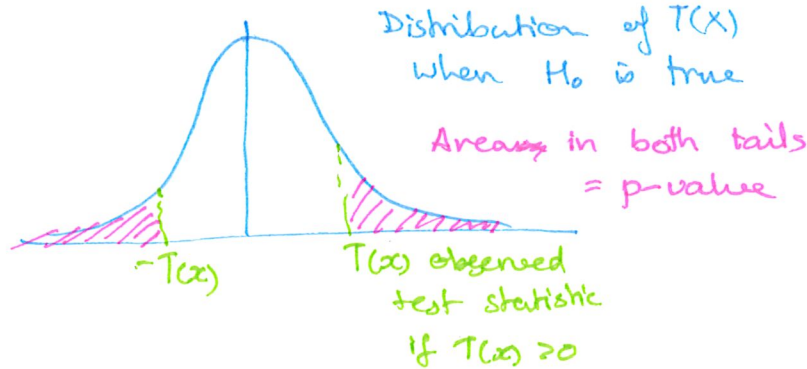


- Two sided alternative ( $H_1 : \theta \neq \theta_0$ ) The  $p$ -value is given by

$$2 \min [\mathbb{P}(T(\mathbf{X}) > T(\mathbf{x})), \mathbb{P}(T(\mathbf{X}) < T(\mathbf{x}))],$$

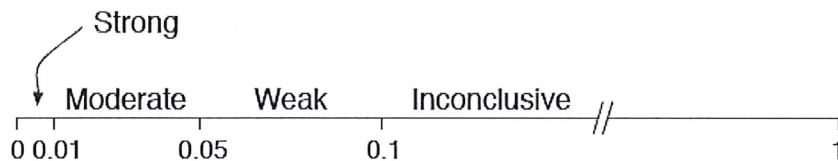
where the probability is evaluated under the null hypothesis.



Like the test statistic, the  $p$ -value is a function data and so it also has a distribution. Under the null hypothesis

$$p\text{-value} \sim \text{Uniform}(0, 1).$$

The strength of evidence against the null hypothesis provided by the  $p$ -value is summarised in the figure below.



We must decide how small the  $p$ -value must be before we reject the null hypothesis. This cut-off point is called the **significance level** and is often denoted by  $\alpha$ . The significance level determines the probability that we reject the null hypothesis when it is in fact true.

**Question:** Suppose you were to toss a coin that you believed was fair several times. How many consecutive heads would need to appear before you begin to doubt that it is really a fair coin?

It is common to use significance levels of 5% or 1%, though sometimes much smaller significance levels are needed.

**Example:** Assume that the measurements from Cavendish's apparatus are a realisation of a simple random sample from  $\text{Normal}(\mu, \sigma^2)$ . We wish to test whether or not Cavendish's apparatus gave unbiased measurements of the density of the earth, that is we are testing

$$H_0 : \mu = 5.517g/cm^3 \quad \text{against} \quad H_1 : \mu \neq 5.517g/cm^3.$$

Cavendish made 23 measurements of the earth's density, with  $\bar{x} = 5.4835g/cm^3$  and  $s = 0.1904g/cm^3$ . The test statistic is

$$T(\mathbf{x}) = \frac{\bar{x} - 5.517}{s/\sqrt{n}} = \frac{5.4835 - 5.517}{0.1904/\sqrt{23}} = -0.8438$$

$$\text{test statistic} = \frac{\text{estimate} - \text{hypothesis}}{\text{s.e. (estimate)}}$$

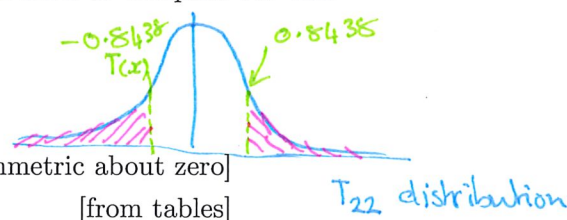
Under  $H_0$ ,  $T(\mathbf{X})$  has a  $t_{n-1}$ -distribution. So in this case we need to compare our test statistic with the  $t_{22}$ -distribution to get the  $p$ -value.

