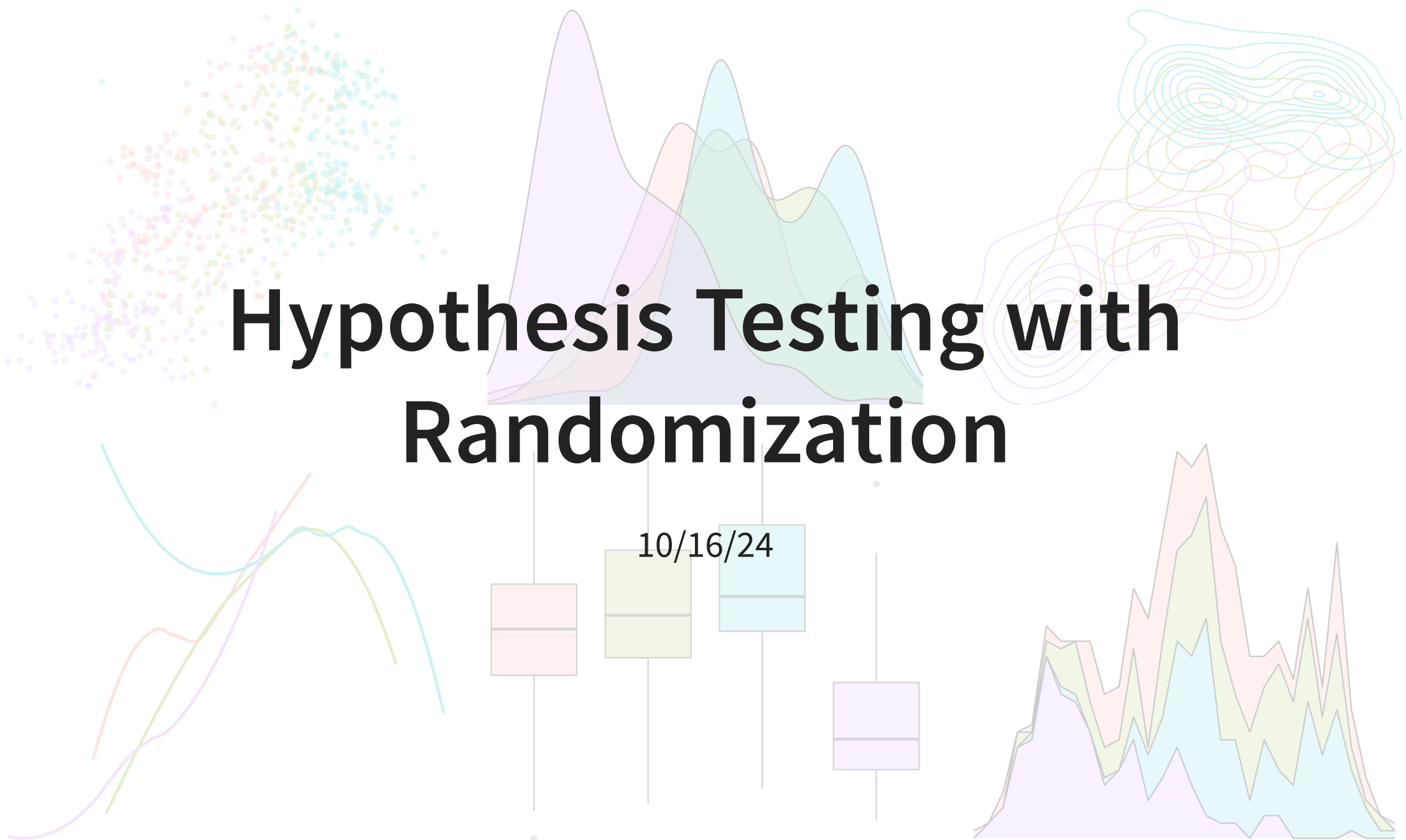# Hypothesis Testing with Randomization

10/16/24

# Housekeeping

- Office hours change this week

# Where we're going today

- We will see another kinds of hypotheses for different types of research questions

- Hypothesis testing framework is the same, but will change how we obtain null distribution

- Try to see the big picture

# Test of independence

# Running example: sex discrimination study

- Note: this study considered sex as binary "male" or "female", and did not take into consideration gender identities
- Participants in the study were 48 bank supervisors who identified as male and were attending a management institute at UNC in 1972
  - Each supervisor was asked to assume the role of personnel director of a bank
  - They were each given a file to judge whether the person in the file should be promoted
  - The files were identical, except half of them indicated that the candidate was male, and the other half were indicated as female
  - Files were randomly assigned to bank managers
  - Experiment or observational study?
- Research question: Are individuals who identify their sex as female discriminated against in promotion decisions made by their managers who identify as male?

