{\sqrt{\sigma_0^2/n}} \right)^2$$

Recall if B, C are independent random variables, then

$$M_{B+C}(t) = E[e^{t(B+C)}] = E[e^{tB+tC}] = E[e^{tB}e^{tC}] = E[e^{tB}]E[e^{tC}] = M_B(t)M_C(t)$$

independence!

$$\sum_{i=1}^n \underbrace{\left(\underbrace{\frac{X_i - \mu_0}{\sigma_0}}_{N(0,1)} \right)^2}_{\text{Gamma}(\frac{1}{2}, \frac{1}{2})} = \underbrace{\frac{(n-1)}{\sigma_0^2} \hat{\sigma}_n^2}_{???} + \underbrace{\left(\frac{\hat{\mu}_n - \mu_0}{\sqrt{\frac{\sigma_0^2}{n}}} \right)^2}_{N(0,1) \text{ Gamma}(\frac{1}{2}, \frac{1}{2})}$$

Gamma($\frac{n}{2}, \frac{1}{2}$)

$$A = B + C$$

$$M_A(t) = M_{B+C}(t) = M_B(t)M_C(t) \Rightarrow M_B(t) = \frac{M_A(t)}{M_C(t)}$$

$$= \frac{\left(\frac{1/2}{1/2 - t} \right)^{n/2}}{\left(\frac{1/2}{1/2 - t} \right)^{1/2}}$$

$$= \left(\frac{1/2}{1/2 - t} \right)^{\frac{n-1}{2}}$$

So $\frac{(n-1)}{\sigma_0^2} \hat{\sigma}_n^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{1}{2}\right)$

$$\Rightarrow \hat{\sigma}_n^2 = \frac{\sigma_0^2}{n-1} \frac{n-1}{\sigma_0^2} \hat{\sigma}_n^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n-1}{2} \frac{1}{\sigma_0^2}\right)$$

What happens if we use n instead of $n-1$?

$$\hat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n-1}{2} \frac{1}{\sigma_0^2}\right)$$

$$(n-1) \hat{\sigma}_n^2 = \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{1}{2} \frac{1}{\sigma_0^2}\right)$$

$$\frac{n-1}{n} \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_n)^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n}{2} \frac{1}{\sigma_0^2}\right)$$

So

$$E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_n)^2\right] = \frac{\frac{n-1}{2}}{\frac{n}{2} \frac{1}{\sigma_0^2}} = \frac{n-1}{n} \sigma_0^2 < \sigma_0^2$$

• biased downward!

• for large n , doesn't matter too much

If you standardize w/ the true variance...

$$\frac{\hat{\mu}_n - \mu_0}{\sqrt{\frac{\sigma_0^2}{n}}} \sim N(0,1)$$

only one source of randomness

if you plug in the estimated variance, which is more realistic...

$$\frac{\hat{\mu}_n - \mu_0}{\sqrt{\frac{\hat{\sigma}_n^2}{n}}} \sim t_{n-1}$$

two sources of randomness make tails heavier

$$\begin{aligned} \frac{\hat{\mu}_n - \mu_0}{\sqrt{\frac{\hat{\sigma}_n^2}{n}}} &= \frac{\hat{\mu}_n - \mu_0}{1/\sqrt{n}} = \frac{1/\sigma_0}{1/\sigma_0} \frac{\hat{\mu}_n - \mu_0}{1/\sqrt{n}} \\ &= \frac{\hat{\mu}_n - \mu_0}{\sqrt{\sigma_0^2/n}} \cdot \frac{1}{\sqrt{\hat{\sigma}_n^2/\sigma_0^2}} \end{aligned}$$

$\frac{\hat{\mu}_n - \mu_0}{\sqrt{\sigma_0^2/n}} \sim N(0,1)$

$\frac{\hat{\sigma}_n^2}{\sigma_0^2} \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n-1}{2}\right)$

independent!

$$\sim t_{n-1}$$

Building the exact CI for the mean of iid normal data

- For any sample size n , we have

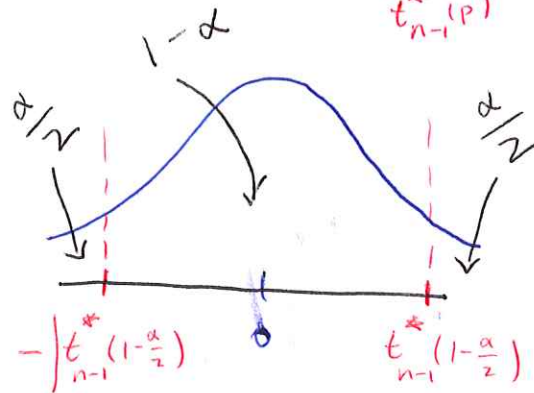
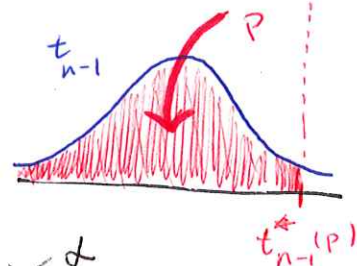
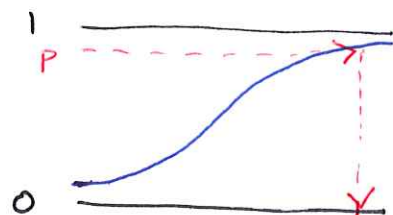
$$\frac{\hat{\mu}_n - \mu_0}{\sqrt{\frac{\hat{\sigma}_n^2}{n}}} \sim t_{n-1}$$

- Let $t_{n-1}^*(p)$ be quantile of t_{n-1} distribution

- We know for $\alpha = 0.01, 0.05, 0.1$, etc ...

$$P\left(-t_{n-1}^*\left(1-\frac{\alpha}{2}\right) < \frac{\hat{\mu}_n - \mu_0}{\sqrt{\frac{\hat{\sigma}_n^2}{n}}} < t_{n-1}^*\left(1-\frac{\alpha}{2}\right)\right)$$

exact! $\Rightarrow 1 - \alpha$



- $P\left(-t_{n-1}^*\left(1-\frac{\alpha}{2}\right) < \frac{\hat{\mu}_n - \mu_0}{\sqrt{\frac{\hat{\sigma}_n^2}{n}}} < t_{n-1}^*\left(1-\frac{\alpha}{2}\right)\right) = 1 - \alpha$

$$P\left(-t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}} < \hat{\mu}_n - \mu_0 < t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}}\right) = 1 - \alpha$$

$$P\left(-\hat{\mu}_n - t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}} < -\mu_0 < -\hat{\mu}_n + t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}}\right) = 1 - \alpha$$

$$P\left(\hat{\mu}_n + t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}} > \mu_0 > \hat{\mu}_n - t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}}\right) = 1 - \alpha$$

- So, if you take
$$L_n = \hat{\mu}_n - t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}}$$
$$U_n = \hat{\mu}_n + t_{n-1}^*\left(1-\frac{\alpha}{2}\right) \sqrt{\frac{\hat{\sigma}_n^2}{n}}$$

then $P(L_n < \mu_0 < U_n) \Rightarrow 1 - \alpha$ For all n !

exact, not approximate