

### 3.3 Confidence Intervals for Comparing Means and Variances of Two Populations

It will be necessary sometimes to compare characteristics of two populations. For that, we will need results on sample functions referring to both collections.

Assume we have two characteristics  $X_{(1)}$  and  $X_{(2)}$ , relative to two populations, with means  $\mu_1 = E(X_{(1)})$ ,  $\mu_2 = E(X_{(2)})$  and variances  $\sigma_1^2 = V(X_{(1)})$ ,  $\sigma_2^2 = V(X_{(2)})$ , respectively.

We draw from both populations random samples of sizes  $n_1$  and  $n_2$ , respectively, that are **independent**. Denote the two sets of random variables by

$$X_{11}, \dots, X_{1n_1} \text{ and } X_{21}, \dots, X_{2n_2}.$$

Then we have two sample means and two sample variances, given by

$$\bar{X}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} X_{1i}, \quad \bar{X}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} X_{2j}$$

and

$$s_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2, \quad s_2^2 = \frac{1}{n_2 - 1} \sum_{j=1}^{n_2} (X_{2j} - \bar{X}_2)^2,$$

respectively. In addition, denote by

$$s_p^2 = \frac{\sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2 + \sum_{j=1}^{n_2} (X_{2j} - \bar{X}_2)^2}{n_1 + n_2 - 2} = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

the **pooled variance** of the two samples, i.e. a variance that considers the sample data from both samples.

In inferential Statistics, when comparing the means of two populations, we estimate their *difference* and when comparing the variances, we estimate their *ratio*.

The formulas for finding confidence intervals for the difference of means  $\mu_1 - \mu_2$  and for the ratio of variances  $\frac{\sigma_1^2}{\sigma_2^2}$  are based on the following results (which follow either from properties of random variables, or are the consequence of some CLT).

**Proposition 3.1.** Assume  $X_{(1)} \in N(\mu_1, \sigma_1)$  and  $X_{(2)} \in N(\mu_2, \sigma_2)$ . Then

$$\text{a) } Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \in N(0, 1);$$

$$\text{b) } T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \in T(n_1 + n_2 - 2);$$

$$\text{c) } T^* = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \in T(n), \text{ where } \frac{1}{n} = \frac{c^2}{n_1 - 1} + \frac{(1 - c)^2}{n_2 - 1} \quad \text{and} \quad c = \frac{\frac{s_1^2}{n_1}}{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}};$$

$$\text{d) } F = \frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} \in F(n_1 - 1, n_2 - 1).$$

**Proposition 3.2.** *If the samples are large enough ( $n_1 + n_2 > 40$ ), then parts a), b) and c) of Proposition 3.1 still hold.*

## CI for the difference of means

### Case $\sigma_1, \sigma_2$ known

If either  $X_{(1)} \in N(\mu_1, \sigma_1)$ ,  $X_{(2)} \in N(\mu_2, \sigma_2)$  or the samples are large enough ( $n_1 + n_2 > 40$ ) and  $\sigma_1, \sigma_2$  are known, then by Propositions 3.1 and 3.2, we can use the pivot

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \in N(0, 1).$$

With the same line of computations as before, we find a  $100(1 - \alpha)\%$  CI for  $\mu_1 - \mu_2$  as

$$\mu_1 - \mu_2 \in \left[ \bar{X}_1 - \bar{X}_2 - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \bar{X}_1 - \bar{X}_2 - z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right], \quad (3.1)$$

or, using symmetry,

$$\left[ \bar{X}_1 - \bar{X}_2 \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right], \quad (3.2)$$

where the quantiles  $z_{\frac{\alpha}{2}}, z_{1-\frac{\alpha}{2}}$  refer to the  $N(0, 1)$  distribution.

