

# Lectures

MVE055 / MSG810

Mathematical statistics and discrete mathematics

---

Moritz Schauer

Last updated September 28, 2022

GU & Chalmers University of Technology

# Central limit theorem/CLT

---

## Recall

---

If  $X \sim N(\mu, \sigma^2)$ , then

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

If  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  independent, then

$$\bar{X}^{(n)} \sim N(\mu, \sigma^2/n).$$

then

$$\frac{\bar{X}^{(n)} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

## Normal approximation of Binomial distribution

If  $X_1 \dots X_n \sim \text{Ber}(p)$ . Then  $X = \sum X_i \sim \text{Bin}(n, p)$ .

$X$  is approximately normally distributed

$$X \overset{\text{approx.}}{\sim} N(np, np(1-p)),$$

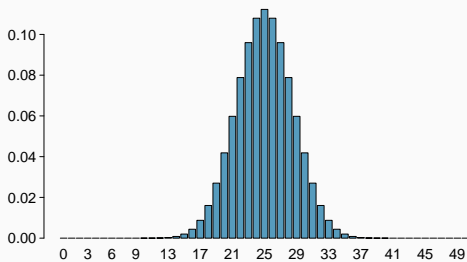
Thus **again** for  $\bar{X}^{(n)} = \frac{1}{n} \sum X_i$ ,

$$\bar{X}^{(n)} \overset{\text{approx.}}{\sim} N(p, p(1-p)/n),$$

or

$$\frac{\bar{X}^{(n)} - p}{\sqrt{p(1-p)/n}} \overset{\text{approx.}}{\sim} N(0, 1)$$

## Normal approximation



$$n = 50, p = 0.5$$

# Central limit theorem

## Central limit theorem (CLT)

If  $X_1, \dots, X_n$  are independent and equally distributed random variables with expected value  $\mu$  and variance  $\sigma^2 < \infty$ , then

$$\mathbb{P} \left( \frac{\bar{X}^{(n)} - \mu}{\sigma/\sqrt{n}} \leq x \right) \rightarrow F(x), \quad \text{for } n \rightarrow \infty.$$

where  $F$  is the distribution function of  $N(0, 1)$ .

This means,

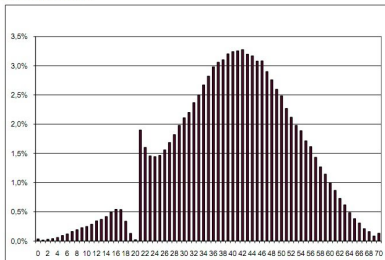
- $\bar{X} = n^{-1} \sum_{i=1}^n X_i$  is approximatively  $N(\mu, \text{SE}^2)$ -distributed, where  $\text{SE} = \sigma/\sqrt{n}$  is the **standard error**

for large  $n$ .

How large is large? Depends on the distribution of the  $X_i$ 's.

# High-school maturity exam in Poland

2.1. Poziom podstawowy



Wykres 1. Rozkład wyników na poziomie podstawowym

Histogram showing the distribution of scores for the obligatory Polish language test. "The dip and spike that occurs at around 21 points just happens to coincide with the cut-off score for passing the exam"

<http://freakonomics.com/2011/07/07/>

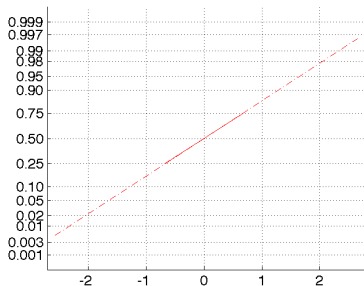
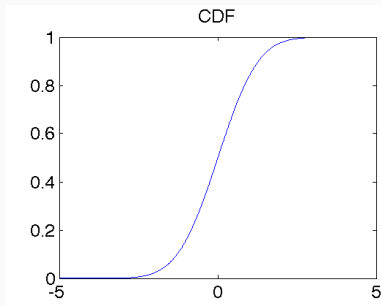
another-case-of-teacher-cheating-or-is-it-just-altruism/

## Normal probability plot

---



# Normal probability plot



The standard normal distribution function (cdf) is

$$F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-y^2/2} dy$$

It is possible to transform the scaling on the y-axis so that  $F$  becomes a straight line in the plot.

