

From Paper to Program: Challenges of Implementing Permutation Tests

Kellie Ottoboni
with Philip B. Stark

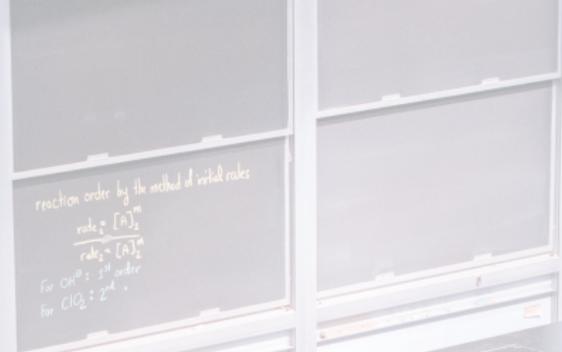
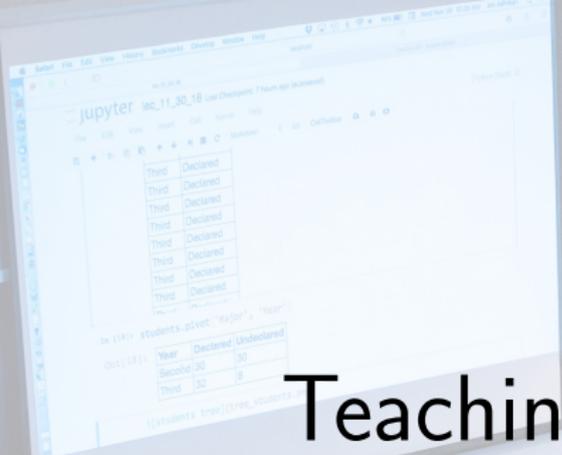
Department of Statistics, UC Berkeley
Berkeley Institute for Data Science

June 12, 2018
ISNPS, Salerno



University of California, Berkeley
DEPARTMENT OF STATISTICS





Teaching Introductory Statistics



Important concepts

- Sampling distributions
- p -values
- Confidence intervals

Introductory Statistics

Important concepts

- Sampling distributions
- p -values
- Confidence intervals

Barriers to learning

- Z tests, t tests
- Assumptions
- Formulas

What if we could teach the concepts without the particular details?

Tools:

- 1 Resampling methods
- 2 Computers



Data Science 8
UC Berkeley

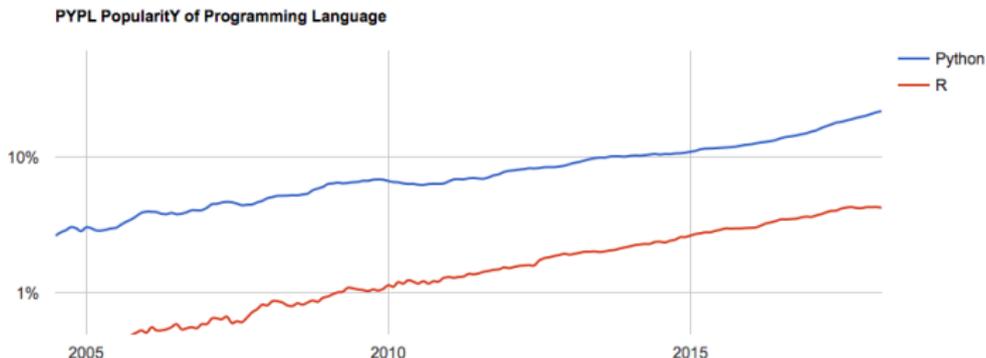
Permutation tests clarify concepts.

- General: it's a procedure, not a formula
- Discrete: counting instead of integration
- Design-based: assumptions come from the data collection

Hesterberg (2015)

Python is gaining popularity for data analysis.

