

Unit 2: Foundations for Inference

1. Randomization and Sampling (2.1-2.2)

2/04/2022

Quiz 2 - Exploratory data analysis

Recap from last time

1. R is an awesome language for rapid prototyping
2. dplyr verbs are a useful framework for transforming data (munging)
3. Every r chunk is a paragraph, every line of code is a sentence, pipes are periods.

Key ideas

1. We generally don't want to make claims about samples, we want to make claims about populations (or the processes that generated the samples)
2. We can use randomization to ask what inferences our sample tells us about the population
3. We are always talking about degrees of evidence.
We can never have certainty.

Case study: Gender discrimination

- In 1972, as a part of a study on gender discrimination, 48 male bank supervisors were each given the same personnel file and asked to judge whether the person should be promoted to a branch manager job that was described as “routine”.
- The files were identical except that half of the supervisors had files showing the person was male while the other half had files showing the person was female.
- It was randomly determined which supervisors got “male” applications and which got “female” applications.

Is this an observational study or an experiment?

Experiment

Rosen & Jerdee (1974, Journal of Applied Psychology)

The results

		<i>Promotion</i>		Total
		Promoted	Not Promoted	
<i>Gender</i>	Male	21	3	24
	Female	14	10	24
	Total	35	13	48

Does it look like there is a relationship between gender and promotion?

87.5% of men promoted (21/24), 58.3% of women promoted (14/24)

Practice question: What can we conclude?

We saw a difference of almost 30% in the proportion of men and women promoted. Based on this information, which of the following is true?

- (a) If we were to repeat the experiment we would definitely see that more women got promoted. This was a fluke.
- (b) Promotion is dependent on gender, males are more likely to be promoted. There was gender discrimination in these promotion decisions.
- (c) The difference in the proportions of promoted men and women is due to chance, this is not evidence of gender discrimination.
- (d) Women were less qualified than men, and this is why fewer women got promoted.

Practice question: What can we conclude?

We saw a difference of almost 30% in the proportion of men and women promoted. Based on this information, which of the following is true?

- (a) If we were to repeat the experiment we would definitely see that more women got promoted. This was a fluke.
- (b) *Promotion is dependent on gender, males are more likely to be promoted. There was gender discrimination in these promotion decisions. **Maybe***
- (c) *The difference in the proportions of promoted men and women is due to chance, this is not evidence of gender discrimination. **Maybe***
- (d) Women were less qualified than men, and this is why fewer women got promoted.

Two competing claims

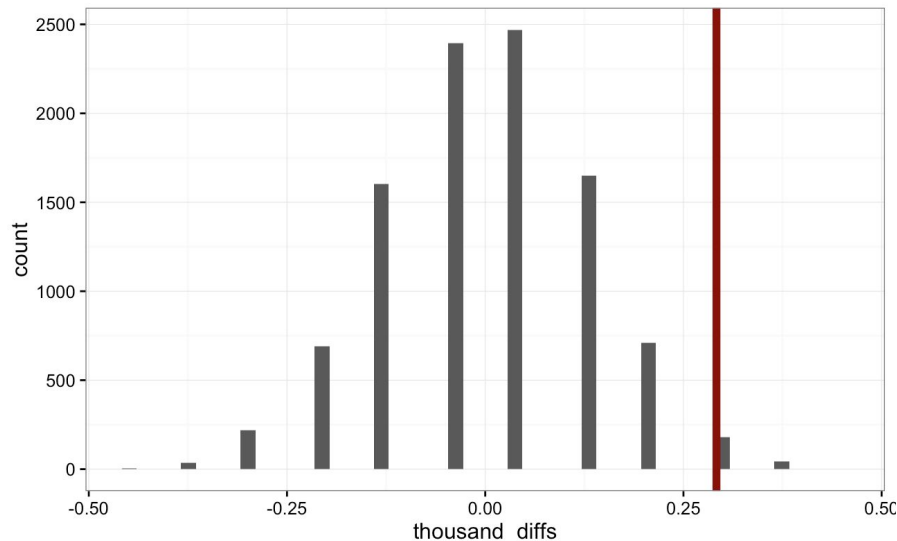
1. “There is nothing going on” (**Null Hypothesis**)
The *process* of promotion is independent of gender
We observed results that *look* dependent due to chance
2. “There is something going on” (**Alternative Hypothesis**)
The *process* of promotion is dependent of gender
We observed results that *look* dependent because they *are dependent*

How can we test the null hypothesis?

What if we generate data from the null hypothesis. What does it look like?

gender	promoted	not_promoted	total
Male	16	8	24
Female	19	5	24
Total	35	13	48

Simulation results



If promotion is independent of gender, we should see a difference like the one we observed *less than 1% of the time*.

Practice question: What can we conclude?

Based on our simulations, what should we conclude?

- (a) Promotion is dependent on gender, males are more likely to be promoted. There was gender discrimination in these decisions.
- (b) The difference in the proportions of promoted men and women is due to chance, this is not evidence of gender discrimination.

Practice question: What can we conclude?

Based on our simulations, what should we conclude?

(a) Promotion is dependent on gender, males are more likely to be promoted. There was gender discrimination in these decisions.

(b) The difference in the proportions of promoted men and women is due to chance, this is not evidence of gender discrimination.

But note we can never be certain!

We can only say that we find a more likely.

Key ideas

1. We generally don't want to make claims about samples, we want to make claims about populations (or the processes that generated the samples)
2. We can use randomization to ask what inferences our sample tells us about the population
3. We are always talking about degrees of evidence.
We can never have certainty.