



# 10

## Statistical Inference about Means and Proportions with Two Populations

### CHAPTER CONTENTS

Statistics in Practice How your name affects your buying behaviour

- 10.1 Inferences about the difference between two population means:  $\sigma_1$  and  $\sigma_2$  known
- 10.2 Inferences about the difference between two population means:  $\sigma_1$  and  $\sigma_2$  unknown
- 10.3 Inferences about the difference between two population means: matched samples
- 10.4 Inferences about the difference between two population proportions

**LEARNING OBJECTIVES** After studying this chapter and doing the exercises, you should be able to:

- |   |   |
|---|---|
| <ul style="list-style-type: none"><li>1 Construct and interpret confidence intervals and hypothesis tests for the difference between two population means, given independent samples from the two populations:<ul style="list-style-type: none"><li>1.1 When the standard deviations of the two populations are known.</li><li>1.2 When the standard deviations of the two populations are unknown.</li></ul></li></ul> | <ul style="list-style-type: none"><li>2 Construct and interpret confidence intervals and hypothesis tests for the difference between two population means, given matched samples from the two populations.</li><li>3 Construct and interpret confidence intervals and hypothesis tests for the difference between two population proportions, given independent samples from the two populations.</li></ul> |
|---|---|

In Chapters 8 and 9 we showed how to construct interval estimates and do hypothesis tests for situations involving a single population mean or a single population proportion. In this chapter we extend our discussion by showing how interval estimates and hypothesis tests can be constructed when the difference between two population means or two population proportions is of prime importance. For example, we may want to construct an interval estimate of the difference between the mean starting salary for a population of men and the mean starting salary for a population of women. Or we may want to conduct a hypothesis test to determine whether there is any difference between the proportion of defective parts in a population of parts produced by supplier A and the proportion of defective parts in a population of parts produced by supplier B.



## STATISTICS IN PRACTICE

### How your name affects your buying behaviour

In an article in the *Journal of Consumer Research* in 2011, two researchers reported results of studies on a phenomenon they called the 'last name effect'. In the consumer behaviour field, *acquisition timing* refers to the speed with which consumers respond to opportunities to acquire goods or services – for instance, opportunities to get discounts or free offers, to acquire new technology or to replace consumer goods with new models. The researchers hypothesized that people with family names starting with a letter near the end of the alphabet would react



more quickly to such opportunities than people with names beginning with a letter near the beginning of the alphabet.

Their reasoning was that, during childhood, people with names near the beginning of the alphabet tend to develop a relatively laid-back approach to 'queuing' opportunities, because their name often gives them an advantage in situations where queuing is arranged on an alphabetical basis. On the other hand, people with names near the end of the alphabet tend to be more proactive, to counteract the disadvantage they experience in alphabetically queued situations.

One of the studies reported in the *Journal of Consumer Research* measured the acquisition timing, or reaction time, of a sample of MBA students to an email offer of free tickets to a basketball game. The mean reaction time of respondents with a family name beginning with one of the last nine letters of the alphabet was 19.38 minutes, compared to 25.08 minutes for respondents whose name began with one of the first nine letters of the alphabet. This difference was found to be statistically significant, using a statistical hypothesis test known as an independent-samples  $t$  test. This result offered support for the researchers' hypothesis.

In this chapter, you will learn how to construct interval estimates and do hypothesis tests about mean and proportions with two populations. The independent-samples  $t$  test used in the consumer behaviour research is an example of such a test.

We begin by showing how to construct interval estimates and do hypothesis tests for the difference between two population means when the population standard deviations are assumed known.

## 10.1 INFERENCES ABOUT THE DIFFERENCE BETWEEN TWO POPULATION MEANS: $\sigma_1$ AND $\sigma_2$ KNOWN

Let  $\mu_1$  denote the mean of population 1 and  $\mu_2$  denote the mean of population 2. We focus on inferences about the difference between the means:  $\mu_1 - \mu_2$ . We select a simple random sample of  $n_1$  units from population 1 and a second simple random sample of  $n_2$  units from population 2. The two samples, taken separately and independently, are referred to as independent simple random samples. In this section, we assume the two population standard deviations,  $\sigma_1$  and  $\sigma_2$ , are known prior to collecting the samples. We refer to this situation as the  $\sigma_1$  and  $\sigma_2$  known case. In the following example we show how to compute a margin of error and construct an interval estimate of the difference between the two population means when  $\sigma_1$  and  $\sigma_2$  are known.

## Interval estimation of $\mu_1 - \mu_2$

Suppose a retailer such as Currys (selling TVs, DVD players, computers, photographic equipment and so on) operates two stores in Dublin, Ireland. One store is in the inner city and the other is in an out-of-town shopping centre. The regional manager noticed that products selling well in one store do not always sell well in the other. The manager believes this may be attributable to differences in customer demographics at the two locations. Customers may differ in age, education, income and so on. Suppose the manager asks us to investigate the difference between the mean ages of the customers who shop at the two stores.

Let us define population 1 as all customers who shop at the inner-city store and population 2 as all customers who shop at the out-of-town store.

$$\begin{aligned}\mu_1 &= \text{mean age of population 1} \\ \mu_2 &= \text{mean age of population 2}\end{aligned}$$

The difference between the two population means is  $\mu_1 - \mu_2$ . To estimate  $\mu_1 - \mu_2$ , we shall select a simple random sample of  $n_1$  customers from population 1 and a simple random sample of  $n_2$  customers from population 2. We then compute the two sample means.

$$\begin{aligned}\bar{x}_1 &= \text{sample mean age for the simple random sample of } n_1 \text{ inner-city customers} \\ \bar{x}_2 &= \text{sample mean age for the simple random sample of } n_2 \text{ out-of-town customers}\end{aligned}$$

The point estimator of the difference between the two populations is the difference between the sample means.

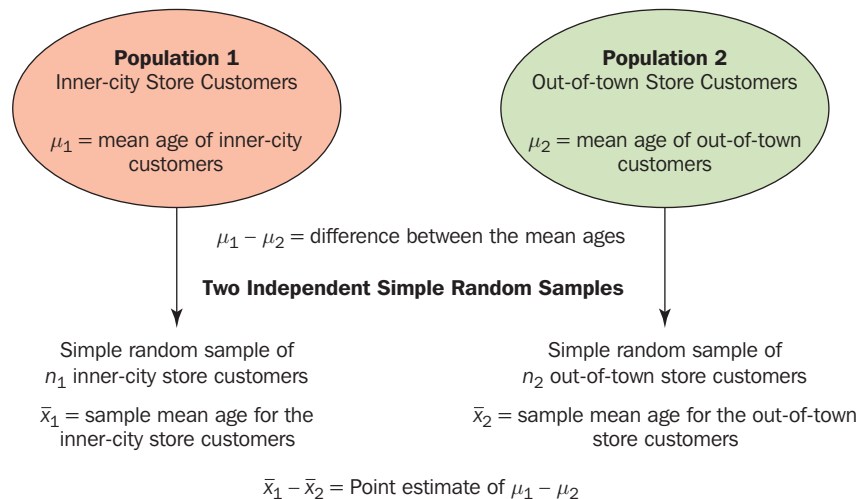
### Point estimator of the difference between two population means

$$\bar{X}_1 - \bar{X}_2 \quad (10.1)$$

Figure 10.1 provides an overview of the process used to estimate the difference between two population means based on two independent simple random samples.

**FIGURE 10.1**

Estimating the difference between two population means



As with other point estimators, the point estimator  $\bar{X}_1 - \bar{X}_2$  has a standard error that describes the variation in the sampling distribution of the estimator. With two independent simple random samples, the standard error of  $\bar{X}_1 - \bar{X}_2$  is as follows:

**Standard error of  $\bar{X}_1 - \bar{X}_2$**

$$\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad (10.2)$$

If both populations have a normal distribution, or if the sample sizes are large enough to use a normal approximation for the sampling distributions of  $\bar{X}_1$  and  $\bar{X}_2$ , the sampling distribution of  $\bar{X}_1 - \bar{X}_2$  will be normal, with a mean of  $\mu_1 - \mu_2$ .

