

Inference and Randomization

EC 607, Set 11

Edward Rubin

27 May 2020

Prologue

Schedule

Last time

An analytical solution to cluster-robust inference

Today

Inference using (re)randomization [†]

Upcoming

The end is near. As is the final.

[†] These notes follow [notes](#) by Kosuke Imai, [Field Experiments](#) by Gerber and Green, and [Causal Inference for Statistics, Social, and Biomedical Sciences](#) by Imbens and Rubin.

Inference and (re)randomization

Inference and (re)randomization

Inference recap

Our inference techniques have focused on (asymptotic) **analytical methods**.

1. Choose (or derive) an estimator
2. Derived the estimator's (asymptotic) distribution[†]
3. Construct confidence intervals or hypothesis tests

[†] And, consequently, standard errors.

Inference and (re)randomization

Resampling

Resampling methods offers a different, more computationally intense (less asymptotically intense) approach.

A **resampling method** involves repeatedly drawing samples (*resampling*) from a dataset and refitting the model of interest on each sample. We can learn about the behavior of the model through its performance across the many iterations.[†]

Common implementations: Bootstrap (and jackknife), cross validation, permutation tests/randomization inference

[†] This approach is very similar to our Monte Carlo simulations, except that we will sample *with replacement* from a single dataset.

The bootstrap

The bootstrap

Basics

Bootstrapping resamples, *with replacement*, from the original dataset.

- In each sample, we apply our estimator.
- Then, we consider the distribution/properties of these estimates.

This resampling helps us better understand the uncertainty associated with our estimator (within the current data setting).

The bootstrap

More formally

Let's formalize the bootstrap a bit.

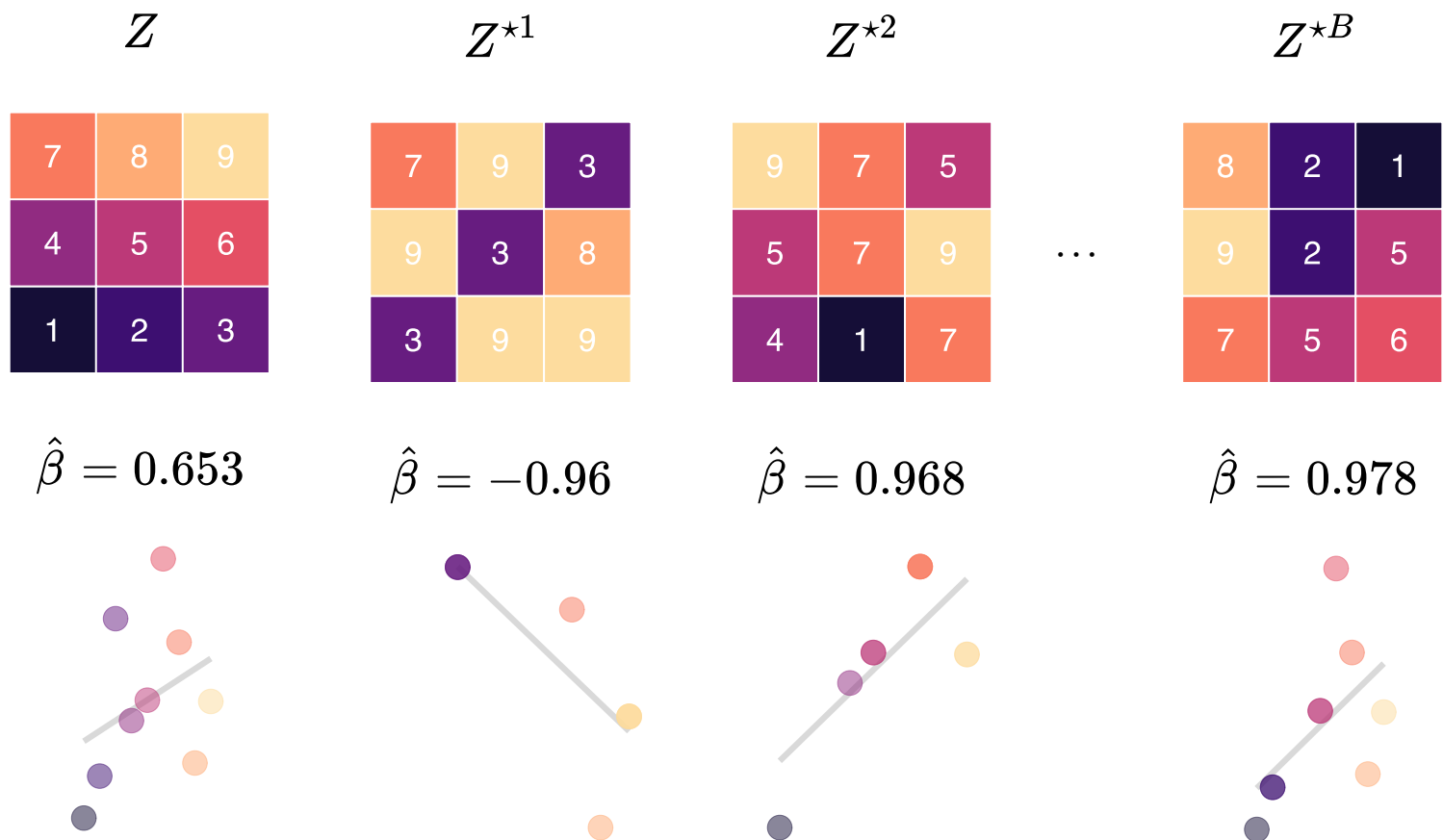
- Z denotes our original dataset (e.g., $Z = [\mathbf{Y} \mid \mathbf{X}]$ in our standard setup).
- $\hat{\alpha}(Z)$ refers to the estimate for α derived from our dataset Z .
- We draw B bootstrap samples $b \in \{1, \dots, B\}$.
- Z^{*1} represents our first bootstrap sample ($b = 1$).
- $\hat{\alpha}^{*1} = \hat{\alpha}(Z^{*1})$ is our estimator evaluated on the first bootstrap sample.