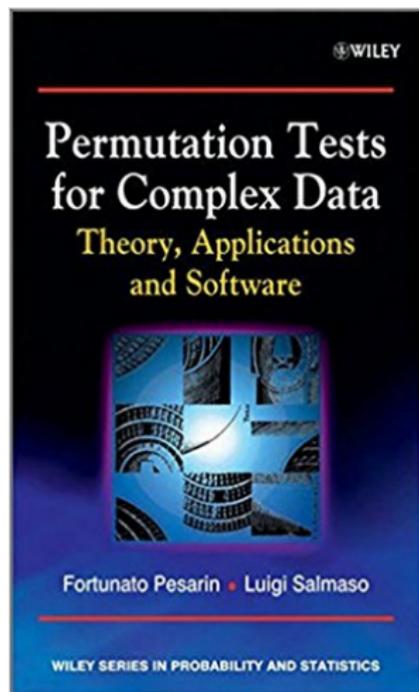
- General purpose language, “batteries included”
- Libraries for a variety of scientific applications





The book

A companion to Pesarin and Salmaso (2010)



- Target audience: late undergraduate to early PhD or domain researchers with some statistics background
- Conceptual: how to determine an appropriate permutation test
- Carefully consider experimental design and assumptions
- Introduce concepts about causal inference

Three chapters

- ① Components of a permutation test
The null invariance and the test statistic
- ② From experiments to observational studies
Randomization, confounding, the “implied experiment”
- ③ Computation for permutation tests
Sampling algorithms and pseudorandom numbers

The book will be open-source online...

Type to search

Introduction

1. Components of a permutation test

1.1 Canonical example: the two sa...

1.2 Exchangeability under the null ...

1.2.1 [Mathematical preliminaries](#)

1.2.2 The null hypothesis

1.2.3 Neyman-Rubin model

1.3 Test statistics

1.4 An example with data

1.5 Exercises

2. From experiments to observational ...

2.1 Randomized experiments

2.1.1 Types of experiments



[Interact](#)

Mathematical preliminaries ¶

Before we talk about hypothesis tests, it's crucial to discuss material that is familiar with concepts from advanced undergraduate math courses.

Sets ¶

A set is a collection of objects, called members or elements of the set. We write $a \in A$, pronounced " a is an element of A ," " a is in A ," or " a belongs to A ," to indicate that a is an element of the set A . This is the same as writing " A contains a ." If a is not an element of A , we write $a \notin A$. Sets may be listed explicitly, or implicitly by specifying a property that they satisfy. The contents of sets are enclosed in curly braces.

with interactive Python code.

Jupyter 01_math_preliminaries Last Checkpoint: 04/12/2018 (autosaved)



File Edit View Insert Cell Kernel Help Python 3.0

CellToolbar

If A is a finite set, the cardinality of the power set of A is $2^{\#A}$. This can be seen as follows: suppose $\#A = n$ is finite. Consider the elements of A to be written in some canonical order. We can specify an element of the power set by an n -digit binary number. The first digit is 1 if the first element of A is in the subset, and 0 otherwise. The second digit is 1 if the second element of A is in the subset, and 0 otherwise, etc. There are 2^n n -digit binary numbers, so there are 2^n subsets. The cardinality of the power set of \mathbb{N} is not \aleph_0 .

If A is a finite set, B is a countable set and $\{A_\beta : \beta \in B\}$ is a partition of A , then

$$\#A = \sum_{\beta \in B} \#A_\beta.$$

Two of the most common data structures in Python are lists and Numpy arrays. These *not* sets, because the items in them have an order. Python also has a data structure for sets, which comes with common set operations we've discussed.

To create a set, pass in a list of items to the `set` function. When you print the set, you'll notice that items are no longer in the same order as when they were passed in.

```
In [1]: set1 = set(['a', 'b', 'c'])
        set2 = set(['a', 'c', 'e', 'g'])
        print(set1)
        print(set2)

{'b', 'c', 'a'}
{'e', 'g', 'c', 'a'}
```

Now, we can do simple set operations with these two sets.

