



# Table of contents

## 1 Inferential Statistics

- Paramater estimation
  - Point estimation
  - Interval estimation
  - Confidence interval for the mean

## 2 Tests of significance

- Statistical hypothesis testing
- Errors, significance level, and power
- Significance level and  $P$  value
- Parametric and non-parametric tests
- Proportion test
- One-tailed and two-tailed tests

## 3 Bibliography

# Inferential Statistics

- Inferential statistics aims to *infer* about a population using probability theory results and statistics computed from one or more samples.
- Without the use of the probability theory (the law of large numbers, the central limit theorem and other similar results) a certain feature may be considered systematic while it is only a random effect, or, on the contrary, some systematic features can remain unnoticed.
- Example of statistical inference: *confidence interval* for parameter estimating, or *significance tests*.
- The inference techniques are all based on sampling distributions.

## Parameter estimation

- The distribution of a certain population can be totally unknown, hence you may want to, at least, find its expectation, its variance and so on.
- These parameters of a population can be estimated using statistics computed from samples.
- There are two types of estimation: *point estimate* and *interval estimate*.

### Definition 1

*Point estimate for a parameter is a number, usually the value of the corresponding sample statistic, designed to estimate the parameter.*

## Parameter estimation

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

### Definition 2

An **interval estimate for a parameter** is an interval whose limits are statistics computed from a sample.

The **level of confidence**  $(1 - \alpha)$  is the proportion of all interval estimates that include the estimated parameter.

A **confidence interval** is an interval estimate with a specified level of confidence.

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

## Point estimation

- Examples of point estimates are the sample mean  $\bar{x}_n$ , the sample variance  $s^2$ , or the sample standard deviation  $s$ .
- For a given parameter we can have more point estimates: the population mean can be estimated by sample mean, median, mode etc.
- There are some questions about the quality of the point estimates.
- How exact is the estimate (i.e., the *confidence*) - it is frequently higher (*over-estimates*) or frequently lower (*under-estimates*) than the estimated parameter?
- In this regard it is preferred an estimate that is *unbiased*.
- What is the variability of the point estimate (when is viewed as a random variable)? This is the *accuracy*.
- For example take the standard deviation of the sample mean which is  $\sigma/\sqrt{n}$ : the larger the sample the smaller the variability of the sample mean.

## Point estimation characteristics

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

### Definition 3

Let  $\theta$  be a certain parameter of a population,  $x$  a sample (of size  $n$ ) from that population, and  $\hat{\theta}_n = \hat{\theta}_n(x)$  a point estimate of  $\theta$ .  $\hat{\theta}_n$  is an **unbiased statistic** if  $\mathbb{E}[\hat{\theta}_n] = \theta$ . Otherwise  $\hat{\theta}_n$  is called **biased**.

A **minimum variance unbiased estimator**, if exists, it is called an **efficient statistic**.

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

## Sample mean distribution

- Consider  $\mathcal{P}$  a population formed with  $N$  individuals those measured attribute values are  $a_1, a_2, \dots, a_N$ .
- The variable having these values is denoted by  $X$  and represents the distribution of the population.
- For this population (therefore for  $X$ ) the expectation (or simply the mean) and variance are

$$\frac{1}{N} \sum_{i=1}^N a_i = \mu, \quad \frac{1}{N} \sum_{i=1}^N (a_i - \mu)^2 = \sigma^2.$$

- When we estimate  $\mu$  and  $\sigma^2$  we use samples of size  $n \ll N$ .



## Sample mean distribution

- For a given sample  $x^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})$ , the sample mean is

$$\bar{x}_n^{(k)} = \frac{1}{n} \sum_{j=1}^n x_j^{(k)}.$$

- Each  $x_j^{(k)}$  is a value of a random variable identically distributed with  $X$ .

### Definition 4

The random variable having as values all possible sample means,  $\bar{x}_n^{(k)}$ , for  $n$  size samples is called the **sample mean distribution**.

## Sample mean distribution

- One can prove the following result

### Theorem 1.1

*The expectation and the variance of the sample mean,  $\bar{x}_n$ , are  $\mu$  and  $\sigma^2/n$ :*

$$\mathbb{E}[\bar{x}_n] = \mu, \text{Var}[\bar{x}_n] = \frac{\sigma^2}{n}.$$

*Moreover, for large values of  $n$  ( $n \geq 30$ ), the sample mean distribution is normal, that is*

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

## Mean and variance estimation

- Obviously, a point estimate for the (true) population mean ( $\mu$ ) is the sample mean  $\bar{x}_n$ .
- An estimate for the (true) population variance ( $\sigma^2$ ) is the sample variance  $s^2$  - which is an unbiased statistic.
- For the standard deviation of the sample mean,  $\sigma/\sqrt{n}$ , an estimate is  $s/\sqrt{n}$  which is called the **standard error of the mean (SEM)**.