The **bootstrapped standard error** of $\hat{\alpha}$ is the standard deviation of the $\hat{\alpha}^{*b}$

$$\text{SE}_B(\hat{\alpha}) = \sqrt{\frac{1}{B} \sum_{b=1}^B \left(\hat{\alpha}^{*b} - \frac{1}{B} \sum_{\ell=1}^B \hat{\alpha}^{*\ell} \right)^2}$$

The bootstrap

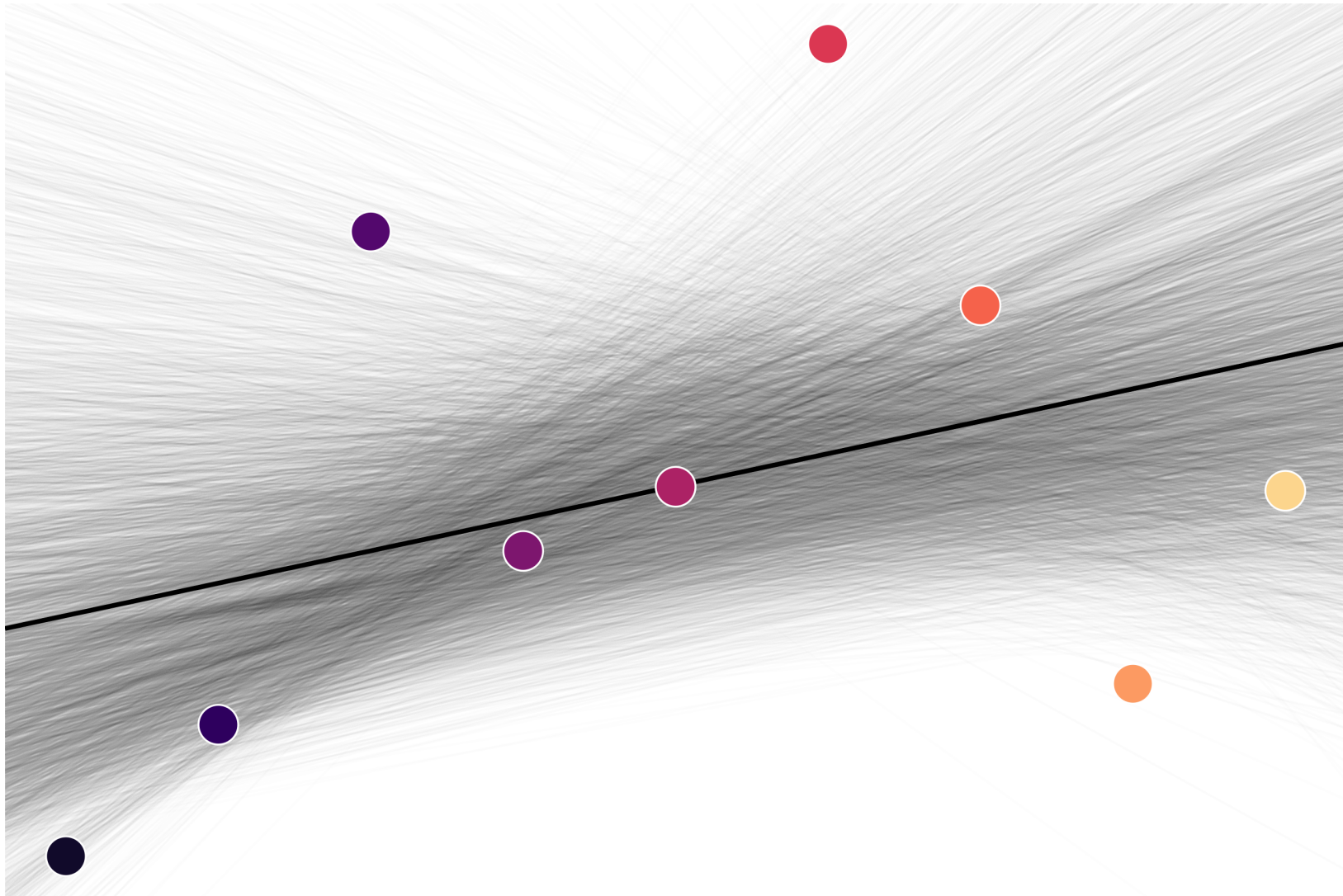
More graphically



The bootstrap

Running this bootstrap 10,000 times

```
plan(multiprocess, workers = 10)
# Set a seed
set.seed(123)
# Run the simulation 1e4 times
boot_df <- future_map_dfr(
  # Repeat sample size 100 for 1e4 times
  rep(n, 1e4),
  # Our function
  function(n) {
    # Estimates via bootstrap
    est <- lm(y ~ x, data = z[sample(1:n, n, replace = T), ])
    # Return a tibble
    data.frame(int = est$coefficients[1], coef = est$coefficients[2])
  },
  # Let furrr know we want to set a seed
  .options = future_options(seed = T)
)
```



The bootstrap

Comparison

In this 10,000-sample bootstrap, we calculate a standard error for $\hat{\beta}_1$ of approximately 0.786.

If we go the old-fashioned OLS route $\left(s^2(\mathbf{X}'\mathbf{X})^{-1}\right)$, we estimate 0.673.

Not bad.

Permutation tests

Permutation tests

Motivation

Consider the null hypothesis of *no average treatment effect*, i.e.,

$$H_0: \bar{Y}_0 = \bar{Y}_1 \quad (\implies \bar{\tau} = 0)$$

We've discussed how randomization avoids the pitfalls of selection bias.

Randomization can also clarify inference—helping quantify uncertainty.

Q How?

A We know exactly how the randomness happened (we assigned it), so we don't need parametric assumptions to derive a distribution under H_0 !

We use the **experimental design**, rather than a probability model.

Permutation tests

Tea drinkers

Classic example Sir R. A. Fisher had a colleague who claimed to be able to tell whether the tea was poured into milk *or* milk was poured into the tea.[†]

Being the friend he was, Fisher designed an experiment to determine whether his colleague was telling the truth.

Fisher randomized the order of 8 cups of tea:

- 4 cups with **m**ilk added first
- 4 cups with **t**ea added first

Vindication! His colleague got all 8 correct.

Q With random guessing, how likely is correctly guessing all 8 cups?

[†] Don't worry, Fisher is known for more than this one experiment.

Permutation tests

Tea drinkers 2

Q With random guessing, how likely is correctly guessing all 8 cups?

This question reflects our understanding of a **p-value**.

If Fisher's colleague had no ability and simply guessed (H_0), what is the probability she would have guessed all 8 cups correctly?

Fisher's H_0 : the answers were unrelated to the cups' actual contents.

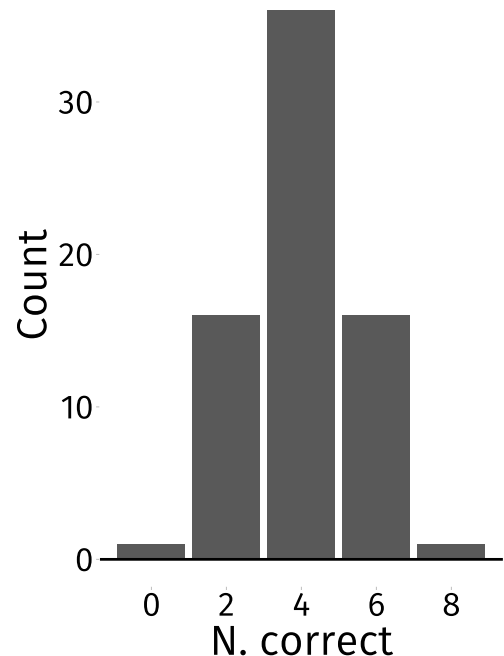
Under this hypothesis, we can re-randomize the cups and see how many times her answer was perfectly correct.