As we showed in Chapter 8, an interval estimate is given by a point estimate  $\pm$  a margin of error. In the case of estimation of the difference between two population means, an interval estimate will take the form  $(\bar{x}_1 - \bar{x}_2) \pm$  margin of error. When the sampling distribution of  $\bar{X}_1 - \bar{X}_2$  is a normal distribution, we can write the margin of error as follows:

$$\text{Margin of error} = z_{\alpha/2} \sigma_{\bar{X}_1 - \bar{X}_2} = z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad (10.3)$$

Therefore the interval estimate of the difference between two population means is as follows:

**Interval estimate of the difference between two population means:  $\sigma_1$  and  $\sigma_2$  known**

$$(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad (10.4)$$

where  $1 - \alpha$  is the confidence coefficient.

We return to the example of the Dublin retailer. Based on data from previous customer demographic studies, the two population standard deviations are known,  $\sigma_1 = 9$  years and  $\sigma_2 = 10$  years. The data collected from the two independent simple random samples of the retailer's customers provided the following results:

	Inner-city store	Out-of-town store
Sample size	$n_1 = 36$	$n_2 = 49$
Sample mean	$\bar{x}_1 = 40$ years	$\bar{x}_2 = 35$ years

Using expression (10.1), we find that the point estimate of the difference between the mean ages of the two populations is  $\bar{x}_1 - \bar{x}_2 = 40 - 35 = 5$  years. We estimate that the customers at the inner-city store have a mean age five years greater than the mean age of the out-of-town customers. We can now use expression (10.4) to compute the margin of error and provide the interval estimate of  $\mu_1 - \mu_2$ . Using 95 per cent confidence and  $z_{\alpha/2} = z_{0.025} = 1.96$ , we have:

$$\bar{x}_1 - \bar{x}_2 \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = (40 - 35) \pm 1.96 \sqrt{\frac{(9)^2}{36} + \frac{(10)^2}{49}} = 5 \pm 4.1$$

The margin of error is 4.1 years and the 95 per cent confidence interval estimate of the difference between the two population means is  $5 - 4.1 = 0.9$  years to  $5 + 4.1 = 9.1$  years.

## Hypothesis tests about $\mu_1 - \mu_2$

Let us consider hypothesis tests about the difference between two population means. Using  $D_0$  to denote the hypothesized difference between  $\mu_1$  and  $\mu_2$ , the three forms for a hypothesis test are as follows:

$$\begin{array}{lll} H_0: \mu_1 - \mu_2 \geq D_0 & H_0: \mu_1 - \mu_2 \leq D_0 & H_0: \mu_1 - \mu_2 = D_0 \\ H_1: \mu_1 - \mu_2 < D_0 & H_1: \mu_1 - \mu_2 > D_0 & H_1: \mu_1 - \mu_2 \neq D_0 \end{array}$$

In most applications,  $D_0 = 0$ . Using the two-tailed test as an example, when  $D_0 = 0$  the null hypothesis is  $H_0: \mu_1 - \mu_2 = 0$ , i.e. the null hypothesis is that  $\mu_1$  and  $\mu_2$  are equal. Rejection of  $H_0$  leads to the conclusion that  $H_1: \mu_1 - \mu_2 \neq 0$  is true, i.e.  $\mu_1$  and  $\mu_2$  are not equal.

The steps for doing hypothesis tests presented in Chapter 9 are applicable here. We must choose a level of significance, compute the value of the test statistic and find the  $p$ -value to determine whether the null hypothesis should be rejected. With two independent simple random samples, we showed that the point estimator  $\bar{X}_1 - \bar{X}_2$  has a standard error  $\sigma_{\bar{X}_1 - \bar{X}_2}$  given by expression (10.2), and the distribution of  $\bar{X}_1 - \bar{X}_2$  can be described by a normal distribution. In this case, the test statistic for the difference between two population means when  $\sigma_1$  and  $\sigma_2$  are known is as follows.

### Test statistic for hypothesis tests about $\mu_1 - \mu_2$ : $\sigma_1$ and $\sigma_2$ known

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad (10.5)$$

Here is an example. As part of a study to evaluate differences in education quality between two training centres, a standardized examination is given to individuals trained at the centres. The difference between the mean examination scores is used to assess quality differences between the centres. The population means for the two centres are as follows:

$$\begin{array}{ll} \mu_1 = & \text{the mean examination score for the population of individuals trained at centre A} \\ \mu_2 = & \text{the mean examination score for the population of individuals trained at centre B} \end{array}$$

We begin with the tentative assumption that no difference exists between the average training quality provided at the two centres. Hence, in terms of the mean examination scores, the null hypothesis is that  $\mu_1 - \mu_2 = 0$ . If sample evidence leads to the rejection of this hypothesis, we shall conclude that the mean examination scores differ for the two populations. This conclusion indicates a quality differential between the two centres and suggests that a follow-up study investigating the reason for the differential may be warranted. The null and alternative hypotheses for this two-tailed test are written as follows:

$$\begin{array}{l} H_0: \mu_1 - \mu_2 = 0 \\ H_1: \mu_1 - \mu_2 \neq 0 \end{array}$$



EXAMSCORES

The standardized examination given previously in a variety of settings always resulted in an examination score standard deviation near ten points. We shall use this information to assume that the population standard deviations are known with  $\sigma_1 = 10$  and  $\sigma_2 = 10$ . An  $\alpha = 0.05$  level of significance is specified for the study.

Independent simple random samples of  $n_1 = 30$  individuals from training centre A and  $n_2 = 40$  individuals from training centre B are taken. The respective sample means are  $\bar{x}_1 = 82$  and  $\bar{x}_2 = 78$ . Do these data suggest a difference between the population means at the two training centres? To help answer this question, we compute the test statistic using equation (10.5):

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{(82 - 78) - 0}{\sqrt{\frac{(10)^2}{30} + \frac{(10)^2}{40}}} = 1.66$$

Next we compute the  $p$ -value for this two-tailed test. Because the test statistic  $z$  is in the upper tail, we first compute the area under the curve to the right of  $z = 1.66$ . Using the standard normal distribution table, the cumulative probability for  $z = 1.66$  is 0.9515, so the area in the upper tail of the distribution is  $1 - 0.9515 = 0.0485$ . Because this test is a two-tailed test, we must double the tail area:  $p$ -value =  $2(0.0485) = 0.0970$ . Following the usual rule to reject  $H_0$  if  $p$ -value  $\leq \alpha$ , we see that the  $p$ -value of 0.0970 does not allow us to reject  $H_0$  at the 0.05 level of significance. The sample results do not provide sufficient evidence to conclude that the training centres differ in quality.

In this chapter we shall use the  $p$ -value approach to hypothesis testing as described in Chapter 9. However, if you prefer, the test statistic and the critical value rejection rule may be used. With  $\alpha = 0.05$  and  $z_{\alpha/2} = z_{0.025} = 1.96$ , the rejection rule using the critical value approach would be to reject  $H_0$  if  $z \leq -1.96$  or if  $z \geq 1.96$ . With  $z = 1.66$ , we reach the same 'do not reject  $H_0$ ' conclusion.

In the preceding example, we demonstrated a two-tailed hypothesis test about the difference between two population means. Lower-tail and upper-tail tests can also be considered. These tests use the same test statistic as given in equation (10.5). The procedure for computing the  $p$ -value and the rejection rules for these one-tailed tests are the same as those presented in Chapter 9.

## Practical advice

In most applications of the interval estimation and hypothesis testing procedures presented in this section, random samples with  $n_1 \geq 30$  and  $n_2 \geq 30$  are adequate. In cases where either or both sample sizes are less than 30, the distributions of the populations become important considerations. In general, with smaller sample sizes, it is more important for the analyst to be satisfied that the distributions of the two populations are at least approximately normal.

## EXERCISES

### Methods

1. Consider the following results for two independent random samples taken from two populations.

Sample 1	Sample 2
$n_1 = 50$	$n_2 = 35$
$\bar{x}_1 = 13.6$	$\bar{x}_2 = 11.6$
$\sigma_1 = 2.2$	$\sigma_2 = 3.0$

- a. What is the point estimate of the difference between the two population means?
  - b. Construct a 90 per cent confidence interval for the difference between the two population means.
  - c. Construct a 95 per cent confidence interval for the difference between the two population means.
2. Consider the following hypothesis test.

$$H_0: \mu_1 - \mu_2 \leq 0$$

$$H_1: \mu_1 - \mu_2 > 0$$



The following results are for two independent samples taken from the two populations.

