

# MATH3014-6027 Design (and Analysis) of Experiments

Dave Woods

2022-01-27



# Contents

<b>Preface</b>	<b>5</b>
<b>1 Motivation, introduction and revision</b>	<b>7</b>
1.1 Motivation . . . . .	8
1.2 Aims of experimentation and some examples . . . . .	12
1.3 Some definitions . . . . .	13
1.4 Principles of experimentation . . . . .	13
1.5 Revision on the linear model . . . . .	15
1.6 Exercises . . . . .	18
<b>2 Completely randomised designs</b>	<b>21</b>
2.1 A unit-treatment linear model . . . . .	23
2.2 The partitioned linear model . . . . .	24
2.3 Generalised inverses . . . . .	26
2.4 Solutions to the reduced normal equations . . . . .	26
<b>3 Blocking</b>	<b>27</b>
<b>4 Factorial experiments</b>	<b>29</b>
<b>5 Blocking in factorial designs</b>	<b>31</b>
<b>6 Fractional factorial designs</b>	<b>33</b>
<b>7 Response surface methodology</b>	<b>35</b>
<b>8 Optimal design of experiments</b>	<b>37</b>



# Preface

These are draft lecture notes for the modules MATH3014 and MATH6027 Design (and Analysis) of Experiments at the University of Southampton for academic year 2021-22. They are very much work in progress.



# Chapter 1

## Motivation, introduction and revision

**Definition 1.1.** An **experiment** is the process through which data are collected to answer a scientific question (physical science, social science, actuarial science ...) by **deliberately** varying some features of the process under study in order to understand the impact of these changes on measureable responses.

In this course we consider only *intervention* experiments, in which some aspects of the process are under the experimenters' control. We do not consider *surveys* or *observational* studies.

**Definition 1.2. Design of experiments** is the topic in Statistics concerned with the selection of settings of controllable variables or factors in an experiment and their allocation to experimental units in order to maximise the effectiveness of the experiment at achieving its aim.

People have been designing experiments for as long as they have been exploring the natural world. Some notable milestones in the history of the design of experiments include:

- prior to the 20th century:
  - Francis Bacon (17th century; pioneer of the experimental methods)
  - James Lind (18th century; experiments to eliminate scurvy)
  - Charles Peirce (19th century; advocated randomised experiments and randomisation-based inference)
- 1920s: agriculture (particularly at the Rothamsted Agricultural Research Station)
- 1940s: clinical trials (Austin Bradford-Hill)
- 1950s: (manufacturing) industry (W. Edwards Deming; Genichi Taguchi)
- 1960s: psychology and economics (Vernon Smith)
- 1980s: in-silico (computer experiments)

- 2000s: online (A/B testing)

See Luca and Bazerman (2020) for further history, anecdotes and examples, especially from psychology and technology.

Figure 1.1 shows the Broadbalk agricultural field experiment at Rothamsted, one of the longest continuous running experiments in the world, which is testing the impact of different manures and fertilizers on the growth of winter wheat.



Figure 1.1: The Broadbalk experiment, Rothamsted (photograph taken 2016)

## 1.1 Motivation

**Example 1.1.** Consider an experiment to compare two treatments (e.g. drugs, diets, fertilisers, ...). We have  $N$  subjects (people, mice, plots of land, ...), each of which can be assigned one of the two treatments. A response (protein measurement, weight, yield, ...) is then measured.

**Question:** How many subjects should be assigned to each treatment to gain the most precise<sup>1</sup> inference about the difference in response from the two treatments?

Consider a linear statistical model<sup>2</sup> for the response (see MATH2010):

$$Y_j = \beta_0 + \beta_1 x_j + \varepsilon_j, \quad j = 1, \dots, n, \quad (1.1)$$

---

<sup>1</sup>Smallest variance.

<sup>2</sup>In this course, we will almost always start with a statistical model which we wish to use to answer our scientific question.



where  $\varepsilon_j \sim N(0, \sigma^2)$  are independent and identically distributed errors and  $\beta_0, \beta_1$  are unknown constants (parameters).