This is the idea behind **permutation testing** and **randomization inference**.

Permutation tests

Tea drinkers with a vengeance

| Cup | Guess | Truth | P_1 | P_2 | P_3 | \dots | P_{70} |
|-----|-------|-------|-------|-------|-------|---------|----------|
| 1 | m | m | m | m | m | | t |
| 2 | t | t | m | m | m | | t |
| 3 | t | t | m | m | m | | t |
| 4 | m | m | m | t | t | | t |
| 5 | m | m | t | m | t | | m |
| 6 | t | t | t | t | m | | m |
| 7 | t | t | t | t | t | | m |
| 8 | m | m | t | t | t | | m |
| | | 8/8 | 4/8 | 4/8 | 2/8 | | 4/8 |



So our permutation-test-based p -value is $1/70 \approx 0.0143$. \implies Reject H_0 .

Permutation tests

Generalization

The procedure for permutation-based hypothesis testing[†] is the same as our "standard" asymptotic-based hypothesis testing.

1. **Define hypotheses**, H_0 and H_a .
2. Choose our **rejection threshold** α (tolerated type-I error rate).
3. Choose a **test statistic** that is a function of our sample.
4. Derive/calculate the **test statistic's distribution under H_0** .
5. **Compute the p -value** by comparing test stat. to its H_0 distribution.
6. **Conclusions**—reject or fail to reject H_0 .

The difference: Permutation tests use the randomization's mechanism to construct the test-statistic's exact distribution under H_0 .

[†] Also called *Fisher's exact test*, as you get exact p -values.

Permutation tests

More generally

Fisher focused on testing a **sharp null hypothesis**—no effect *for anyone*, i.e.,

$$H_0: Y_{1i} - Y_{0i} = 0 \quad \forall i \quad (\implies \tau_i = 0 \quad \forall i)$$

against an alternative hypothesis that someone has a non-zero effect

$$H_a: Y_{1i} - Y_{0i} \neq 0 \text{ for some } i \quad (\implies \exists i \text{ s.t. } \tau_i \neq 0)$$

A **sharp null hypothesis** is specified *for all individuals*, e.g.,

$$H_0: Y_{1i} - Y_{0i} = C \quad \forall i$$

which differs from the ATE-based nulls that we normally consider, e.g.,

$$H_0: E[Y_{1i} - Y_{0i}] = C.$$

Permutation tests

Key insight

Our estimate (or test statistic) is a function of

1. individuals' responses (\mathbf{Y}_i)
2. individuals' treatment assignments (\mathbf{D}_i)

Under the sharp null $H_0: \tau_i = 0 \ (\forall i)$

- $\mathbf{Y}_{0i} = \mathbf{Y}_{1i} = \mathbf{Y}_i \ \forall i$ (i.e., changing \mathbf{D}_i will not affect observed \mathbf{Y}_i)
- Permutations of \mathbf{D} construct the *exact* null distribution (unchanged \mathbf{Y}).

The number of possible permutations can get big—e.g., 500 treated and 500 control has 2.7×10^{299} options. Approximate the distribution by sampling.

Permutation tests

Different inference

In his 2019 paper [Channeling Fisher: Randomization Tests and the Statistical Insignificance of Seemingly Significant Experimental Results](#), Alwyn Young 'updates' inference from 53 experimental papers by using randomization-based inference.

In the average paper, randomization tests of the significance of individual treatment effects find 13% to 22% fewer significant results than are found using authors' methods.

Young (2019)

Permutation tests

Different inference?

It's certainly possible authors and methods can be wrong.

However, permutation-based inference itself may generate differences relative to the more standard, derived, asymptotics-based estimators.

Why?

1. We are testing **different null hypotheses** (sharp vs. non-sharp).
2. The two estimators have **different asymptotic properties**.[†]

[†] Thanks go to Alberto Abadie for this point.

Permutation tests

On average

The sharp null was central to Fisher's interpretation.

Neyman *et al.* (1935) extended[†] this idea of permutation-based tests to the average treatment effect (testing $H_0: E[Y_{1i}] - E[Y_{0i}] = 0$).

Neyman and others also added standard errors and confidence intervals.

These extensions have come to be known as **randomization inference**.^{††}

[†] Fisher, paraphrased: 🤖

^{††} *Permutation tests* and *Randomization inference* are not the most strictly defined terms.