### Case $\sigma_1 = \sigma_2$ unknown

Assume that either  $X_{(1)} \in N(\mu_1, \sigma_1)$ ,  $X_{(2)} \in N(\mu_2, \sigma_2)$  or the samples are large enough ( $n_1 + n_2 > 40$ ). The population variances are *not* known anymore, but they are known to be equal. Then each is approximated by the pooled variance  $s_p^2$ . Then by Propositions 3.1 and 3.2, we use the pivot

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \in T(n_1 + n_2 - 2).$$

A  $100(1 - \alpha)\%$  CI for  $\mu_1 - \mu_2$  is given by

$$\mu_1 - \mu_2 \in \left[ \bar{X}_1 - \bar{X}_2 - t_{1-\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \bar{X}_1 - \bar{X}_2 - t_{\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right], \quad (3.3)$$

where the quantiles  $t_{\frac{\alpha}{2}}, t_{1-\frac{\alpha}{2}}$  refer to the  $T(n_1 + n_2 - 2)$  distribution. Again, by symmetry we can write the CI in short as

$$\left[ \bar{X}_1 - \bar{X}_2 \pm t_{\frac{\alpha}{2}} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right]. \quad (3.4)$$

### Case $\sigma_1, \sigma_2$ unknown

Assuming that either  $X_{(1)} \in N(\mu_1, \sigma_1)$ ,  $X_{(2)} \in N(\mu_2, \sigma_2)$  or the samples are large enough ( $n_1 + n_2 > 40$ ), by Propositions 3.1 and 3.2, we use the pivot

$$T^* = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \in T(n),$$

where  $\frac{1}{n} = \frac{c^2}{n_1 - 1} + \frac{(1 - c)^2}{n_2 - 1}$  and  $c = \frac{\frac{s_1^2}{n_1}}{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$

We find a  $100(1 - \alpha)\%$  CI for  $\mu_1 - \mu_2$  as

$$\mu_1 - \mu_2 \in \left[ \bar{X}_1 - \bar{X}_2 - t_{1-\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}, \bar{X}_1 - \bar{X}_2 - t_{\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right], \quad (3.5)$$

or, by symmetry,

$$\left[ \bar{X}_1 - \bar{X}_2 \pm t_{\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right] \quad (3.6)$$

where the quantile  $t_{\frac{\alpha}{2}}, t_{1-\frac{\alpha}{2}}$  refer to the  $T(n)$  distribution, with  $n$  given above.

### CI for the ratio of variances

Assume the two independent samples were drawn from approximately Normal distributions  $N(\mu_1, \sigma_1)$  and  $N(\mu_2, \sigma_2)$ , respectively. By Proposition 3.2, we use the pivot

$$F = \frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} \in F(n_1 - 1, n_2 - 1).$$

A  $100(1 - \alpha)\%$  CI for  $\frac{\sigma_1^2}{\sigma_2^2}$  is given by

$$\frac{\sigma_1^2}{\sigma_2^2} \in \left[ \frac{1}{f_{1-\frac{\alpha}{2}}} \cdot \frac{s_1^2}{s_2^2}, \frac{1}{f_{\frac{\alpha}{2}}} \cdot \frac{s_1^2}{s_2^2} \right] \quad (3.7)$$

and, from here, a  $100(1 - \alpha)\%$  CI for  $\frac{\sigma_1}{\sigma_2}$  is

$$\frac{\sigma_1}{\sigma_2} \in \left[ \sqrt{\frac{1}{f_{1-\frac{\alpha}{2}}}} \cdot \frac{s_1}{s_2}, \sqrt{\frac{1}{f_{\frac{\alpha}{2}}}} \cdot \frac{s_1}{s_2} \right], \quad (3.8)$$

where the quantiles  $f_{\frac{\alpha}{2}}, f_{1-\frac{\alpha}{2}}$  refer to the  $F(n_1 - 1, n_2 - 1)$  distribution.