## Interval estimation

- An estimate of the population mean using intervals may be given by Tchebychev inequality:

$$P(|\bar{x}_n - \mathbb{E}[\bar{x}_n]| \geq k \cdot \sqrt{\text{Var}[\bar{x}_n]}) \leq \frac{1}{k^2} \Leftrightarrow$$

$$P\left(|\bar{x}_n - \mu| \geq k \cdot \frac{\sigma}{\sqrt{n}}\right) \leq \frac{1}{k^2} \Leftrightarrow$$

$$P\left(|\bar{x}_n - \mu| < k \cdot \frac{\sigma}{\sqrt{n}}\right) \geq 1 - \frac{1}{k^2} \Leftrightarrow$$

$$P\left(\bar{x}_n - k \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x}_n + k \cdot \frac{\sigma}{\sqrt{n}}\right) \geq 1 - \frac{1}{k^2}.$$

- An estimate of  $\mu$  is the interval

$$\left(\bar{x}_n - k \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x}_n + k \cdot \frac{\sigma}{\sqrt{n}}\right).$$

## Confidence interval

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

### Definition 5

*A confidence interval for a parameter  $\theta$  with  $(1 - \alpha)$  level of confidence is defined by two statistics  $L$  and  $U$  such that*

$$P(L \leq \theta \leq U) = (\geq) 1 - \alpha.$$

- $L$  and  $U$  are random variables and their values are statistics: for different samples they have different values.
- Usually the level of confidence is a probability close to one: 0.90, 0.95, or 0.99 (which give  $\alpha \in \{0.10, 0.05, 0.01\}$ ).

## Confidence intervals for the mean of a population with known variance

- Suppose we have a sample of size  $n$  and a given level of confidence  $(1 - \alpha)$  and we have to build a confidence interval for  $\mu$ .
- We know that the sample mean  $\bar{x}_n$  has a distribution which is (almost, or exactly if the population is normally distributed) normal 
$$N\left(\mu, \frac{\sigma^2}{n}\right).$$
- After standardization the variable  $Z = \frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}}$  is a  $N(0, 1)$  distributed variable.
- We look for the value  $z^*$ , called *the critical value*, such that the normal distribution covers an area  $(1 - \alpha)$  over the interval centered in the mean (which is 0 for standard normal distribution) and having the length equal with  $2z^*$  standard deviations.

## Confidence intervals for the mean of a population with known variance

- Let  $Z : N(0, 1)$ , the critical value,  $z^*$ , is such that

$$P(-z^* \leq Z \leq z^*) = 1 - \alpha.$$

- Equivalent definitions of  $z^*$ :

$$P(Z \leq -z^*) = \alpha/2 \text{ or } P(Z \geq z^*) = \alpha/2.$$

- The distribution function of a standard normal variable is  $\Phi(a) = P(Z \leq a)$ .
- Therefore  $z^* = -\Phi^{-1}(\alpha/2)$  - value that can be taken from tables or estimated with statistical software analysis packages (R, MiniTab, SPSS etc).
- Now, having determined this value, we know that

$$P\left(-z^* \leq \frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}} \leq z^*\right) = 1 - \alpha.$$

## Confidence intervals for the mean of a population with known variance

- Equivalently

$$P\left(\bar{x}_n - z^* \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x}_n + z^* \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

- In this way we have proved

### Theorem 1.2

*A confidence interval with  $(1 - \alpha)$  confidence level for the mean of a population with known variance is*

$$\left(\bar{x}_n - z^* \frac{\sigma}{\sqrt{n}}, \bar{x}_n + z^* \frac{\sigma}{\sqrt{n}}\right),$$

*where  $z^*$  is the critical value corresponding to  $\alpha/2$ . Furthermore, this interval is an exact one for a normal distributed population and approximative for large samples ( $n \geq 30$ ) otherwise.*



## Confidence intervals for the mean of a population with known variance

Level of confidence	$\alpha/2$	$z^*$
90%	0.05	1.645
95%	0.025	1.960
99%	0.005	2.576

- $z^* \frac{\sigma}{\sqrt{n}}$  is called the *marginal error*.
- The length of a confidence interval is  $2z^* \frac{\sigma}{\sqrt{n}}$ .
- If we want a prescribed length for this interval,  $w$ , then we must sample with a size  $n = \frac{(2z^* \sigma)^2}{w^2}$ . This size could be unrealistic (if  $n$  is too big).
- For a larger  $n$  the interval length (or the marginal error) shrinks which is ok, but can be impractical in many situations (we cannot always sample very large samples).

