Statistical Power

1

Overview

- 1. The Central Limit Theorem revisited
- 2. From Null worlds to True effect worlds
- 3. Statistical power
- 4. A power simulation

The Central Limit Theorem revisited

Are action movies better than comedies?

In the previous session on hypothesis testing, we invented a null world:

We simulated a population of 1'000'000 movies with no difference.

```
1 set.seed(1234) # For reproducibility
2
3 imaginary_movies_null <- tibble(
4 movie_id = 1:1000000,
5 rating = sample(seq(1, 10, by = 0.1), size = 1000000, replace = TRUE),
6 genre = sample(c("Comedy", "Action"), size = 1000000, replace = TRUE)
7 )</pre>
```

Sampling distribution

We randomly drew 1,000 samples from this population and calculated the difference between action movies and comedies for each.

We called the distribution of the differences from the different samples the **sampling distribution**

Our sample size was always the same: 20,000.

```
1 n simulations <- 1000
 2 differences <- c() # make an empty vector
 3 sample size <- 20000
 4
 5 for (i in 1:n simulations) {
     # draw a sample of 20'000 films
 6
     imaginary sample <- imaginary movies |>
 7
       sample n(sample size)
 8
     # compute rating difference in the sample
 9
     estimate <- imaginary sample |>
10
       group by(genre) |>
11
12
       summarize(avg rating = mean(rating)) |>
       summarise(diff = avg rating[genre == "Action"] - avg rating[genre == "Comedy"]) %>%
13
       pull(diff)
14
15
     differences[i] <- estimate
16
17 }
```

The Central Limit Theorem (part I)

The sampling distribution approximates the shape of a normal distribution



The Central Limit Theorem (part II)

This is what the sampling distribution looks like with samples of size 10,000



The Central Limit Theorem (part II)

And with samples of size 1,000



The Central Limit Theorem (part II)

The smaller the sample, the larger the standard deviation of the sampling distribution



How is this relevant for hypothesis testing?

Imagine we find an effect of -0.2 in our sample

In a sample based on 20'000 movies, that seems reasonably unlikely in a Null world



But the same effect in a sample of 1000 seems not so unlikely...



The Central Limit Theorem

(part I)

• The sampling distribution approximates the shape of a normal distribution

(part II)

• The smaller the sample, the larger the standard deviation of the sampling distribution

From Null worlds to True effect worlds

For example, let's imagine a world where the true difference between action and comedy movies is -0.2

```
1 # Generate Comedy and Action movie ratings using truncated normal distributions
 2 imaginary movies true <- tibble(</pre>
     movie id = 1:1000000,
 3
     genre = sample(c("Comedy", "Action"), size = 1000000, replace = TRUE),
 4
     rating = ifelse(
 5
       genre == "Comedy",
 6
       rtruncnorm(1000000, a = 1, b = 10, mean = 6.0, sd = 2),
 7
       rtruncnorm(1000000, a = 1, b = 10, mean = 5.8, sd = 2)
 8
 9
     )
10 )
```

Let's plot the population data

► Code



It worked, we see our imagined effect of ~ -0.2

► Code



Imagine we look at a sample of 1000 movie ratings.

```
1 n simulations <- 1000
 2 differences <- c() # make an empty vector
  sample size <- 1000</pre>
 3
 4
 5 for (i in 1:n simulations) {
     # draw a sample of 20'000 films
 6
 7
     imaginary sample <- imaginary movies |>
       sample n(sample size)
 8
     # compute rating difference in the sample
 9
     estimate <- imaginary sample |>
10
11
       group by(genre) |>
       summarize(avg rating = mean(rating)) |>
12
       summarise(diff = avg rating[genre == "Action"] - avg rating[genre == "Comedy"]) %>%
13
14
       pull(diff)
15
16
     differences[i] <- estimate
17 }
```

This is what our sampling distribution would look like



Distribution of Rating Differences by Sample Size

Now imagine we do a hypothesis test.

We can simulate both a null world and a true effect world (so far, nothing new)

Null world

```
1 # Null world population
2 imaginary_movies_null <- tibble(
3 movie_id = 1:1000000,
4 rating = sample(seq(1, 10, by = 0.1), size = 100
5 genre = sample(c("Comedy", "Action"), size = 100
6 )</pre>
```

True effect world

1	# True effect world population			
2	<pre>imaginary_movies_true <- tibble(</pre>			
3	movie_id = 1:1000000,			
4	<pre>genre = sample(c("Comedy", "Action"), size = 1000000, repla</pre>			
5	<pre>rating = ifelse(</pre>			
6	<pre>genre == "Comedy",</pre>			
7	rtruncnorm(1000000, a = 1, b = 10, mean = 6.0, sd = 2),			
8	rtruncnorm(1000000, a = 1, b = 10, mean = 5.8, sd = 2)			
9)			
10)			
		J .		