**Example 3.3.** An account on server A is more expensive than an account on server B. However, server A is faster. To see if it's optimal to go with the faster but more expensive server, a manager needs to know how much faster it is. A certain computer algorithm is executed 30 times on server A and 20 times on server B with the following results:

| Server A                      | Server B                      |
|-------------------------------|-------------------------------|
| $n_1 = 30$                    | $n_2 = 20$                    |
| $\bar{X}_1 = 6.7 \text{ min}$ | $\bar{X}_2 = 7.5 \text{ min}$ |
| $s_1 = 0.6 \text{ min}$       | $s_2 = 1.2 \text{ min}$       |

- Construct a 95% confidence interval for the difference  $\mu_1 - \mu_2$  between the mean execution times on server A and server B.
- Assuming that the observed times are approximately Normal, find a 95% confidence interval for the ratio of the two population standard deviations.

**Solution.**

a) The samples are large enough ( $n_1 + n_2 = 50$ ), that we can use Proposition 3.2. Nothing is said about the population variances (that they might be known, or known to be equal). Also, the second sample standard deviation is twice as large as the first one, therefore, equality of population variances can hardly be assumed. We use the general case for unknown, unequal variances and use formula (3.6).

We want confidence level  $1 - \alpha = 0.95$ , so  $\alpha = 0.05$  and  $\alpha/2 = 0.025$ .

The parameter  $n$  in (3.6) is found to be

$$n = 25.3989 \approx 25.$$

For the  $T(25)$  distribution, we find the quantile

$$t_{0.025} = -2.0595.$$

Then the 95% CI for the difference of means is

$$\left[ \bar{X}_1 - \bar{X}_2 \pm t_{\frac{\alpha}{2}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right] = \left[ 6.7 - 7.5 \pm 2.06 \sqrt{\frac{0.6^2}{30} + \frac{1.2^2}{20}} \right] = [-0.8 \pm 0.505],$$

so,

$$\mu_1 - \mu_2 \in [-1.305, -0.295]$$

with probability 0.95. Since *all* values in the CI are negative, with high probability, it seems that  $\mu_1 - \mu_2 < 0$ , so indeed the first server seems to be faster, on average.

b) Since now the times are assumed to be approximately Normal, we can use formula (3.8). For the  $F(29, 19)$  distribution, the quantiles are

$$f_{0.025} = 0.4482, \quad f_{0.975} = 2.4019.$$

Now,

$$\begin{aligned} \frac{s_1}{s_2} &= \frac{0.6}{1.2} = 0.5, \\ \frac{s_1^2}{s_2^2} &= \frac{0.36}{1.44} = 0.25. \end{aligned}$$

Then, the 95% CI for the ratio of variances is

$$\frac{\sigma_1^2}{\sigma_2^2} \in \left[ \frac{1}{f_{1-\frac{\alpha}{2}}} \cdot \frac{s_1^2}{s_2^2}, \frac{1}{f_{\frac{\alpha}{2}}} \cdot \frac{s_1^2}{s_2^2} \right] = \left[ \frac{1}{2.4019} \cdot 0.25, \frac{1}{0.4482} \cdot 0.25 \right] = [0.104, 0.558]$$

and the 95% CI for the ratio of standard deviations is

$$\frac{\sigma_1}{\sigma_2} \in \left[ \sqrt{\frac{1}{f_{1-\frac{\alpha}{2}}}} \cdot \frac{s_1}{s_2}, \sqrt{\frac{1}{f_{\frac{\alpha}{2}}}} \cdot \frac{s_1}{s_2} \right] = \left[ \sqrt{\frac{1}{2.4019}} \cdot 0.5, \sqrt{\frac{1}{0.4482}} \cdot 0.5 \right] = [0.323, 0.747].$$

It seems that  $\frac{\sigma_1^2}{\sigma_2^2} < 1$ , with high probability, so  $\sigma_1 < \sigma_2$ . ■

## 4 Hypothesis Testing

In the previous sections we have considered the basic ideas of parameter estimation in some detail. We attempted to approximate the value of some population parameter  $\theta$ , based on a sample, *without* having any predetermined notion concerning the actual value of this parameter. We simply tried to ascertain its value, to the best of our ability, from the information given by a random sample. In contrast, **statistical hypothesis testing** is a method of making statistical inferences on some unknown population characteristic, when *there is* a preconceived notion concerning its value or its properties.

