

# MS-C1620 Statistical inference

## 4 Inference for binary data

Jukka Kohonen

Department of Mathematics and Systems Analysis  
School of Science  
Aalto University

Academic year 2020–2021  
Period III–IV

# Contents

- 1 Binary data
- 2 Single binary sample
- 3 Two binary samples
- 4 Lecture quiz

## Binary observations

In many applications the observations are binary.

- Something is true/false.
- Something happened/did not happen.
- Someone belongs/does not belong to a group.

In such a case the observations are most conveniently coded as 0/1.

Recall that if we have a iid sample of binary observations, their distribution is necessarily the *Bernoulli distribution*.

## Bernoulli distribution

The random variable  $x$  is said to obey the Bernoulli distribution with the probability of success  $p$  if,

$$\mathbb{P}(x = 1) = p \quad \text{and} \quad \mathbb{P}(x = 0) = 1 - p.$$

The expected value and variance of  $x$  are,

$$\mathbb{E}(x) = p$$

$$\text{Var}(x) = p(1 - p).$$

That is, the Bernoulli distribution has only a single parameter to estimate.

The sum of  $n$  i.i.d. Bernoulli random variables with the success probability  $p$  has the binomial distribution with the parameters  $n$  and  $p$ .

# Contents

- 1 Binary data
- 2 Single binary sample**
- 3 Two binary samples
- 4 Lecture quiz

## Approximate confidence interval

Central limit theorem can be used to obtain a confidence interval for the success probability  $p$  of a Bernoulli distribution.

Let  $x_1, x_2, \dots, x_n$  be an i.i.d. sample from the Bernoulli distribution with the success probability/expected value  $p$ .

For large  $n$ , a level  $100(1 - \alpha)\%$  confidence interval for the success probability  $p$  is obtained as

$$\left( \hat{p} - z_{\alpha/2} \frac{\sqrt{\hat{p}(1 - \hat{p})}}{\sqrt{n}}, \hat{p} + z_{\alpha/2} \frac{\sqrt{\hat{p}(1 - \hat{p})}}{\sqrt{n}} \right),$$

where  $\hat{p}$  is the observed proportion of successes and  $z_{\alpha/2}$  is the  $(1 - \alpha/2)$ -quantile of the standard normal distribution.

# One-sample proportion test

To test whether the success probability of a Bernoulli distribution equals some pre-specified value, we employ **one-sample proportion test**.

## One-sample proportion test, assumptions

Let  $x_1, x_2, \dots, x_n$  be an i.i.d. sample from a Bernoulli distribution with the success probability  $p$ .

## One-sample proportion test, hypotheses

$$H_0 : p = p_0 \quad H_1 : p \neq p_0.$$

# One-sample proportion test

## One-sample proportion test, test statistic

- The test statistic,

$$C = \sum_{i=1}^n x_i,$$

follows the binomial distribution with parameters  $n$  and  $p_0$  under  $H_0$ .

- Under  $H_0$ , the test statistic has  $E[C] = np_0$  and  $\text{Var}(C) = np_0(1 - p_0)$  and both large and both **large** and **small** values of the test statistic suggest that the null hypothesis  $H_0$  is false.

The distribution of the test statistic  $C$  is tabulated and statistical software calculate exact  $p$ -values of the test.



## Asymptotic one-sample proportion test

If the sample size is large, then under the null hypothesis  $H_0$  the standardized test statistic,

$$Z = \frac{\hat{p} - p_0}{\sqrt{p_0(1 - p_0)/n}},$$

where  $\hat{p}$  is the unbiased estimator  $\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i$  of the parameter  $p$ , follows approximately the standard normal distribution.

The approximation is usually accurate enough if  $n\hat{p} > 10$  and  $n(1 - \hat{p}) > 10$ . For smaller sample sizes one should rely on the exact distribution of the test statistic  $C$ .

