







We face the problem that we want to investigate, whether some universally quantified statement holds, while we only have access to a subset of the overall population of entities the statement is quantifying over. This subset of the population we have access to is called the sample.

 \Rightarrow Inferential statistics is about what we can reasonably say about the population given a sample.





Sampling Distribution of the Sample Mean



The Gist

The sample mean will be approximately normally distributed for large sample sizes, regardless of the distribution from which we are sampling.

Nebel, Lindner, Engesser, Kuhnert, Wächter - Social Robotics

6 / 27

UNI FREIBURG

Mean of the Sampling Distribution of the Sample Mean

Let $X_1, ..., X_N$ be *N* independently drawn observations from a distribution with mean μ and variance σ^2 . Thus, $E[X_i] = \mu$ for all *i*. Let's derive $E[\overline{X}]$, which we call the mean of the sampling distribution of the sample mean (also written as $\mu_{\overline{X}}$):

$$\mathbb{E}[\overline{X}] = \mathbb{E}[\frac{1}{N}\sum_{i}^{N}X_{i}] = \frac{1}{N}\mathbb{E}[\sum_{i}^{N}X_{i}] = \frac{1}{N}\sum_{i}^{N}\mathbb{E}[X_{i}] = \frac{1}{N}N\mu = \mu$$

Variance of the Sampling Distribution of the Sample Mean

Let $X_1, ..., X_N$ be *N* independently drawn observations from a distribution with mean μ and variance σ^2 . Thus, $Var[X_i] = \sigma^2$ for all *i*. Let's derive $Var[\overline{X}]$, which we call the variance of the sampling distribution of the sample mean (also written as $\sigma_{\overline{X}}^2$):

$$Var[\overline{X}] = Var[\frac{1}{N}\sum_{i}^{N}X_{i}] = (\frac{1}{N})^{2}Var[\sum_{i}^{N}X_{i}] = (\frac{1}{N})^{2}\sum_{i}^{N}Var[x_{i}] = (\frac{1}{N})^{2}N\sigma^{2} = \frac{\sigma^{2}}{N}$$

- Hence, the standard deviation of the sampling distribution of the sample mean is $\sigma_{\overline{X}} = \frac{\sigma}{\sqrt{N}}$.
- $\sigma_{\overline{X}}$ is also called the Standard Error.

Nebel, Lindner, Engesser, Kuhnert, Wächter – Social Robotics



- Suppose we know the population mean μ and standard deviation σ.
- Can we find boundaries within which we believe the mean of a sample of size N will fall with 95% probability?
- We know how our sample means are distributed, viz., $\overline{X} \sim \mathcal{N}(\mu, \frac{\sigma^2}{M})$
 - The lower boundary \overline{X}_{low} will be 1.96 standard errors below μ , and the upper boundary \overline{X}_{up} will be 1.96 standard errors above μ .

$$\mu - \overline{X}_{low} = 1.96 \times \frac{\sigma}{\sqrt{N}} \Rightarrow \overline{X}_{low} = \mu - 1.96 \times \frac{\sigma}{\sqrt{N}}$$
$$\overline{X}_{up} - \mu = 1.96 \times \frac{\sigma}{\sqrt{N}} \Rightarrow \overline{X}_{up} = \mu + 1.96 \times \frac{\sigma}{\sqrt{N}}$$





9 / 27

UNI FREIBURG

UNI FREIBURG



Application: Hypothesis Testing



- This is useful, because given μ and σ , we can compute the probability that some sample of size *N* with mean \overline{X} stems from that population!
- We already know how we can judge whether some value from a normal distribution is 'usual' or rather 'extreme': z-Scores!
- Hence, we can judge a sample mean as 'usual' or 'extreme' by computing its z-Score.
- Let's see how we can use this for hypothesis testing!

UNI FREIBURG





Report

We recorded the number of interactions with our robot per day for nine days (N = 9). The number of interactions ranged from 35 to 150 (\overline{X} = 65.11, *s* = 33.59, 95% CI [43.16, 87.05]).

Remember the data 35, 50, 50, 50, 56, 60, 60, 75, 150.

Nebel, Lindner, Engesser, Kuhnert, Wächter - Social Robotics

14 / 27

BURG

FREI

Very First Hypothesis Test

Suppose you have been deploying a robot (Robo-One) in your museum. You have recorded the number of interaction for a very long time, such that you can assume the collected mean and variance of the number of interactions to be the population mean $\mu_0 = 40$ and standard deviation $\sigma_0 = 4$. You have now bought a fancy new version of the robot, viz., Robo-Two. Your Hypothesis is that Robo-Two will generate much more interactions compared to Robo-One.