## Confidence intervals for the mean of a population with known variance

- Remember that  $\frac{\sigma}{\sqrt{n}}$  is the standard deviation of the sample mean.

Now a confidence interval can be viewed as

$$\text{estimatedMean} \pm z^* \text{StdDevofMean}$$

- Example.** A certain pharmaceutical drug is analyzed by measuring three time its active substance, the results are 0.8403, 0.8636, and 0.8447 g/l. It is known that the concentration of this substance follows a normal law with standard deviation  $\sigma = 0.0068$  g/l. Find a 99% confidence interval for the real mean concentration,  $\mu$ .

$$\bar{x}_3 = 0.8495, \alpha = 0.01, \alpha/2 = 0.005, z^* = 2.576, z^* \frac{\sigma}{\sqrt{n}} = 0.0101$$

- The confidence interval is (0.8394, 0.8596).

## Confidence intervals for the mean of a population with unknown variance

- Remember that, when the variance of the population,  $\sigma^2$ , is known, an interval of  $(1 - \alpha)$  level of confidence for the true mean,  $\mu$ , is

$$\left( \bar{x}_n - z^* \frac{\sigma}{\sqrt{n}}, \bar{x}_n + z^* \frac{\sigma}{\sqrt{n}} \right),$$

- When the variance is unknown we can estimate the standard deviation of the sample by the standard error of the mean:  $s/\sqrt{n}$ .
- We will get a new statistic

$$T = \frac{\bar{x}_n - \mu}{s/\sqrt{n}}$$

known as **Student's t-statistics**.

- This is because the new statistic follows the Student distribution with  $(n - 1)$  *degrees of freedom*,  $t(n - 1)$  - if the population follows a normal distribution.

## Confidence intervals for the mean of a population with unknown variance

- Let  $T : t(n - 1)$ , the critical value,  $t^*$ , for  $(1 - \alpha)$  level of confidence is such that

$$P(-t^* \leq T \leq t^*) = 1 - \alpha, P(T \leq -t^*) = \alpha/2 \text{ or } P(T \geq t^*) = \alpha/2.$$

### Proposition 1

*A confidence interval with  $(1 - \alpha)$  confidence level for the mean of a normal population with unknown variance is*

$$\left( \bar{x}_n - t^* \frac{s}{\sqrt{n}}, \bar{x}_n + t^* \frac{s}{\sqrt{n}} \right).$$

*where  $t^*$  is the critical value corresponding to  $\alpha/2$ .*

## Confidence intervals for the mean of a population with unknown variance - Example

### Example

- A certain town has 10,000 rental units. A local real estate office does a survey of these units: 250 are chosen at random, and the occupants are interviewed. The sample rent average is 568\$ and the standard deviation of the sample is 385\$. (We know that the rent follows a normal distribution.)
- Find a 95%-confidence interval for the mean rent of all 10,000 units.

### Solution

- The collected data gives

$$\bar{x}_{250} = 568, s = 385, \alpha = 0.05, \alpha/2 = 0.025,$$

$$t^* = 1.9695, t^* \frac{s}{\sqrt{n}} = 48.0535$$

- The confidence interval is (616.0535, 519.9465).

## Confidence intervals for the mean - Exercises

- I. As part of an opinion survey, a simple random sample of 400 persons of age 25 and over is taken in a certain town in Appalachia. The total years of schooling completed by the sample persons is 4,635, and the standard deviation of the sample is 4.1 years. Find a 95%-confidence interval for the average educational level of all persons of age 25 and over in this town. (Assume the normality of the data.)
- II. Suppose that a city manager wants to know the average income of the 25000 families living in his town. He hires a survey organization to take a simple random sample of 1000 families. The total income of the sample turns out to be 62,396,714\$, and the standard deviation of the sample is 53,000\$. Find an 99%-confidence interval for the average income of a family in this town. (Assume that population is normally distributed.)