<i>Sample 1</i>	<i>Sample 2</i>
$n_1 = 40$	$n_2 = 50$
$\bar{x}_1 = 25.2$	$\bar{x}_2 = 22.8$
$\sigma_1 = 5.2$	$\sigma_2 = 6.0$

- What is the value of the test statistic?
  - What is the  $p$ -value?
  - With  $\alpha = 0.05$ , what is your hypothesis testing conclusion?
3. Consider the following hypothesis test.

$$H_0: \mu_1 - \mu_2 = 0$$

$$H_1: \mu_1 - \mu_2 \neq 0$$

The following results are for two independent samples taken from the two populations.

<i>Sample 1</i>	<i>Sample 2</i>
$n_1 = 80$	$n_2 = 70$
$\bar{x}_1 = 104$	$\bar{x}_2 = 106$
$\sigma_1 = 8.4$	$\sigma_2 = 7.6$

- What is the value of the test statistic?
- What is the  $p$ -value?
- With  $\alpha = 0.05$ , what is your hypothesis testing conclusion?

### Applications

4. A study of wage differentials between men and women reported that one of the reasons wages for men are higher than wages for women is that men tend to have more years of work experience than women. Assume that the following sample summaries show the years of experience for each group.

<i>Men</i>	<i>Women</i>
$n_1 = 100$	$n_2 = 85$
$\bar{x}_1 = 14.9$ years	$\bar{x}_2 = 10.3$ years
$\sigma_1 = 5.2$ years	$\sigma_2 = 3.8$ years

- What is the point estimate of the difference between the two population means?
  - At 95 per cent confidence, what is the margin of error?
  - What is the 95 per cent confidence interval estimate of the difference between the two population means?
5. The Dublin retailer age study (used as an example above) provided the following data on the ages of customers from independent random samples taken at the two store locations.

<i>Inner-city store</i>	<i>Out-of-town store</i>
$n_1 = 36$	$n_2 = 49$
$\bar{x}_1 = 40$ years	$\bar{x}_2 = 35$ years
$\sigma_1 = 9$ years	$\sigma_2 = 10$ years

- a. State the hypotheses that could be used to detect a difference between the population mean ages at the two stores.
  - b. What is the value of the test statistic?
  - c. What is the  $p$ -value?
  - d. At  $\alpha = 0.05$ , what is your conclusion?
6. Consider the following results from a survey looking at how much people spend on gifts on Valentine's Day (14 February). The average expenditure of 40 males was €135.67, and the average expenditure of 30 females was €68.64. Based on past surveys, the standard deviation for males is assumed to be €35, and the standard deviation for females is assumed to be €20. Do male and female consumers differ in the average amounts they spend?
- a. What is the point estimate of the difference between the population mean expenditure for males and the population mean expenditure for females?
  - b. At 99 per cent confidence, what is the margin of error?
  - c. Construct a 99 per cent confidence interval for the difference between the two population means.

## 10.2 INFERENCES ABOUT THE DIFFERENCE BETWEEN TWO POPULATION MEANS: $\sigma_1$ AND $\sigma_2$ UNKNOWN

In this section we extend the discussion of inferences about  $\mu_1 - \mu_2$  to the case when the two population standard deviations,  $\sigma_1$  and  $\sigma_2$ , are unknown. In this case, we use the sample standard deviations,  $s_1$  and  $s_2$ , to estimate the unknown  $\sigma_1$  and  $\sigma_2$ . The interval estimation and hypothesis testing procedures are based on the  $t$  distribution rather than the standard normal distribution.

### Interval estimation of $\mu_1 - \mu_2$

The Union Bank is conducting a study designed to identify differences between cheque account practices by customers at two of its branches. A simple random sample of 28 cheque accounts is selected from the Northern Branch and an independent simple random sample of 22 cheque accounts is selected from the Eastern Branch. The current cheque account balance is recorded for each of the accounts. A summary of the account balances follows:

	Northern	Eastern
Sample size	$n_1 = 28$	$n_2 = 22$
Sample mean	$\bar{x}_1 = €1025$	$\bar{x}_2 = €910$
Sample standard deviation	$s_1 = €150$	$s_2 = €125$



CHEQACCT

The Union Bank would like to estimate the difference between the mean cheque account balances maintained by the population of Northern customers and the population of Eastern customers. In Section 10.1, we provided the following interval estimate for the case when the population standard deviations,  $\sigma_1$  and  $\sigma_2$ , are known:

$$(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

With  $\sigma_1$  and  $\sigma_2$  unknown, we shall use the sample standard deviations  $s_1$  and  $s_2$  to estimate  $\sigma_1$  and  $\sigma_2$  and replace  $z_{\alpha/2}$  with  $t_{\alpha/2}$ . As a result, the interval estimate of the difference between two population means is given by the following expression:



**Interval estimate of the difference between two population means:  $\sigma_1$  and  $\sigma_2$  unknown**

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (10.6)$$

where  $1 - \alpha$  is the confidence coefficient.

In this expression, the use of the  $t$  distribution is an approximation, but it provides excellent results and is relatively easy to use. The only difficulty in using expression (10.6) is determining the appropriate degrees of freedom for  $t_{\alpha/2}$ . Statistical software packages compute the appropriate degrees of freedom automatically. The formula used is as follows:

**Degrees of freedom for the  $t$  distribution using two independent random samples**

$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\left(\frac{1}{n_1 - 1}\right)\left(\frac{s_1^2}{n_1}\right)^2 + \left(\frac{1}{n_2 - 1}\right)\left(\frac{s_2^2}{n_2}\right)^2} \quad (10.7)$$

We return to the Union Bank example. The sample data show  $n_1 = 28$ ,  $\bar{x}_1 = \text{€}1025$  and  $s_1 = \text{€}150$  for the Northern Branch, and  $n_2 = 22$ ,  $\bar{x}_2 = \text{€}910$  and  $s_2 = \text{€}125$  for the Eastern Branch. The calculation for degrees of freedom for  $t_{\alpha/2}$  is as follows:

$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\left(\frac{1}{n_1 - 1}\right)\left(\frac{s_1^2}{n_1}\right)^2 + \left(\frac{1}{n_2 - 1}\right)\left(\frac{s_2^2}{n_2}\right)^2} = \frac{\left(\frac{(150)^2}{28} + \frac{(125)^2}{22}\right)^2}{\left(\frac{1}{28 - 1}\right)\left(\frac{(150)^2}{28}\right)^2 + \left(\frac{1}{22 - 1}\right)\left(\frac{(125)^2}{22}\right)^2} = 47.8$$

We round the non-integer degrees of freedom *down* to 47 to provide a larger  $t$ -value and a more conservative interval estimate. Using the  $t$  distribution table with 47 degrees of freedom, we find  $t_{0.025} = 2.012$ . Using expression (10.6), we construct the 95 per cent confidence interval estimate of the difference between the two population means as follows:

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = (1025 - 910) \pm 2.012 \sqrt{\frac{(150)^2}{28} + \frac{(125)^2}{22}} = 115 \pm 78$$

The point estimate of the difference between the population mean cheque account balances at the two branches is €115. The margin of error is €78, and the 95 per cent confidence interval estimate of the difference between the two population means is  $115 - 78 = \text{€}37$  to  $115 + 78 = \text{€}193$ .

The computation of the degrees of freedom (equation (10.7)) is cumbersome if you are doing the calculation by hand, but it is easily implemented with a computer software package. Note that the terms  $s_1^2/n_1$  and  $s_2^2/n_2$  appear in both expression (10.6) and in (10.7). These need to be computed only once in order to evaluate both (10.6) and (10.7).

## Hypothesis tests about $\mu_1 - \mu_2$

Let us now consider hypothesis tests for  $\mu_1 - \mu_2$  when the population standard deviations  $\sigma_1$  and  $\sigma_2$  are unknown. Letting  $D_0$  denote the hypothesized value for  $\mu_1 - \mu_2$ , Section 10.1 showed that the test statistic used for the case where  $\sigma_1$  and  $\sigma_2$  are known is as follows. The test statistic,  $z$ , follows the standard normal distribution:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

When  $\sigma_1$  and  $\sigma_2$  are unknown, we use  $s_1$  as an estimator of  $\sigma_1$  and  $s_2$  as an estimator of  $\sigma_2$ . Substituting these sample standard deviations for  $\sigma_1$  and  $\sigma_2$  gives the following test statistic when  $\sigma_1$  and  $\sigma_2$  are unknown.