$$2 \min [\mathbb{P}(T_{22} > -0.8438), \mathbb{P}(T_{22} < -0.8438)]$$

$$= 2\mathbb{P}(T_{22} > 0.8438) \quad [\text{as } t\text{-distribution is symmetric about zero}]$$

$$= 2 \times (0.1, 0.25)$$

$$= (0.2, 0.5)$$



The  $p$ -value is between 0.2 and 0.5. This is inconclusive evidence against  $H_0$ . In other words, there is no evidence of bias in Cavendish's apparatus. At the 5% significance level, we retain the null hypothesis.

### Connection to confidence intervals

In the previous chapter, we saw how to construct a confidence interval for the mean. Let's now construct a confidence interval for the mean density reading from Cavendish's apparatus.

95% CI      estimate  $\pm$  (critical value)  $\times$  SE (estimate)

$$\bar{x} \pm t_{1-\alpha/2; n-1} \times \frac{s}{\sqrt{n}}$$

$\bar{x} = 5.4835 \text{ g/cm}^3$        $s = 0.1904 \text{ g/cm}^3$        $n = 23$

$t_{0.975; 22} = 2.074$

$$5.4835 \pm 2.074 \times \frac{0.1904}{\sqrt{23}}$$

$$5.4835 \pm 0.0823 \text{ g/cm}^3 \quad (5.401, 5.5658) \text{ g/cm}^3$$

We can be 95% confident that the mean value of the density measurements made by his apparatus was between 5.401 and 5.566  $\text{g/cm}^3$ . Note that this interval contains the hypothesised true value of 5.517  $\text{g/cm}^3$ . Is it just a coincidence that 5.517 was accepted in our hypothesis test?

There is a nice duality between confidence intervals and hypothesis testing. In fact, confidence intervals can be defined as the "inverse" of hypothesis tests:

An alternative definition of a  $(1 - \alpha)100\%$  confidence interval for a parameter  $\theta$  is

$$\{\theta \mid \theta \text{ is accepted at } \alpha \text{ significance level (two-sided test)}\}.$$

This is the set of all hypothesised parameter values that would be accepted in a two-sided hypothesis test at significance level  $\alpha$ .

## Type I and II errors

Whenever we make decisions, we run the risk of making errors. If we reject the null hypothesis when it is in fact true, we have made a **Type I error**. The probability of making a Type I error is precisely the significance level  $\alpha$  that we choose for making decisions. For example, if we think a  $p$ -value less than  $0.05 = 5\%$  is too rare to accept  $H_0$ , then we will accidentally reject  $H_0$  precisely 5% of the time.

On the other hand, if we accept  $H_0$  when it is false, then we make a **Type II error**. Related to the notion of Type II errors is the **power** of a statistical test. The power of a statistical test is the probability of detecting an effect when there is indeed an effect. If  $\beta$  is the probability of making a Type II error, then the power is given by  $1 - \beta$ .

We can think of these errors in terms of a court case:

- A Type I error is accidentally finding someone guilty when they are in fact innocent.
- A Type II error is accidentally finding someone innocent when they are in fact guilty.
- Power is the probability of finding a guilty person guilty.

To summarise, we have the following probabilities for all four scenarios:

|                | Decision                     |                              |
|----------------|------------------------------|------------------------------|
|                | Retain $H_0$                 | Reject $H_0$                 |
| $H_0$ is true  | Correct<br>( $1 - \alpha$ )  | Type I Error<br>( $\alpha$ ) |
| $H_0$ is false | Type II Error<br>( $\beta$ ) | Correct<br>( $1 - \beta$ )   |

## Comparing two means

**Example:** A real estate agency wants to compare the appraised values of studio apartments in Toowong and Dutton Park. The following results were obtained from random samples:

|                           | Toowong    | Dutton Park |
|---------------------------|------------|-------------|
| Sample Size               | 25         | 30          |
| Sample Mean               | \$ 226 716 | \$ 206 634  |
| Sample Standard Deviation | \$ 32 338  | \$ 13 464   |

Do the two regions have the same (population) mean value for studio apartments?

