

CSE P 590

Building Data Analysis Pipelines

Fall 2024

Significance and confidence



Today

- Housekeeping: Group work and grading
- Recap: Terminology
- Live demo: Statistical significance
- In-class exercise 4

Housekeeping

Canvas groups

How we see them

People > Groups

Everyone **In-class-1-R** In-class-2-data-wrangling In-class-3-stats-modeling In-class-4-stats-nhst In-class-5-big-data [+ Group Set](#)

How you see them

Collaborations

UW Libraries

Panopto Recordings

Zoom

People

UW Resources

Poll Everywhere

Ed Discussion

Files

▶ In-class-1-R 49 In-class-1-R	2 students	🔒
In-class-1-R 50 In-class-1-R	0 students	🔒
▶ In-class-2-data-wrangling 1 In-class-2-data-wrangling	2 students	🔒
▶ In-class-2-data-wrangling 2 In-class-2-data-wrangling	2 students	🔒

Groups and group sets on Canvas



Self-assign to a group by Friday EOD

Grading

- Holistic grading – reasoning and justifications
- Fine-grained grading breakdown (now) on each assignment
 - Completion
 - Questions
 - Optional questions
- Reach out with questions/concerns

Recap: Terminology

Population vs. Sample vs. Sampling distribution

Population

- All possible individuals
- Parameters
 - μ : mean, σ : standard deviation

Sample

- A subset of the population
- Sample statistics
 - \bar{x} : mean, s : standard deviation

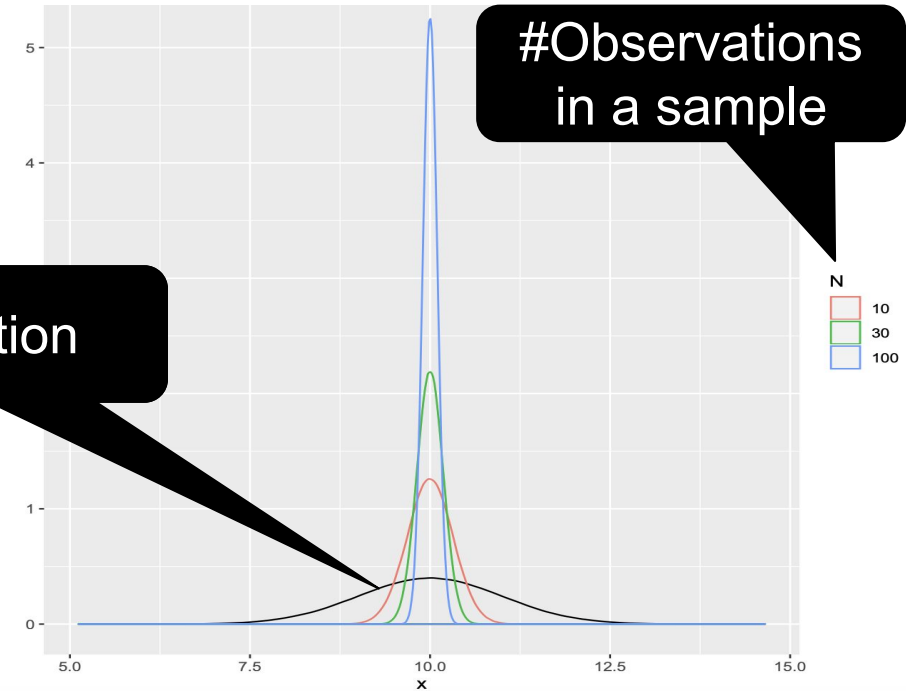
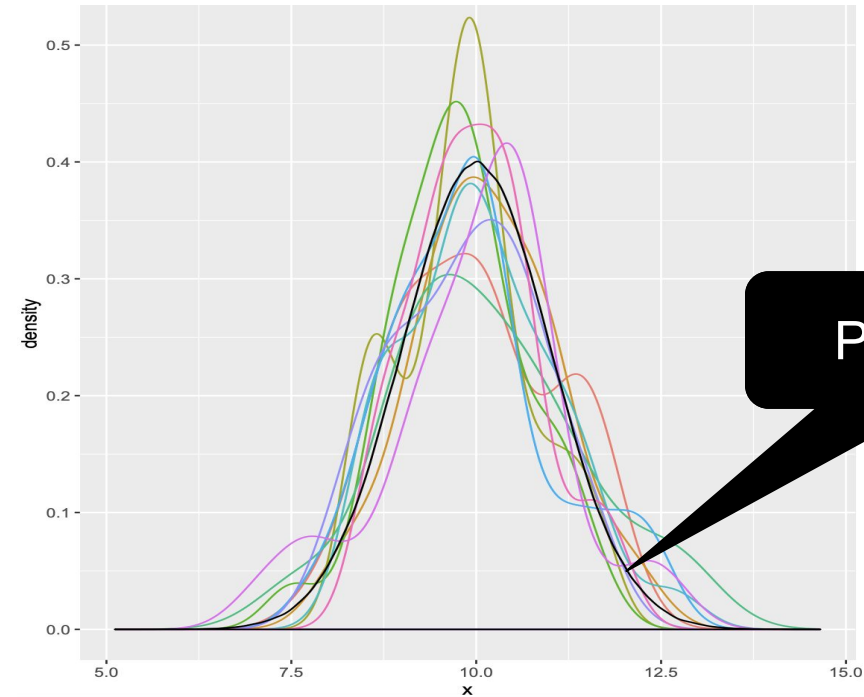
Sampling distribution