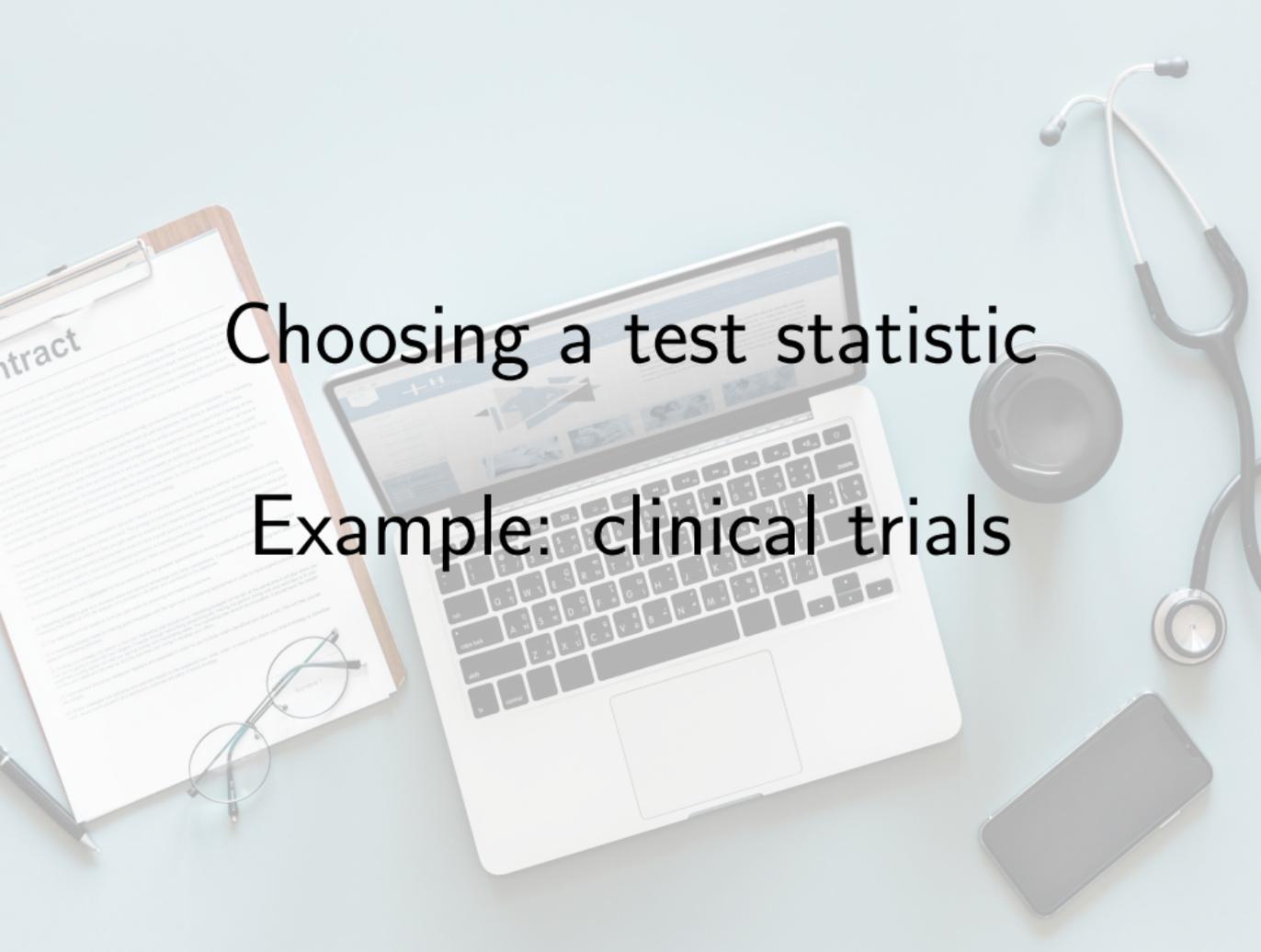
```
In [2]: len(set1)
```

```
Out[2]: 3
```

```
In [3]: set1.union(set2)
        # Or equivalently
        set1 | set2
```

```
Out[3]: {'a', 'b', 'c', 'e', 'g'}
```



A top-down view of a light blue desk with various items: a laptop in the center, a clipboard with a document on the left, a stethoscope on the right, a smartphone at the bottom right, and a pair of glasses on the clipboard. The text is overlaid in the center.