## Normal probability plot

---

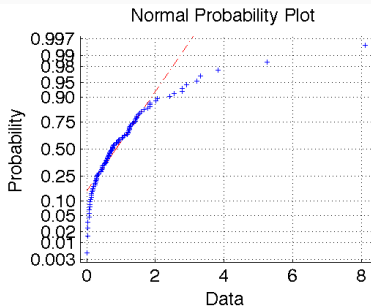
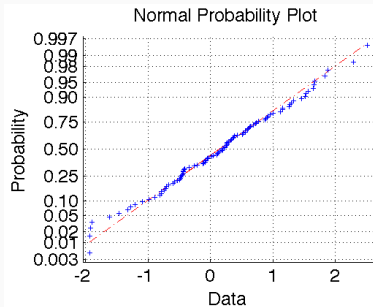
Suppose we have the data  $x_1, \dots, x_n$  and want to see if a normal distribution is a reasonable model for the data. We can use the normal probability plot for this.

First we compute the *empirical distribution function*

$$F^*(x) = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{I}(x_i \leq x)}_{\text{proportion of values smaller than } x}$$

We plot the points  $F^*(x_j)$  in the normal probability diagram, and if the data is normally distributed, these points should lie along a straight line.

# Normal probability plot



**Example:** left normally distributed data and right exponentially distributed data in normal probability diagram. In Matlab: `normplot`.

## Confidence interval

---

## Confidence interval

If  $X_1, \dots, X_n$  i.i.d random variables with distribution depending on a parameter  $\theta$ , with  $\theta_0$  being the unknown value. A  $100(1 - \alpha)\%$  confidence interval for  $\theta$  with confidence level  $1 - \alpha$  is an interval  $I_\theta = [A, B]$  computed from the data such that

$$P(A \leq \theta_0 \leq B) = 1 - \alpha.$$

## Confidence interval for parameter $\mu$ of a normal distribution

Let  $X_1, \dots, X_n$  be independent  $N(\mu, \sigma^2)$ .

**Known variance  $\sigma^2$**

$$I_\mu = (A, B) = \left( \bar{X}^{(n)} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{X}^{(n)} + 1.96 \frac{\sigma}{\sqrt{n}} \right)$$

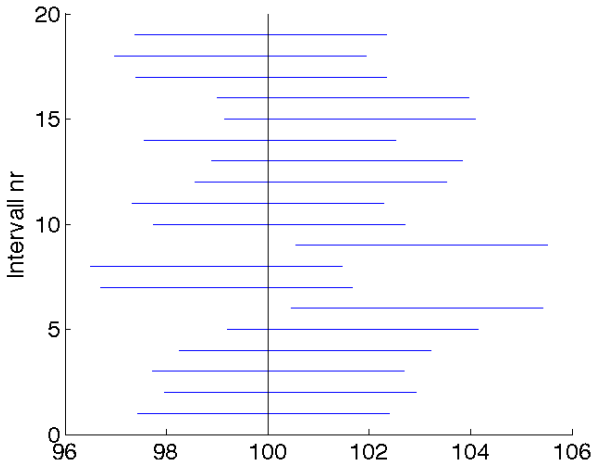
is a confidence interval for  $\mu$  with confidence level 95%.

Here 1.96 is the  $0.975 = (100 - 2.5)\%$  quantile of  $Z \sim N(0, 1)$ :

$$P(-1.96 < Z < 1.96) = 0.95.$$

$$P\left(-1.96 < \frac{\bar{X}^{(n)} - \mu}{\sigma/\sqrt{n}} < 1.96\right) = 0.95.$$

$$P(A \leq \mu \leq B) = 0.95$$



20 confidence intervals for  $\mu$ , that where each constructed from 20 different samples of 10  $N(100, 16)$ -observations.

- $[A, B]$  is a random interval, because  $A$  and  $B$  are random variables (transformations of the random variables  $X_1, \dots, X_n$ ).
- Interpretation. Let  $\mathbf{x}_1 = (x_{11}, \dots, x_{n1})$ ,  $\mathbf{x}_2 = (x_{12}, \dots, x_{n2})$ ,  $\dots$  be repeated measurements of  $X_1, \dots, X_n$ . If we make the confidence interval for  $\theta$  based on every  $\mathbf{x}_i$ , then  $100(1 - \alpha)\%$  of these intervals cover the true value  $\theta_0$ .

