

# TO HYPOTESIS TESTING AND PARAMETRIC TESTS

Marianna Bolla, DSc

Institute of Mathematics, BME

2020. március 29.

## 1. Continuous distributions

### 1.1. One-sample, two-sided z-test

There are complaints that the bread (written 1 kg on it) weights less. We randomly select  $n$  breads, their weights are  $X_1, \dots, X_n$ ,  $n = 25$ . The sample average is 0.98 kg. What to do? The 0.02 kg difference can be caused by randomness. Even if  $X_1, \dots, X_n \sim \mathcal{N}(\mu_0, \sigma_0)$  with  $\mu_0 = 1$  kg and  $\sigma_0 = 0.05$ , there are fluctuations. Investigate the alternative:

$$H_0 : \mu = \mu_0 (= 1 \text{ kg}) \quad \text{versus} \quad H_1 : \mu \neq \mu_0.$$

By the assumption of innocence, the jury assumes  $H_0$  and we must have enough evidence to prove the contrary:  $H_1$ .

If we construct, say, a 95% confidence interval for the population mean, the hypothetical  $\mu_0 = 1$  kg should be in it with high probability. If not, then

- either the complementary event happens, but this has small 5% probability;
- we rather suspect, that the population mean is not 1 kg, and reject  $H_0$ .

We can only rule the Type I error probability (now 0.05): we reject  $H_0$  if it is true (the jury sentences an innocent). The Type II error probability is opposite: we accept  $H_0$  if not true (the jury acquits someone who is guilty). Its probability increases if we decrease the Type I error probability, so selecting the Type I error probability (called significance, and denoted by  $\alpha$ ) raises ethical issues.

To simplify things, first we construct the test statistic:

$$Z = \frac{\bar{X} - \mu_0}{\sigma_0} \sqrt{n}$$

and select a significance (now  $\alpha = 0.05$ ), then find the *critical value*

$$z_{\alpha/2} = \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right)$$

from the standard normal distribution table. When we constructed confidence interval last week, we saw that

$$\mathbb{P}\left(\mu_0 \in \left(\bar{X} - \frac{z_{\alpha/2}\sigma_0}{\sqrt{n}}, \bar{X} + \frac{z_{\alpha/2}\sigma_0}{\sqrt{n}}\right)\right) = \mathbb{P}(|Z| < z_{\alpha/2}) = 1 - \alpha.$$

Therefore, our decision is: if  $|z| < z_{\alpha/2}$ , then we accept, and if  $|z| \geq z_{\alpha/2}$ , then we reject  $H_0$  with significance  $\alpha$ . The region

$$R = \{\mathbf{x} : |z(\mathbf{x})| \geq z_{\alpha/2}\}$$

is called *rejection or critical region*.

In our numerical example:  $\bar{x} = 0.98$ ,  $\mu_0 = 1$ ,  $n = 25$  and let  $\sigma_0 = 0.05$ . So  $z = -2$ . With sign.  $\alpha = 0.05$   $z_{\alpha/2} = 1.96$ , therefore, with sign. 0.05 we reject  $H_0$  and the decision of the lower court: the shop is guilty. They work with sign. 0.05, i.e., give 0.05 prob. to the event that an innocent is convicted (jury basically defends the innocents, and they must have enough evidence to convict someone).

Then the shop goes to the higher court. They work with sign.  $\alpha = 0.01$ , because they are more strict and better defend the innocents (give only 0.01 prob. to the event that an innocent is convicted). Now  $z_{\alpha/2} = 2.58$ , our  $|z| = 2 < 2.58$ , so the higher court cannot reject  $H_0$  and acquits the shop.

From the table we can see that with  $\alpha = 0.0456$ ,  $z_{\alpha/2} = 2$ , so 0.0456 is the smallest possible sign., at which we can reject  $H_0$ . Program packages output this sometimes called *P-value* and if it is small enough, we can reject  $H_0$ .