Based on a random sample, we can use Statistics to verify a various number of statements, such as:

- the average connection speed is as claimed by the internet service provider,
- a system has not been infected,
- the proportion of defective products is at most a certain percentage, as promised by the manufacturer,
- a hardware upgrade was efficient,
- service times have a certain distribution,
- the average number of customers has increased by a certain number this year, etc.

Testing statistical hypotheses has wide applications far beyond Mathematics or Computer Science. These methods can be used to prove efficiency of a new medical treatment, safety of a new automobile brand, innocence of a defendant, authorship of a document; to establish cause-and-effect relationships; to identify factors that can significantly improve performance; to detect information leaks; and so forth.

## 4.1 Basic Concepts

So, we will work with **statistical hypotheses**, about some characteristic  $X$  (relative to a population), whose pdf  $f(x; \theta)$  depends on the parameter  $\theta$ , which is to be estimated.

The method(s) used to decide whether a hypothesis is true or not (in fact, to decide whether to *reject* a hypothesis or not) make up the **hypothesis test**. To begin with, we need to state *exactly* what we are testing. Any hypothesis test will involve two theories, two hypotheses,

- the **null hypothesis**, denoted by  $H_0$  and
- the **alternative (research) hypothesis**, denoted by  $H_1$  (or  $H_a$ ).

A null hypothesis is always an equality, showing absence of an effect or relation, some “normal” usual statement that people have believed in for years. The alternative is the opposite (in some way) of the null hypothesis, a “new” theory proposed by the researcher to “challenge” the old one. In order to overturn the common belief and to reject the null hypothesis, *significant* evidence is needed. Such evidence can only be provided by data. Only when such evidence is found, and when it *strongly* supports the alternative  $H_1$ , can the hypothesis  $H_0$  be rejected in favor of  $H_1$ . The purpose of each test is to determine whether the data provides sufficient evidence *against*  $H_0$  in favor of  $H_1$ . This is similar to a criminal trial. The jury are required to determine if the presented evidence against the defendant is sufficient and convincing. By default, the *presumption of innocence*, insufficient evidence leads to acquittal.

To determine the truth value of a hypothesis, we use a sample function called

- the **test statistic (TS)**.

The set of values of the test statistic for which we decide to *reject*  $H_0$  is called

- the **rejection region (RR)** or **critical region (CR)**.

The purpose of the experiment is to decide if the evidence (the data from a sample) tends to rebut the null hypothesis (if the value of the test statistic is in the rejection region) or not (if that value falls outside the rejection region).

If the statistical hypothesis refers to the parameter(s) of the distribution of the characteristic  $X$ , then we have a **parametric** test, otherwise, a **nonparametric** test. For parametric tests, we will consider that the target parameter

$$\theta \in A = A_0 \cup A_1, \quad A_0 \cap A_1 = \emptyset,$$

and then the two hypotheses will be set as

$$\begin{aligned} H_0 : \quad & \theta \in A_0 \\ H_1 : \quad & \theta \in A_1. \end{aligned}$$

If the set  $A_0$  consists of one single value,  $A_0 = \{\theta_0\}$ , which completely specifies the population distribution, then the hypothesis is called **simple**, otherwise, it is called a **composite** hypothesis (and the same is true for  $A_1$  and the alternative hypothesis). The null hypothesis will *always* be taken to be simple. Then the null hypothesis

$$H_0 : \theta = \theta_0$$

will have one of the alternatives

$$H_1 : \theta < \theta_0 \text{ (left-tailed test),}$$

$$H_1 : \theta > \theta_0 \text{ (right-tailed test),}$$

$$H_1 : \theta \neq \theta_0 \text{ (two-tailed test).}$$

**Remark 4.1.** The first and one of the most important tasks in a hypothesis testing problem is to state the *relevant* null and alternative hypotheses to be tested. The null hypothesis is usually taken to be a simple hypothesis, but the *appropriate* alternate has to be *understood from the context*. We mentioned that  $H_1$  is the opposite “in some way” of  $H_0$ . Let us clarify this.