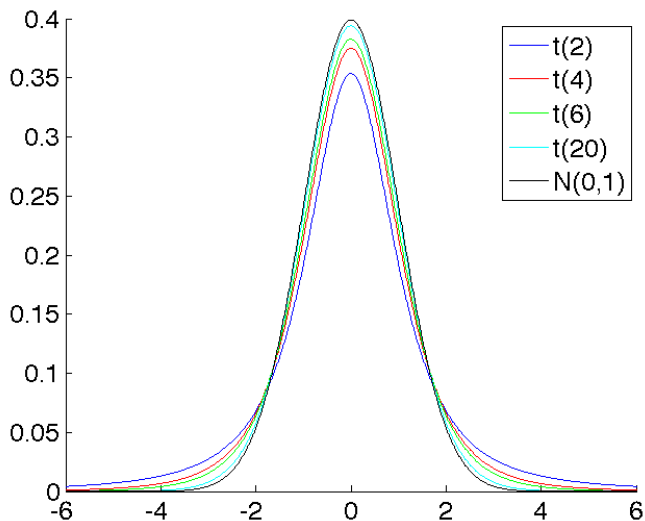


## Table 2: Quantiles of the normal distribution

Table gives  $P(X > \lambda_\alpha) = \alpha$  for  $X \sim N(0, 1)$

$\alpha$	.1	.05	.025	.01	.005	.001	...	.00001
$\lambda_\alpha$	1.2816	1.6449	1.9600	2.3263	2.5758	3.0902	...	4.2649

## $t(n)$ -distribution



**Table 3: Quantiles of the  $t$ -distribution**

Table gives  $P(X > t_\alpha(f)) = \alpha$  for  $X \sim t(f)$ .

$\alpha$	.1	.05	.025	.01	.001
$t_\alpha(1)$	3.0777	6.3138	12.706	31.820	318.31
$t_\alpha(2)$	1.8856	2.9200	4.3027	6.9646	22.327
$t_\alpha(3)$	1.6377	2.3534	3.1824	4.5407	10.215
$t_\alpha(4)$	1.5332	2.1318	2.7764	3.7469	7.1732
$t_\alpha(5)$	1.4759	2.0150	2.5706	3.3649	5.8934
$t_\alpha(6)$	1.4398	1.9432	2.4469	3.1427	5.2076
$t_\alpha(7)$	1.4149	1.8946	2.3646	2.9980	4.7853
$t_\alpha(8)$	1.3968	1.8595	2.3060	2.8965	4.5008
$t_\alpha(9)$	1.3830	1.8331	2.2622	2.8214	4.2968
$t_\alpha(10)$	1.3722	1.8125	2.2281	2.7638	4.1437
$t_\alpha(15)$	1.3406	1.7531	2.1314	2.6025	3.7328
$t_\alpha(20)$	1.3253	1.7247	2.0860	2.5280	3.5518
$t_\alpha(30)$	1.3104	1.6973	2.0423	2.4573	3.3852
$t_\alpha(40)$	1.3031	1.6839	2.0211	2.4233	3.3069
$t_\alpha(60)$	1.2958	1.6706	2.0003	2.3901	3.2317
$t_\alpha(\infty)$	1.2816	1.6449	1.9600	2.3263	3.0902

## Confidence interval for $\mu$ of a normal distribution

Let  $X_1, \dots, X_n$  be independent  $N(\mu, \sigma^2)$ .

### Known variance $\sigma^2$

$$I_\mu = \left( \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$$

is a confidence interval for  $\mu$  with confidence level  $1 - \alpha$ .

### Unknown variance $\sigma^2$

$$I_\mu = \left( \bar{X} - t_{\alpha/2}(n-1) \frac{s}{\sqrt{n}}, \bar{X} + t_{\alpha/2}(n-1) \frac{s}{\sqrt{n}} \right)$$

is a confidence interval for  $\mu$  with confidence level  $1 - \alpha$ . Here  $s^2$  is the sample variance and  $t_{\alpha/2}(n-1)$  are the  $(1 - \alpha/2)$ -quantiles of the  $t(n-1)$ -distribution.