- Distribution of the sample statistic (e.g., mean)

Population vs. Sample vs. Sampling distribution

Population and Samples

Sampling distribution



Population

#Observations
in a sample

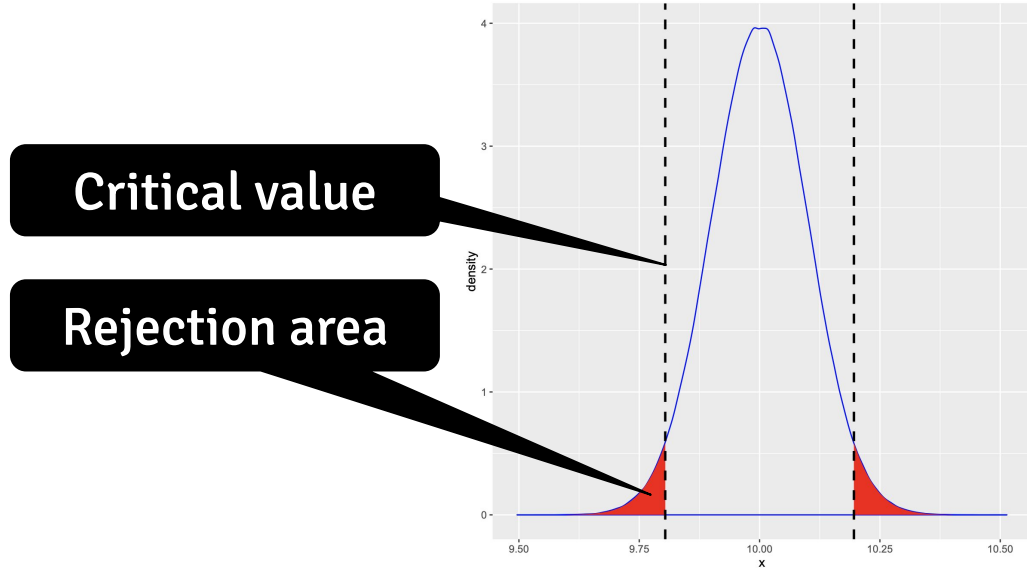
Given a sample mean, what is its p value?

P-value, critical value, and rejection area

Zooming in on the Sampling Distribution

Sampling distribution for $N(10, 1)$ and **sample size 100**

- $\mu = 10, \sigma = 0.1 \rightarrow$ critical values ($\alpha = 0.05$, two-tailed) = 9.804 and 10.196