1. Consider a problem in which a medicine which is believed to have the side effect of increasing the body temperature above normal, is tested. If the temperature values of a number of patients taking this medicine are considered, then for the mean temperature the relevant hypotheses would be

$$H_0 : \mu = 37$$

$$H_1 : \mu > 37,$$

since an average lower than or equal to  $37^\circ\text{C}$  would mean the same thing in this context, the patients are fine. A problem would be a mean temperature *greater* than  $37^\circ\text{C}$ . In this sense,  $H_0$  and  $H_1$  are “opposites” of each other.

2. To verify that the average broadband internet connection speed is 100 Mbps, we test the hypothesis

$$H_0 : \mu = 100$$

$$H_1 : \mu \neq 100.$$

However, if we worry about a *low* connection speed only, we can conduct a one-sided test of

$$H_0 : \mu = 100$$

$$H_1 : \mu < 100.$$

In this case, we only measure the amount of evidence supporting the one-sided alternative  $H_1 : \mu <$



100. In the absence of such evidence, we gladly accept the null hypothesis.

Designing a hypothesis test means constructing the rejection region  $RR$ , such that for a given  $\alpha \in (0, 1)$ , the conditional probability, conditioned by  $H_0$  being true,

$$P(TS \in RR \mid H_0) = \alpha. \quad (4.1)$$

The value  $\alpha$  is called **significance level** or **risk probability**.

For any given hypothesis testing problem, we have the following possibilities:

| Decision         | Actual situation                  |                                   |
|------------------|-----------------------------------|-----------------------------------|
|                  | $H_0$ true                        | $H_1$ true                        |
| Reject $H_0$     | Type I error<br>(prob. $\alpha$ ) | Right<br>decision                 |
| Not reject $H_0$ | Right<br>decision                 | Type II error<br>(prob. $\beta$ ) |

Table 1: Decisions and errors

In two of the cases, we make the right decision, in the other two, we make an error.

A **type I error** occurs when we reject a true null hypothesis and by (4.1), the probability of making such an error is the significance level

$$P(\text{type I error}) = P(\text{reject } H_0 \mid H_0) = P(TS \in RR \mid H_0) = \alpha, \quad (4.2)$$

while a **type II error** happens when we fail to reject a false null hypothesis, and its probability is denoted by  $\beta$ ,

$$P(\text{type II error}) = P(\text{not reject } H_0 \mid H_1) = P(TS \notin RR \mid H_1) = \beta. \quad (4.3)$$

**Remark 4.2.**

1. The rejection region and hence, the hypothesis test, are *not* uniquely determined by (4.1), as was the case with confidence intervals.
2. Since both  $\alpha$  and  $\beta$  represent risks of making an error, we would like to design tests such that both of their values are small. Unfortunately, making one of them very small will result in the other being unreasonably large. But, for almost all statistical tests,  $\alpha$  and  $\beta$  will both decrease as the

sample size increases.

3. In general,  $\alpha$  is preset and a procedure is given for finding an appropriate rejection region.

## 4.2 General Framework, $Z$ -Tests

Just like with confidence intervals, we start with the case where the test statistic has a  $N(0, 1)$  distribution, so we can better understand the ideas.

Let  $\theta$  be a target parameter and let  $\bar{\theta}$  be an unbiased estimator for  $\theta$  ( $E(\bar{\theta}) = \theta$ ), with standard error  $\sigma_{\bar{\theta}}$ , such that, under certain conditions, it is known that

$$Z = \frac{\bar{\theta} - \theta}{\sigma_{\bar{\theta}}} \left( = \frac{\bar{\theta} - E(\bar{\theta})}{\sigma(\bar{\theta})} \right) \quad (4.4)$$

has an approximately Standard Normal  $N(0, 1)$  distribution. We design a hypothesis testing procedure for  $\theta$  the following way: for a given level of significance  $\alpha \in (0, 1)$ , consider the hypotheses

$$H_0 : \theta = \theta_0,$$

with one of the alternatives

$$H_1 : \begin{cases} \theta < \theta_0 \\ \theta > \theta_0 \\ \theta \neq \theta_0. \end{cases} \quad (4.5)$$

We will use the test statistic  $TS = Z$  given by (4.4).