# Quiz

---

$x_1, \dots, x_n$  are a sample of i.i.d observations with distribution depending on a parameter  $\theta$ .

Winnie computes a 95 % confidence interval for  $\theta$ .

Piglet computes a 90 % confidence interval for  $\theta$  using the same data.

Which interval is smallest? Piglet's 90 % confidence interval.

## Confidence interval for $\mu$ from central limit theorem

- By the CLT the sample mean  $\bar{X}^{(n)}$  is approximatively  $N(\mu, \sigma^2/n)$ -distributed for large  $n$ .
- If we have a sample with known variance  $\sigma^2$ ,

$$I_\mu = \left( \bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$$

is a confidence interval for the mean  $\mu$  with confidence level  $1 - \alpha$ .

- If  $\sigma$  is not known we can estimate it by  $S$ . For the estimate to be good, it is important that  $n$  is large and the distribution for  $X_i$  is not too heavy tailed.
- Since  $n$  is big, we use  $t_{\alpha/2}(n-1) \approx z_{\alpha/2}$ , so if  $\sigma$  is unknown, we use

$$I_\mu = \left( \bar{X} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right).$$

## Confidence interval for $\sigma^2$ for the normal distribution

### Confidence interval for $\sigma$

If  $X_1, \dots, X_n$  are independent  $N(\mu, \sigma^2)$  then a confidence interval with confidence level  $1 - \alpha$  for  $\sigma$  is

$$I_\sigma = \left( \sqrt{\frac{(n-1)s^2}{\chi_{\alpha/2}^2(n-1)}}, \sqrt{\frac{(n-1)s^2}{\chi_{1-\alpha/2}^2(n-1)}} \right).$$

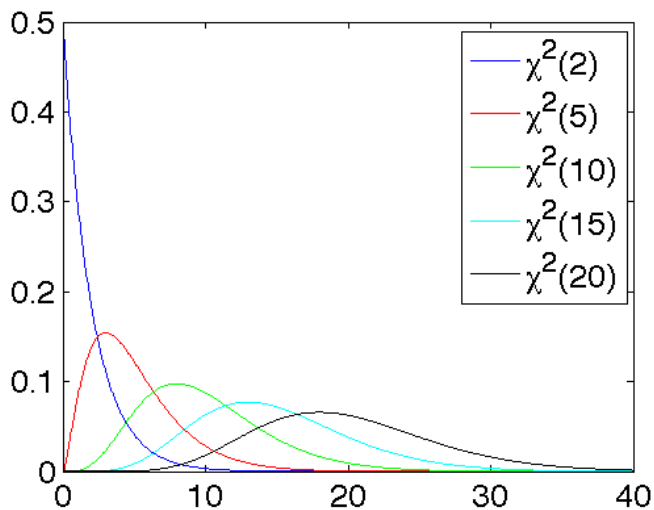
Here  $\chi_{\alpha/2}^2(n-1)$  are the  $(1 - \alpha/2)$ -quantiles of the  $\chi^2(n-1)$  distribution.

If  $Z_i$  are independent  $N(0, 1)$ , it holds

$$\sum_{i=1}^n Z_i^2$$

is  $\chi^2(n)$ -distributed

## $\chi^2(n)$ -distribution





## Confidence interval for $\sigma^2$ for the normal distribution

### Confidence interval for $\sigma$

If  $X_1, \dots, X_n$  are independent  $N(\mu, \sigma^2)$  then a confidence interval with confidence level  $1 - \alpha$  for  $\sigma$  is

$$I_\sigma = \left( \sqrt{\frac{(n-1)s^2}{\chi_{\alpha/2}^2(n-1)}}, \sqrt{\frac{(n-1)s^2}{\chi_{1-\alpha/2}^2(n-1)}} \right).$$

Important: In contrast to the confidence interval for the expected value, the confidence interval for the variance is very sensitive to deviations from the normal distribution.