## Confidence intervals for the mean - Exercises

- III. An university has 30,000 registered students; as part of a survey, 900 of these students are chosen at random. The average age of the sample is 22.3 years with a sample standard deviation 4.5 years. Find 90%- and 95%-confidence intervals for the average age of all the students from this university. (Assume the normality of the data.)

## Statistical hypothesis testing

- Decisions are made in every day life; some of them are more significant than others, but all decisions follow the same pattern.
- We have two or more alternatives and we have to choose one of them based on evidence/beliefs/context/information etc.
- A statistical test of significance follows a similar pattern except that the decision is made using statistical information.
- That is, in this process we compute some statistics and based on these statistics we formulate a decision.
- In the process the first step is to define a situation having some level of uncertainty, and than to formulate two *hypothesis* about it.



## Statistical hypothesis testing

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

### Definition 6

A **statistical hypothesis test** is a process by which a decision is made between two opposing hypothesis.

The **statistical hypothesis** are formulated such that one of them is always true and the other is false.

One of the hypothesis is tested hoping that it can be shown to be of very improbable occurrence, thereby implying that the other hypothesis is likely to be true.

- The two hypothesis are known as the **null hypothesis** and the **alternative hypothesis**.
- A statistical test tries to find that the null hypothesis is likely to be false.

## Statistical hypothesis testing

Probabilități și Statistică

Probabilități și Statistică

Probabilități și Statistică

### Definition 7

**Null hypothesis**,  $H_0$ , is the hypothesis of the status-quo concerning of a population; formally it could be an assertion about a population: that it has a certain expectation or variance, or a certain distribution etc.

**Alternative hypothesis**,  $H_a$ , is the research hypothesis, and it says a different thing about the object of the null hypothesis.

- The null hypothesis is the starting point of the research and conservatively says that "there is no difference" or "nothing's happening" (it tends to oppose any change).
- In some way the alternative hypothesis attacks the null hypothesis and says a different thing than the statement of the null hypothesis.

## Statistical hypothesis testing - Example I

- A dairy company producing cheese, sour cream, ice cream and the likes has many milk distributors. There is a doubt about the quality of the milk they distribute.
- The freezing point of milk follows a normal law with mean  $\mu_0 = -0.545^\circ C$  and standard deviation  $\sigma = 0.008^\circ C$ . Adding water to milk changes this normal distribution and slightly increases the freezing point.
- For a given milk distributor five different milk batches are measured and the mean freezing point is found to be  $\bar{x}_5 = -0.538$ . Is this an evidence that the distributor adds water to milk? Or is just due to the chance?

## Statistical hypothesis testing - Example I

- The work hypothesis, that is, the null hypothesis, is that the expectation of the population is  $\mu = \mu_0$ .
- The research hypothesis, i.e., the alternative hypothesis, is that the expectation of the population is larger:  $\mu > \mu_0$ .
- We can ask a more formal question related to the alternative hypothesis: in normal conditions what are the chances that  $\mu > \mu_0$ ?
- Let  $X$  be the random variable associated to the milk freezing point;  $X : N(-0.545, (0.008)^2)$ .

## Statistical hypothesis testing - Example I

- Then

$$P(X > -0.538) = P\left(\frac{X - (-0.545)}{0.008/\sqrt{5}} > \frac{(-0.538) - (-0.545)}{0.008/\sqrt{5}}\right) =$$

$$= P\left(\frac{X - (-0.545)}{0.008/\sqrt{5}} > 1.95655948\right).$$

- This probability is  $P(Z > 1.95655948)$ , where  $Z : N(0, 1)$ .
- For a standard normal variable  $Z$ ,  $P(Z > 1.95655948) \sim 0.0250$ . That is, with probability 0.025, we can say that the difference observed in the sample is due to the usual random variations.

## Statistical hypothesis testing

- If, after supposing that a particular (null) hypothesis is true, we find that the statistics computed from a random sample greatly differs from the the results expected under our hypothesis, then we would say that the difference is *significant*.
- In this case and we will find tempting to *reject the null hypothesis*.
- If the statistics doesn't differ from the expected results under the null hypothesis, it is said that we *fail to reject the null hypothesis*.
- A *test of hypothesis* is a procedure that enable us to determine whether computed statistics significantly differs from the expected results under the null hypothesis.

