# Revision of Friday's Lecture

A total of 80 city council buses were randomly selected. The arrival times of each bus at their last stop were compared to the published bus timetable to determine if they were late. Sixteen buses were observed to be late. Construct a 90% confidence interval for the true proportion for council buses that arrived late at their last stop.

Let $p$ be the proportion of buses that are late (all city council buses, not just those in the sample).

general form of a CI: estimate $\pm$ critical value $\times$ s.e. (estimate)

estimate $\hat{p} = \frac{16}{80} = 0.2$, s.e. $(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{0.2 \times 0.8}{80}} = 0.04472\ldots$

90% CI: $1 - \alpha = 0.9 \rightarrow \alpha = 0.1$

critical value $Z_{1-\alpha/2} = Z_{0.95} = 1.645$

$$0.2 \pm 1.645 \times 0.04472$$

$$0.2 \pm 0.0736 \iff (0.126, 0.274)$$

We are 90% confident that the true proportion of late buses is between 0.126 and 0.274.

Having received many complaints from upset passengers that buses were always late, the council decides to implement an electronic ticketing system. Three months after introducing the electronic ticketing system, another random sample of 80 city council buses was selected and 10 were observed to be late at their last stop. Construct a 99% confidence interval for the change in the proportion of buses that arrive late at their last stop.

Assume new sample is independent of old sample.

Let $p_{new}$ be the proportion of late buses with new system

$p_{old}$ be the " " " " " " old system

$\hat{p}_{old} = 0.2$ $\qquad$ $\hat{p}_{new} = \frac{10}{80} = 0.125$

s.e. $(\hat{p}_{new} - \hat{p}_{old}) = \sqrt{\frac{\hat{p}_{new}(1-\hat{p}_{new})}{n_{new}} + \frac{\hat{p}_{old}(1-\hat{p}_{old})}{n_{old}}} = \sqrt{\frac{0.125 \times 0.875}{80} + \frac{0.2 \times 0.8}{80}}$

$$= 0.05803$$

99% CI $\qquad$ $0.99 = 1 - \alpha \iff \alpha = 0.01$ $\qquad$ $Z_{1-\alpha/2} = Z_{0.995} = 2.576$

$$\hat{p}_{new} - \hat{p}_{old} \pm Z_{0.995} \times s.e. (\hat{p}_{new} - \hat{p}_{old})$$

$$-0.075 \pm 2.576 \times 0.05803$$

$$-0.075 \pm 0.1495$$

We are 99% confident that the true difference in the proportion of late buses (new - old) is between -0.2245 and 0.0745.

# 7

# Hypothesis Testing

By the end of this chapter you should:

- Know how to specify null and alternative hypotheses.

- Be able to apply basic statistical tests.

- Know how to interpret a test statistic and a p-value.

- Be able to understand the types of errors that occur in hypothesis testing.

- Know what factors controls the probablity of these errors in hypothesis testing.

In the previous chapter we saw how to estimate basic quantities such as a mean or proportion and how to quantify our uncertainty about those estimates. Another problem that arises in the analysis of data is how to make a decision about our model. This arises naturally in a number of settings:

- Do video games increase aggressive behaviour in children?

- Does the new website design get more hits than the old?

- Does eating chilli cause memory problems?

- Note: questions of causation require other tools (in addition to hypothesis testing).

## Null and Alternative Hypotheses

In statistics, this problem is called a *hypothesis test*. In hypothesis testing, given data, we wish to determine which of two competing hypotheses: the **null hypothesis** ($H_0$) and the **alternative hypothesis** ($H_1$).

We begin with a model for the process generating our data. For example, suppose our data is a realisation of a simple random sample (that is, a realisation of a collection of

independent random variables all having the same distribution) from a $\text{Normal}(\mu, \sigma^2)$ distribution. In general, we will denote the parameter(s) of the model by $\theta$ and the set of all possible parameter values by $\Theta$. We can now specify the null and alternative hypothese in terms of the parmeter $\theta$.

Let $\Theta_0$ and $\Theta_1$ form a partition of the parmeter space $\Theta$. That is, $\boxed{\Theta_0 \cup \Theta_1 = \Theta}$ and $\boxed{\Theta_0 \cap \Theta_1 = \emptyset}$ . The null and alternative hypotheses are then specified as

$$H_0 : \theta \in \Theta_0, \qquad\qquad H_1 : \theta \in \Theta_1.$$

**Example:** The currently accepted value for the mean density of the Earth is $5.517 g/cm^3$. In 1798 Henry Cavendish presented some observations for the mean density of the Earth. Suppose Cavendish's apparatus produced measurements from a $\text{Normal}(\mu, \sigma^2)$ distribution. Potential hypotheses to test would be

- Test $H_0 : \mu = 5.517 g/cm^3$ versus $H_1 : \mu \neq 5.517 g/cm^3$ (two-sided alternative)

  $H_0$: measurements from the apparatus are unbiased.
  $H_1$: measurements from the apparatus are biased.

- Test $H_0 : \mu = 5.517 g/cm^3$ versus $H_1 : \mu > 5.517 g/cm^3$ (one-sided alternative)

  $H_1$: measurements from the apparatus over estimate the density of the earth.

- Test $H_0 : \mu = 5.517 g/cm^3$ versus $H_1 : \mu < 5.517 g/cm^3$ (one-sided alternative)

  $H_1$: measurements from the apparatus under estimate the density of the earth.

These two hypotheses are not treated symmetrically. The null hypothesis $H_0$ is taken as a statement of the "status quo" and we examine the data looking for evidence against $H_0$.