To address problems like this we follow the same argument that we used to construct the test of a single mean. Suppose we have a simple random sample  $X_1, \dots, X_m$  from a  $\text{Normal}(\mu_X, \sigma^2)$  distribution and another simple random sample from  $Y_1, \dots, Y_n$  from a  $\text{Normal}(\mu_Y, \sigma^2)$ . We want to test the null hypothesis  $H_0 : \mu_X - \mu_Y = d$ , for some given value  $d$  against an alternative hypothesis  $H_1$ . The alternative hypothesis is usually one of the following forms:

- One sided alternative:  $H_1 : \mu_X - \mu_Y > d$ .
- One sided alternative:  $H_1 : \mu_X - \mu_Y < d$ .
- Two sided alternative:  $H_1 : \mu_X - \mu_Y \neq d$ .

**Example:** For the real estate example, we formulate the null and alternative hypothesis as follows: Let  $\mu_T$  be the mean appraised value of a studio apartment in Toowong and let  $\mu_D$  be the mean appraised value of a studio apartment in Dutton Park.

|   |   |
|---|---|
| $( \mu_T - \mu_D = 0 )$ $H_0 : \mu_T = \mu_D$ | $( \mu_T - \mu_D \neq 0 )$ $H_1 : \mu_T \neq \mu_D$ |
|---|---|

The test statistic for this hypothesis test is

$$T(\mathbf{X}, \mathbf{Y}) = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{m} + \frac{1}{n}}}, \quad \text{estimate} - \text{hypothesis} \\ \text{se. (estimate)}$$

where  $S_p^2$  is the sample pooled variance estimator

$$\begin{aligned} S_p^2 &= \frac{(m-1)S_X^2 + (n-1)S_Y^2}{n+m-2} \\ &= \frac{\sum_{i=1}^m (X_i - \bar{X})^2 + \sum_{j=1}^n (Y_j - \bar{Y})^2}{n+m-2}. \end{aligned}$$

As we saw in the previous chapter on confidence intervals, under  $H_0$ ,

$$T(\mathbf{X}, \mathbf{Y}) \sim t_{n+m-2}.$$

When our test statistic computed from the sample data  $T(\mathbf{x}, \mathbf{y})$  is 'large' in an appropriate sense, this will indicate evidence against the null hypothesis. The  $p$ -value is given by:



- One sided alternative ( $H_1 : \mu_X - \mu_Y > d$ ) The  $p$ -value is given by

$$\mathbb{P}(T(\mathbf{X}, \mathbf{Y}) > T(\mathbf{x}, \mathbf{y})),$$

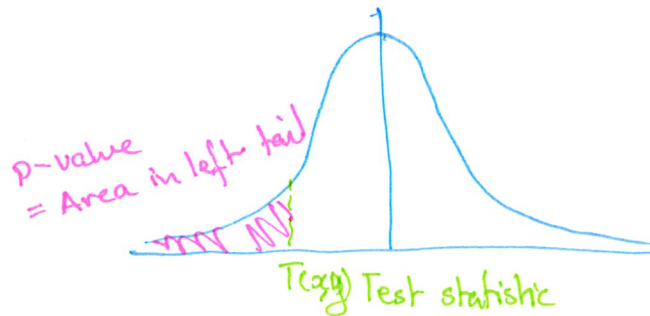
where the probability is evaluated under the null hypothesis.