# Defining hypotheses

Research question: Are individuals who identify their sex as female discriminated against in promotion decisions made by their managers who identify as male?

- What is/are the variables(s) here? What types of variables are they?

- We need to construct hypotheses where $H_0$ is "status quo" and $H_A$ is the claim researchers have

- $H_0$: the variables `sex` and `decision` are independent.
  - i.e. any observed difference in promotion rates is due to variability

- $H_A$: the variables `sex` and `decision` are *not* independent, and equally-qualified female personnel are less likely to be promoted than male personnel

# Data

For each of the 48 supervisors, the following were recorded:

- The sex of the candidate in the file (male/female)

- The decision (promote/not promote)

| sex | not promote | promote | total |
|---|---|---|---|
| female | 10 | 14 | 24 |
| male | 3 | 21 | 24 |
| total | 13 | 35 | 48 |

- Are we prepared to answer our research question: Are individuals who identify their sex as female discriminated against in promotion decisions made by their managers who identify as male?

  - What evidence do we have?

# Data (cont.)

Conditional probability of getting promoted by sex:

```
1  discrimination |>
2    count(sex, decision) |>
3    group_by(sex) |>
4    mutate(cond_prob = n/sum(n)) |>
5    filter(decision == "promote") |>
6    select(-n)
```

```
# A tibble: 2 × 3
# Groups:    sex [2]
  sex     decision cond_prob
  <chr>   <chr>        <dbl>
1 female  promote      0.583
2 male    promote      0.875
```

- Is the observed difference -0.2916667 convincing evidence? We need to examine variability in the data, assuming $H_0$ true.

- Let's set $\alpha = 0.05$

# Simulate under null

- Simulating under $H_0$ means operating in a hypothetical word where `sex` and `decision` are independent.

    - This means that knowing the `sex` of the candidate should have no bearing on the `decision` to promote or not

- We will perform a simulation called a **randomization test** where we randomly re-assign `decision` and `sex` outcome pairs to see what would have happened if the bankers' `decision` had been independent of candidate's `sex` (i.e. if $H_0$ true)

# Randomization test

| sex | not promote | promote | total |
|---|---|---|---|
| female | 10 | 14 | 24 |
| male | 3 | 21 | 24 |
| total | 13 | 35 | 48 |

- Write down "promote" on 35 cards and "not promote" on 13 cards. Repeat the following:

  - Thoroughly shuffle these 48 cards.

  - Deal out a stack of 24 cards to represent males, and the remaining 24 cards to represent females

    - This is how we simulate independence under $H_0$

  - Calculate the proportion of "promote" cards in each stack, $\hat{p}_{m,sim}$ and $\hat{p}_{f,sim}$

  - Calculate and record the difference $\hat{p}_{f,sim} - \hat{p}_{m,sim}$ (order of difference doesn't matter so long as you are consistent)

# Randomization test (activity)

Try it!

