

# t-tests

## Recap: $t$ distribution

If  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then

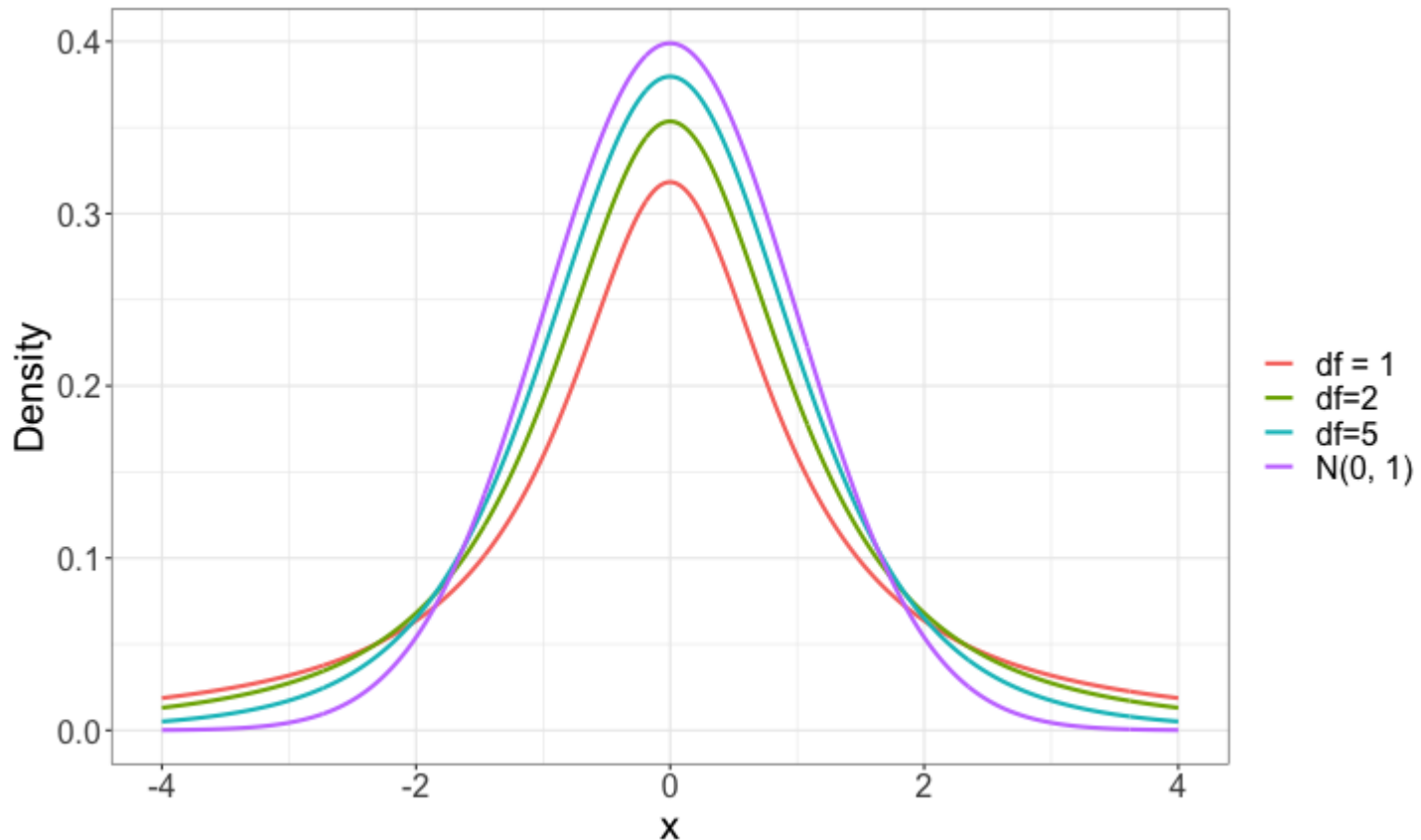
$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{s} \sim t_{n-1}$$

**Definition:** Let  $Z \sim N(0, 1)$  and  $V \sim \chi_d^2$  be independent. Then

$$T = \frac{Z}{\sqrt{V/d}} \sim t_d$$

$$\text{As } n \rightarrow \infty, t_n \approx N(0,1)$$

## t-distribution



Intuition:  $\frac{\sqrt{n}(\bar{X}_n - \mu)}{S}$  is more variable than  $\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}$   
 As  $n \rightarrow \infty$ ,  $\frac{S}{\sigma} \rightarrow 1$ , so  $\frac{\sqrt{n}(\bar{X}_n - \mu)}{S} \approx \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}$

Derivation: WTS if  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$  then  $\frac{\sqrt{n}(\bar{X} - \mu)}{S} \sim t_{n-1}$

$$\textcircled{1} \quad \frac{\sqrt{n}(\bar{X} - \mu)}{S} = \underbrace{\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}}_{N(0,1)} \cdot \sqrt{\frac{(n-1)\sigma^2}{(n-1)S^2}}$$

$$\Rightarrow \text{WTS} \quad \frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1} \quad \text{and} \quad \text{WTS} \quad \frac{(n-1)S^2}{\sigma^2} \perp \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}$$

$$\textcircled{2} \quad \frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \left( \frac{X_i - \bar{X}}{\sigma} \right)^2$$

$$\text{we know} \quad \sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2_n$$

$$\text{Now,} \quad \underbrace{\sum_{i=1}^n \left( \frac{X_i - \mu}{\sigma} \right)^2}_{\chi^2_n} = \underbrace{\sum_{i=1}^n \left( \frac{X_i - \bar{X}}{\sigma} \right)^2}_{\chi^2_{n-1} ?} + \underbrace{n \left( \frac{\bar{X} - \mu}{\sigma} \right)^2}_{\chi^2_1}$$