Choosing a test statistic

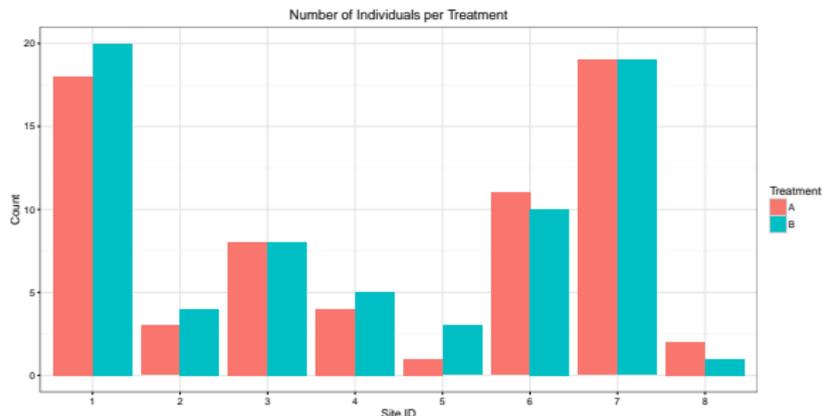
Example: clinical trials

Clinical trial data

Randomized experiment to treat GERD

- 136 patients
- 8 sites
- 7 clinical endpoints, baseline X and follow-up Y measures

Did treatment A reduce symptoms compared to treatment B?

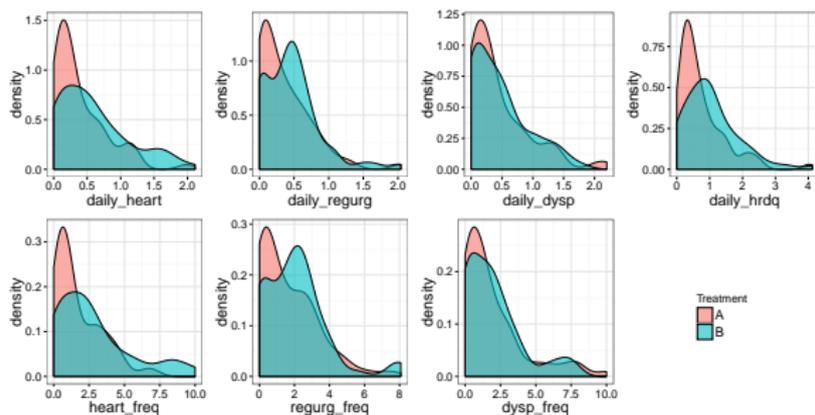


Ottoboni, Lewis, and Salmaso (2018)

The permutation test

Test the **sharp null hypothesis** that treatment has no effect:

$$H_0 : Y_{ij}(A) = Y_{ij}(B), \forall \text{ individuals } i \text{ and sites } j$$