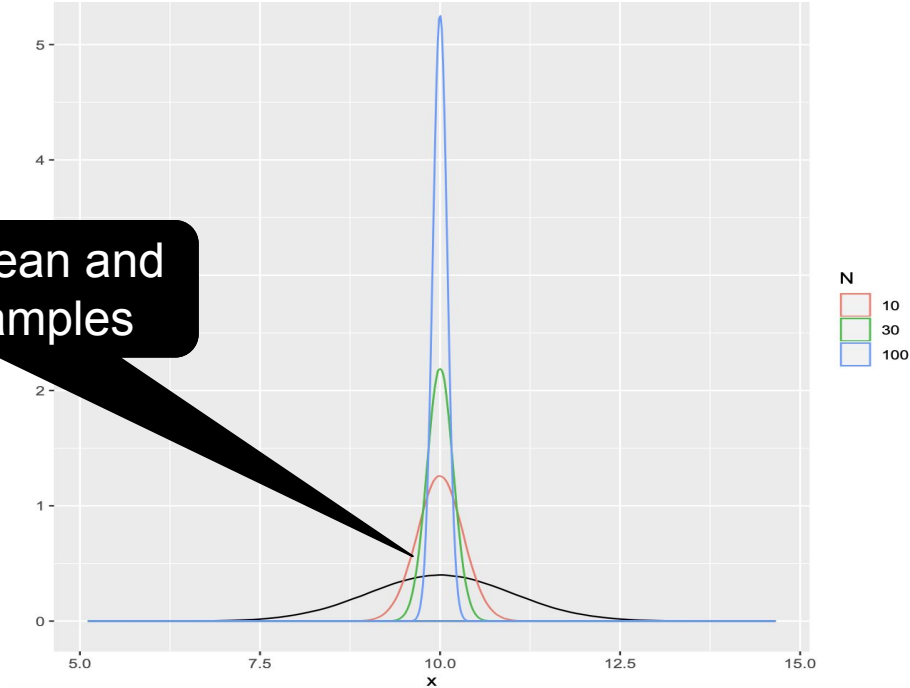
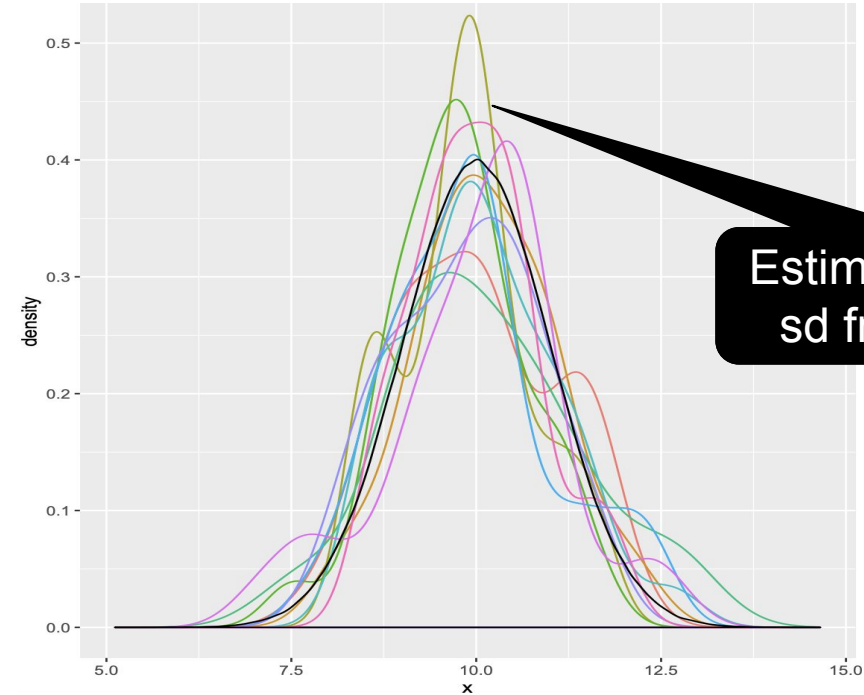


Great, but how do we know the sampling distribution?

Estimating parameters from samples

Population and Samples

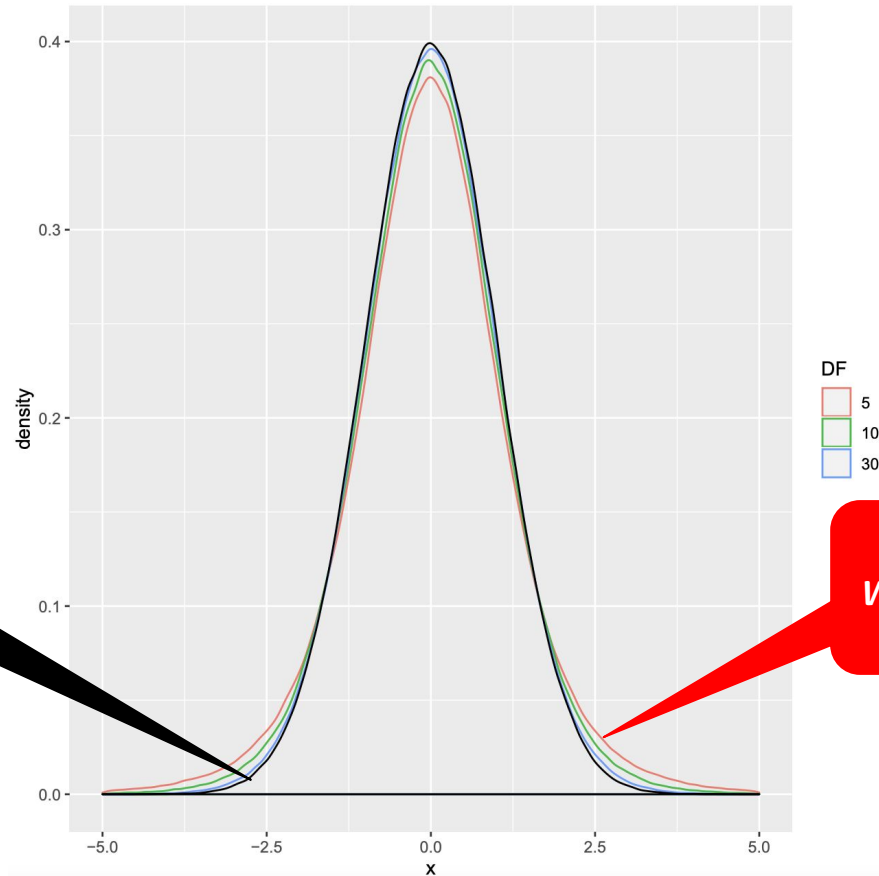
Sampling distribution



Estimate mean and sd from samples

How do we account for uncertainty in those estimates?

z distribution vs. t distribution (small samples)



z distribution:
standard normal
 $N(0,1)$

t distribution
with 5 DF (degrees of
freedom)

NHST: live demo

In-class exercise