### Test statistic for hypothesis tests about $\mu_1 - \mu_2$ : $\sigma_1$ and $\sigma_2$ unknown

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (10.8)$$

The degrees of freedom for  $t$  are given by equation (10.7).

Consider an example involving a new computer software package developed to help systems analysts reduce the time required to design, develop and implement an information system. To evaluate the benefits of the new software package, a random sample of 24 systems analysts is selected. Each analyst is given specifications for a hypothetical information system. Then 12 of the analysts are instructed to produce the information system by using current technology. The other 12 analysts are trained in the use of the new software package and then instructed to use it to produce the information system.

This study involves two populations: a population of systems analysts using the current technology and a population of systems analysts using the new software package. In terms of the time required to complete the information system design project, the population means are as follows:

- $\mu_1$  = the mean project completion time for systems analysts using the current technology
- $\mu_2$  = the mean project completion time for systems analysts using the new software package

The researcher in charge of the new software evaluation project hopes to show that the new software package will provide a shorter mean project completion time, i.e. the researcher is looking for evidence to conclude that  $\mu_2$  is less than  $\mu_1$ . In this case,  $\mu_1 - \mu_2$  will be greater than zero. The research hypothesis  $\mu_1 - \mu_2 > 0$  is stated as the alternative hypothesis. The hypothesis test becomes:

$$H_0: \mu_1 - \mu_2 \leq 0$$

$$H_1: \mu_1 - \mu_2 > 0$$

We shall use  $\alpha = 0.05$  as the level of significance. Suppose that the 24 analysts complete the study with the results shown in Table 10.1.

Using the test statistic in equation (10.8), we have:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{(325 - 286) - 0}{\sqrt{\frac{(40)^2}{12} + \frac{(44)^2}{12}}} = 2.27$$



SOFTWARE  
TEST

**TABLE 10.1** Completion time data and summary statistics for the software testing study

	Current technology	New software
	300	274
	280	220
	344	308
	385	336
	372	198
	360	300
	288	315
	321	258
	376	318
	290	310
	301	332
	283	263
<b>Summary statistics</b>		
Sample size	$n_1 = 12$	$n_2 = 12$
Sample mean	$\bar{x}_1 = 325$ hours	$\bar{x}_2 = 286$ hours
Sample standard deviation	$s_1 = 40$	$s_2 = 44$

Computing the degrees of freedom using equation (10.7), we have:

$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\left(\frac{1}{n_1 - 1}\right)\left(\frac{s_1^2}{n_1}\right)^2 + \left(\frac{1}{n_2 - 1}\right)\left(\frac{s_2^2}{n_2}\right)^2} = \frac{\left(\frac{(40)^2}{12} + \frac{(44)^2}{12}\right)^2}{\left(\frac{1}{12 - 1}\right)\left(\frac{(40)^2}{12}\right)^2 + \left(\frac{1}{12 - 1}\right)\left(\frac{(44)^2}{12}\right)^2} = 21.8$$

Rounding down, we shall use a  $t$  distribution with 21 degrees of freedom. This row of the  $t$  distribution table is as follows:

Area in upper tail	0.20	0.10	0.05	0.025	0.01	0.005
$t$ value (21 df)	0.859	1.323	1.721	2.080	2.518	2.831

$t = 2.27$

With an upper-tail test, the  $p$ -value is the area in the upper tail to the right of  $t = 2.27$ . From the above results, we see that the  $p$ -value is between 0.025 and 0.01. Hence, the  $p$ -value is less than  $\alpha = 0.05$  and  $H_0$  is rejected. The sample results enable the researcher to conclude that  $\mu_1 - \mu_2 > 0$  or  $\mu_1 > \mu_2$ . The research study supports the conclusion that the new software package provides a smaller population mean completion time.

## Practical advice

The interval estimation and hypothesis testing procedures presented in this section are robust and can be used with relatively small sample sizes. In most applications, equal or nearly equal sample sizes such that the total sample size  $n_1 + n_2$  is at least 20 can be expected to provide very good results even if the populations are not normal. Larger sample sizes are recommended if the distributions of the populations are highly skewed or contain outliers. Smaller sample sizes should only be used if the analyst is satisfied that the distributions of the populations are at least approximately normal.

Another approach sometimes used to make inferences about the difference between two population means when  $\sigma_1$  and  $\sigma_2$  are unknown is based on the assumption that the two population standard

deviations are equal. You will find this approach as an option in MINITAB, IBM SPSS and EXCEL. Under the assumption of equal population variances, the two sample standard deviations are combined to provide the following 'pooled' sample variance  $s^2$ :

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

The  $t$  test statistic becomes:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

and has  $n_1 + n_2 - 2$  degrees of freedom. At this point, the computation of the  $p$ -value and the interpretation of the sample results are identical to the procedures discussed earlier in this section. A difficulty with this procedure is that the assumption of equal population standard deviations is usually difficult to verify. Unequal population standard deviations are frequently encountered. Using the pooled procedure may not provide satisfactory results especially if the sample sizes  $n_1$  and  $n_2$  are quite different. The  $t$  procedure that we presented in this section does not require the assumption of equal population standard deviations and can be applied whether the population standard deviations are equal or not. It is a more general procedure and is recommended for most applications.

## EXERCISES

### Methods

7. Consider the following results for independent random samples taken from two populations.

Sample 1	Sample 2
$n_1 = 20$	$n_2 = 30$
$\bar{x}_1 = 22.5$	$\bar{x}_2 = 20.1$
$s_1 = 2.5$	$s_2 = 4.8$

- What is the point estimate of the difference between the two population means?
- What are the degrees of freedom for the  $t$  distribution?
- At 95 per cent confidence, what is the margin of error?
- What is the 95 per cent confidence interval for the difference between the two population means?

8. Consider the following hypothesis test.

$$H_0: \mu_1 - \mu_2 = 0$$

$$H_1: \mu_1 - \mu_2 \neq 0$$

The following results are from independent samples taken from two populations.

Sample 1	Sample 2
$n_1 = 35$	$n_2 = 40$
$\bar{x}_1 = 13.6$	$\bar{x}_2 = 10.1$
$s_1 = 5.2$	$s_2 = 8.5$



**COMPLETE  
SOLUTIONS**

- a. What is the value of the test statistic?
  - b. What are the degrees of freedom for the  $t$  distribution?
  - c. What is the  $p$ -value?
  - d. At  $\alpha = 0.05$ , what is your conclusion?
9. Consider the following data for two independent random samples taken from two normal populations.

Sample 1	10	7	13	7	9	8
Sample 2	8	7	8	4	6	9

- a. Compute the two sample means.
- b. Compute the two sample standard deviations.
- c. What is the point estimate of the difference between the two population means?
- d. What is the 90 per cent confidence interval estimate of the difference between the two population means?

### Applications

10. The International Air Transport Association surveyed business travellers to determine ratings of various international airports. The maximum possible score was ten. Suppose 50 business travellers were asked to rate airport L and 50 other business travellers were asked to rate airport M. The rating scores follow.

#### Airport L

10 9 6 7 8 7 9 8 10 7 6 5 7 3 5 6 8 7 10 8 4 7 8 6 9  
9 5 3 1 8 9 6 8 5 4 6 10 9 8 3 2 7 9 5 3 10 3 5 10 8

#### Airport M

6 4 6 8 7 7 6 3 3 8 10 4 8 7 8 7 5 9 5 8 4 3 8 5 5  
4 4 4 8 4 5 6 2 5 9 9 8 4 8 9 9 5 9 7 8 3 10 8 9 6

Construct a 95 per cent confidence interval estimate of the difference between the mean ratings of the airports L and M.

11. Suppose independent random samples of 15 unionized women and 20 non-unionized women in a skilled manufacturing job provide the following hourly wage rates (€).