# Randomization test (code)

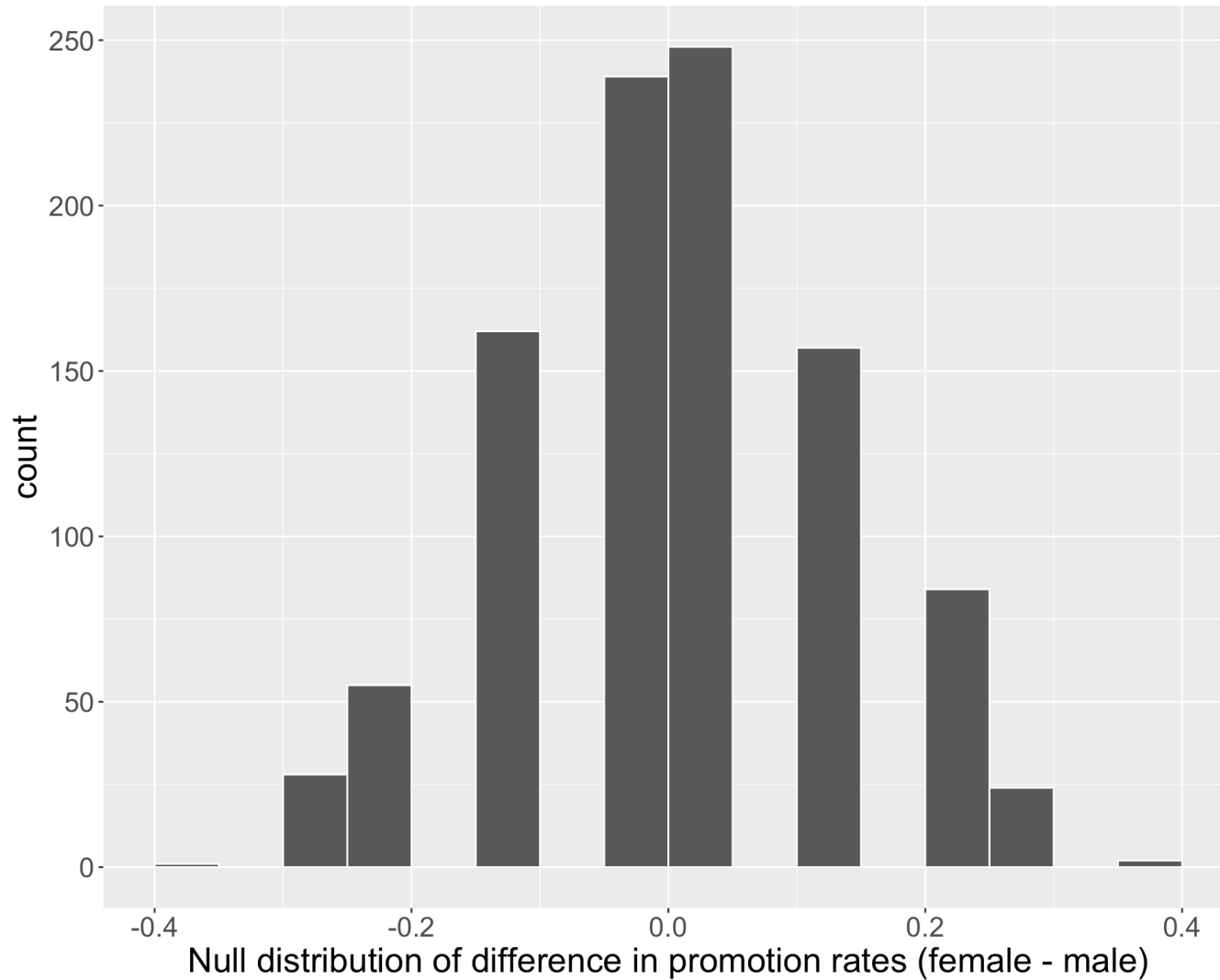Let's perform one iteration of the simulation.

```r
1   # reproducibility
2   set.seed(1)
3   n_female <- sum(discrimination$sex == "female")
4   n_male <- sum(discrimination$sex == "male")
5
6   # create cards
7   cards <- discrimination$decision
8
9   # shuffle cards
10  shuffled <- sample(cards)
```

| sex | not promote | promote | total |
|---|---|---|---|
| female | 8 | 16 | 24 |
| male | 5 | 19 | 24 |
| total | 13 | 35 | 48 |

- Under this simulation, 0.6666667 of females were promoted, and 0.7916667 of males were promoted. Simulated difference: -0.125
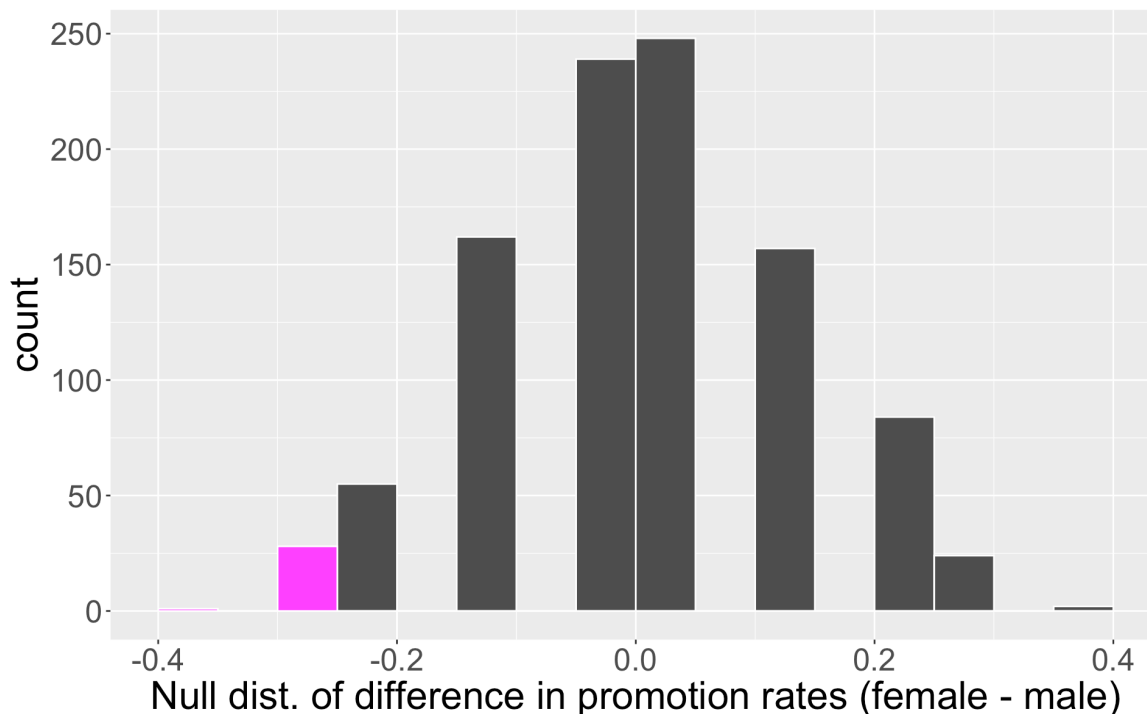
# Null distribution