## Statistical hypothesis testing

- If the decision is to reject  $H_0$ , then the conclusion should be: *"There is sufficient evidence that at the  $\alpha$  level of significance to conclude that..."* (the meaning of the alternative hypothesis).
- If the decision is the failure to reject  $H_0$ , then the conclusion should be: *"There is not sufficient evidence that at the  $\alpha$  level of significance to conclude that ..."* (the meaning of the alternative hypothesis).

## Type I and II errors

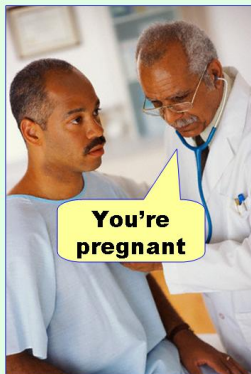
- Obviously, since we make decisions based on probabilities, our procedures we can be wrong. This is the place where errors occur.
- There are two type of errors: the first type occurs when we reject a hypothesis that should be accepted and the second type occurs when we accept a hypothesis when it should be rejected.

	$H_0$ validity	
	true	false
reject	Type I error (false positive)	Correct (true positive)
fail to reject	Correct (true negative)	Type II error (false negative)



## Type I and II errors

**Type I error**  
(false positive)



**Type II error**  
(false negative)

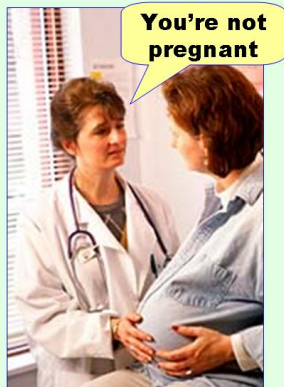


Figure: Type I and II errors. ( $H_0$  : You're not pregnant.)

## Statistical hypothesis testing - Example II

- You suspect that a brand-name detergent outperforms the store's brand of detergent, and you wish to test the two detergents (anyway, you would prefer to buy the better brand, but for the same quality you will buy the cheaper store's brand). State the null and alternative hypotheses.
- Your suspicion, "The brand-name detergent performs better than the store brand," is the reason for the test and therefore becomes the alternative hypothesis.

$H_0$  : "There is no difference in detergent performance."

$H_a$  : "The brand-name detergent outperforms the store brand."

- The test is designed hoping to reject the null hypothesis, but as a consumer, you are hoping not to reject the null hypothesis for budgetary reasons.

## Statistical hypothesis testing - Example II

	$H_0$ is true	$H_0$ is false
reject	<p><b>Type I error</b></p> <p>Truth of the situation: There is no difference.</p> <p>Conclusion: The brand name detergent is better.</p> <p>Action: The consumer spends more for no better results.</p>	<p><b>Correct decision</b></p> <p>Truth of the situation: The brand name is better.</p> <p>Conclusion: The brand name that detergent is better.</p> <p>Action: The consumer spends more and gets better results.</p>
fail to reject	<p><b>Correct decision</b></p> <p>Truth of the situation: There is no difference.</p> <p>Conclusion: there is no difference.</p> <p>Action: The consumer spends less and gets same results.</p>	<p><b>Type II error</b></p> <p>Truth of the situation: The brand name is better.</p> <p>Conclusion: there is no difference.</p> <p>Action: The consumer spends less and gets worse results.</p>

## Statistical hypothesis testing - Example II

- We described the four possible outcomes and the resulting actions that would occur for the hypothesis test.
- The truth of the situation is not known before the decision is made, the conclusion reached, and the resulting actions take place. The truth of  $H_0$  may never be known.
- The type II error often results in what represents a "lost opportunity"; lost in this situation is the chance to use a product that yields better results.

## Effect size, significance level, and power of the test

- The decision rule should minimize errors of decision.
- There are four items that influence the decision made from a statistical test (leaving apart the sampling itself): *the effect size*, the *sample size*, the *level of significance*, and the *power of the test*.

### Definition 8

The **effect size** is the magnitude of the difference (if any) discovered in the random sample.

The **level of significance**,  $\alpha$ , is the maximum (conditional) probability that someone is willing to risk a Type I error.

The **power of the test**, is one minus the probability of making a Type II error.

## Effect size, significance level, and power of the test

- The significance level is generally specified before sampling so that the results will not influence our sample choice.
- Usually the level of significance is 0.05 or 0.01; for example 0.01 (or 1% level of significance) is used for a decision rule which gives at most 1 in 100 chances that we would reject the null hypothesis when in fact it is true.
- The power of the test is the likelihood that our test detects a difference when there is a difference to detect.
- When the statistical power is high the probability of making Type II errors (the probability that the test finds no difference when there is such a difference) is low.