The **observed value of the test statistic** from the sample data is

$$TS_0 = TS(\theta = \theta_0). \quad (4.6)$$

In our case, this is

$$Z_0 = TS(\theta = \theta_0) = \frac{\bar{\theta} - \theta_0}{\sigma_{\bar{\theta}}}.$$

How to design the rejection region RR? Let us start with the left-tailed case. We need to determine the RR such that (4.1) holds. Intuitively, we reject  $H_0$  if the observed value of the test statistic is *far* from the value specified in  $H_0$ , “far” in the sense of the alternative  $H_1$ , in this case *far to the*

left of  $\theta_0$ . So, we determine a rejection region of the form

$$RR = \{Z_0 \mid Z_0 \leq k_1\} = (-\infty, k_1].$$

We have

$$\begin{aligned}\alpha &= P(Z_0 \in RR \mid H_0) \\ &= P(Z_0 \leq k_1 \mid \theta = \theta_0) \\ &= P(Z_0 \leq k_1 \mid Z_0 \in N(0, 1)).\end{aligned}$$

Now, we know that if  $Z_0 \in N(0, 1)$ ,  $P(Z_0 \leq z_\alpha) = \alpha$ , where  $z_\alpha$  is the quantile of order  $\alpha$  for the  $N(0, 1)$  distribution. Thus, we choose  $k_1 = z_\alpha$  and

$$RR_{\text{left}} = \{Z_0 \leq z_\alpha\}. \quad (4.7)$$

Similarly, for a right-tailed test, we want to find a rejection region of the form

$$RR = \{Z_0 \mid Z_0 \geq k_2\} = [k_2, \infty),$$

so that

$$\begin{aligned}\alpha &= P(Z_0 \in RR \mid H_0) \\ &= P(Z_0 \geq k_2 \mid \theta = \theta_0) \\ &= P(Z_0 \geq k_2 \mid Z_0 \in N(0, 1)) \\ &= 1 - P(Z_0 < k_2 \mid Z_0 \in N(0, 1)).\end{aligned}$$

Since  $P(Z_0 < z_{1-\alpha}) = 1 - \alpha$ , then  $P(Z_0 \geq z_{1-\alpha}) = \alpha$  and so we choose  $k_2 = z_{1-\alpha}$ , the quantile of order  $1 - \alpha$  for the  $N(0, 1)$  distribution and

$$RR_{\text{right}} = \{Z_0 \geq z_{1-\alpha}\}. \quad (4.8)$$

Finally, for a two-tailed test, we reject the null hypothesis if the observed value of the test statistic is far away from  $\theta_0$  *on either side*. That is, the rejection region should be of the form  $RR = \{Z_0 \mid Z_0 \leq k_1 \text{ or } Z_0 \geq k_2\} = (-\infty, k_1] \cup [k_2, \infty)$ . The rejection region should be chosen such that

$$P(Z_0 \leq k_1 \text{ or } Z_0 \geq k_2 \mid \theta = \theta_0) = \alpha,$$

or, equivalently,

$$P(k_1 < Z_0 < k_2 \mid Z_0 \in N(0, 1)) = 1 - \alpha.$$

We encountered such problems before in the previous section, when finding (two-sided) confidence intervals. As we did then, we will choose  $k_1 = z_{\frac{\alpha}{2}}$  and  $k_2 = z_{1-\frac{\alpha}{2}}$ , so

$$RR_{\text{two}} = \{Z_0 \leq z_{\frac{\alpha}{2}} \text{ or } Z_0 \geq z_{1-\frac{\alpha}{2}}\}, \quad (4.9)$$

or, since the distribution of  $Z$  is symmetric and  $z_{1-\frac{\alpha}{2}} > 0$ ,

$$\begin{aligned} RR_{\text{two}} &= \{Z_0 \leq -z_{1-\frac{\alpha}{2}} \text{ or } Z_0 \geq z_{1-\frac{\alpha}{2}}\} \\ &= \{|Z_0| \geq z_{1-\frac{\alpha}{2}}\}. \end{aligned}$$