Repeat the previous step 1000 times:

# Obtain p-value

Recall, the observed difference in our data was $\hat{p}_f - \hat{p}_m = $ -0.2916667.

- p-value is probability of observing data as or more extreme than our original data, given $H_0$ true.

- Where does "as or more extreme" correspond to on our plot?



- Out of 1000 simulations under $H_0$, 29 resulted in a difference in promotion rates as or more extreme than our observed

- So the p-value is 0.029

# Making decision and conclusion

Our research question: Are individuals who identify their sex as female discriminated against in promotion decisions made by their managers who identify as male?

- $H_0$: `sex` and `decision` are independent
- $H_A$: `sex` and `decision` are not independent and equally-qualified female personnel are less likely to get promoted than male personnel by male supervisors
- $\alpha = 0.05$

- Interpret our p-value in context.
- Make a decision and conclusion in response to the research question.

# Making decision and conclusion (answer)

- Assuming that `sex` and `decision` are independent, the probability of observing a difference in promotion rates as or more extreme as we did is 0.029.

- Because the observed p-value of 0.029 is less than our significant level 0.05, we reject $H_0$. The data provide strong evidence of sex discrimination against female candidates by the male supervisors.

- What kind of error could we have made?

# Comparing two proportions

# Running example: CPR

An experiment was conducted, consisting of two treatments on 90 patients who underwent CPR for a heart attack and subsequently went to the hospital. Each patient was randomly assigned to either:

- treatment group: received a blood thinner

- control group: did not receive a blood thinner

For each patient, the outcome recorded was whether they survived for at least 24 hours.

| group | died | survived | total |
|-----------|------|----------|-------|
| control | 39 | 11 | 50 |
| treatment | 26 | 14 | 40 |
| total | 65 | 25 | 90 |

- What is/are the variables(s) here? What types of variables are they?

# Defining hypotheses

The researchers are interested in learning if the blood thinner treatment was effective.

In words, try to determine $H_0$ and $H_A$.

- Let $p_T$ and $p_C$ denote the proportion of patients who survive when receiving the thinner (Treatment) and when not receiving the treatment (Control), respectively

Option 1

- $H_0: p_T \leq p_C$
- $H_A: p_T > p_C$

Option 2 (preferred)

- $H_0: p_T - p_C \leq 0$
- $H_A: p_T - p_C > 0$

# Collect data

Using the data, obtain the observed difference in sample proportions.

| group | died | survived | total |
|---|---|---|---|
| control | 39 | 11 | 50 |
| treatment | 26 | 14 | 40 |
| total | 65 | 25 | 90 |

```
1  p_hat_c <- cpr |>
2    filter(group == "control") |>
3    summarise(p = mean(outcome == "survived")) |>
4    pull()
5  p_hat_t <- cpr |>
6    filter(group == "treatment") |>
7    summarise(p = mean(outcome == "survived")) |>
8    pull()
9  obs_diff <- p_hat_t - p_hat_c
```

- $\hat{p}_C = \frac{11}{50} = 0.22$
- $\hat{p}_T = \frac{14}{40} = 0.35$
- Observed difference:
  $\hat{p}_T - \hat{p}_C = 0.13$