Naive, but correct, statistic: difference in mean outcomes

$$\overline{Y(A)} - \overline{Y(B)}$$

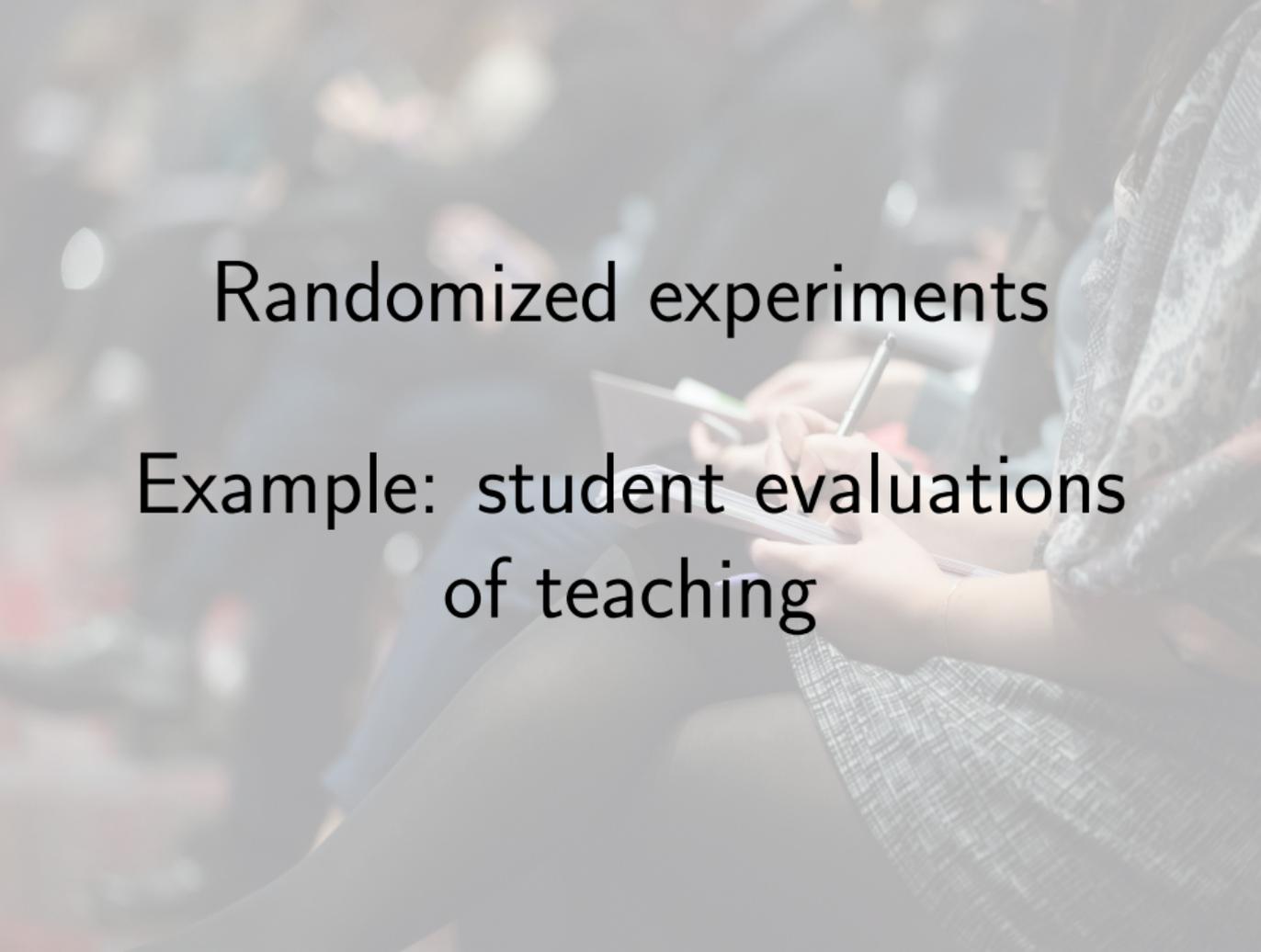
Controlling for covariates

- Difference scores: $\left(\overline{Y(A)} - \overline{X(A)}\right) - \left(\overline{Y(B)} - \overline{X(B)}\right)$

- t statistic from a linear model

$$Y_{ij} = \alpha_j + \beta X_{ij} + \gamma Z_{ij} + \varepsilon_{ij}$$

- Two ways to find a permutation distribution:
 - Permute Z
 - Permute estimates $\hat{\varepsilon}$ from reduced model (*Freedman and Lane (1983)*)



Randomized experiments

Example: student evaluations
of teaching

Teaching evaluations

Are SET a valid measure of teaching effectiveness?

Data from MacNell et al. (2014)

- Students were randomized to 4 online sections of a course.
- In two sections, the instructors swapped identities.
- Was the instructor who identified as female rated lower on average?

Boring, Ottoboni, and Stark (2016)

Neyman-Rubin model, generalized

Rating r_{ijk} = SET given by student i to instructor j
when they appear to have gender k

If gender doesn't matter, $r_{ij\text{male}} = r_{ij\text{female}}$.

