

Eltecon Data Science Course by Emarsys

Measuring uncertainty

András Bérczi

October 14, 2020

Homeworks from last week

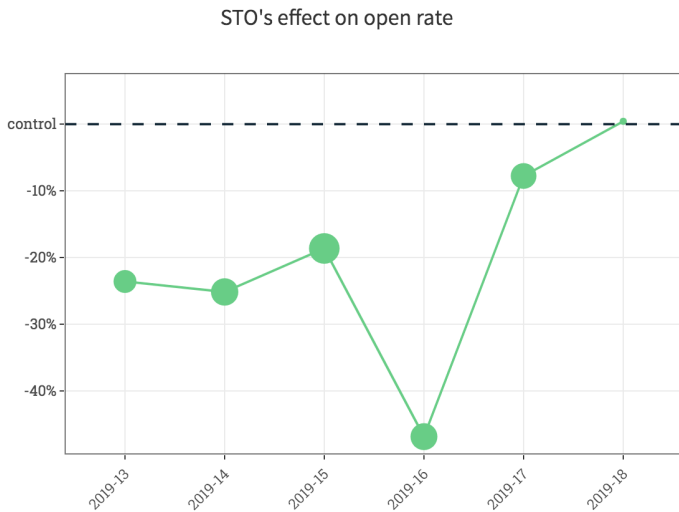
Any questions about final project?

Measuring uncertainty

We can always measure something from our data...

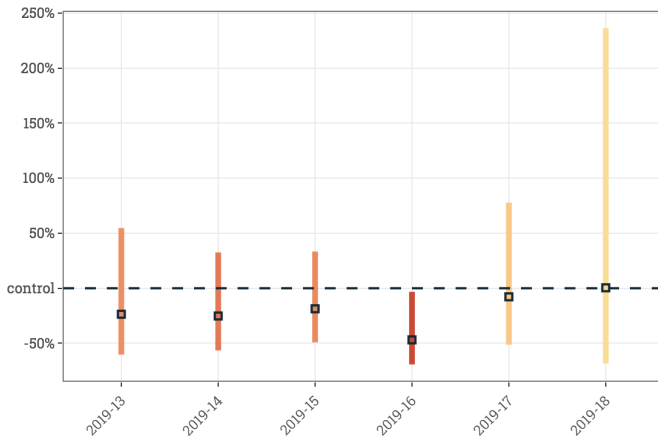
... but how sure can we be about our measurement?

We can always measure something from our data...



But not necessarily significant!