Randomization inference

Randomization inference

Setup

In order to **generalize our null hypothesis to the average treatment effect**,

$$H_0: \bar{\tau} = 0 \implies E[Y_{1i} - Y_{0i}] = 0$$

we have to give up something.

1. If we want an exact null distribution, then we must **assume a uniform treatment effect**. (Assuming our way back to a sharp null.)
2. If we want to avoid assuming $\tau_i = \bar{\tau} \ \forall i$, then we have to **accept a non-exact null distribution**. (We don't observe Y_{0i} for $D_i = 1$.)

If we don't like either option, then we need to go back to deriving asymptotic properties via probability modeling assumptions.

Randomization inference

Implementation

Once we decide which simplification we're willing to accept, we proceed similarly to permutation tests:

- shuffle \mathbf{D} in a way that mimics treatment assignment
- collect test statistics from each iteration

Note Monte Carlo simulations, bootstrap, permutation tests, and randomization all apply very similar processes.

Randomization inference

(Which) Test statistics

We still need to choose a test statistic on which we base the p -value.

- The **actual estimate**—difference in means or coefficient
- **Transformed estimates**
- **Quantiles**, *e.g.*, the median
- **t statistic**
- **Rank** statistics

We can also extend this idea to **confidence intervals**.

E.g., Use the point estimates associated with the 2.5th and 95th percentiles to construct a 95% confidence interval.

Randomization inference

Example

Back to the LaLonde NSW dataset. We previously estimated

- the NSW increased real earnings by $\hat{\beta}_1 \approx \$886.30$
- (het.-robust) standard error of $\$488.20$
- t statistic $t_{\text{stat}} \approx 1.82$ with p -value ≈ 0.0699

Let's re-randomize treatment 10,000 times. In each **iteration** r , calculate

1. $\hat{\beta}_1^r$, the **point estimate** (the regression coefficient)
2. t_{stat}^r , the **t statistic**

Then calculate the implied p -values using the location of $\hat{\beta}_1$ and t_{stat} in the distributions of $\hat{\beta}_1^r$ and t_{stat}^r , respectively.[†]

[†] Very similar exercise for confidence intervals.

Randomization inference

Example: Re-randomization

The main decision is how to generate treatment.

Q Should we permute \mathbf{D} or draw \mathbf{D}_i for each individual?[†]

A How was the original randomization conducted?

We'll assume the NSW started with a set number of treatments to disperse.

[†] The difference is in whether we hold the number of treated individuals constant.