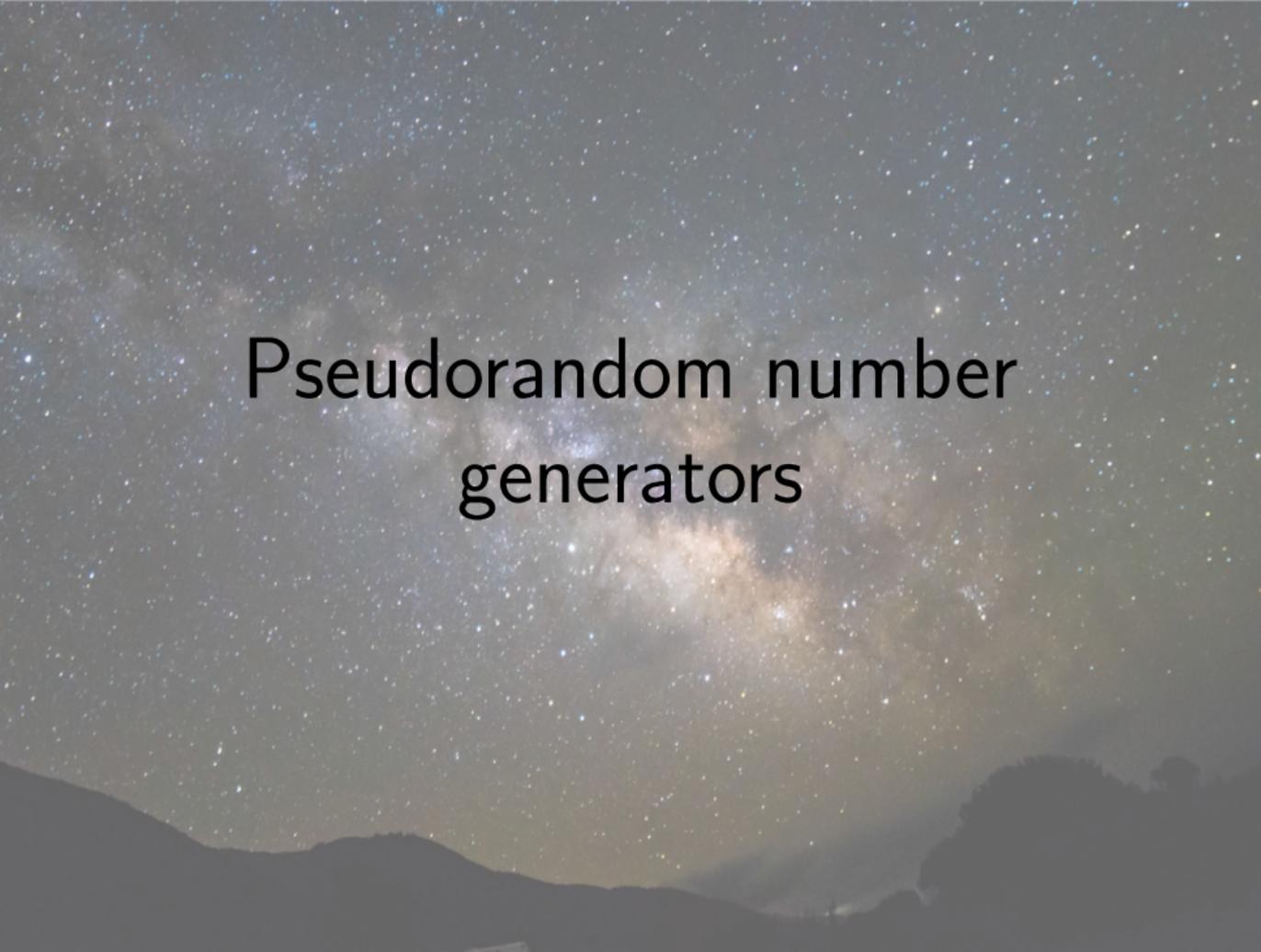
Stratified two-sample test

- For each instructor, permute perceived gender assignments
- Use difference in mean ratings for female-identified minus male-identified

Results

In all categories, the male-identified instructor was rated higher.