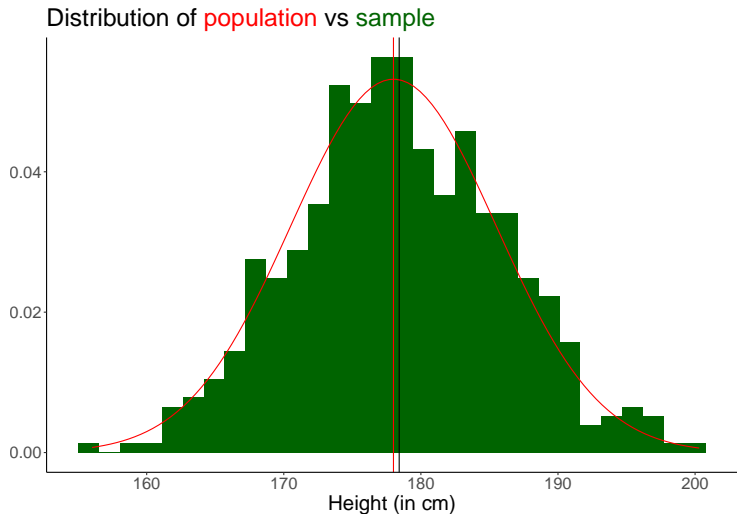
STO's effect on open rate



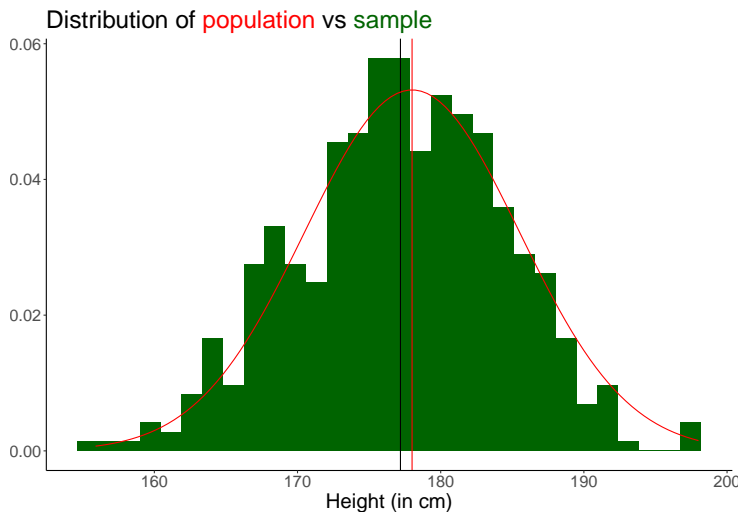
Why do we have uncertainty in the measurement?

- If you knew the whole population, there wouldn't be uncertainty in your measurement
- But we only see 1 'segment' of the data = we have a sample of the population

Sampling from a population



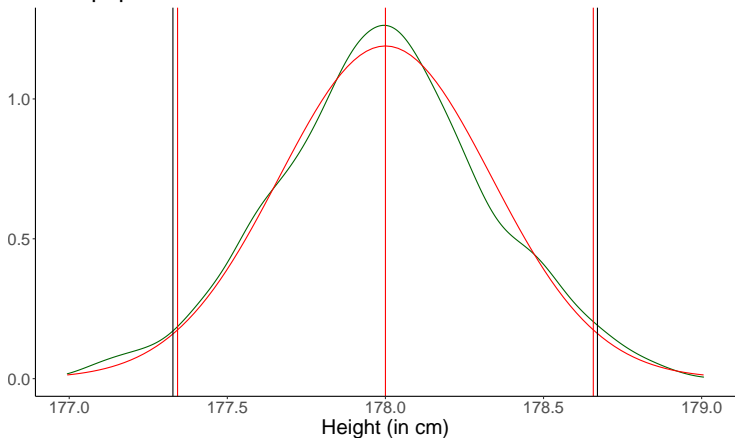
Sampling from a population



Sampling from a population

Distribution of sample means

Distribution of **sample means**
compared to **normal distribution** with 'true' parameters
from population



Distribution of sample means

Normal distribution with parameters:

\bar{x} is the sample mean,

s is the standard deviation of the sample distribution,

n is the sample size,

add 95% 'CI' interval as:

$$\bar{x} \pm 1.96 * \frac{s}{\sqrt{n}}$$

Law of Large Numbers

The average of the results obtained from a large number of trials should be close to the expected value and will tend to become closer to the expected value as more trials are performed. - Wikipedia

Central Limit Theorem

When independent random variables are added, their properly normalized sum tends toward a normal distribution (informally a bell curve) even if the original variables themselves are not normally distributed. - Wikipedia

What are Confidence Intervals?

- The normal table gives us the fact that $P[-1.96 < Z < 1.96] = 0.95$.
- With a sample of n values from a population with mean μ and standard deviation σ , the Central Limit theorem gives us the result that $Z = \sqrt{n} \frac{\bar{x} - \mu}{\sigma}$ is approximately normally distributed with mean 0 and with standard deviation 1.

What are Confidence Intervals?

Start from $P[-1.96 < Z < 1.96] = 0.95$ and then substitute for Z the expression $\sqrt{n} \frac{\bar{X} - \mu}{\sigma}$.

This will give us

$$P\left[-1.96 < \sqrt{n} \frac{\bar{X} - \mu}{\sigma} < 1.96\right] = 0.95$$

We can rewrite this as

$$P\left[-1.96 \frac{\sigma}{\sqrt{n}} < \bar{X} - \mu < 1.96 \frac{\sigma}{\sqrt{n}}\right] = 0.95$$