# Cochran's theorem

Let  $Z_1, \dots, Z_n \stackrel{iid}{\sim} N(0, 1)$ , and let  $Z = [Z_1, \dots, Z_n]^T$ . Let  $A_1, \dots, A_k \in \mathbb{R}^{n \times n}$  be symmetric matrices such that

$Z^T Z = \sum_{i=1}^k Z^T A_i Z$ , and let  $r_i = \text{rank}(A_i)$ . Then the following

are equivalent:

- +  $r_1 + \dots + r_k = n$
- + The  $Z^T A_i Z$  are independent
- + Each  $Z^T A_i Z \sim \chi_{r_i}^2$

## Application to t-tests

Let  $z_i = \frac{x_i - \mu}{\sigma}$ ,  $Z = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix}$   $\sum_i z_i^2 = Z^T Z$

$$\Rightarrow Z^T Z = \sum_{i=1}^n \left( \frac{x_i - \mu}{\sigma} \right)^2 = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma} \right)^2 + n \left( \frac{\bar{x} - \mu}{\sigma} \right)^2$$

want to find  $A_1, A_2$  such that  $Z^T A_1 Z = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma} \right)^2$

and  $Z^T A_2 Z = n \left( \frac{\bar{x} - \mu}{\sigma} \right)^2$

$$\sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma} \right)^2 = \sum_{i=1}^n (z_i - \bar{z})^2 = \sum_{i=1}^n \left( z_i - \frac{1}{n} \sum_j z_j \right)^2$$

$$n \left( \frac{\bar{x} - \mu}{\sigma} \right)^2 = n (\bar{z})^2 = \frac{1}{n} \left( \sum_j z_j \right)^2$$

$$= \frac{1}{n} [z_1 \dots z_n] \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix}$$

$$= Z^T \left( \frac{1}{n} J_n \right) Z \quad J_n = n \times n \text{ matrix of all 1s}$$

$$A_2 = \frac{1}{n} J_n \quad \text{rank}(A_2) = 1$$

$$Z^T Z = \sum_i (Z_i - \frac{1}{n} \sum_j Z_j)^2 + Z^T (\frac{1}{n} J_n) Z$$

$$\Rightarrow \sum_i (Z_i - \frac{1}{n} \sum_j Z_j)^2 = Z^T Z - Z^T (\frac{1}{n} J_n) Z$$

$$= Z^T (I_n - \frac{1}{n} J_n) Z$$

$$\Rightarrow \underbrace{Z^T Z}_{\sim \chi_n^2} = \underbrace{Z^T (I_n - \frac{1}{n} J_n) Z}_{\sum_i (\frac{x_i - \bar{x}}{\sigma})^2} + \underbrace{Z^T (\frac{1}{n} J_n) Z}_{n(\frac{\bar{x} - \mu}{\sigma})}$$

$$\text{rank}(I_n - \frac{1}{n} J_n) = n-1 \quad \text{rank}(\frac{1}{n} J_n) = 1$$

Cochran's theorem;

$$\underbrace{\sum_i (\frac{x_i - \bar{x}}{\sigma})^2}_{\frac{(n-1)s^2}{\sigma^2}} = Z^T (I_n - \frac{1}{n} J_n) Z \sim \chi_{n-1}^2$$

$$\Rightarrow \frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$$

$$n\left(\frac{\bar{x} - \mu}{\sigma}\right)^2 \sim \chi_1^2$$

, and

$$\frac{(n-1)s^2}{\sigma^2}$$

$$\parallel \left( \frac{\sqrt{n}(\bar{x} - \mu)}{\sigma} \right)^2 \parallel$$

# Global F tests for linear regression

# Test for a population <sup>proportion</sup> ~~mean~~

Suppose  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$ . We want to test

$$H_0 : p = p_0 \quad H_A : p \neq p_0$$

Wald test: 
$$Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \quad \text{or} \quad Z = \frac{\hat{p} - p_0}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}}$$

Why is a  $t$ -test not appropriate?

- We don't have to separately estimate the variance (only estimating  $p$ )

Linear regression:  $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$   $\text{var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$

## Test for logistic regression

$$Y_i \sim \text{Bernoulli}(p_i) \quad \log\left(\frac{p_i}{1-p_i}\right) = \beta^T X_i$$

We want to test

$$H_0 : C\beta = \gamma_0 \quad H_A : C\beta \neq \gamma_0$$
$$(C\hat{\beta} - \gamma_0)^T (C X^{-1} (X^T X)^{-1} C^T)^{-1} (C\hat{\beta} - \gamma_0) \approx \chi^2_{\varepsilon} \text{ under } H_0$$

Why is a  $t$ -test not appropriate?

- No separate variance term ( $\sigma^2$ ) to estimate

## Philosophical question

- + If  $X_1, \dots, X_n$  are iid from a population with mean  $\mu$  and variance  $\sigma^2$ , then  $\frac{\sqrt{n}(\bar{X}_n - \mu)}{s} \xrightarrow{d} N(0, 1)$
- + If  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , then  $\frac{\sqrt{n}(\bar{X}_n - \mu)}{s} \sim t_{n-1}$
- + **Position 1:** For any reasonable sample size, the test statistic is approximately normal. And we never really have data from a normal distribution, so the  $t$  distribution is an approximation anyway. So always use the normal distribution
- + **Position 2:** We always have a finite sample size, so our test statistic is never truly normal. And the  $t$  distribution is more conservative than the normal (heavier tails). So always use the  $t$  distribution

With which position do you agree?