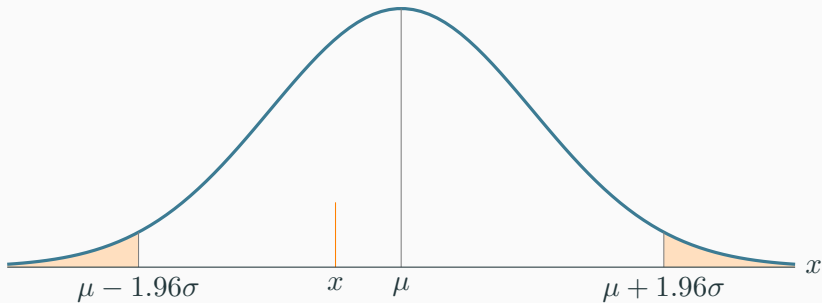
# Summary

---

For a confidence interval

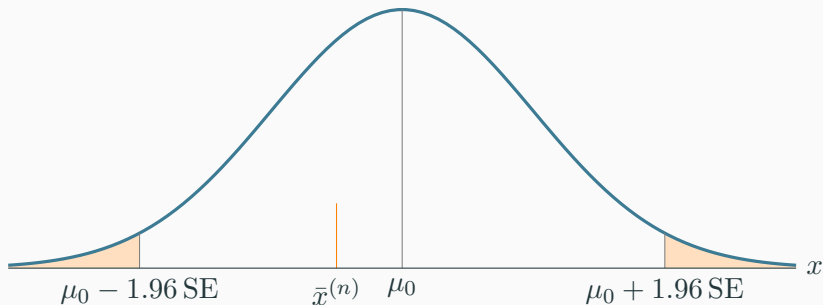
- for the expected value  $\mu$ 
  - of the normal distribution: Slide: confidence interval for  $\mu$  of a normal distribution
    - Known  $\sigma$  or large  $n$ : use confidence interval based on normal quantiles.
    - Small  $n$  and unknown  $\sigma$ : use quantiles based on  $t$ -distribution.
  - of a general distribution
    - Large  $n$ : use confidence interval based on normal quantiles (valid approximation by CLT). Slide: Confidence interval for  $\mu$  from central limit theorem.
- for the variance  $\sigma^2$ 
  - of the normal distribution: Slide: Confidence interval for  $\sigma^2$  for the normal distribution.

$$X \sim N(\mu, \sigma^2)$$



$$\mathbb{P}(X \in [\mu - 1.96\sigma, \mu + 1.96\sigma]) = 0.95$$

CLT:  $\bar{X}^{(n)} \overset{\text{approx}}{\sim} N(\mu, \text{SE}^2), \quad \mu = \mathbb{E}[X_i], \sigma^2 = \text{V}(X_i).$



$$\mathbb{P}(\bar{X}^{(n)} \in [\mu - 1.96 \text{SE}, \mu + 1.96 \text{SE}]) = 0.95$$

# Hypothesis tests

---

An important problem in statistics is to test whether a theory or a *research hypothesis* is right or wrong.

Examples of such problems include:

- Does a new drug have any effect?  $\text{Mean effect} > 0$
- Do smokers die sooner than non-smokers?  $\text{Mean life time difference} < 0$
- Does the measuring device have a systematic error?  $\text{Mean measurement error} \neq 0$

# Hypothesis tests

---

Answers the statistical analysis could give are

1. that the research hypothesis is supported by the data (and possibly a quantification of the degree of support),
2. that the data doesn't support the hypothesis,
3. a decision rule.

## Example

---

The length of a certain lumber from a national home building store is supposed to be 2.5 m.

A builder wants to check whether the lumber cut by the lumber mill has a mean length different smaller than 2.5 m.

A statistical formulation of this problem is that we want to test the **null hypothesis**

$$H_0: \text{mean length} = 2.5 \text{ m}$$

against the **alternative/research hypothesis**

$$H_1: \text{mean length} < 2.5 \text{ m}$$

$H_1$  is actionable knowledge. If  $H_1$  is true she needs to write an angry letter.

## Example

---

- You have new laboratory equipment to measure the chlorine content in water and want to check it. You mix water with true chlorine content 60 (you can do that very precisely), and take 6 measurements.
- Results of the measurement are  $\bar{x} = 59.62$  and  $s^2 = 4.6920$ .
- Assume that the measurements are samples of a random variable  $X \sim N(\mu, \sigma^2)$ .
- The question now is whether we can claim that the new equipment has systematic measurement error,  $\mu \neq 60$ .