#### Union workers

22.40 18.90 16.70 14.05 16.20 20.00 16.10 16.30 19.10 16.50  
18.50 19.80 17.00 14.30 17.20

#### Non-union workers

17.60 14.40 16.60 15.00 17.65 15.00 17.55 13.30 11.20 15.90  
19.20 11.85 16.65 15.20 15.30 17.00 15.10 14.30 13.90 14.50

- a. What is the point estimate of the difference between mean hourly wages for the two populations?
  - b. Develop a 95 per cent confidence interval estimate of the difference between the two population means.
  - c. Does there appear to be any difference in the mean wage rate for these two groups? Explain.
12. The Scholastic Aptitude Test (SAT) is a commonly used entrance qualification for university. Consider the research hypothesis that students whose parents had attained a higher level of



AIRPORTS



UNION

education would on average score higher on the SAT. SAT verbal scores for independent samples of students follow. The first sample shows the SAT verbal test scores for students whose parents are college graduates with a bachelor's degree. The second sample shows the SAT verbal test scores for students whose parents are high school graduates but do not have a college degree.

<i>Students' parents</i>			
<i>College grads</i>		<i>High school grads</i>	
485	487	442	492
534	533	580	478
650	526	479	425
554	410	486	485
550	515	528	390
572	578	524	535
497	448		
592	469		

- Formulate the hypotheses that can be used to determine whether the sample data support the hypothesis that students show a higher population mean verbal score on the SAT if their parents attained a higher level of education.
  - What is the point estimate of the difference between the means for the two populations?
  - Compute the  $p$ -value for the hypothesis test.
  - At  $\alpha = 0.05$ , what is your conclusion?
- 13.** Periodically, Merrill Lynch customers are asked to evaluate Merrill Lynch financial consultants and services. Higher ratings on the client satisfaction survey indicate better service, with 7 the maximum service rating. Independent samples of service ratings for two financial consultants in the Dubai office are summarized here. Consultant A has ten years of experience while consultant B has one year of experience. Use  $\alpha = 0.05$  and test to see whether the consultant with more experience has the higher population mean service rating.

<i>Consultant A</i>	<i>Consultant B</i>
$n_1 = 16$	$n_2 = 10$
$\bar{x}_1 = 6.82$	$\bar{x}_2 = 6.25$
$s_1 = 0.64$	$s_2 = 0.75$

- State the null and alternative hypotheses.
  - Compute the value of the test statistic.
  - What is the  $p$ -value?
  - What is your conclusion?
- 14.** Safegate Foods is redesigning the checkouts in its supermarkets throughout the country and is considering two designs. Tests on customer checkout times conducted at two stores where the two new systems have been installed result in the following summary of the data.

<i>System A</i>	<i>System B</i>
$n_1 = 120$	$n_2 = 100$
$\bar{x}_1 = 4.1$ minutes	$\bar{x}_2 = 3.4$ minutes
$s_1 = 2.2$ minutes	$s_2 = 1.5$ minutes



**COMPLETE  
SOLUTIONS**

Test at the 0.05 level of significance to determine whether the population mean checkout times of the two systems differ. Which system is preferred?

15. Samples of final examination scores for two statistics classes with different instructors provided the following results.

<i>Instructor A</i>	<i>Instructor B</i>
$n_1 = 12$	$n_2 = 15$
$\bar{x}_1 = 72$	$\bar{x}_2 = 76$
$s_1 = 8$	$s_2 = 10$

With  $\alpha = 0.05$ , test whether these data are sufficient to conclude that the population mean grades for the two classes differ.

16. Educational testing companies provide tutoring, classroom learning and practice tests in an effort to help students perform better on tests such as the Scholastic Aptitude Test (SAT). The test preparation companies claim that their courses will improve SAT score performances by an average of 120 points. A researcher is uncertain of this claim and believes that 120 points may be an overstatement in an effort to encourage students to take the test preparation course. In an evaluation study of one test preparation service, the researcher collects SAT score data for 35 students who took the test preparation course and 48 students who did not take the course.

	<i>Course</i>	<i>No course</i>
Sample mean	1058	983
Sample standard deviation	90	105

- Formulate the hypotheses that can be used to test the researcher's belief that the improvement in SAT scores may be less than the stated average of 120 points.
- Use  $\alpha = 0.05$  and the data above. What is your conclusion?
- What is the point estimate of the improvement in the average SAT scores provided by the test preparation course? Provide a 95 per cent confidence interval estimate of the improvement.
- What advice would you have for the researcher after seeing the confidence interval?

### 10.3 INFERENCES ABOUT THE DIFFERENCE BETWEEN TWO POPULATION MEANS: MATCHED SAMPLES

Suppose employees at a manufacturing company can use two different methods to perform a production task. To maximize production output, the company wants to identify the method with the smaller population mean completion time. Let  $\mu_1$  denote the population mean completion time for production method 1 and  $\mu_2$  denote the population mean completion time for production method 2. With no preliminary indication of the preferred production method, we begin by tentatively assuming that the two production methods have the same population mean completion time. The null hypothesis is  $H_0: \mu_1 - \mu_2 = 0$ . If this hypothesis is rejected, we can conclude that the population mean completion

times differ. In this case, the method providing the smaller mean completion time would be recommended. The null and alternative hypotheses are written as follows:

$$H_0: \mu_1 - \mu_2 = 0$$

$$H_1: \mu_1 - \mu_2 \neq 0$$

In choosing the sampling procedure that will be used to collect production time data and test the hypotheses, we consider two alternative designs. One is based on **independent samples** and the other is based on **matched samples**.

- 1 Independent sample design:** A simple random sample of workers is selected and each worker in the sample uses method 1. A second independent simple random sample of workers is selected and each worker in this sample uses method 2. The test of the difference between population means is based on the procedures in Section 10.2.
- 2 Matched sample design:** One simple random sample of workers is selected. Each worker first uses one method and then uses the other method. The order of the two methods is assigned randomly to the workers, with some workers performing method 1 first and others performing method 2 first. Each worker provides a pair of data values, one value for method 1 and another value for method 2.

In the matched sample design the two production methods are tested under similar conditions (i.e. with the same workers). Hence this design often leads to a smaller sampling error than the independent sample design. The primary reason is that in a matched sample design, variation between workers is eliminated because the same workers are used for both production methods.

Let us demonstrate the analysis of a matched sample design by assuming it is the method used to test the difference between population means for the two production methods. A random sample of six workers is used. The data on completion times for the six workers are given in Table 10.2. Note that each worker provides a pair of data values, one for each production method. Also note that the last column contains the difference in completion times  $d_i$  for each worker in the sample.

The key to the analysis of the matched sample design is to realize that we consider only the column of differences. Therefore, we have six data values (0.6, -0.2, 0.5, 0.3, 0.0, 0.6) that will be used to analyze the difference between population means of the two production methods.

Let  $\mu_d$  = the mean of the *difference* values for the population of workers. With this notation, the null and alternative hypotheses are rewritten as follows:

$$H_0: \mu_d = 0$$

$$H_1: \mu_d \neq 0$$

If  $H_0$  is rejected, we can conclude that the population mean completion times differ. The  $d$  notation is a reminder that the matched sample provides *difference* data. The sample mean and sample standard deviation for the six difference values in Table 10.2 follow.



**TABLE 10.2** Task completion times for a matched sample design

Worker	Completion time for Method 1 (minutes)	Completion time for Method 2 (minutes)	Difference in completion times ( $d_i$ )
1	6.0	5.4	0.6
2	5.0	5.2	-0.2
3	7.0	6.5	0.5
4	6.2	5.9	0.3
5	6.0	6.0	0.0
6	6.4	5.8	0.6



Other than the use of the  $d$  notation, the formulae for the sample mean and sample standard deviation are the same ones used previously in the text.

$$\bar{d} = \frac{\sum d_i}{n} = \frac{1.8}{8} = 0.30$$

$$s_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}} = \sqrt{\frac{0.56}{5}} = 0.335$$

With the small sample of  $n = 6$  workers, we need to make the assumption that the population of differences has a normal distribution. This assumption is necessary so that we may use the  $t$  distribution for hypothesis testing and interval estimation procedures. Sample size guidelines for using the  $t$  distribution were presented in Chapters 8 and 9. Based on this assumption, the following test statistic has a  $t$  distribution with  $n - 1$  degrees of freedom.