To summarize, the rejection regions for the three alternatives (4.11) are given by

$$RR : \begin{cases} \{Z_0 \leq z_{\alpha}\} \\ \{Z_0 \geq z_{1-\alpha}\} \\ \{Z_0 \leq z_{\frac{\alpha}{2}} \text{ or } Z_0 \geq z_{1-\frac{\alpha}{2}}\} = \{|Z_0| \geq z_{1-\frac{\alpha}{2}}\}. \end{cases} \quad (4.10)$$

**Remark 4.3.**

1. Since a test statistic  $Z \in N(0, 1)$  was used, these are commonly known as **Z-tests**.
2. We will derive hypothesis tests for all the common parameters (mean, variance, difference of means, ratio of variances). The test statistics and their distributions will change, but the ideas and the principles will remain the same, as for the case we just described.
3. Notice from our derivation of the rejection region for a two-tailed test, that there is a strong relationship between confidence intervals and rejection regions: The values  $\theta_0$  of a target parameter  $\theta$  in a  $100(1 - \alpha)\%$  CI ( $\alpha \in (0, 1)$ ), are precisely the values for which the test statistic falls *outside* the RR, and hence, for which the null hypothesis  $\theta = \theta_0$  is not rejected at the significance level  $\alpha$ . We say that the  $100(1 - \alpha)\%$  two-sided CI consists of all the *acceptable* values of the parameter, at the significance level  $\alpha$ .
4. **Caution!** This is **not** saying that the rejection region is the complement of the confidence interval! The RR contains values for the *test statistic* TS, while the CI consists of values of the *parameter*  $\theta$ .

**Example 4.4.** The number of monthly sales at a firm is known to have a mean of 20 and a standard deviation of 4 and all salary, tax and bonus figures are based on these values. However, in times of economical recession, a sales manager fears that his employees do not average 20 sales per month, but less, which could seriously hurt the company. For a number of 36 randomly selected salespeople,

it was found that in one month they averaged 19 sales. At the 5% significance level, does the data confirm or contradict the manager's suspicion?

**Solution.** Recall that for  $n$  large, we have that

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

has an approximately  $N(0, 1)$  distribution. Since the sample size  $n = 36 > 30$  and  $\sigma$  is known, we can use a  $Z$ -test. The test is on the *average* number of sales per month, so for the mean  $\mu$ . The manager's suspicion is that the average is *less* than 20, which is supposed to be, so the two relevant hypotheses for this problem are

$$H_0 : \mu = 20$$

$$H_1 : \mu < 20,$$

a left-tailed test.

A type I error would mean concluding that the average number of monthly sales is less than 20, when in fact, it is not; a type II error would be deciding that the average number of monthly sales is 20 (or higher), but it actually is not. We allow for the probability of a type I error (the significance level) to be  $\alpha = 0.05$ . The population standard deviation is known,  $\sigma = 4$  and the sample mean is  $\bar{X} = 19$ .

The value of the test statistic is

$$Z_0 = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} = \frac{19 - 20}{\frac{4}{6}} = -1.5.$$

The rejection region is, by (4.10),

$$RR = (-\infty, z_\alpha] = (-\infty, -1.645].$$

Since  $Z_0 \notin RR$ , we *do not reject*  $H_0$ . The evidence obtained from the data is not sufficient to reject it. In the absence of sufficient evidence, by default, we accept the null hypothesis. So, at the 5% significance level, the data *does not* confirm the manager's suspicion.