# Contents

- 1 Binary data
- 2 Single binary sample
- 3 Two binary samples**
- 4 Lecture quiz

## Two-sample proportion test

The one-sample proportion test can be seen as the equivalent of  $t$ -test when the normal distribution is replaced by the Bernoulli distribution.

As with  $t$ -test, a two-sample version readily follows and in **two-sample proportion test** parameters of two independent Bernoulli-distributed samples are compared.

### Two-sample proportion test, assumptions

Let  $x_1, x_2, \dots, x_n$  be an i.i.d. sample from a Bernoulli distribution with the success probability  $p_x$  and let  $y_1, y_2, \dots, y_m$  be an i.i.d. sample from a Bernoulli distribution with the success probability  $p_y$ . Furthermore, let the two samples be independent.

### Two-sample proportion test, hypotheses

$$H_0 : p_x = p_y \quad H_1 : p_x \neq p_y.$$

## Two-sample proportion test

### One-sample proportion test, test statistic

- Calculate the sample proportions

$$\hat{p}_x = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{p}_y = \frac{1}{m} \sum_{i=1}^m y_i, \quad \hat{p} = \frac{n\hat{p}_x + m\hat{p}_y}{n + m}.$$

- The test statistic,

$$Z = \frac{\hat{p}_x - \hat{p}_y}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n} + \frac{1}{m}\right)}},$$

follows *for large n* under  $H_0$  the standard normal distribution.

- Both **large** and **small** values of the test statistic suggest that the null hypothesis  $H_0$  is false.

The normal approximation is usually good enough if  $n\hat{p}_x > 5$ ,  $n(1 - \hat{p}_x) > 5$ ,  $m\hat{p}_y > 5$  and  $m(1 - \hat{p}_y) > 5$ .

## Frequency tables

Assuming a “paired binary sample”, the previous test is no longer valid.

id	X	Y
1	0	1
2	0	0
3	0	1
4	1	1
$\vdots$	$\vdots$	$\vdots$

This kind of data is most conveniently represented in a **frequency table**.

	Y = 0	Y = 1
X = 0	173	40
X = 1	65	53

Inference for frequency tables is discussed next time.

# Contents

- 1 Binary data
- 2 Single binary sample
- 3 Two binary samples
- 4 Lecture quiz**

## Lecture quiz

A lecture quiz to determine what you have learned thus far!

Answer the following questions on your own or in small groups.

# Lecture quiz

## Question 1

Consider the following random sample: 5, -4, -2, 2. Calculate the following sample quantities:

- 1 Sample mean
- 2 Sample standard deviation
- 3 Sample median
- 4 Sample median absolute deviation
- 5 Sample range
- 6 Signs of the sample points
- 7 Ranks of the sample points
- 8 Signed ranks of the sample points with respect to distance to 0.



# Lecture quiz

## Question 2

Give concrete examples when you would/would not use the following measures of location:

- 1 Sample mean
- 2 Sample median
- 3 Mode

## Question 3

Give concrete examples when you would/would not use the following measures of scatter:

- 1 Standard deviation
- 2 Median absolute deviation
- 3 Sample range

# Lecture quiz

## Question 4

What does it mean in practice if:

- The confidence interval of a parameter is narrow
- The significance level of a test is set low
- The  $p$ -value of a test is high
- Type I error occurs in a statistical test
- Type II error occurs in a statistical test

# Lecture quiz

## Question 5

How would you visualize the following samples:

- The heights of the male and female students attending a course.
- The exam points (0-24) on a large course.
- The proportions of faulty products produced by 5 different production lines.
- Stock prices of 3 companies.
- The monthly salaries and postal codes of adults living in Helsinki area.

# Lecture quiz

## Question 6

The following plots show the distributions of the test statistics of  $t$ -test, sign test and signed rank test for the null hypothesis of zero location when the data is a sample of size  $n = 10$  from the standard normal distribution. Which plot corresponds to which test?