**Test statistic for hypothesis test involving matched samples**


$$t = \frac{\bar{d} - \mu_d}{s_d/\sqrt{n}} \quad (10.9)$$

Let us use equation (10.9) to test the hypotheses  $H_0: \mu_d = 0$  and  $H_1: \mu_d \neq 0$ , using  $\alpha = 0.05$ . Substituting the sample results  $\bar{d} = 0.30$ ,  $s_d = 0.335$  and  $n = 6$  into equation (10.9), we compute the value of the test statistic.

$$t = \frac{\bar{d} - \mu_d}{s_d/\sqrt{n}} = \frac{0.30 - 0}{0.335/\sqrt{6}} = 2.20$$

Now let us compute the  $p$ -value for this two-tailed test. Because  $t = 2.20 > 0$ , the test statistic is in the upper tail of the  $t$  distribution. With  $t = 2.20$ , the area in the upper tail to the right of the test statistic can be found by using the  $t$  distribution table with degrees of freedom  $= n - 1 = 6 - 1 = 5$ . Information from the five degrees of freedom row of the  $t$  distribution table is as follows:

Area in upper tail	0.20	0.10	0.05	0.025	0.01	0.005
$t$ value (5 df)	0.920	1.476	2.015	2.571	3.365	4.032


  
 $t = 2.20$

We see that the area in the upper tail is between 0.05 and 0.025. Because this test is a two-tailed test, we double these values to conclude that the  $p$ -value is between 0.10 and 0.05. This  $p$ -value is greater than  $\alpha = 0.05$ , so the null hypothesis  $H_0: \mu_d = 0$  is not rejected. MINITAB, EXCEL and IBM SPSS show the  $p$ -value as 0.080.

In addition we can obtain an interval estimate of the difference between the two population means by using the single population methodology of Chapter 8. At 95 per cent confidence, the calculation follows:

$$\bar{d} \pm t_{0.025} \frac{s_d}{\sqrt{n}} = 0.30 \pm 2.527 \left( \frac{0.335}{\sqrt{6}} \right) = 0.30 \pm 0.35$$

The margin of error is 0.35 and the 95 per cent confidence interval for the difference between the population means of the two production methods is  $-0.05$  minutes to  $0.65$  minutes.

In the example presented in this section, workers performed the production task with first one method and then the other method. This example illustrates a matched sample design in which each sampled element (worker) provides a pair of data values. It is also possible to use different but 'similar' elements to

provide the pair of data values. For example, a worker at one location could be matched with a similar worker at another location (similarity based on age, education, gender, experience, etc.). The pairs of workers would provide the difference data that could be used in the matched sample analysis. A matched sample procedure for inferences about two population means generally provides better precision than the independent samples approach, therefore it is the recommended design. However, in some applications matching is not feasible, or perhaps the time and cost associated with matching are excessive. In such cases, the independent samples design should be used.

## EXERCISES

### Methods

17. Consider the following hypothesis test.

$$H_0: \mu_d \leq 0$$

$$H_1: \mu_d > 0$$

The following data are from matched samples taken from two populations.

Element	Population	
	1	2
1	21	20
2	28	26
3	18	18
4	20	20
5	26	24

- Compute the difference value for each element.
- Compute  $\bar{d}$ .
- Compute the standard deviation  $s_d$ .
- Conduct a hypothesis test using  $\alpha = 0.05$ . What is your conclusion?

18. The following data are from matched samples taken from two populations.

Element	Population	
	1	2
1	11	8
2	7	8
3	9	6
4	12	7
5	13	10
6	15	15
7	15	14

- Compute the difference value for each element.
- Compute  $\bar{d}$ .
- Compute the standard deviation  $s_d$ .
- What is the point estimate of the difference between the two population means?
- Provide a 95 per cent confidence interval for the difference between the two population means.



**COMPLETE  
SOLUTIONS**

## Applications

19. In recent years, a growing array of entertainment options has been competing for consumer time. Researchers used a sample of 15 individuals and collected data on the hours per week spent watching cable television and hours per week spent listening to the radio.

<i>Individual</i>	<i>Television</i>	<i>Radio</i>	<i>Individual</i>	<i>Television</i>	<i>Radio</i>
1	22	25	9	21	21
2	8	10	10	23	23
3	25	29	11	14	15
4	22	19	12	14	18
5	12	13	13	14	17
6	26	28	14	16	15
7	22	23	15	24	23
8	19	21			

- a. What is the sample mean number of hours per week spent watching cable television? What is the sample mean number of hours per week spent listening to radio? Which medium has the greater usage?
- b. Use a 0.05 level of significance and test for a difference between the population mean usage for cable television and radio. What is the  $p$ -value?

20. A market research firm used a sample of individuals to rate the purchase potential of a particular product before and after the individuals saw a new television commercial about the product. The purchase potential ratings were based on a 0 to 10 scale, with higher values indicating a higher purchase potential. The null hypothesis stated that the mean rating 'after' would be less than or equal to the mean rating 'before'. Rejection of this hypothesis would show that the commercial improved the mean purchase potential rating. Use  $\alpha = 0.05$  and the following data to test the hypothesis and comment on the value of the commercial.

<i>Purchase rating</i>			<i>Purchase rating</i>		
<i>Individual</i>	<i>After</i>	<i>Before</i>	<i>Individual</i>	<i>After</i>	<i>Before</i>
1	6	5	5	3	5
2	6	4	6	9	8
3	7	7	7	7	5
4	4	3	8	6	6

21. Figures on profit margins (%) for 2010 and 2011 are given below for a sample of large French companies. Use the data to comment on differences between profit margins in the two years.

<i>Profit margin (%)</i>		
<i>Company</i>	<i>2010</i>	<i>2011</i>
BNP Paribas	29.74	23.43
Carrefour	1.29	-1.50
Danone	14.64	12.59
Lafarge	8.83	4.42
L'Oréal	16.17	17.03
Michelin	8.69	9.63
Pernod-Ricard	16.43	17.94
Renault	8.61	6.17
Thales	-2.90	4.61
Vinci	8.04	7.87



TVRADIO

- a. Use  $\alpha = 0.05$  and test for any difference between the population mean profit margins in 2010 and 2011. What is the  $p$ -value? What is your conclusion?
- b. What is the point estimate of the difference between the two mean profit margins?
- c. At 95 per cent confidence, what is the margin of error for the estimate in part (b)?

22. A survey was made of Book-of-the-Month Club members to ascertain whether members spend more time watching television than they do reading. Assume a sample of 15 respondents provided the following data on weekly hours of television watching and weekly hours of reading. Using a 0.05 level of significance, can you conclude that Book-of-the-Month Club members spend more hours per week watching television than reading?

<i>Respondent</i>	<i>Television</i>	<i>Reading</i>	<i>Respondent</i>	<i>Television</i>	<i>Reading</i>
1	10	6	9	4	7
2	14	16	10	8	8
3	16	8	11	16	5
4	18	10	12	5	10
5	15	10	13	8	3
6	14	8	14	19	10
7	10	14	15	11	6
8	12	14			



PROFITS



TVREAD

## 10.4 INFERENCES ABOUT THE DIFFERENCE BETWEEN TWO POPULATION PROPORTIONS

Let  $\pi_1$  denote the proportion for population 1 and  $\pi_2$  denote the proportion for population 2. We next consider inferences about the difference between the two population proportions:  $\pi_1 - \pi_2$ . We shall select two independent random samples consisting of  $n_1$  units from population 1 and  $n_2$  units from population 2.

### Interval estimation of $\pi_1 - \pi_2$

An accountancy firm specializing in the preparation of income tax returns is interested in comparing the quality of work at two of its regional offices. The firm will be able to estimate the proportion of erroneous returns by randomly selecting samples of tax returns prepared at each office and verifying their accuracy. The difference between these proportions is of particular interest:

$\pi_1$  = proportion of erroneous returns for population 1 (office 1)

$\pi_2$  = proportion of erroneous returns for population 2 (office 2)

$P_1$  = sample proportion for a simple random sample from population 1

$P_2$  = sample proportion for a simple random sample from population 2

The difference between the two population proportions is given by  $\pi_1 - \pi_2$ . The point estimator of  $\pi_1 - \pi_2$  is as follows:

**Point estimator of the difference between two population proportions**

$$P_1 - P_2$$

**(10.10)**

The point estimator of the difference between two population proportions is the difference between the sample proportions of two independent simple random samples.