# Setup

A statistical formulation of this problem is that we want to test the **null hypothesis**

$$H_0: \mu = 60$$

against the **alternative hypothesis** or **research hypothesis**

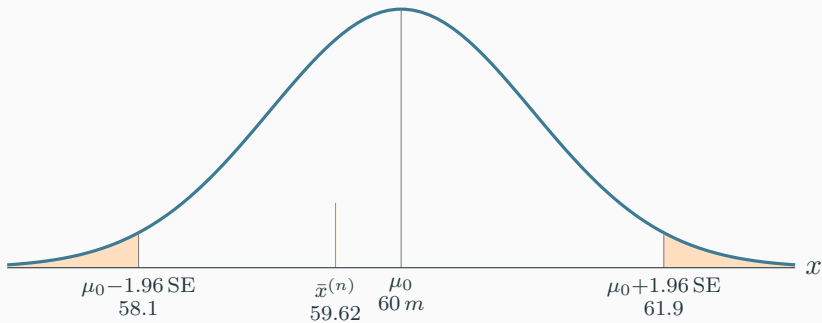
$$H_1: \mu \neq 60.$$

If the test we perform finds that there is a systematic error,  $H_0$  is rejected in favour of  $H_1$ .

Is  $H_1$  actionable knowledge?

## Choosing the alternative $H_1$

Choose  $H_1$  such if someone would tell you it is true, you can do something useful with that knowledge!



$$\text{SE} \approx \frac{\sqrt{4.6920}}{\sqrt{5}}$$

The **outcome** of a hypothesis test can be:

- Reject  $H_0$  (accept  $H_1$ .)
  - Action!
- Do not reject  $H_0$ 
  - Could be lack of data, or  $H_0$  being correct. The question of  $H_0$  or  $H_1$  is truly left open. Meh. Should still report it though.

## Decision errors

		Decision	
		fail to reject $H_0$	reject $H_0$
Truth	$H_0$ true	✓	Type 1 Error
	$H_1$ true	Type 2 Error	✓

- A **Type 1 Error** is rejecting the null hypothesis when  $H_0$  is true. We want to avoid that, control the probability for this error.
- A Type 2 Error is failing to reject the null hypothesis when  $H_1$  is true.

## Burden of proof

If we again think of a hypothesis test as a criminal trial then it makes sense to frame the verdict in terms of the null and alternative hypotheses:

$H_0$  : Defendant is innocent

$H_1$  : Defendant is guilty

Which type of error is being committed in the following circumstances?

- Declaring the defendant innocent when they are actually guilty

Type 2 error

- Declaring the defendant guilty when they are actually innocent

Type 1 error

Which error do you think is the worse error to make?

# Statistical reasoning

*Classical logic:* If the null hypothesis is correct, then **these data can not occur**.

These data have occurred.

Therefore, the null hypothesis is **false**.

*Tweak the language, so that it becomes **probabilistic**...      Statistical reasoning:*

If the null hypothesis is correct, then **these data are highly unlikely**.

These data have occurred.

Therefore, the null hypothesis is **unlikely**.

## Definition

In statistical hypothesis testing, a **result has statistical significance** when it is very unlikely to have occurred under the null hypothesis. So significance corresponds to "statistical evidence against the null".

The **significance level**  $\alpha$  is the (tolerated) probability of making a type I error:

## About failure to reject $H_0$

---

If you want to take a decision in the case the test fails to reject  $H_0$ , you should compute the type II error probability first. This is typically difficult.

Therefore we should avoid far reaching decisions if our tests fail to reject  $H_0$ .

## Tests from confidence intervals

**Data** (samples from a distribution with unknown parameter  $\mu$ ).

**Hypothesis** about parameter. Here  $H_0 : \mu = \mu_0$  and  $H_1 : \mu \neq \mu_0$ .

**Significance level**  $\alpha$ , e.g  $\alpha = 5\%$ .

**Decision rule:** Compute a  $(1 - \alpha)(= 95\%)$ -confidence interval  $[A, B]$  for the parameter  $\mu$ . If the  $\mu_0 \notin [A, B]$ , reject  $H_0$ .

**Type 1 error:** This rule has type 1 error of 5 %, so this is a valid test for level  $\alpha = 5\%$ .



## Tests with test statistics

**Data** (samples with unknown population parameter  $\mu$ ).

**Hypothesis** about parameter. Here  $H_0 : \mu = \mu_0$  and  $H_1 : \mu \begin{matrix} \neq \\ \geq \\ < \end{matrix} \mu_0$ .