Now subtract \bar{X} from all items to get

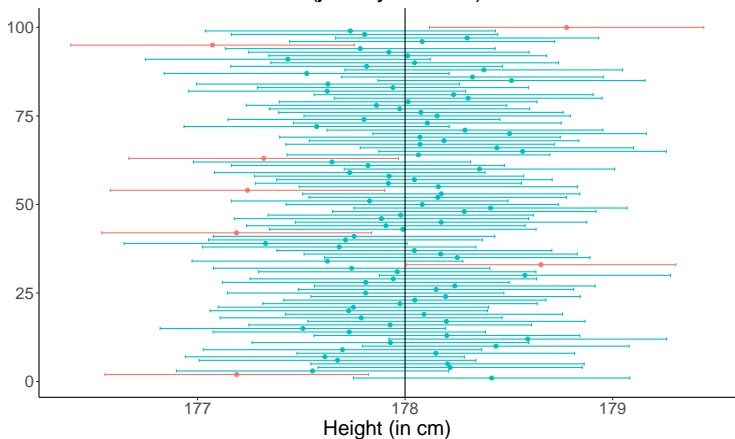
$$P\left[-\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} < -\mu < -\bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right] = 0.95$$

Multiply by -1 (which requires reversing inequality direction) to obtain

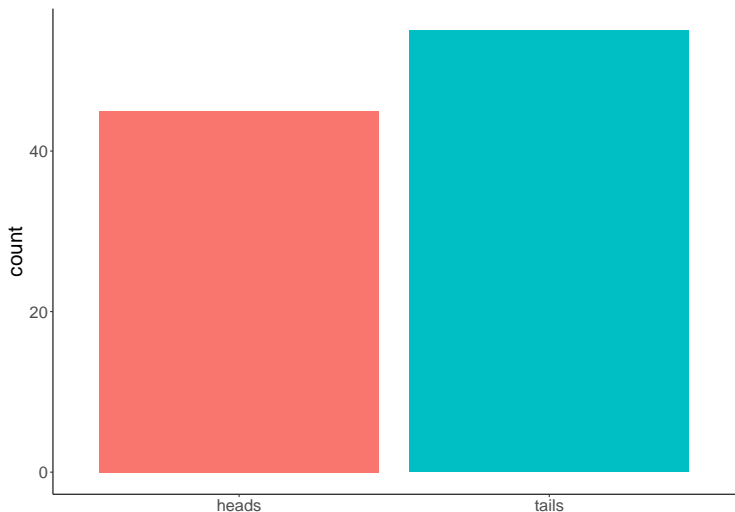
$$P\left[\bar{X} + 1.96 \frac{\sigma}{\sqrt{n}} > \mu > \bar{X} - 1.96 \frac{\sigma}{\sqrt{n}}\right] = 0.95$$

What are Confidence Intervals?

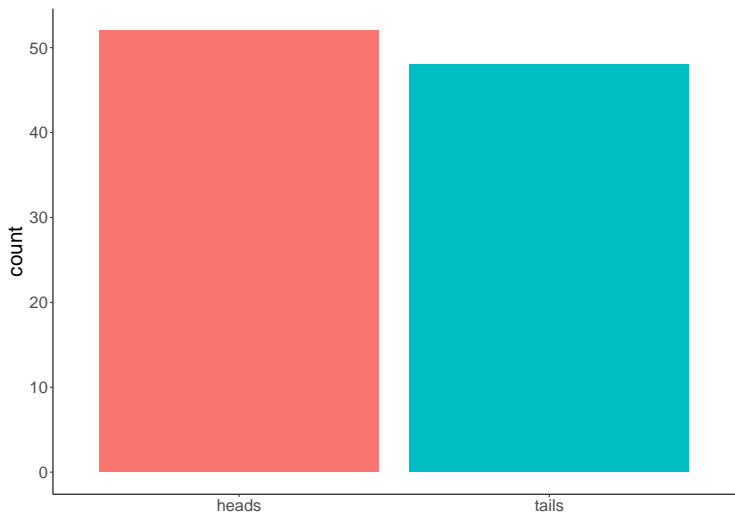
Mean and CI from different samples:
About 95% of the CIs **contains** the true mean,
but 5% **does not contain** (just by chance)