As with other point estimators, the point estimator  $P_1 - P_2$  has a sampling distribution that reflects the possible values of  $P_1 - P_2$  if we repeatedly took two independent random samples. The mean of this sampling distribution is  $\pi_1 - \pi_2$  and the standard error of  $P_1 - P_2$  is as follows:

**Standard error of  $P_1 - P_2$**

$$\sigma_{P_1 - P_2} = \sqrt{\frac{\pi_1(1 - \pi_1)}{n_1} + \frac{\pi_2(1 - \pi_2)}{n_2}} \quad (10.11)$$

If the sample sizes are large enough that  $n_1\pi_1$ ,  $n_1(1 - \pi_1)$ ,  $n_2\pi_2$  and  $n_2(1 - \pi_2)$  are all greater than or equal to five, the sampling distribution of  $P_1 - P_2$  can be approximated by a normal distribution.

As we showed previously, an interval estimate is given by a point estimate  $\pm$  a margin of error. In the estimation of the difference between two population proportions, an interval estimate will take the form  $p_1 - p_2 \pm$  margin of error. With the sampling distribution of  $P_1 - P_2$  approximated by a normal distribution, we would like to use  $z_{\alpha/2}\sigma_{P_1 - P_2}$  as the margin of error. However,  $\sigma_{P_1 - P_2}$  given by equation (10.11) cannot be used directly because the two population proportions,  $\pi_1$  and  $\pi_2$ , are unknown. Using the sample proportion  $p_1$  to estimate  $\pi_1$  and the sample proportion  $p_2$  to estimate  $\pi_2$ , the margin of error is as follows:

$$\text{Margin of error} = z_{\alpha/2} \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}} \quad (10.12)$$

The general form of an interval estimate of the difference between two population proportions is as follows:

**Interval estimate of the difference between two population proportions**

$$(p_1 - p_2) \pm z_{\alpha/2} \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}} \quad (10.13)$$

where  $1 - \alpha$  is the confidence coefficient.

Returning to the tax returns example, we find that independent simple random samples from the two offices provide the following information:

Office 1	Office 2
$n_1 = 250$	$n_1 = 300$
Number of returns with errors = 35	Number of returns with errors = 27

The sample proportions for the two offices are:

$$p_1 = \frac{35}{250} = 0.14 \quad p_2 = \frac{27}{300} = 0.09$$

The point estimate of the difference between the proportions of erroneous tax returns for the two populations is  $p_1 - p_2 = 0.14 - 0.09 = 0.05$ . We estimate that Office 1 has a 0.05, or 5 percentage points, greater error rate than Office 2.



Expression (10.13) can now be used to provide a margin of error and interval estimate of the difference between the two population proportions. Using a 90 per cent confidence interval with  $z_{\alpha/2} = z_{0.05} = 1.645$ , we have:

$$\begin{aligned} (p_1 - p_2) \pm z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \\ = (0.14 - 0.09) \pm 1.645 \sqrt{\frac{0.14(1-0.14)}{250} + \frac{0.09(1-0.09)}{300}} = 0.05 \pm 0.045 \end{aligned}$$

The margin of error is 0.045, and the 90 per cent confidence interval is 0.005 to 0.095.

## Hypothesis tests about $\pi_1 - \pi_2$

Let us now consider hypothesis tests about the difference between the proportions of two populations. The three forms for a hypothesis test are as follows:

$$\begin{array}{lll} H_0: \pi_1 - \pi_2 \geq 0 & H_0: \pi_1 - \pi_2 \leq 0 & H_0: \pi_1 - \pi_2 = 0 \\ H_1: \pi_1 - \pi_2 < 0 & H_1: \pi_1 - \pi_2 > 0 & H_1: \pi_1 - \pi_2 \neq 0 \end{array}$$

When we assume  $H_0$  is true as an equality, we have  $\pi_1 - \pi_2 = 0$ , which is the same as saying that the population proportions are equal,  $\pi_1 = \pi_2$ . The test statistic is based on the sampling distribution of the point estimator  $P_1 - P_2$ .

In expression (10.11), we showed that the standard error of  $P_1 - P_2$  is given by:

$$\sigma_{P_1 - P_2} = \sqrt{\frac{\pi_1(1-\pi_1)}{n_1} + \frac{\pi_2(1-\pi_2)}{n_2}}$$

Under the assumption that  $H_0$  is true as an equality, the population proportions are equal and  $\pi_1 = \pi_2 = \pi$ . In this case,  $\sigma_{P_1 - P_2}$  becomes:

**Standard error of  $P_1 - P_2$  when  $\pi_1 = \pi_2 = \pi$**

$$\sigma_{P_1 - P_2} = \sqrt{\frac{\pi(1-\pi)}{n_1} + \frac{\pi(1-\pi)}{n_2}} = \sqrt{\pi(1-\pi) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} \quad (10.14)$$

With  $\pi$  unknown, we pool, or combine, the point estimates from the two samples ( $p_1$  and  $p_2$ ) to obtain a single point estimate of  $\pi$  as follows:

**Pooled estimate of  $\pi$  when  $\pi_1 = \pi_2 = \pi$**

$$p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \quad (10.15)$$

This **pooled estimate of  $\pi$**  is a weighted average of  $p_1$  and  $p_2$ .

Substituting  $p$  for  $\pi$  in equation (10.14), we obtain an estimate of  $\sigma_{P_1 - P_2}$ , which is used in the test statistic. The general form of the test statistic for hypothesis tests about the difference between two population proportions is the point estimator divided by the estimate of  $\sigma_{P_1 - P_2}$ .

**Test statistic for hypothesis tests about  $\pi_1 - \pi_2$** 

$$z = \frac{(p_1 - p_2)}{\sqrt{p(1-p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \quad (10.16)$$

This test statistic applies to large sample situations where  $n_1\pi_1$ ,  $n_1(1 - \pi_1)$ ,  $n_2\pi_2$  and  $n_2(1 - \pi_2)$  are all greater than or equal to five.

Let us return to the tax returns example and assume that the firm wants to use a hypothesis test to determine whether the error proportions differ between the two offices. A two-tailed test is required. The null and alternative hypotheses are as follows:

$$H_0: \pi_1 - \pi_2 = 0$$

$$H_1: \pi_1 - \pi_2 \neq 0$$

If  $H_0$  is rejected, the firm can conclude that the error rates at the two offices differ. We shall use  $\alpha = 0.10$  as the level of significance.

The sample data previously collected showed  $p_1 = 0.14$  for the  $n_1 = 250$  returns sampled at Office 1 and  $p_2 = 0.09$  for the  $n_2 = 300$  returns sampled at Office 2. The pooled estimate of  $\pi$  is:

$$p = \frac{n_1p_1 + n_2p_2}{n_1 + n_2} = \frac{250(0.14) + 300(0.09)}{250 + 300} = 0.1127$$

Using this pooled estimate and the difference between the sample proportions, the value of the test statistic is as follows:

$$z = \frac{(p_1 - p_2)}{\sqrt{p(1-p)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} = \frac{(0.14 - 0.09)}{\sqrt{0.1127(1 - 0.1127)\left(\frac{1}{250} + \frac{1}{300}\right)}} = 1.85$$

To compute the  $p$ -value for this two-tailed test, we first note that  $z = 1.85$  is in the upper tail of the standard normal distribution. Using the standard normal distribution table, we find the area in the upper tail for  $z = 1.85$  is  $1 - 0.9678 = 0.0322$ . Doubling this area for a two-tailed test, we find the  $p$ -value =  $2(0.0322) = 0.0644$ . With the  $p$ -value less than  $\alpha = 0.10$ ,  $H_0$  is rejected at the 0.10 level of significance. The firm can conclude that the error rates differ between the two offices. This hypothesis test conclusion is consistent with the earlier interval estimation results that showed the interval estimate of the difference between the population error rates at the two offices to be 0.005 to 0.095, with Office 1 having the higher error rate.

**EXERCISES****Methods**

- 23.** Consider the following results for independent samples taken from two populations.

Sample 1	Sample 2
$n_1 = 400$	$n_2 = 300$
$p_1 = 0.48$	$p_2 = 0.36$

- a. What is the point estimate of the difference between the two population proportions?
- b. Construct a 90 per cent confidence interval for the difference between the two population proportions.
- c. Construct a 95 per cent confidence interval for the difference between the two population proportions.

**24.** Consider the hypothesis test

$$H_0: \pi_1 - \pi_2 \leq 0$$

$$H_1: \pi_1 - \pi_2 > 0$$

The following results are for independent samples taken from the two populations.