**Significance level**  $\alpha$ , e.g  $\alpha = 5\%$ .

**Test statistic**  $T$ : Typically,  $T$  comes from an estimator for our parameter with known distribution under  $H_0$ .

$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \quad (\text{example})$$

**Decision rule:** Reject  $H_0$  if the  $p$ -value is less than the significance level  $\alpha$ .

or: Reject  $H_0$  if the  $T_{obs}$  is in the critical region/rejection region (see next slide).

**Type I error:** The type I error for this test is  $\leq \alpha$ .

## Critical region

The **critical region**  $C_\alpha$  of a test are those values of the test statistic  $T$  for which  $H_0$  can be rejected while obeying significance level  $\alpha$ . Typically represented by one or two critical values.

We compute rejection region for the data. We reject  $H_0$  if  $T_{obs}$  is in the rejection region.

## Example: critical region for mean of normal population

---

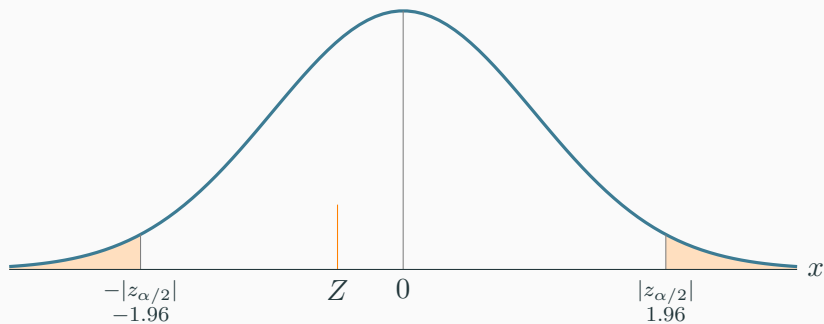
We want to use a quantity  $T$  that we know the distribution of under  $H_0$ , so that we can calculate the critical region.

In case of the normal distribution with known variance

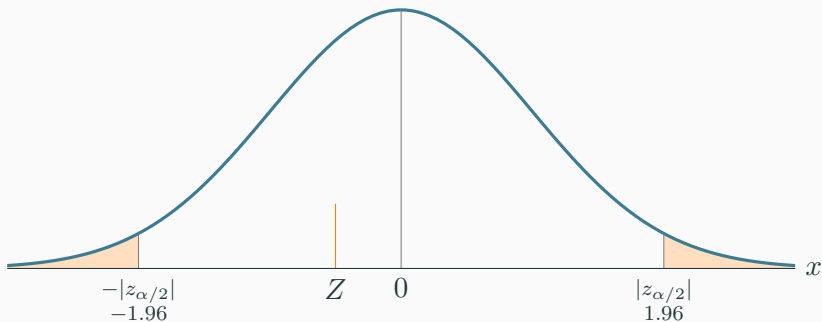
$$(T =) Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

we know that  $Z$  under  $H_0$  is  $N(0, 1)$ -distributed and

Reject  $H_0$  at level  $\alpha$  if  $|Z| > z_{\alpha/2}$ .



Rejection region for  $\alpha = 0.05$ .



Rejection region for  $\alpha = 0.05$  (on the  $x$ -axis below the yellow area).

Rule: Reject  $H_0$  (yeah) if  $Z$  is in the rejection region.

## Example: $p$ -value for mean of normal population

---

### $p$ -value

The  $p$ -value is the probability **under the null hypothesis  $H_0$**  to obtain a test statistic  $T$  with more evidence for the alternative (more “extreme”) than the one we observed,  $t_{obs}$ .

## Example: $p$ -value for normal distribution (two-sided)

Again we want to use a quantity  $T$  that we know the distribution of under  $H_0$ , so that we can calculate the  $p$ -value.