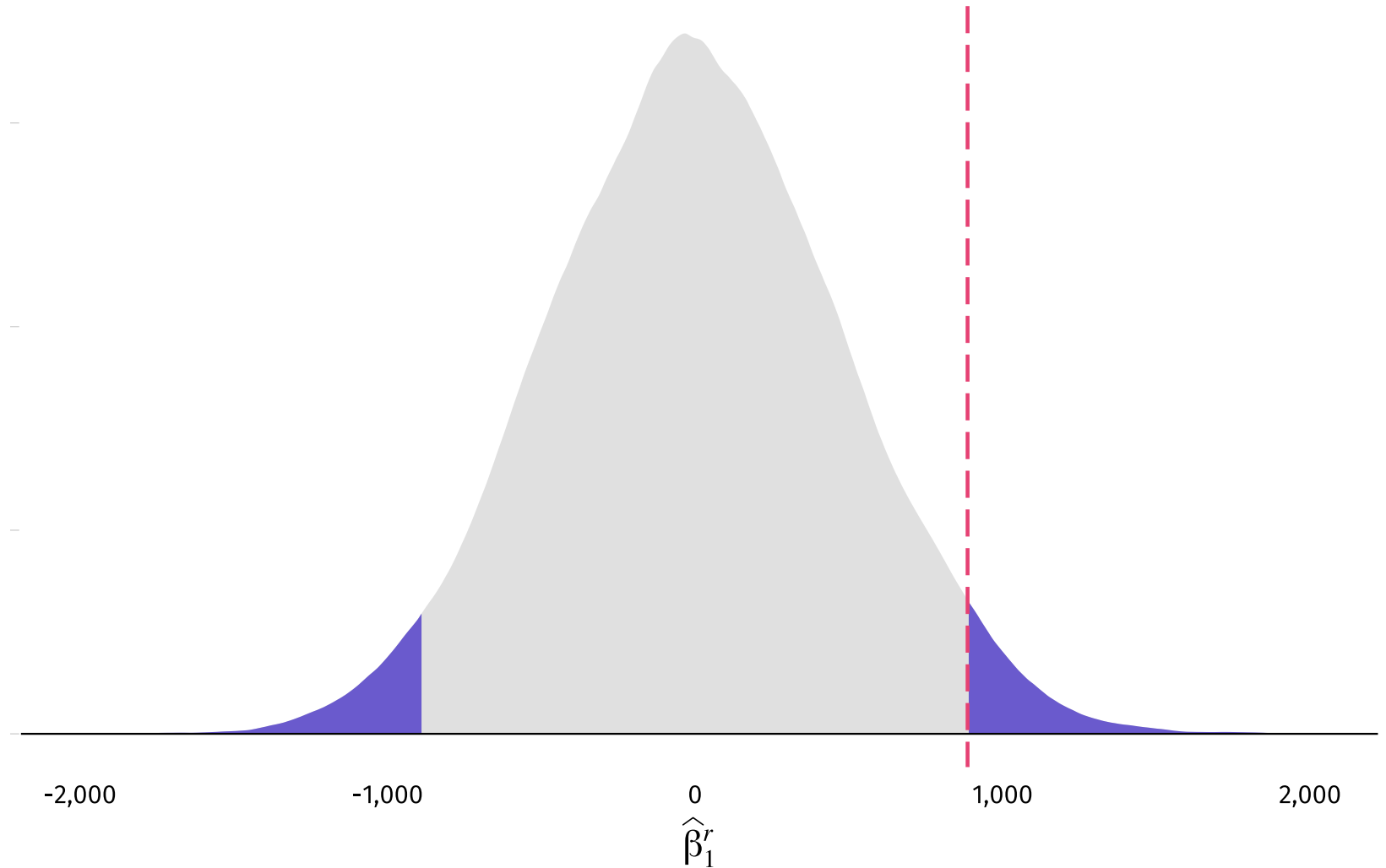
First, we'll write a function that performs one iteration.

```
# Arguments: 'i' (iteration), 'n_t' (# of trt)
fun_randomization <- function(i) {
  # Sample the treatment vector. NOTE: Sampling WITHOUT replacement
  t_i <- sample(nsw_df$treat, size = nrow(nsw_df), replace = F)
  # Regression using our re-randomized treatment
  est_i <- lm_robust(re78 ~ t_i, data = nsw_df) %>% tidy()