Investigate the Type I error prob. ( $\alpha$ ) and the Type II error prob. ( $\beta$ ), which also depends on the true value of  $\mu$ :

$$\alpha = \mathbb{P}_{\mu_0}(|Z| \geq z_{\alpha/2})$$

per def, and if  $\mu \neq \mu_0$ :

$$\beta(\mu) = \mathbb{P}(|Z| < z_{\alpha/2} | \mu).$$

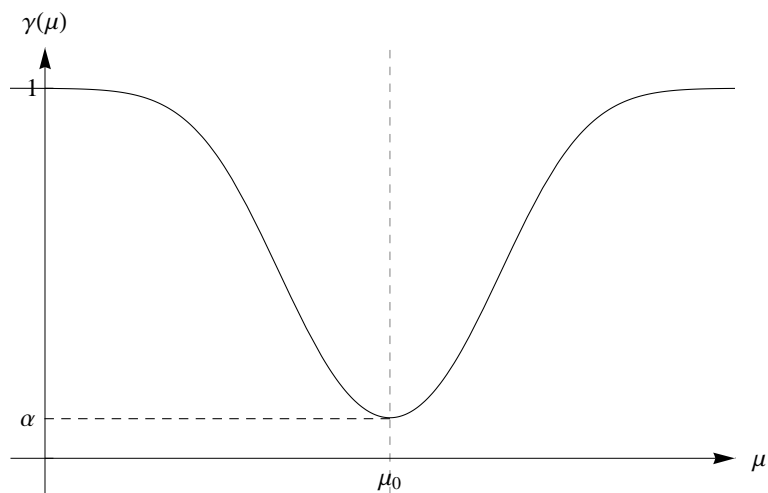
Sometimes the so-called *power function*  $\gamma$  is used, which is

$$\gamma(\mu) = 1 - \beta(\mu) = \mathbb{P}(|Z| \geq z_{\alpha/2} | \mu).$$

The theory guarantees the existence of a Uniformly Most Powerful (UMP) test in this situation (Neyman-Pearson): given  $\alpha$ , the UMP test has the largest possible power (so the smallest possible Type II error) for any  $\mu$ . This means, that if the jury convicts innocents with given prob, then they acquit any criminal with the smallest possible prob.

In our example,

$$\begin{aligned} \gamma(\mu) &= 1 - \mathbb{P}\left(-z_{\alpha/2} < \frac{\bar{X} - \mu_0}{\sigma_0} \sqrt{n} < z_{\alpha/2} | \mu\right) = \\ &= 1 - \left(-z_{\alpha/2} - \Delta_n < \frac{\bar{X} - \mu}{\sigma_0} \sqrt{n} < z_{\alpha/2} - \Delta_n\right) = \\ &= 1 - \Phi(z_{\alpha/2} - \Delta_n) + \Phi(-z_{\alpha/2} - \Delta_n) = \\ &= 2 - \Phi(z_{\alpha/2} - \Delta_n) - \Phi(z_{\alpha/2} + \Delta_n), \end{aligned}$$



1. ábra. Power function of the one-sample two-sided z-test

where

$$\Delta_n = \frac{\mu - \mu_0}{\sigma_0} \sqrt{n}$$

and  $\frac{\bar{X} - \mu}{\sigma_0} \sqrt{n} \sim \mathcal{N}(0, 1)$ , if  $\mu$  is the true population mean.

From this form, it is easy to see that if  $n \rightarrow \infty$  (more and more witnesses) or  $\mu$  gets farther and farther from  $\mu_0$  (more and more guilty), then  $\gamma(\mu) \rightarrow 1$  and  $\beta(\mu) \rightarrow 0$ , see Fig. 1.