- Is this "convincing evidence" that blood thinner usage after CPR is effective?
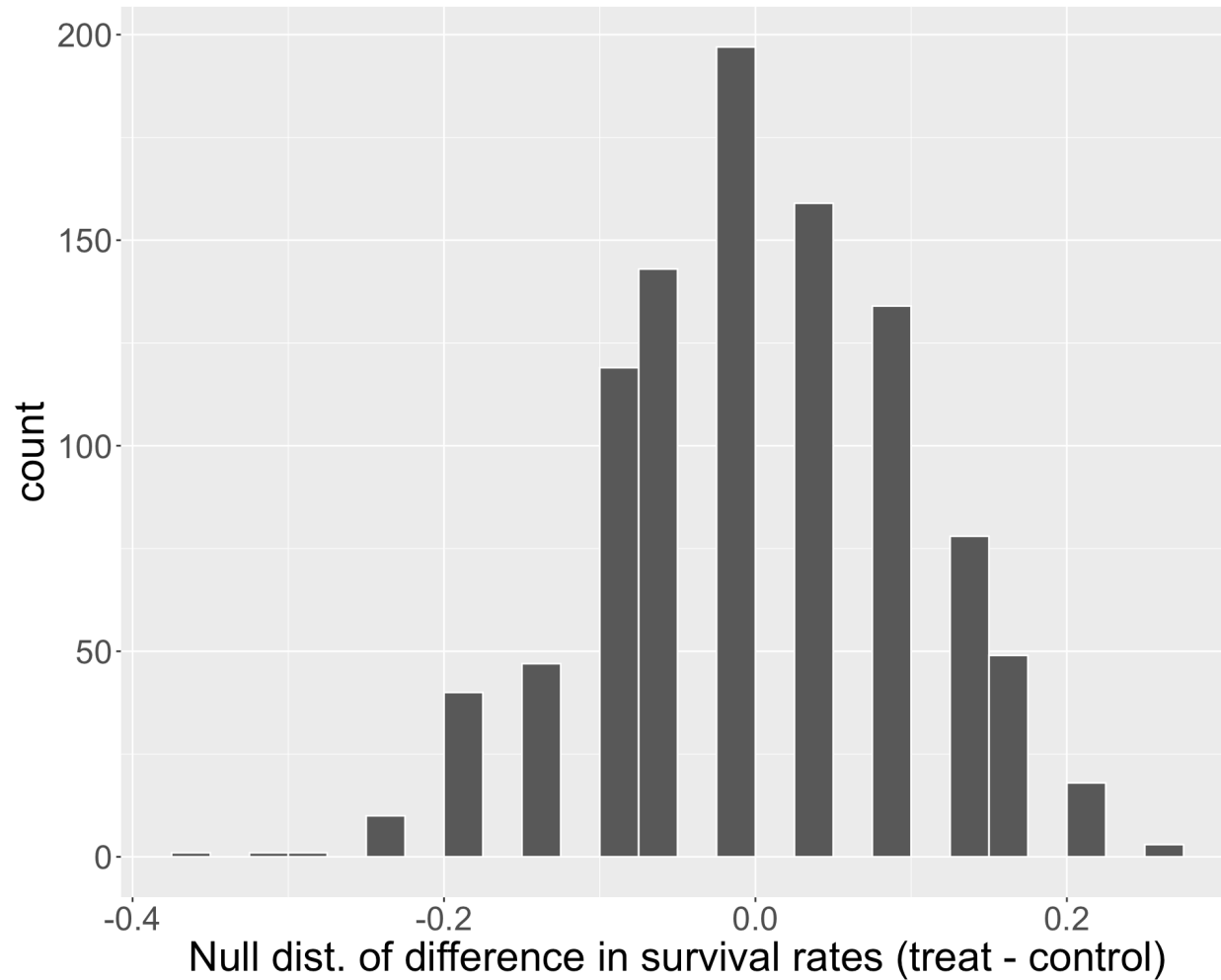- Set $\alpha = 0.05$

# Simulate under null

- We will once again perform a *randomization* test to try and simulate the difference in proportions under $H_0$

  - Under $H_0$, treatment group is no better than control group, so let's simulate assuming that outcome and treatment are independent

- Write down `died` on 65 cards, and `survived` on 25 cards. Then repeat several times:

  - Shuffle cards well

  - Deal out 50 to be Control group, and remaining 40 to be Treatment group

  - Calculate proportions of survival $\hat{p}_{C,sim}$ and $\hat{p}_{T,sim}$

  - Obtain and record the simulated difference $\hat{p}_{T,sim} - \hat{p}_{C,sim}$
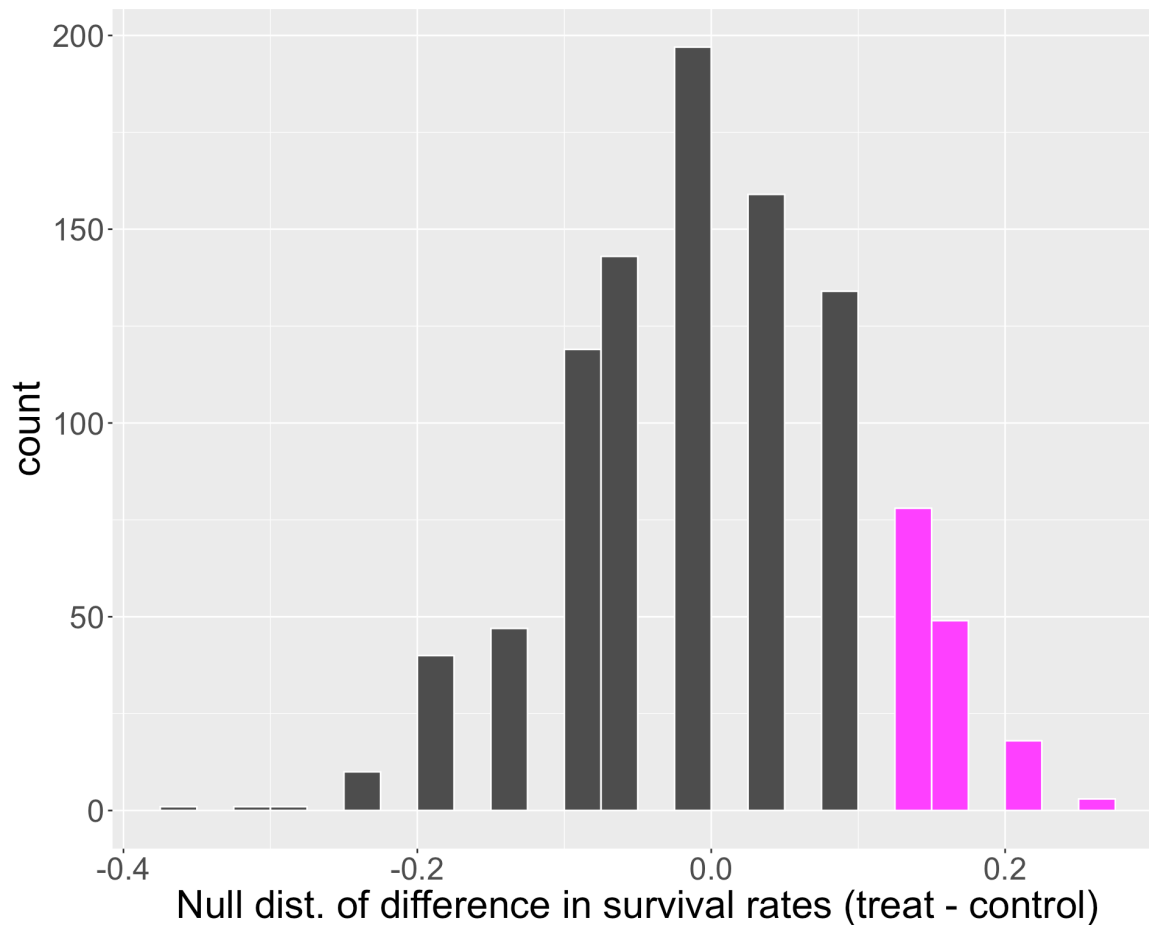
# Simulate under null (code)

Live code

# Visualizing null distribution



How would we obtain the p-value in this problem?

# Calculate p-value



- We observed 148 out of 1000 simulations where the difference in proportions under $H_0$ was as or more extreme than our observed difference of 0.13

- So p-value is 0.148

# Interpret and make conclusion

The researchers are interested in learning if the blood thinner treatment was effective.

Our p-value is 0.148 and our selected significance level was $\alpha = 0.05$.

- Make a decision and conclusion about the research question in context.

# Comprehension questions

- What were the similarities and differences between:
    - hypothesis test for independence
    - hypothesis test for two proportions
- How do the randomization tests today differ from the test for one proportion that we learned last class?