  # Return tibble with iteration, point estimate, and test statistic
  tibble(i, est = est_i[2,"estimate"], t_stat = est_i[2,"statistic"])
}
```

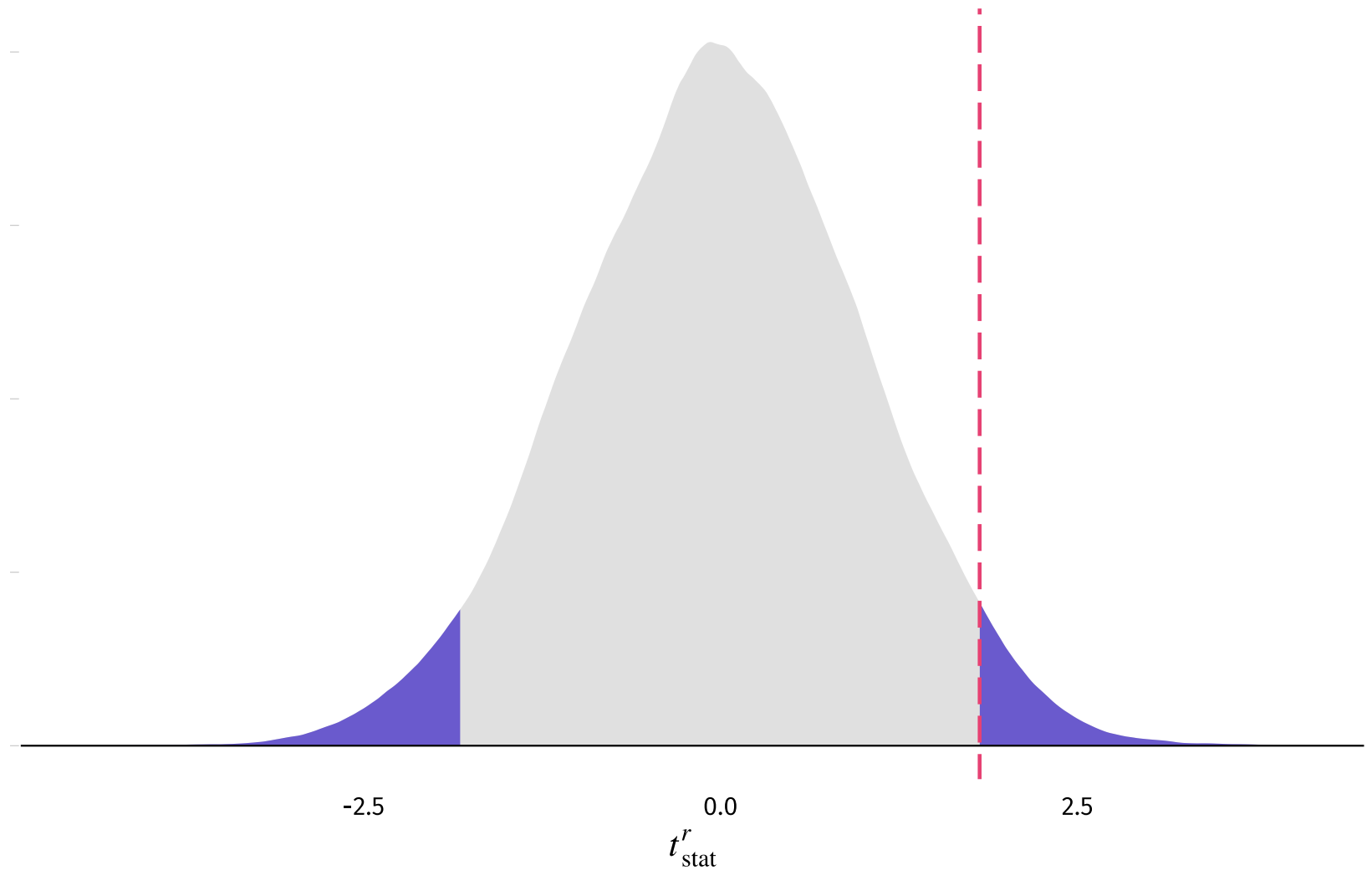
And now run the re-randomization function 10,000 times.

```
# Set up parallelization and seed
plan(multiprocess, workers = 4); set.seed(1234)
# Run the simulation 1e4 times
random_df <- future_map_dfr(
  1:1e4,
  fun_randomization,
  .options = future_options(seed = T)
)
```