■

### 4.3 Significance Testing, $P$ -Values

There is a problem that might occur in hypothesis testing: We preset  $\alpha$ , the probability of a type I error and henceforth determine a rejection region. We get a value of the test statistic that *does not belong* to it, so we cannot reject the null hypothesis  $H_0$ , i.e. we accept it as being true. However, when we compute the probability of getting that value of the test statistic under the assumption that  $H_0$  is true, we find it is *very small*, comparable with our preset  $\alpha$ . So, we accept  $H_0$ , yet considering it to be true, we find that it is *very unlikely* (very improbable) that the test statistic takes the observed value we found for it. That makes us wonder if we set our RR right and if we didn't "accept"  $H_0$  too easily, by hastily dismissing values of the test statistic that did not fall into our RR. So we should take a look at how "far-fetched" does the value of the test statistic seem, under the assumption that  $H_0$  is true. If it seems really implausible to occur by chance, i.e. if its probability is *small*, then maybe we should reject the null hypothesis  $H_0$ .

To avoid this situation, we perform what is called a **significance test**: for a given random sample (i.e. sample variables  $X_1, \dots, X_n$ ), we still set up  $H_0$  and  $H_1$  as before and we choose an appropriate test statistic. Then, we compute the probability of observing a value *at least as extreme* (in the sense of the test conducted) of the test statistic  $TS$  as the value observed from the sample,  $TS_0$ , under the assumption that  $H_0$  is true. This probability is called the critical value, the descriptive significance level, the probability of the test, or, simply the  **$P$ -value** of the test. If it is small, we reject  $H_0$ , otherwise we do not reject it. The  $P$ -value is a numerical value assigned to the test, it depends only on the sample data and its distribution, but *not* on  $\alpha$ .

In general, for the three alternatives (4.1), if  $TS_0$  is the value of the test statistic  $TS$  under the assumption that  $H_0$  is true and  $F$  is the cdf of  $TS$ , the  $P$ -value is computed by

$$P = \begin{cases} P(TS \leq TS_0 \mid H_0) & = F(TS_0) \\ P(TS \geq TS_0 \mid H_0) & = 1 - F(TS_0) \\ 2 \cdot \min\{P(TS \leq TS_0 \mid H_0), P(TS \geq TS_0 \mid H_0)\} & = 2 \cdot \min\{F(TS_0), 1 - F(TS_0)\}. \end{cases} \quad (4.11)$$

Then the decision will be

$$\begin{aligned} & \text{if } P \leq \alpha, \text{ reject } H_0, \\ & \text{if } P > \alpha, \text{ do not reject } H_0. \end{aligned} \quad (4.12)$$

So, more precisely, the  $P$ -value of a test is the smallest level at which we could have preset  $\alpha$  and still have been able to reject  $H_0$ , or the lowest significance level that *forces* rejection of  $H_0$ , i.e. the *minimum rejection level*.

**Remark 4.5.**

1. Thus, we can avoid the costly computation of the rejection region (costly because of the quantiles) and compute the  $P$ -value instead. Then, we simply compare it to the significance level  $\alpha$ . If  $\alpha$  is above the  $P$ -value, we reject  $H_0$ , but if it is below that minimum rejection level, we can no longer reject the null hypothesis.
2. Hypothesis testing (determining the rejection region) and significance testing (computing the  $P$ -value) are two methods for testing *the same* thing (the same two hypotheses), so, of course, the outcome (the decision of rejecting or not  $H_0$ ) will be *the same*, for the same data.

**Example 4.6.** Let us perform a significance test for the problem in Example 4.4.

**Solution.** We tested a left-tailed alternative for the mean

$$H_0 : \mu = 20$$

$$H_1 : \mu < 20.$$

The population standard deviation was given,  $\sigma = 4$ , and for a sample of size  $n = 36$ , the sample mean was  $\bar{X} = 19$ . For the test statistic

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \in N(0, 1),$$

the observed value was

$$Z_0 = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}} = \frac{19 - 20}{\frac{4}{\sqrt{36}}} = -1.5.$$

Now, we compute the  $P$ -value

$$P = P(Z \leq Z_0) = P(Z \leq -1.5) = \Phi(-1.5) = 0.0668.$$

Since

$$\alpha = 0.05 < 0.0668 = P,$$

(is below the minimum rejection level), we do not reject  $H_0$ , so, at the 5% significance level, we conclude that the data contradicts the manager's suspicion. But, for example, at the 7% significance level, we would have rejected it.

■