Distribution of sample means with different distributions

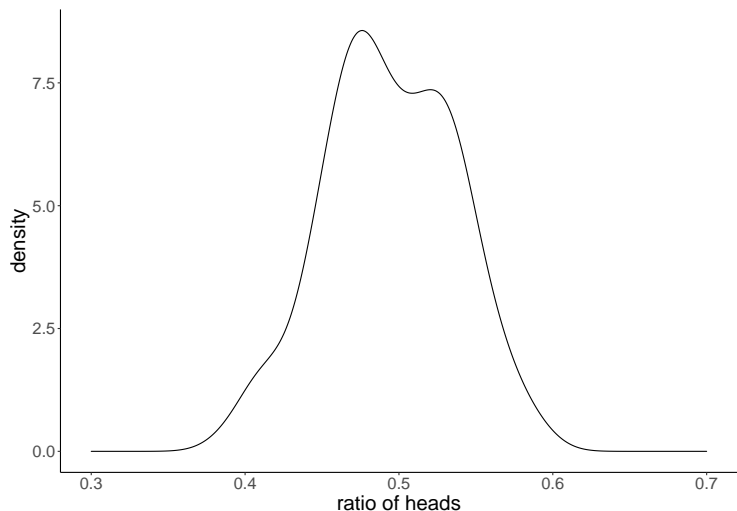


Distribution of sample means with different distributions

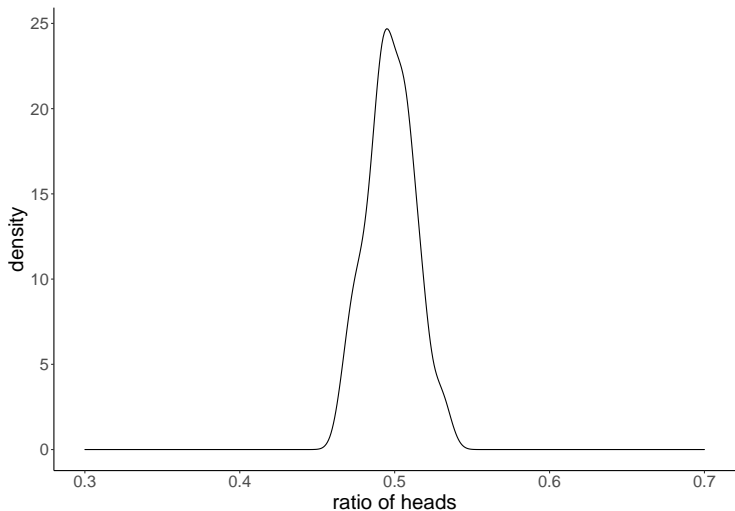


Distribution of sample means with different distributions

Distribution of sample means with different distributions



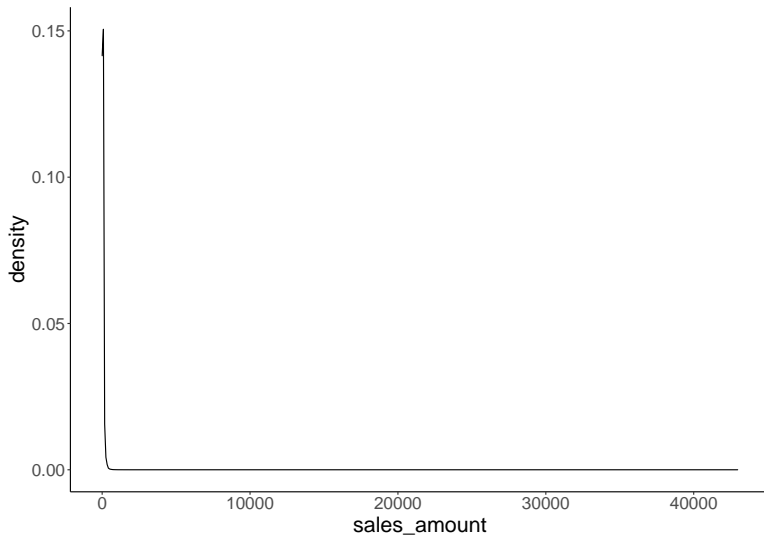
Why does sample size matter?