Result 1 Share $|\hat{\beta}_1^r| > \hat{\beta}_1 = 0.0615$. (Original p -value = 0.0699)



Result 2 Share $|t_{\text{stat}}^r| > t_{\text{stat}} = 0.0703$. (Original p -value = 0.0699)



Randomization inference

Confidence intervals

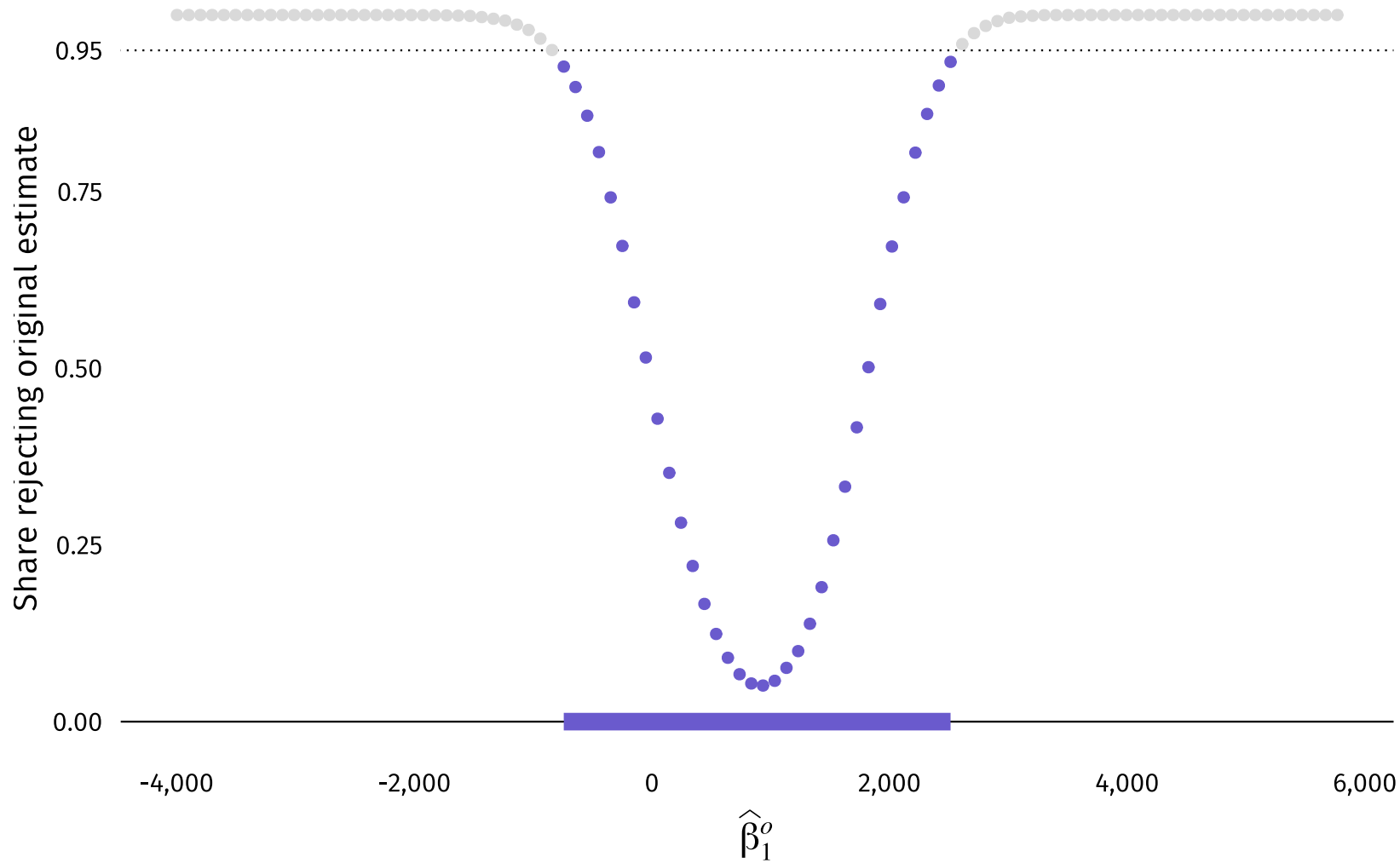
To construct confidence intervals, we **invert** randomization-based **hypothesis tests**, imposing a range of null hypotheses.