- One sided alternative ( $H_1 : \mu_X - \mu_Y < d$ ) The  $p$ -value is given by

$$\mathbb{P}(T(\mathbf{X}, \mathbf{Y}) < T(\mathbf{x}, \mathbf{y})),$$

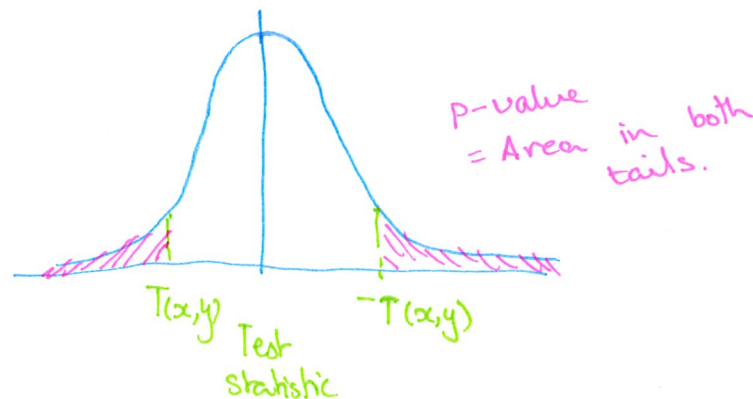
where the probability is evaluated under the null hypothesis.



- Two sided alternative ( $H_1 : \mu_X - \mu_Y \neq d$ ) The  $p$ -value is given by

$$2 \min [\mathbb{P}(T(\mathbf{X}, \mathbf{Y}) > T(\mathbf{x}, \mathbf{y})), \mathbb{P}(T(\mathbf{X}, \mathbf{Y}) < T(\mathbf{x}, \mathbf{y}))],$$

where the probability is evaluated under the null hypothesis.



**Example:** Lets now perform the test of  $H_0 : \mu_T = \mu_D$  against  $H_1 : \mu_T \neq \mu_D$ . Recall the sample data

|                           | Toowong    | Dutton Park |
|---------------------------|------------|-------------|
| Sample Size               | 25         | 30          |
| Sample Mean               | \$ 226 716 | \$ 206 634  |
| Sample Standard Deviation | \$ 32 338  | \$ 13 464   |

To compute the test statistic we need the pooled variance estimator of  $\sigma^2$ .

$$\begin{aligned}
 s_p^2 &= \frac{(n_T - 1)s_T^2 + (n_D - 1)s_D^2}{n_T + n_D - 2} \\
 &= \frac{24 \times 32338^2 + 29 \times 13464^2}{25 + 30 - 2} \\
 &= 5.7274 \times 10^8
 \end{aligned}$$

The test statistic is

$$\begin{aligned}
 T(\mathbf{x}_T, \mathbf{x}_D) &= \frac{(\bar{x}_T - \bar{x}_D) - (\mu_T - \mu_D)}{s_p \sqrt{\frac{1}{n_T} + \frac{1}{n_D}}} \\
 &= \frac{(226716 - 206634) - 0}{\sqrt{5.7274 \times 10^8} \sqrt{1/25 + 1/30}} \\
 &= 3.0987
 \end{aligned}$$

Under the null hypothesis, the test statistic has a  $t_{53}$ -distribution. The  $p$ -value is

$$\begin{aligned}
 &2 \min [\mathbb{P}(T_{53} > 3.0987), \mathbb{P}(T_{53} < -3.0987)] \\
 &= 2\mathbb{P}(T_{53} > 3.0987) \\
 &= 2 \times (0.001, 0.005) \quad \text{[from tables]} \quad \text{(round down to 50 degrees of freedom)} \\
 &= (0.002, 0.01)
 \end{aligned}$$

This is strong evidence against the null hypothesis in favour of the alternative hypothesis that the mean appraisal value for studio apartments is different for the two regions.

### Paired $t$ -test

There are situations where we have two samples  $(X_1, \dots, X_n)$  and  $(Y_1, \dots, Y_n)$  and although  $(X_1, \dots, X_n)$  are independent and  $(Y_1, \dots, Y_n)$  are independent,  $X_i$  and  $Y_i$  are dependent for all  $i$ . To compare the means of the two populations in this setting, we first take difference  $D_i = X_i - Y_i$  and then test the mean of  $D_i$ . This often arises when we have two measurements on a single subject before and after some treatment.