In case of the normal distribution with known variance

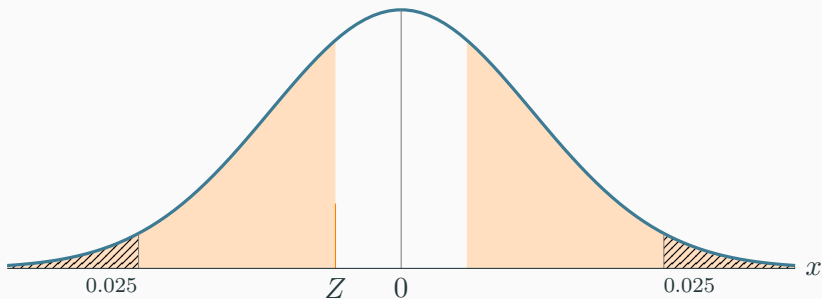
$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

we know that  $T$  under  $H_0$  is  $N(0, 1)$ -distributed and

$$p = P(|T| \geq |T_{obs}|) = 2 \cdot P(T \geq |T_{obs}|) = 2(1 - \Phi(|T_{obs}|)).$$

We compute  $p$  for the data. We reject  $H_0$  if  $p \leq \alpha$

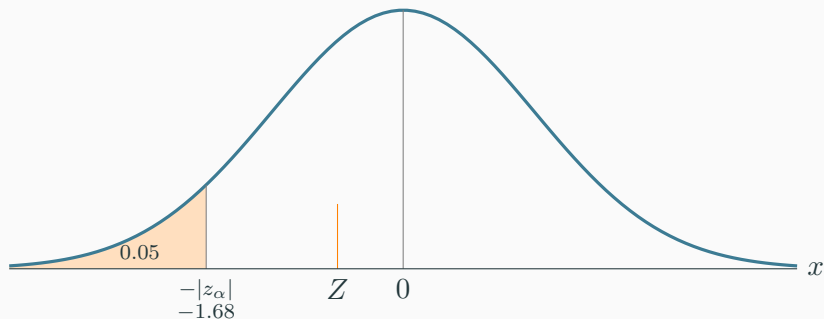
We compute rejection region for the data. We reject  $H_0$  if  $T_{obs}$  is in the rejection region.



Yellow area:  $p$ -value, dashed area:  $\alpha = 0.05$ .

Rule: Reject  $H_0$  if  $p \leq \alpha$





One-sided rejection region for  $\alpha = 0.05$ .

Rule: Reject  $H_0$  if  $Z$  is inside the rejection region.

## Example: $p$ -value for normal distribution (one-sided)

Again we want to use a quantity  $T$  that we know the distribution of under  $H_0$ , so that we can calculate the  $p$ -value.

In case of the normal distribution with known variance

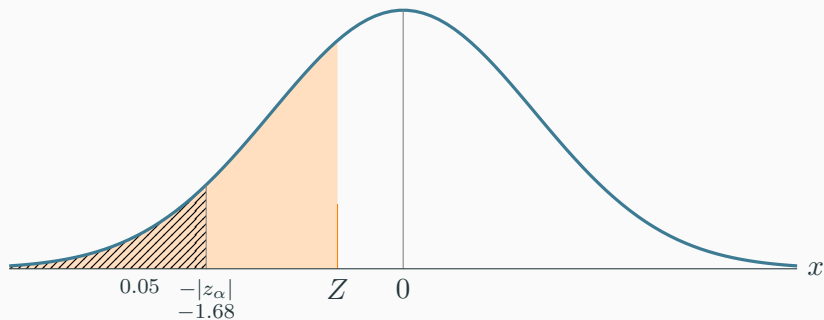
$$T = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

we know that  $T$  under  $H_0$  is  $N(0, 1)$ -distributed.

- 1.) Check if  $T$  is on the right side to give evidence in favour of  $H_1$ .
- 2.)  $p = P(T \text{ more extreme than } T_{obs})$  on the right side  
 $1 - \Phi(|T_{obs}|).$

We compute  $p$  for the data. We reject  $H_0$  if  $p \leq \alpha$

We compute rejection region for the data. We reject  $H_0$  if  $T_{obs}$  is in the orange rejection region.



Yellow area:  $p$  value, dashed area:  $\alpha = 0.05$ .

Rule: Reject  $H_0$  if  $p \leq \alpha$ .

## How many observations are needed?

---

A test detects a deviation of  $\mu - \mu_0$  more easily if:

- If the significance level  $\alpha$  is not very small.
- The number of observations  $n$  is large.
- The population variance relatively  $\sigma^2$  is small.