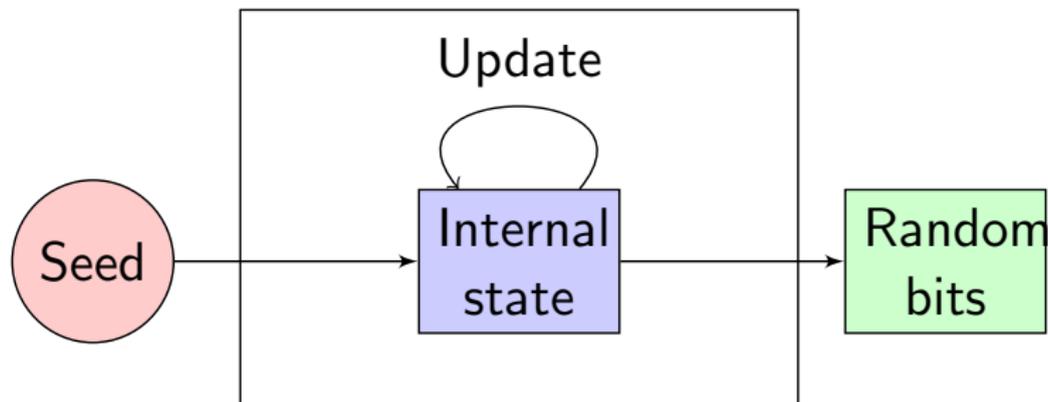
Characteristic	M-F	perm <i>P</i>	t-test <i>P</i>
Overall	0.47	0.12	0.128
Caring	0.52	0.10	0.071
Consistent	0.47	0.21	0.045
Enthusiastic	0.57	0.06	0.112
Fair	0.76	0.01	0.188
Feedback	0.47	0.16	0.054
Helpful	0.46	0.17	0.049
Knowledgeable	0.35	0.29	0.038
Praise	0.67	0.01	0.153
Professional	0.61	0.07	0.124
Prompt	0.80	0.01	0.191
Respectful	0.61	0.06	0.124
Responsive	0.22	0.48	0.013

A night sky photograph featuring the Milky Way galaxy, which appears as a bright, hazy band of light stretching across the frame. The sky is filled with numerous stars of varying brightness. In the foreground, the dark silhouettes of mountains and hills are visible against the starry background.

Pseudorandom number generators

What is a PRNG?

- A user-supplied **seed value** used to set the internal state
- A function that maps the **internal state to pseudorandom bits**
- A function that **updates the internal state**



Pigeonhole principle in action

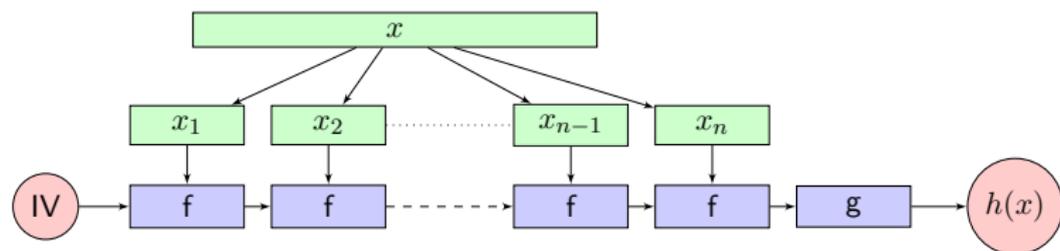
Theorem (Too few pigeons)

If $\binom{n}{k}$ is greater than the size of a PRNG's state space, then the PRNG cannot possibly generate all samples of size k from a population of n .

PRNG	# Internal states	# Possibilities	Proportion of attainable possibilities
32-bit linear congruential generators	4 billion	Samples of 10 out of 50 items ≈ 10 billion	0.4
Mersenne Twister	$\approx 2 \times 10^{6010}$	Permutations of 2084 items $\approx 3 \times 10^{6013}$	0.0001

Cryptographically secure PRNGs

Hash functions take in a message x of arbitrary length and return a value $h(x)$ of fixed size (e.g. 256 bits)



Cryptographic hash functions:

- computationally infeasible to invert
- difficult to find two inputs that map to the same output
- small input changes produce large, unpredictable changes to output
- resulting bits are uniformly distributed

From paper to program

- Teach the concepts of hypothesis testing using permutation tests
- Open-source online book with interactive Python examples
- Packages on GitHub: `permute`, `cryptorandom`, `pscore_match`