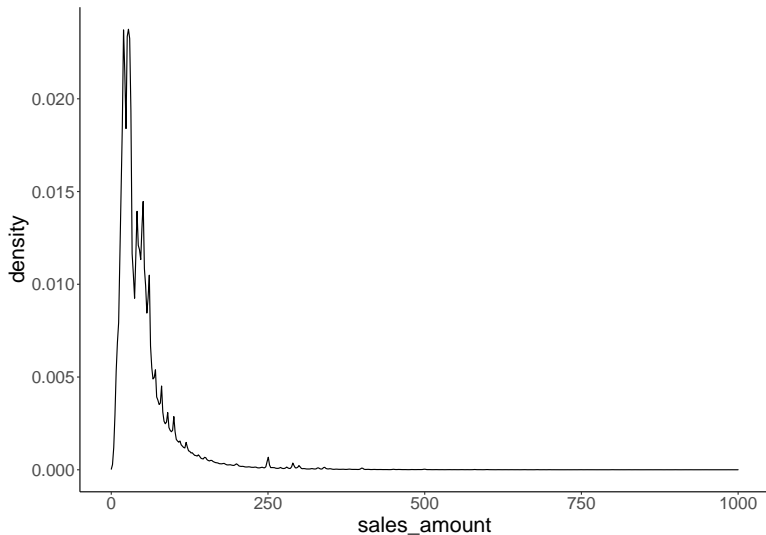
What are the key assumptions?

- we have a random sample from the distribution (= independent and identically distributed random variables are drawn)
- finite variance / distribution is not 'long tailed'

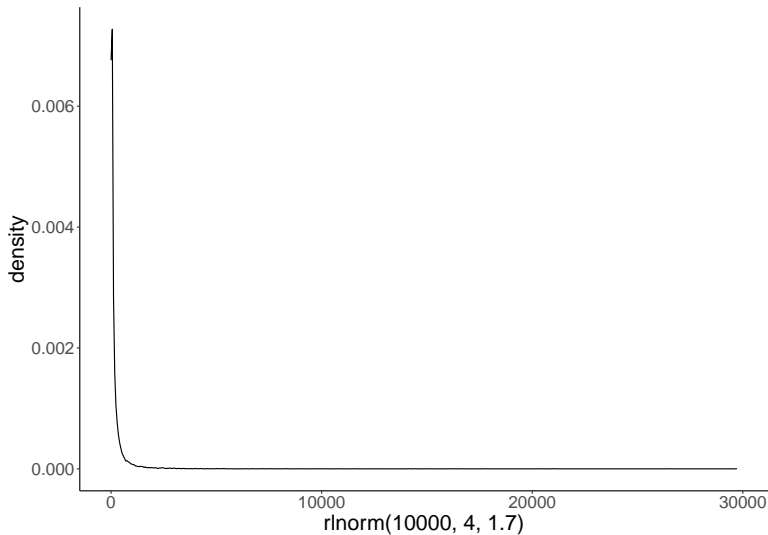
What if variance is infinite?



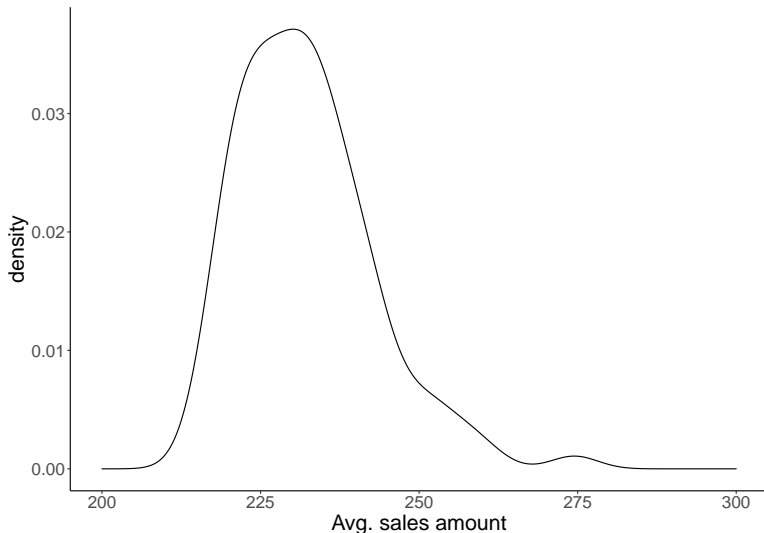
It stays very skewed even if we zoom in



What if variance is infinite?