## Significance level and $P$ value

- There are two ways to conduct a statistical significance test: using the *score* and *critical value* or using the  *$P$ -value*.

### Definition 9

*The  $P$ -value is probability of getting an outcome as extreme or more extreme than the sample outcome. When the  $P$ -value is small "enough" we can reject  $H_0$ .*

- Return now to the first example. We computed  $P(X > -0.538) \sim 0.0250$ ; this is the  $P$ -value of the corresponding test.
- 2.5% may be considered a small probability and, therefore, we can reject the null hypothesis.

## Significance level and $P$ value

- But if someone consider that 2.5% is not small enough, then we fail to reject  $H_0$ .
- For example  $2.5\% < 5\%$ , hence, with 5% level of significance we can reject  $H_0$  and accept  $H_a$ .
- With 1% level of significance we fail to reject the null hypothesis, our data is not significant enough in order to accept the alternative hypothesis,  $H_a$ .



Significance level and  $P$  value

## Definition 10

The **score** of the test is the statistic corresponding to the  $P$ -value; the **critical value** corresponds to the level of significance. The conclusion of the test follows comparing these two values.

- We compute the statistic

$$z = \frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}} = \frac{(-0.538) - (-0.545)}{0.008/\sqrt{5}} = 1.95655948.$$

- $z$  is the score of the test.
- The critical value,  $z^*$ , for 1% level of significance is computed such that

$$P(Z > z^*) = 1 - \alpha = 0.99 \Rightarrow z^* = 2.32,$$

where  $Z : N(0, 1)$ .

## Significance level and $P$ value

- With 1% level of significance  $z < z^*$ , therefore  $H_0$  cannot be rejected
- For 5% level of significance, the critical value,  $z^*$ , is computed such that

$$P(Z > z^*) = 1 - \alpha = 0.95 \Rightarrow z^* = 1.64,$$

- With 5% level of significance  $z > z^*$ , therefore  $H_0$  can be rejected; we accept  $H_a$ .
- This example shows us that with different levels of significance we can get different conclusions.

## Parametric and non-parametric tests

- In our previous example we supposed that the population follows a normal distribution.
- Sometimes is difficult to asses whether a population follows a certain distribution (usually a normal one).
- Based on information concerning the population distribution we have two types of significance tests: parametric and non parametric tests.

### Definition 11

A **parametric test** *assumes that a population follows a specific distribution and infers about the parameters of that distribution.*

## Parametric and non-parametric tests

- Non-parametric tests infer about the population's distribution (which is unknown) rather than the parameters of the population.
- Usually the parametric tests are those that assume that the population follows a normal distribution. (Some authors consider that anything else is non-parametric).

### Definition 12

A **non-parametric test** also called **distribution-free** or **parameter-free** is based on fewer assumptions - the distribution and its parameters (expectation, variance) are not known.

## Proportion test

- One of the most common inference concerns the *binomial parameter*  $p$ , the probability of success.
- In many situations we are concerned about something either "happening" or "not happening"; there are only two possible outcomes of concern, and that is the fundamental property of a binomial experiment.
- The proportion of the individuals in a population that have a certain feature can be considered the probability of success.
- Examples: the proportion of people that smoke, of people that vote for a certain politician, of car drivers that speed up to beat the yellow light, of young adults that start out on their own etc.

## Proportion test

- If  $X : B(n, p)$  (the number of successes), then the population parameters are

$$\mu = np, \sigma^2 = np(1 - p).$$

- If we have a sample of size  $n$ , and  $x$  is the number of successes between them, then the **sample binomial probability** is

$$p' = \frac{x}{n}.$$

- For  $n \geq 20$  and  $np \geq 5$ ,  $p'$  can be approximated by the normal distribution.

## Proportion test

- The parameters of  $p'$  are

$$\mathbb{E}[p'] = \frac{\mathbb{E}[X]}{n} = p, \text{Var}[p'] = \frac{\text{Var}[X]}{n^2} = \frac{p(1-p)}{n}.$$

- Hence,

$$p' \sim N\left(p, \frac{p(1-p)}{n}\right).$$

- Therefore the following statistic follows an almost standard normal distribution

$$z = \frac{p' - p}{\sqrt{p(1-p)/n}}.$$

- This will be the statistic or the score of the test (where  $p' = x/n$ ).