And make a sampling distribution with sample size 1000 for both worlds (also nothing new)

```
1 n simulations <- 1000
 2 differences <- c() # make an empty vector</pre>
 3 sample size <- 1000
 4
   data <- imaginary movies true # replace with imaginary movies null</pre>
 5
 6
   for (i in 1:n simulations) {
 7
     # draw a sample of films
 8
     imaginary sample <- data |>
 9
        sample n(sample size)
10
     # compute rating difference in the sample
11
     estimate <- imaginary sample |>
12
       group by(genre) |>
13
       summarize(avg rating = mean(rating)) |>
14
       summarise(diff = avg rating[genre == "Action"] - avg rating[genre == "Comedy"]) %>%
15
       pull(diff)
16
17
      differences[i] <- estimate
18
19 }
```

We can plot both simulated worlds together.



Sampling distributions for worlds with and without true difference

Sample Size = 1000

Let's bring both worlds on the scale of a standard normal distribution, dividing their respective standard deviation



Sample Size = 1000

With a sample size of 1000, in many cases, we would say: "This could have occurred in a Null World"



Two errors

Type I error

(false positive)

Vou're

Hypothesis Testing

Type II error (false negative)

You're not pregnant

Statistical Power

Two errors

	Hypothesis Testing	Power Analysis
Aim	Rule out that we observe something just by chance.	Ensure that we would find an effect.
Error question	"What are the chances that we find an effect at least this large in our sample, given that there is no effect in the population?"	"What are the chances that we do not find a statistically significant effect in our sample, although there is a certain effect in the population?"
Typical threshold for acceptable	<i>α</i> = 5 %	β=20%

error

Statistical Power

Statistical Power

Statistical power is the probability of detecting an effect with a hypothesis test, given a certain effect size

 $({\rm Or}\ 1-\beta)$

Your turn #1: Calculating power

- 1. Create your true effect world: Simulate a population with a true difference of -0.5 between action and comedy movies.
- 2. Get your sampling distribution: Simulate 1000 random samples of size n = 1000 and store the results.
- 3. Plot your Sampling distribution (use data.frame() to turn your vector into a data frame that can be read by ggplot)
- 4. Prepare for hypothesis testing: bring the results on scale of the standard normal distribution (divide by the standard deviation of your distribution)
- 5. Check the transformation: plot the new, standardized values
- 6. Calculate your power: Check how many (standardized) differences are below the 5% threshold, i.e. smaller than or equal to -1.96 (hint: use mutate() in combination with ifelse to create a new variable significant that takes the values of TRUE or FALSE. Then use summarize() and sum() to calculate the share). Are you above the 80% power threshold?



1. Create your true effect world: Simulate a population with a true difference of -0.5 between action and comedy movies.

```
1 # True effect world population
 2 imaginary movies true <- tibble(</pre>
     movie id = 1:1000000,
 3
     genre = sample(c("Comedy", "Action"), size = 1000000, replace = TRUE),
 4
 5
     rating = ifelse(
       genre == "Comedy",
 6
       rtruncnorm(1000000, a = 1, b = 10, mean = 6.0, sd = 2),
 7
       rtruncnorm(1000000, a = 1, b = 10, mean = 5.5, sd = 2)
 8
 9
     )
10 )
```

2. Get your sampling distribution: Simulate 1000 random samples of size n = 1000 and store the results.

```
1 n simulations <- 1000
 2 differences <- c() # make an empty vector
   sample size <- 1000</pre>
 3
 4
   data <- imaginary movies true # replace with imaginary movies null</pre>
 5
 6
   for (i in 1:n simulations) {
 7
     # draw a sample of films
 8
     imaginary sample <- data |>
 9
       sample n(sample size)
10
     # compute rating difference in the sample
11
     estimate <- imaginary sample |>
12
       group by(genre) |>
13
       summarize(avg rating = mean(rating)) |>
14
       summarise(diff = avg rating[genre == "Action"] - avg rating[genre == "Comedy"]) %>%
15
       pull(diff)
16
17
18
      differences[i] <- estimate
19 }
```

3. Plot your Sampling distribution (use data.frame() to turn your vector into a data frame that can be read by ggplot)





4. Prepare for hypothesis testing: bring the results on scale of the standard normal distribution (divide by the standard deviation of your distribution)

1 differences_standardized <- differences/sd(differences)</pre>

5. Check the transformation: plot the new, standardized values

```
1 ggplot(data.frame(differences_standardized), aes(x = differences_standardized)) +
2 geom_histogram() +
3 labs(title = "Sampling Distribution of Rating Differences",
4 x = "Mean Rating Difference (Action - Comedy)",
5 y = "Frequency") +
6 theme_minimal()
```



Mean Rating Difference (Action - Comedy)

6. Calculate your power: Check how many (standardized) differences are below the 5% threshold, i.e. smaller than or equal to -1.96 (hint: use mutate() in combination with ifelse to create a new variable significant that takes the values of TRUE or FALSE. Then use summarize() and sum() to calculate the share). Are you above the 80% power threshold?

sum_significant share_significant
 975
 0.975

1

A power simulation

The central limit theorem (again)

The larger the sample size, the more statistical power



An example of a power analysis

For an example, head over to the guide on power analysis on the course website.

That's it for today :)