Distribution of avg. sales amount from samples

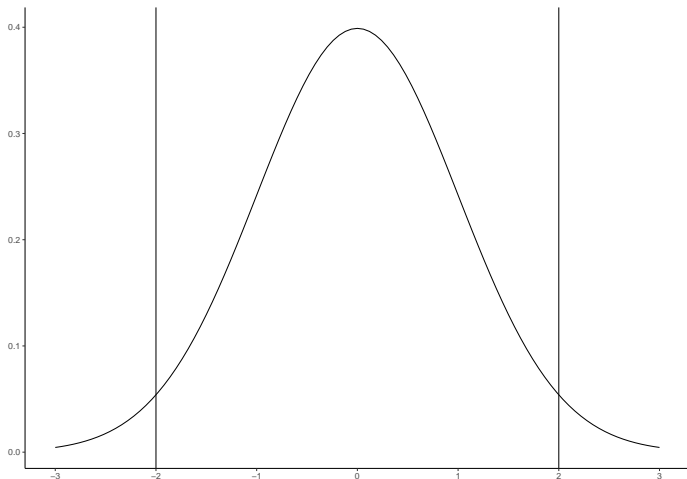


How can we calculate the uncertainty of our measurement?

- **Based on variance of known distribution**
- Monte-Carlo method
- Bootstrapping
- (and other methods as well of course)

Calculate uncertainty based on variance

Calculate interval to show how wide range we need in order to be 'sure' that we have the real value included.



How to calculate uncertainty from sampling distribution

By CLT + LLN, you can add uncertainty to your point estimate (for a 95% Confidence Interval), such as:

$$\bar{x} \pm 1.96 * \frac{\sigma}{\sqrt{n}}$$

\bar{x} is the sample mean,

σ standard deviation of the population,

n is the sample size

Standard Error

We want to calculate the variance of means of different samples.

Since the standard deviation for the population is rarely known, we estimate SE as:

$$\text{Standard Error} = \frac{s}{\sqrt{n}}$$

s standard deviation of the sample,

n is the sample size

Standard Error vs Standard Deviation

Looks similar, but it's not the same!

Derivation for Standard Error:

This comes from the fact that $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + \text{Cov}(X, Y)$ and for a constant a , $\text{Var}(aX) = a^2 \text{Var}(X)$.

Since we are assuming that the individual observations are independent the $\text{Cov}(X, Y)$ term is 0 and since we assume that the observations are identically distributed all the variances are σ^2 . So

$$\text{Var}\left(\frac{1}{n} \sum X_i\right) = \frac{1}{n^2} \sum \text{Var}(X_i) = \frac{1}{n^2} \times \sum \sigma^2 = \frac{n}{n^2} \sigma^2 = \frac{\sigma^2}{n}$$

And when we take the square root of that (because it is harder to think on the variance scale) we get $\frac{\sigma}{\sqrt{n}}$.

Derivation is from here.

Calculating uncertainty for proportion

Distribution of sample means of Bernoulli distr. is normally distributed with parameters:

$$\bar{x} = p$$

$$\sigma = \sqrt{\frac{p * (1 - p)}{n}}$$

Calculating uncertainty for absolute uplift

Distribution of difference of sample means is also normally distributed with parameters:

$$\bar{x}_{diff} = \bar{x}_1 - \bar{x}_2$$

$$\sigma_{diff} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Calculating uncertainty

We can use the same formula:

$$\bar{x} \pm 1.96 * \frac{\sigma}{\sqrt{n}}$$

Just use the right equations for \bar{x} and σ , depending on the distribution!

Calculating uncertainty for relative uplift

Distribution of ratio of sample means is ??

How to calculate it? -> Next time :)

Calculate uncertainty for a point estimate

contact_id	group	num_send	num_open	num_click	sales_amou
1	treatment	0	0	0	
2	treatment	3	0	0	
3	treatment	2	1	0	
4	treatment	3	0	0	
5	treatment	0	0	0	
6	treatment	0	0	0	

How to calculate the open rate absolute uplift with uncertainty!

Let me show that!