## 1.2. One-sample, one-sided z-test

The shop example is better formulated as

$$H_0 : \mu \geq \mu_0 \quad \text{versus} \quad H_1 : \mu < \mu_0$$

(the accuse is only about smaller breads, if larger, no problem). Then the rejection region is

$$R = \{\mathbf{x} : z(\mathbf{x}) \leq -z_\alpha\},$$

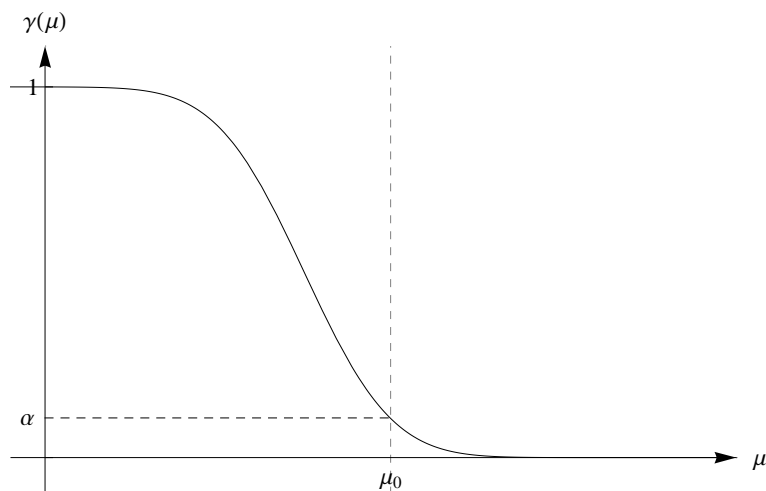
see the formulas to hypothesis testing.

Then with sign. 0.05 and 0.025 we reject  $H_0$ . From  $-z_\alpha = -2$  we get that the smallest possible  $\alpha$  (P-value) at which we can reject  $H_0$  is  $\alpha = 0.0228$ . (In the two-sided situation the P-value was the double of this: 0.0456.) So, in the one-sided, modified accuse case, we can sooner convict the shop than in the two-sided situation.

The power function, fixing  $n$  and  $\alpha$  is

$$\gamma(\mu) = \mathbb{P}(Z \leq -z_\alpha | \mu) = \mathbb{P}_\mu\left(\frac{\bar{X} - \mu}{\sigma_0} \sqrt{n} + \Delta_n \leq -z_\alpha\right) = \Phi(-z_\alpha - \Delta_n).$$

Here  $\gamma(\mu)$  decreases in  $\mu$ , see Fig. 2.



2. ábra. Power function of the one-sample one-sided z-test

## 2. Discrete distributions

In the reference book (Chapter 6), the Problem (cure-rate, p. 167) deals with the one-sided and Ex. 6.1 (cat-food) with the two sided problem (underlying Bernoulli distribution with binomial test statistic). Please, read the different strategies (rejection regions). These are small sample exercises. In the large sample case, we can use  $z$ -test for the population proportions, and get surprisingly other results.

In case of 15 cats, Fig. 3 shows the power curves under strategies

$$a : \mathcal{X}_k = \{X \leq 4 \text{ or } X \geq 11\}, \quad b : \mathcal{X}_k = \{X \leq 3 \text{ or } X \geq 12\},$$

and the smallest rejection region, containing our evidence that 5 cats eat A:

$$c : \mathcal{X}_k = \{X \leq 5 \text{ or } X \geq 10\}.$$

Now  $X_1, \dots, X_{15}$  is i.i.d. Bernoulli sample with parameter  $p$  ( $0 < p < 1$ ), where  $X_i = 1$  if cat  $i$  eats A, and 0, otherwise. The test statistic is  $X = \sum_{i=1}^{15} X_i \sim \mathcal{B}_{15}(p)$ , and the underlying alternative:

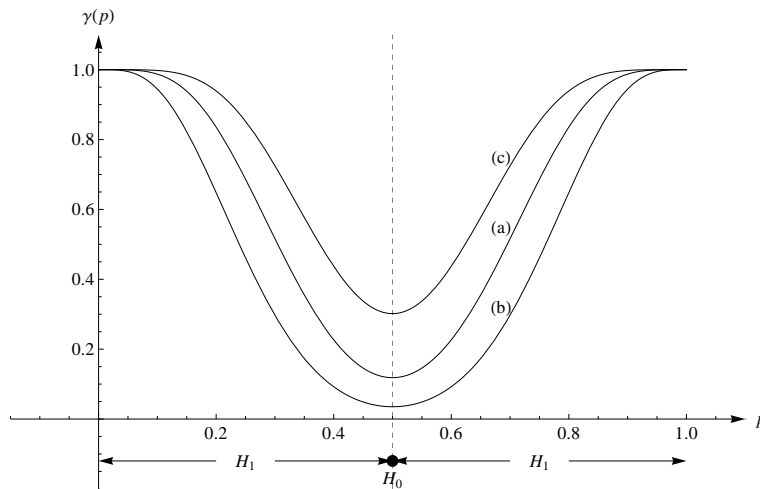
$$H_0 : p = 0.5 \quad \text{versus} \quad H_1 : p \neq 0.5.$$

Here, strategy (b) has the smallest significance (Type I error):  $\alpha = \gamma(0.5) = 0.036$  and strategy (c) has the largest: 0.302. However, for  $p \neq 0.5$ , the Type II error  $\beta(p) = 1 - \gamma(p)$  is the smallest in strategy (c).

So the division 5:10 of the cats is not enough evidence to reject  $H_0$ , but in case of 150 cats, this is a strong evidence as follows.

### 2.1. Testing the population proportion for large samples, two-sided alternative

Now  $n = 150 \geq 30$ , and so, by the special case of the CLT (Moivre–Laplace theorem),  $X = \sum_{i=1}^n X_i \sim \mathcal{B}_n(p)$  is approximately  $\mathcal{N}(np, \sqrt{np(1-p)})$ . Therefore, the population proportion  $\bar{X}$  is approximately  $\mathcal{N}(p, \sqrt{\frac{p(1-p)}{n}})$ , where for



3. ábra. The power function based on 15 cats under the three strategies,  $\alpha = \gamma(0.5)$ .

$p$ ,  $r = \bar{X}$ , and for the standard deviation,  $\sqrt{\frac{r(1-r)}{n}}$  are efficient estimators. Therefore, for the alternative

$$H_0 : p = 0.5 \quad \text{versus} \quad H_1 : p \neq 0.5$$

the test statistic

$$Z = \frac{r - 0.5}{\sqrt{r(1-r)}} \sqrt{n}$$

is approximately standard normal under  $H_0$ . The rejection region is the same as that of the two-sided  $z$ -test:

$$R = \{|z| \geq z_{\alpha/2}\}.$$

If 50 out of the 150 cats eat A, then  $r = \frac{1}{3}$  and  $z = -4.33$ . This is on the boundary of  $R$  such that  $z_{\alpha/2} = |-4.33|$ . From the standard normal table, the corresponding  $\alpha$  is near 0 (with many 0 decimals). So the P-value is practically 0, we can reject  $H_0$  with a very small significance (Type I error, that we state the difference of the foods without any reason, is very small, indeed). Also, we need not worry about the possibly large Type II error, as it is always 'small' if  $n$  is 'large'. So 50 out of 150 is a strong evidence that the foods appeal differently to cats. Of course, 40:110 or 30:120 are much stronger.

Note that accepting  $H_0$  is equivalent that the hypothetical  $p_0 = 0.5$  is within the confidence interval of level  $1 - \alpha$  constructed for  $p$ :

$$r \pm z_{\alpha/2} \sqrt{\frac{r(1-r)}{n}}. \quad (1)$$

## 2.2. Testing the population proportion for large samples, one-sided alternative

We revisit the patient recovery exercise with  $n = 200$  patients trying the new pill.