E.g., To construct a 95% C.I. for $\hat{\tau}$

1. Impose the null hypothesis $H_0: \tau = \tau_0$ for many values of τ_0 .
2. Find all values of τ_0 that do not reject $\hat{\tau}$ at the 5% level.

Note We must to be able to clearly impose the null in our "model".

Constructing a 95% confidence interval



Athey and Imbens (2016) **on regression and randomization inference:**[†]

Although these methods [regression] remain the most popular way of analyzing data from randomized experiments, **we suggest caution in using them.**

... In particular there is a disconnect between the way the conventional assumptions in regression analyses are formulated and the implications of randomization. As a result it is easy for the researcher using regression methods to go beyond analyses that are justified by randomization, and end up with analyses that rely on a **difficult-to-assess mix of randomization assumptions, modeling assumptions, and large sample approximation.**

[†] Specifically in the context of experiments, though the concerns should remain in other contexts.

Athey and Imbens (2016) **on regression and randomization inference:**[†]

Ultimately we recommend that researchers wishing to use regression or other model-based methods rather than the randomization-based methods we prefer, do so with care. For example, using only indicator variables based on partitioning the covariate space, rather than using multi-valued variables as covariates in the regression function preserves many of the finite sample properties that simple comparisons of means have, and leads to regression estimates with clear interpretations. In addition, in many cases the potential gains from regression adjustment can also be captured by careful ex ante design, that is, through stratified randomized experiments to be discussed in the next section, without the potential costs associated with ex post regression adjustment.

[†] Specifically in the context of experiments, though the concerns should remain in other contexts.

Randomization and clustering

Randomization and clustering

The plot thickens

Permutation tests and randomization inference both work because we know[†] the process through which treatment was randomly assigned.

If treatment is correlated within groups, then our bootstraps, permutations, and re-randomizations need to reflect this dependence.

[†] Or claim to understand.

Further reading

Papers

Bootstrap-Based Improvements for Inference with Clustered Errors

Cameron, Gelbach, and Miller (2008)

The Econometrics of Randomized Experiments Athey and Imbens (2016)

Randomization Inference With Natural Experiments

Ho and Imai (2012)

Also: [Notes](#) by Kosuke Imai

Further reading

Books: Resampling methods and the bootstrap

An Introduction to Statistical Learning

James, Witten, Hastie, and Tibshirani

Elements of Statistical Learning

Hastie, Tibshirani, and Friedman

Books: Permutation tests and randomization inference

Causal Inference for Statistics, Social, and Biomedical Sciences

Imbens and Rubin

Field Experiments

Gerber and Green

Table of contents

Admin

1. Schedule
2. Further reading

Inference and randomization

1. Resampling
2. The bootstrap
 - Basics
 - Semi-formally
 - Graphically
3. Permutation tests
 - Motivation
 - Tea tests
 - Different inference
 - Generalization
 - Basics
4. Randomization inference
 - Setup
 - Example
 - Confidence intervals
5. Clustering