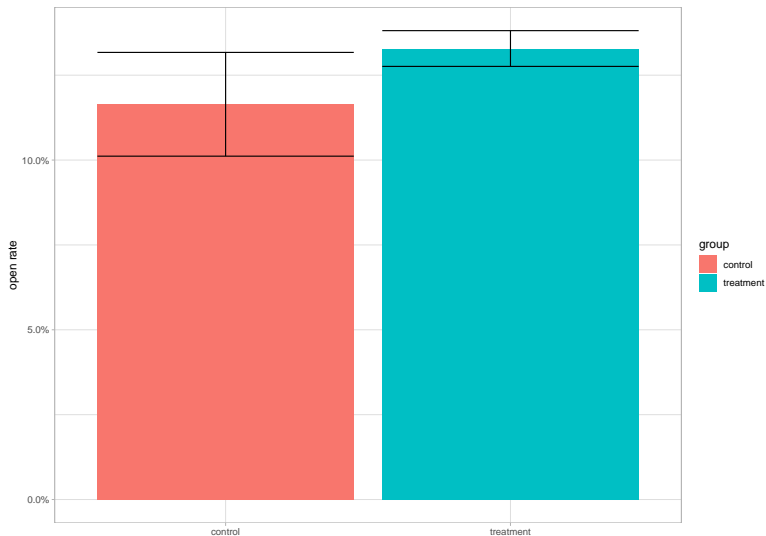
Now your turn!

Calculate the open rate and the uncertainty for both groups!

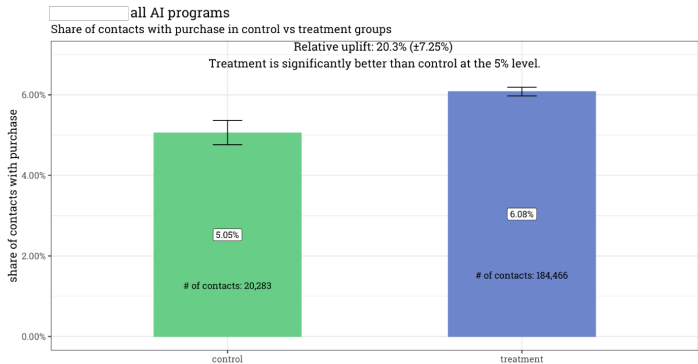
Results from an experiment - How to plot uncertainty?



How to plot uncertainty



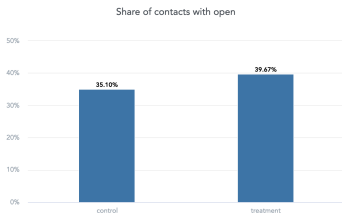
How to plot uncertainty - some examples



Contact behaviour is measured for 7 days from entering the program (currently until May 22, 2019)

How to plot uncertainty - some examples

Share of contacts with open



Relative uplift

13.03%

Likelihood that treatment group outperforms control

100.0%

Control contacts

117,571

Treatment contacts

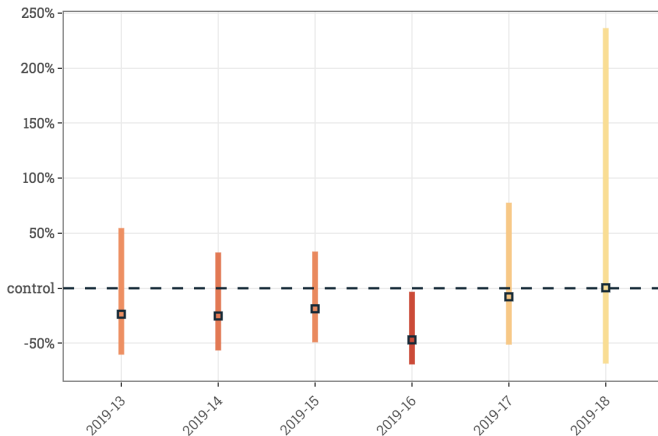
1,057,814

This result means

strong evidence for the program bringing added value

How to plot uncertainty - some examples

STO's effect on open rate



Now your turn!

1. Calculate the click rate and the uncertainty!
2. Plot the results! What do you see on the plots? Are the results significant?

Bayesian uncertainty

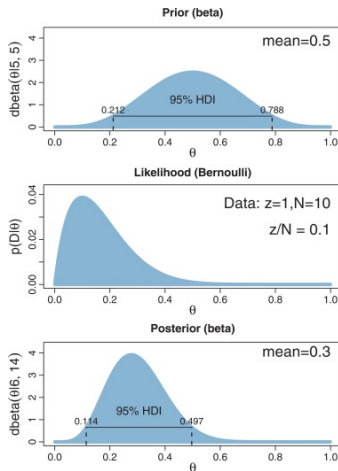
Confidence Interval: If we would resample from our population, 95% of times the Confidence Interval will contain the true, unknown parameter.