Now  $X_1, \dots, X_n$  is i.i.d. Bernoulli sample with parameter  $p$  ( $0 < p < 1$ );  $X_i = 1$  if patient  $i$  recovers from the pill, and 0 if not. Since  $n$  is large, by the CLT (Moivre–Laplace theorem),  $X = \sum_{i=1}^n X_i \sim \mathcal{B}_n(p)$  and so,  $\bar{X}$  is approximately normal, as before. Here, for the alternative

$$H_0 : p \leq 0.6 \quad \text{versus} \quad H_1 : p > 0.6$$

the test statistic is

$$Z = \frac{r - 0.6}{\sqrt{r(1-r)}} \sqrt{n},$$

that is approximately standard normal if  $p = 0.6$  (boundary of  $H_0$ ). The rejection region is:

$$R = \{z \geq z_\alpha\}.$$

Because of the monotonic nature of the power function, this is good for composite  $H_0$  too.

If 140 out of 200 patients recover, then  $r = \frac{140}{200}$  and  $z = 5.09$ . This is on the boundary of  $R$  with  $z_\alpha = 5.09$ ; so,  $\alpha$  is again 0, practically. Therefore, 140 recover out of 200 is a strong evidence to prove that the new pill is more efficient than the old one. Note that 14 out of 20 was not strong enough, but this is the law of large numbers.

### 2.3. Comparing two population proportions for large samples

Now, we have two independent Bernoulli samples, with sizes  $n_1 \geq 30$ ,  $n_2 \geq 30$ , and population proportions  $r_1, r_2$ . For example, we want to compare the recovery rate in two patient groups. Then

$$Z = \frac{r_1 - r_2}{\sqrt{\frac{r_1(1-r_1)}{n_1} + \frac{r_2(1-r_2)}{n_2}}}$$

under  $H_0$  (that  $p_1 = p_2$ ) is approximately  $\mathcal{N}(0, 1)$ .

Note that in this two-sample, two-sided case, the acceptance of  $H_0$  is equivalent to the fact that  $p_1 - p_2$  is within the confidence interval

$$r_1 - r_2 \pm z_{\alpha/2} \sqrt{\frac{r_1(1-r_1)}{n_1} + \frac{r_2(1-r_2)}{n_2}}.$$

Remark: instead of the standard deviation,  $\sqrt{\frac{\hat{r}(1-\hat{r})}{n}}$  can as well be used, where  $\hat{r} = \frac{n_1 r_1 + n_2 r_2}{n}$  and  $n = n_1 + n_2$ . With this pooled s.d.,

$$Z' = \frac{r_1 - r_2}{\sqrt{\hat{r}(1-\hat{r})} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

resembling the formula of the independent sample  $t$ -test in the next section.

### 3. Two-sample $t$ -test

Let  $X_1, \dots, X_{n_1} \sim \mathcal{N}(\mu_1, \sigma)$  be i.i.d. sample, and independently of it, let  $Y_1, \dots, Y_{n_2} \sim \mathcal{N}(\mu_2, \sigma)$  be another i.i.d. sample. Here  $n_1, n_2 < 30$  and the unknown s.d.  $\sigma$  is assumed to be the same in the two samples. For this, first we perform an  $F$ -test, see the next section.

First test the following two-sided alternative:

$$H_0 : \mu_1 = \mu_2 \quad \text{vers.} \quad H_1 : \mu_1 \neq \mu_2 \quad (2)$$

We construct a statistic the distribution of which under  $H_0$  is Student  $t$ . Indeed, under  $H_0$

$$\bar{X} - \bar{Y} \sim \mathcal{N}\left(0, \sqrt{\frac{n_1 + n_2}{n_1 n_2}} \sigma\right),$$

standardize it, and put it into the numerator of the  $t$ -statistic to be constructed.

By the Lukács's theorem,

$$(n_1 - 1)S_X^{*2}/\sigma^2 \sim \chi^2(n_1 - 1)$$

and independently of this,

$$(n_2 - 1)S_Y^{*2}/\sigma^2 \sim \chi^2(n_2 - 1),$$

therefore,

$$\frac{(n_1 - 1)S_X^{*2} + (n_2 - 1)S_Y^{*2}}{\sigma^2} \sim \chi^2(n_1 + n_2 - 2).$$