- Hypothesis H₁: Robo-Two generates more interactions than Robo-One.
- *H*₁ is of type (difference, directional)
- Can be written as H₁: μ > μ₀, i.e., the population mean for interactions with Robo-Two (μ) is bigger than the population mean for interactions with Robo-One (μ₀), i.e., people generally interact more with Robo-Two than with Robo-One.

Very First Hypothesis Test



17/27

- The trick of inferential statistics is to first assume that the negation of H₁ is the case, which is called the Null-Hypothesis, written H₀.
- Then, we collect the data (viz., our sample)
- Subsequently, we show that our sample is so unlikely under H_0 that we are allowed to reject H_0 in favor of H_1 .
 - In the example: $H_1: \mu > \mu_0, H_0: \mu \le \mu_0$.

Nebel, Lindner, Engesser, Kuhnert, Wächter – Social Robotics





■ We compute the z-Score to assess how far our sample mean 42 is from the mean of the sampling distribution of the sample mean, 40: $z = (42 - 40)/\frac{4}{4} = (42 - 40) = 2$.

Nebel, Lindner, Engesser, Kuhnert, Wächter – Social Robotics

18/27



Report

The number of interactions with Robo-Two is significantly higher than the number of interactions with Robo-One (z = 2.0, p = 0.0228).

Because the hypothesis was directional, we checked if the z-Score of \overline{X} was $z_{.95} = 1.65$ or higher. The is called a one-tailed test. The p-Value is just the probability $P(z \ge 2.0) = 0.0228$. This is below the significance level $\alpha = 0.05$.

20 / 27

Second Hypothesis Test



- This time, our H_1 hypothesis was that there is a difference between Robo-One and Robo-Two: $H_1 : \mu \neq \mu_0$.
- The null-hypothesis then is $H_0: \mu = \mu_0$.
- We will reject *H*₀, if *µ* is too low or too high. Thus, we split our 5% significance level into two (2.5% at the lower end, and 2.5% at the higher end).
- We thus check if the z-Value is below z_{.025} = -1.96 or above z_{.975} = 1.96. This is a two-tailed test.
- As our z-Score was 2, we will also reject H_0 this time.

Nebel, Lindner, Engesser, Kuhnert, Wächter – Social Robotics

21 / 27







Report

The number of interactions with Robo-Two and with Robo-One differ significantly (z = 2.0, p = 0.044).

Because the hypothesis was non-directional, we compute the probability to observe a z-Score at least as extreme as 2.0 (in both directions). The probability is thus $P(z \ge 2.0) + P(z \le -2.0) = 0.0228 + 0.0228 = 0.0456$. This is below the significance level $\alpha = 0.05$.

Nebel, Lindner, Engesser, Kuhnert, Wächter – Social Robotics

22 / 27



Report

The hypothesis H_1 stating that the number of interactions with Robo-Two will be less than with Robo-One was not supported (z = 2.0, p = 0.9772).

This time we look only at the lower end, thus, we compute the probability $P(z \le 2.0) = 0.9772$, which clearly is above the significance level $\alpha = 0.05$.

Type-I and Type-II Errors



25 / 27

- Our decisions to reject H_0 or not are based on probabilities! We see that our sample would be rather unusual if H_0 were true, thus we reject H_0 . But it could be that we just had an unusual sample by chance. If we decide to reject H_0 although H_0 is actually true, then we commit a Type-I Error. Using the 5% significance level, we have a 5% chance per rejected H_0 hypothesis that we were wrong.
- If we instead reject H₁ although H₀ is wrong, then we commit a Type-II Error. This can happen, when there is an effect in the population, but our sample size was too small to detect that effect.



- Note that we have assumed that μ and σ² are known to us a-priori, or can be reasonably be approximated in case of a sufficiently big sample size.
- In many applications, we will not be able to enjoy this luxury.
- Therefore, we will learn about other test statistics, as well. But the main idea is the same, most of the time.

Nebel, Lindner, Engesser, Kuhnert, Wächter – Social Robotics

26 / 27

Nebel, Lindner, Engesser, Kuhnert, Wächter – Social Robotics