Credible Interval: There is a 95% chance that this interval contains the true parameter.

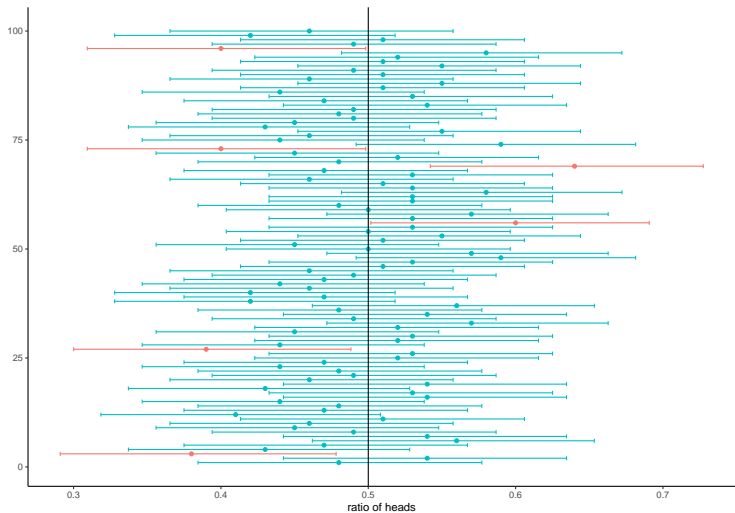
Difference in calculating uncertainty

Confidence Interval: Based on sampling distribution of means

Credible Interval: Based on data and prior belief



95% Credible interval



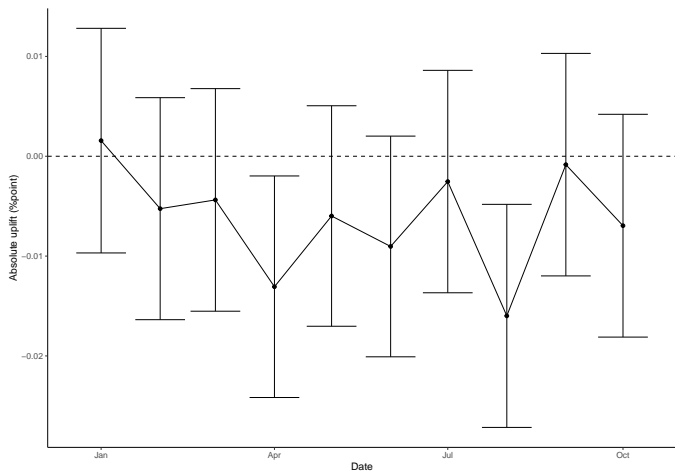
Confidence Interval vs Credible Interval

- Credible Interval is easier to understand
- Credible Interval gives smaller interval if we have some prior knowledge
 - eg.: Use group averages as prior for contact level data
- With wrong prior provided, our posterior distribution is going to be wrong as well!
- With a lot of data or with non-informative priors, the two intervals are about the same

head_ratio	cred_int_lower	cred_int_higher
0.48	0.4320385	0.5480755

head_ratio	conf_int_lower	conf_int_higher
0.48	0.3820784	0.5779216

Calculate uncertainty over time



Your turn!

Try to re-create the plot seen before, using
`experiment_results_over_time.csv`!

Takeaways

- Always show uncertainty
- Think about your audience

Homework for next week and presenters

Homework for next week

- ❶ Figure out your research question / hypothesis!
- ❷ Show (plot) that if there is a difference between two/multiple groups, add uncertainty as well!
 - ❶ If it makes sense for your research question, do it for that! Eg.: Is there a difference in house prices on Buda compared to Pest?
 - ❷ If not, use `experiment_result_HW.csv`: calculate the the point estimates and add uncertainty to one of your KPIs (defined by you in last week's homework) for both periods! Include the *absolute* uplift in the subtitle with confidence intervals!

Presenting next week

Both Márton - Kamenár Gyöngyvér

Emerson, Ian - Ralbovski Judit

Bat-Erdene, Boldmaa - Kashirin, Andrey