The squareroot of this divided by the d.f.  $n_1 + n_2 - 2$  is put into the denominator. The numerator and denominator are independent r.v.'s by the Lukács's theorem. Therefore, the test statistic is

$$\frac{\frac{\bar{X} - \bar{Y} - 0}{\sqrt{\frac{n_1 + n_2}{n_1 n_2}} \sigma}}{\sqrt{\frac{(n_1 - 1)S_X^{*2} + (n_2 - 1)S_Y^{*2}/\sigma^2}{n_1 + n_2 - 2}}} \sim t(n_1 + n_2 - 2),$$

where the unknown (but same)  $\sigma$  cancels, and we get that

$$t = \frac{\bar{X} - \bar{Y}}{s_{pooled} \sqrt{\frac{n_1 + n_2}{n_1 n_2}}} = \frac{\bar{X} - \bar{Y}}{s_{pooled} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

where

$$s_{pooled} = \sqrt{\frac{(n_1 - 1)S_X^{*2} + (n_2 - 1)S_Y^{*2}/\sigma^2}{n_1 + n_2 - 2}}$$

is the *pooled s.d.*, and the pooled variance  $s_{pooled}^2$  is unbiased estimator of  $\sigma^2$ .

For the two-sided alternative (4) the rejection region defining the  $\alpha$ -sign. test is:

$$R = \{|t| \geq t_{\alpha/2}(n_1 + n_2 - 2)\}.$$

. Likewise, for the one-sided alternative

$$H_0 : \mu_1 \leq \mu_2 \quad \text{vers.} \quad H_1 : \mu_1 > \mu_2, \quad (3)$$

the test statistic is the same, but the rejection region defining the  $\alpha$ -sign. test is:

$$R = \{t \geq t_\alpha(n_1 + n_2 - 2)\}.$$

If in (3) we want to test the opposite direction, we simply interchange the rolocast of the  $X - Y$  samples.

Note that the acceptance of  $H_o$  in (4) is equivalent that 0 is within the confidence interval of level  $1 - \alpha$ :

$$\bar{x} - \bar{y} \pm t_{\alpha/2}(n_1 + n_2 - 2)s_{pooled}\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}.$$

This was the independent sample  $t$ -test. IMPORTANT: in case of PARED (MATCHED) SAMPLES, one-sample test should be used for the differences  $X_i - Y_i, i = 1, \dots, n$ .

Note tat when the equality of variances is rejected (see the upcoming  $F$ -test), then a modified  $t$ -test, e.g., the Welch-test, should be used. In case of the paired sample case, it is not needed as we have only one sample  $D_i = X_i - Y_i, i = 1, \dots, n$ .

#### 4. $F$ -test

Let  $X_1, \dots, X_{n_1} \sim \mathcal{N}(\mu_1, \sigma_1)$  be i.i.d. sample, and independently of it, let  $Y_1, \dots, Y_{n_2} \sim \mathcal{N}(\mu_2, \sigma_2)$  be another i.i.d. sample.

We test the following two-sided alternative:

$$H_0 : \sigma_1 = \sigma_2 \quad \text{vers.} \quad H_1 : \sigma_1 \neq \sigma_2 \quad (4)$$

We construct a statistic the distribution of which under  $H_0$  is Fischer  $F$  (by definition, this is the distribution of the ratio of two independent  $\chi^2$ -distributed r.v.'s, each divided with its own d.f.).

Because of the above considerations,

$$F = \frac{S_X^{*2}}{S_Y^{*2}}$$

follows  $\mathcal{F}(n_1 - 1, n_2 - 1)$  distribution under  $H_0$ . However in the  $F$ -table only values at least 1 can be seen. So our test statistic is actually

$$F^* = \max\left\{\frac{S_X^{*2}}{S_Y^{*2}}, \frac{S_Y^{*2}}{S_X^{*2}}\right\} \geq 1$$

and the rejection region is

$$R = \{F^* \geq F_{\alpha/2}(f_1 - 1, f_2 - 1)\},$$

where  $f_1$  is the sample size of the sample having the largest, while  $f_2$  is the size of the sample having the smallest empirical variance.