## Proportion test - Example

- According to a nationwide Harris Poll during August 2008, 68% of American adults own a library card. Suppose you conduct a survey of 1000 randomly chosen adults in order to test  $H_0 : p = 0.68$  versus  $H_a : p < 0.68$ , where  $p$  represents the proportion of adults who currently have a library card.
- 651 of the 1,000 sampled had a library card. Use  $\alpha = 0.01$ .
  - a. Calculate the value of the test (the statistic).
  - b. Solve using the  $P$ -value approach.
  - c. Solve using the classical approach (with critical value).



## Proportion test - Example

$$z = \frac{p' - p}{\sqrt{p(1-p)/n}} = \frac{0.651 - 0.68}{\sqrt{0.68 * 0.32/1000}} = -1.965$$

- With  $P$ -value:

$$P = P(Z < z | H_0) = P(Z < -1.965) = 0.0246$$

- With  $\alpha = 1\%$  we can't reject  $H_0$ , the data sample is not significant.
- With critical value:

$$z^* = -2.326, (P(Z < z^*) = 0.01).$$

- Since  $z \not< z^*$  we can't reject  $H_0$ , there is not sufficient evidence to accept the alternative hypothesis.

## Proportion test - Example

- Let's rework the exercise with  $\alpha = 5\%$ .
- With  $P$ -value:  $P = 0.0246 < 0.05$ , hence we can reject the null hypothesis and accept  $H_a$ , that is the proportion of adults owning a library card is smaller than 0.68.
- With critical value:  $z^* = -1.644 > z$  which means that  $H_0$  can be rejected.

## Proportion test - Exercises

- I. A coin is tossed 10,000 times, and it lands heads 5,167 times. Is the chance of heads equal to 50%? Or are there too many heads for that? Perform appropriate proportion tests in order to answer to these questions.
- II. Repeat exercise I if the coin lands heads 5,067 times.
- III. In 2009, an USA Today article titled "On road, it's do as I say, not as I do" reported that 58% of U.S. adults speed up to beat a yellow light. Suppose you conduct a survey in your hometown of 150 randomly selected adults and find that 71 out of the 150 admit to speeding up to beat a yellow light. Does your hometown have a lower rate for speeding up to beat a yellow light than the nation as a whole? Use a 0.05 level of significance.

## One-tailed and two-tailed tests

- There are three types of alternative hypothesis concerning a given parameter, for example the binomial probability  $p$ .  
 $H_a : p < p_0$ ,  $H_a : p > p_0$  or  $H_a : p \neq p_0$ .
- The first two are "one-tailed" alternative hypotheses, and the third is "two-tailed" alternative hypothesis.
- For each of these the critical values are computed using different formulas.

## Proportion one-tailed test

- For  $H_0 : p = p_0$  and  $H_a : p > p_0$  - *right tail*.

$$P\text{-value } P(Z > z),$$

the critical value  $z^* > 0$ , s. t.  $P(Z > z^*) = \alpha = 1 - P(Z < -z^*)$ .

- For  $H_0 : p = p_0$  and  $H_a : p < p_0$  - *left tail*.

$$P\text{-value } P(Z < z),$$

the critical value  $z^* < 0$ , s. t.  $P(Z < z^*) = \alpha = 1 - P(Z > -z^*)$ .

## Proportion two-tailed test

- For  $H_0 : p = p_0$  and  $H_a : p \neq p_0$ .






$$P\text{-value } P(Z > |z|) + P(Z < -|z|),$$

the critical value  $z^* > 0$ , s. t.  $P(Z < -|z^*|) = \alpha/2 =$

$$= 1 - P(Z > |z^*|).$$



## Bibliography

-  Freedman, D., R. Pisani, R. Purves, *Statistics*, W. W. Norton & Company, 4th edition, 2007.
-  Johnson, R., P. Kuby, *Elementary Statistics*, Brooks/Cole, Cengage Learning, 11th edition, 2012.
-  Shao, J., *Mathematical Statistics*, Springer Verlag, 1998.
-  Spiegel, M. R., L. J. Stephens, *Theory and Problems of Statistics*, Schaum's Outline Series, McGraw Hill, 3rd edition, 1999.
-  <http://effectsizefaq.com/2010/05/31/i-always-get-confused-about-type-i-and-ii->