- If no evidence against $H_0$ is found, then we accept $H_0$.

- On the other hand, if evidence against $H_0$ is found (in the direction of $H_1$), then we will reject $H_0$ in favour of the alternative hypothesis $H_1$.

## Test statistics and $p$-values

So before we can decided whether or not to accept the null hypothesis, we need to be able to quantify the evidence against the null hypothesis. We do this using be constructing a **test statistic** and a **$p$-value**.

> A test statistic is a function of the data whose distribution under the null hypothesis is known.

**Example:** Suppose $X_1, \ldots, X_n$ be a simple random sample from Normal$(\mu, \sigma^2)$ with $\bar{X}$ and $S^2$ be the usual estimators of $\mu$ and $\sigma^2$ constructed from the $X_1, \ldots, X_n$. Under the null hypothesis $H_0 : \mu = 5.517 g/cm^3$, the test statistic

$$T(\mathbf{X}) = \frac{\bar{X} - 5.517}{S/\sqrt{n}}$$

has a $t_{n-1}$-distribution.

When our test statistic computed from the sample data $T(\mathbf{x})$ is 'large' in an appropriate sense, this will indicate evidence against the null hypothesis. This evidence against the null hypothesis is summarised more clearly through the use of a $p$-value.

*p-value — probability of observing data "more extreme" than what we observed if the null hypothesis is true.*

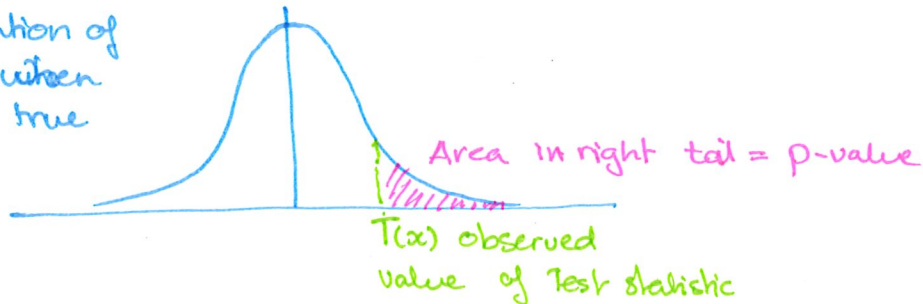- One sided alternative $(H_1 : \theta > \theta_0)$ The $p$-value is given by

*$H_0 : \Theta = \Theta_0$*

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

*computed from data*

where the probability is evaluated under the null hypothesis.

*Distribution of $T(X)$ when $H_0$ is true*



*Area in right tail = p-value*

*$T(x)$ observed value of Test statistic*

- One sided alternative $(H_1 : \theta < \theta_0)$ The $p$-value is given by
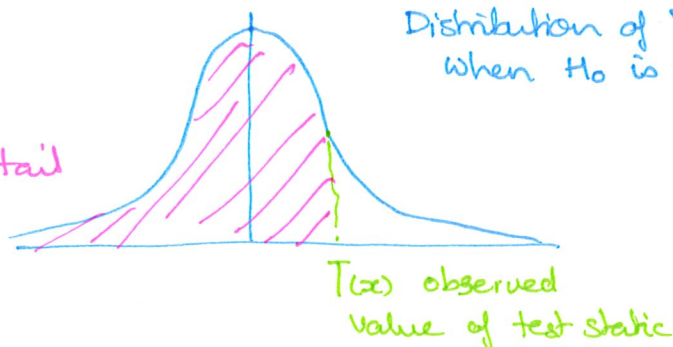
*$H_0 : \Theta = \Theta_0$*

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

where the probability is evaluated under the null hypothesis.
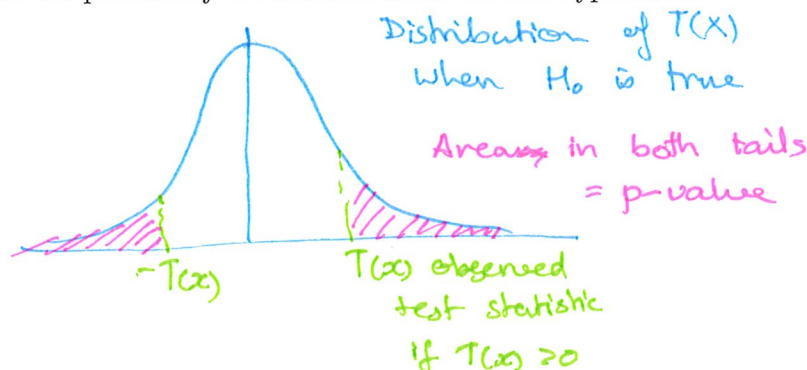
*Distribution of $T(X)$ when $H_0$ is true*

*Area in left tail = p-value*



*$T(x)$ observed value of test static*

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

$$2 \min \left[ \mathbb{P}(T(\mathbf{X}) > T(\mathbf{x})), \mathbb{P}(T(\mathbf{X}) < T(\mathbf{x})) \right],$$
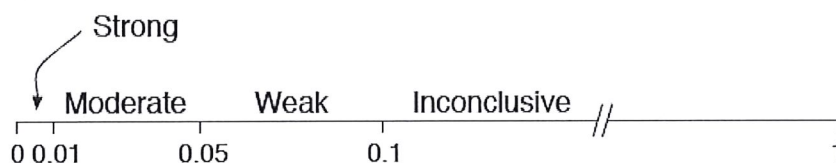
where the probability is evaluated under the null hypothesis.



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

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

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 hypotheis. 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 appartus 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

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

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

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