<i>Sample 1</i>	<i>Sample 2</i>
$n_1 = 200$	$n_2 = 300$
$p_1 = 0.22$	$p_2 = 0.10$

- a. What is the  $p$ -value?
- b. With  $\alpha = 0.05$ , what is your hypothesis testing conclusion?

### Applications

- 25.** In November and December 2008, research companies affiliated to the Worldwide Independent Network of Market Research carried out polls in 17 countries to assess people's views on the economic outlook. In the Canadian survey, conducted by Léger Marketing, 61 per cent of the sample of 1511 people thought the economic situation would worsen over the next three months. In the UK survey, conducted by ICM Research, 78 per cent of the sample of 1050 felt that economic conditions would worsen over that period. Provide a 95 per cent confidence interval estimate for the difference between the population proportions in the two countries. What is your interpretation of the interval estimate?
- 26.** In the results of the NUS 2011/12 Student Experience Research, it was reported that 34.3 per cent of students studying Business ( $n = 2171$ ) said a main reason for choosing their course was that the course was well-regarded by potential employers. The corresponding figure amongst students studying Maths and Computer Science ( $n = 1180$ ) was 28.1 per cent. Construct a 95 per cent confidence interval for the difference between the proportion of Business students who gave this as main reason and the proportion of Maths and Computer Science students who did likewise.
- 27.** In a test of the quality of two television commercials, each commercial was shown in a separate test area six times over a one-week period. The following week a telephone survey was conducted to identify individuals who had seen the commercials. Those individuals were asked to state the primary message in the commercials. The following results were recorded.

	<i>Commercial A</i>	<i>Commercial B</i>
Number who saw commercial	150	200
Number who recalled message	63	60

- a. Use  $\alpha = 0.05$  and test the hypothesis that there is no difference in the recall proportions for the two commercials.
- b. Compute a 95 per cent confidence interval for the difference between the recall proportions for the two populations.



**COMPLETE  
SOLUTIONS**





### COMPLETE SOLUTIONS

28. In the UNITE 2007 *Student Experience Report*, it was reported that 49 per cent of 1600 student respondents in UK universities considered the academic reputation of the university an important factor in their choice of university. In the 2012 *Student Experience Report*, 343 out of 488 respondents considered academic reputation to be important. Test the hypothesis  $\pi_1 - \pi_2 = 0$  with  $\alpha = 0.05$ . What is the  $p$ -value. What is your conclusion?
29. A large car insurance company selected samples of single and married male policyholders and recorded the number who made an insurance claim over the preceding three-year period.

Single policyholders	Married policyholders
$n_1 = 400$	$n_2 = 900$
Number making claims = 76	Number making claims = 90

- Use  $\alpha = 0.05$  and test to determine whether the claim rates differ between single and married male policyholders.
- Provide a 95 per cent confidence interval for the difference between the proportions for the two populations.



### ONLINE RESOURCES

For the data files, online summary, additional questions and answers, and software section for Chapter 10, go to the online platform.

### SUMMARY

In this chapter we discussed procedures for constructing interval estimates and doing hypothesis tests involving two populations. First, we showed how to make inferences about the difference between two population means when independent simple random samples are selected. We considered the case where the population standard deviations,  $\sigma_1$  and  $\sigma_2$ , could be assumed known. The standard normal distribution  $z$  was used to develop the interval estimate and served as the test statistic for hypothesis tests. We then considered the case where the population standard deviations were unknown and estimated by the sample standard deviations  $s_1$  and  $s_2$ . In this case, the  $t$  distribution was used to develop the interval estimate and served as the test statistic for hypothesis tests.

Inferences about the difference between two population means were then discussed for the matched sample design. In the matched sample design each element provides a pair of data values, one from each population. The difference between the paired data values is then used in the statistical analysis. The matched sample design is generally preferred to the independent sample design, when it is feasible, because the matched-samples procedure often improves the precision of the estimate.

Finally, interval estimation and hypothesis testing about the difference between two population proportions were discussed. Statistical procedures for analyzing the difference between two population proportions are similar to the procedures for analyzing the difference between two population means.

## KEY TERMS

Independent samples

Matched samples

Pooled estimator of  $\pi$

## KEY FORMULAE

Point estimator of the difference between two population means

$$\bar{X}_1 - \bar{X}_2 \quad (10.1)$$

Standard error of  $\bar{X}_1 - \bar{X}_2$

$$\sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad (10.2)$$

Interval estimate of the difference between two population means:  $\sigma_1$  and  $\sigma_2$  known

$$(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad (10.4)$$

Test statistic for hypothesis tests about  $\mu_1 - \mu_2$ :  $\sigma_1$  and  $\sigma_2$  known

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad (10.5)$$

Interval estimate of the difference between two population means:  $\sigma_1$  and  $\sigma_2$  unknown

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (10.6)$$

Degrees of freedom for the  $t$  distribution using two independent random samples

$$df = \frac{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\left( \frac{1}{n_1 - 1} \right) \left( \frac{s_1^2}{n_1} \right)^2 + \left( \frac{1}{n_2 - 1} \right) \left( \frac{s_2^2}{n_2} \right)^2} \quad (10.7)$$

Test statistic for hypothesis tests about  $\mu_1 - \mu_2$ :  $\sigma_1$  and  $\sigma_2$  unknown

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (10.8)$$

Test statistic for hypothesis test involving matched samples

$$t = \frac{\bar{d} - \mu_d}{s_d / \sqrt{n}} \quad (10.9)$$

Point estimator of the difference between two population proportions

$$P_1 - P_2 \quad (10.10)$$

Standard error of  $P_1 - P_2$

$$\sigma_{P_1 - P_2} = \sqrt{\frac{\pi_1(1 - \pi_1)}{n_1} + \frac{\pi_2(1 - \pi_2)}{n_2}} \quad (10.11)$$

Interval estimate of the difference between two population proportions

$$(p_1 - p_2) \pm z_{\alpha/2} \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}} \quad (10.13)$$

Standard error of  $P_1 - P_2$  when  $\pi_1 = \pi_2 = \pi$

$$\sigma_{P_1 - P_2} = \sqrt{\frac{\pi(1 - \pi)}{n_1} + \frac{\pi(1 - \pi)}{n_2}} = \sqrt{\pi(1 - \pi) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} \quad (10.14)$$

Pooled estimate of  $\pi$  when  $\pi_1 = \pi_2 = \pi$

$$p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \quad (10.15)$$

Test statistic for hypothesis tests about  $\pi_1 - \pi_2$

$$z = \frac{(p_1 - p_2)}{\sqrt{p(1 - p) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad (10.16)$$

## CASE PROBLEM



### Par Products

Par Products is a major manufacturer of golf equipment. Management believes that Par's market share could be increased with the introduction of a cut-resistant,

longer-lasting golf ball. Therefore, the research group at Par has been investigating a new golf ball coating designed to resist cuts and provide a more durable ball. The tests with the coating have been promising.

One of the researchers voiced concern about the effect of the new coating on driving distances. Par would like the new cut-resistant ball to offer driving distances comparable to those of the current-model

golf ball. To compare the driving distances for the two balls, 40 balls of both the new and current models were subjected to distance tests. The testing was performed with a mechanical hitting machine so that any difference between the mean

distances for the two models could be attributed to a difference in the two models. The results of the tests, with distances measured to the nearest metre, are available on the online platform, in the file 'Golf'.



Model		Model		Model		Model	
Current	New	Current	New	Current	New	Current	New
264	277	270	272	263	274	281	283
261	269	287	259	264	266	274	250
267	263	289	264	284	262	273	253
272	266	280	280	263	271	263	260
258	262	272	274	260	260	275	270
283	251	275	281	283	281	267	263
258	262	265	276	255	250	279	261
266	289	260	269	272	263	274	255
259	286	278	268	266	278	276	263
270	264	275	262	268	264	262	279



### Managerial report

1. Formulate and present the rationale for a hypothesis test that Par could use to compare the driving distances of the current and new golf balls.
2. Analyze the data to provide the hypothesis test conclusion. What is the  $p$ -value for your test? What is your recommendation for Par Products?
3. Provide descriptive statistical summaries of the data for each model.
4. What is the 95 per cent confidence interval for the population mean of each model, and what is the 95 per cent confidence interval for the difference between the means of the two populations?
5. Do you see a need for larger sample sizes and more testing with the golf balls? Discuss.