Let<sup>3</sup>

$$x_j = \begin{cases} -1 & \text{if treatment 1 is applied to the } j\text{th subject} \\ +1 & \text{if treatment 2 is applied to the } j\text{th subject,} \end{cases}$$

for  $j = 1, \dots, n$ .<sup>4</sup>

The difference in expected response from treatments 1 and 2 is

$$\begin{aligned} E[Y_j | x_j = +1] - E[Y_j | x_j = -1] &= \beta_0 + \beta_1 - \beta_0 + \beta_1 \\ &= 2\beta_1. \end{aligned} \quad (1.2)$$

Therefore, we require the the most precise estimator of  $\beta_1$  possible. That is, we wish to make the variance of our estimator of  $\beta_1$  as small as possible.

Parameters  $\beta_0$  and  $\beta_1$  can be estimated using least squares (see MATH2010). For  $Y_1, \dots, Y_n$ , we can write the model down in matrix form:

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}.$$

Or, by defining some notation:

$$Y = X\beta + \varepsilon \quad (1.3)$$

where

- $Y$  -  $n \times 1$  vector of responses;
- $X$  -  $n \times p$  model matrix;
- $\beta$  -  $p \times 1$  vector of parameters;
- $\varepsilon$  -  $n \times 1$  vector of errors.

The **least squares estimators**,  $\hat{\beta}$ , are chosen such that the quadratic form

$$(Y - X\beta)^T(Y - X\beta)$$

is minimised (recall that  $E(\mathbf{Y}) = X\beta$ ). Therefore

$$\hat{\beta} = \operatorname{argmin}_{\beta} (Y^T Y + \beta^T X^T X \beta - 2\beta^T X^T Y).$$

---

<sup>3</sup>Other codings can be used: e.g. 0,1; see later in the module. It makes no difference for our current purpose.

<sup>4</sup>We will discuss the choice of *coding* -1, +1 later.

If we differentiate with respect to  $\beta^5$ ,

$$\frac{\partial}{\partial \beta} = 2X^T X \beta - 2X^T Y,$$

and equate to 0, we get the estimators

$$\hat{\beta} = (X^T X)^{-1} X^T Y. \quad (1.4)$$

These are the least squares estimators.

For Example 1.1,

$$X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad X^T X = \begin{bmatrix} n & \sum x_j \\ \sum x_j & \sum x_j^2 \end{bmatrix},$$

$$(X^T X)^{-1} = \frac{1}{n \sum x_j^2 - (\sum x_j)^2} \begin{bmatrix} \sum x_j^2 & -\sum x_j \\ -\sum x_j & n \end{bmatrix}, \quad X^T Y = \begin{bmatrix} \sum Y_j \\ \sum x_j Y_j \end{bmatrix}.$$

Then,

$$\begin{aligned} \hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} &= \frac{1}{n \sum x_j^2 - (\sum x_j)^2} \begin{bmatrix} \sum x_j^2 & -\sum x_j \\ -\sum x_j & n \end{bmatrix} \begin{bmatrix} \sum Y_j \\ \sum x_j Y_j \end{bmatrix} \\ &= \frac{1}{n \sum x_j^2 - (\sum x_j)^2} \begin{bmatrix} \sum Y_j \sum x_j^2 - \sum x_j \sum x_j Y_j \\ n \sum x_j Y_j - \sum x_j \sum Y_j \end{bmatrix}. \end{aligned} \quad (1.5)$$

We don't usually work through the algebra in such detail; the matrix form is often sufficient for theoretical and numerical calculations and software, e.g. R, can be used.

The precision of  $\hat{\beta}$  is measured via the variance-covariance matrix, given by

$$\text{Var}(\hat{\beta}) = \text{Var}\{(X^T X)^{-1} X^T Y\} \quad (1.6)$$

$$= (X^T X)^{-1} X^T \text{Var}(Y) X (X^T X)^{-1} \quad (1.7)$$

$$= (X^T X)^{-1} \sigma^2, \quad (1.8)$$

where  $Y \sim N(X\beta, I_n \sigma^2)$ , where  $I_n$  is an  $n \times n$  identity matrix.

---

<sup>5</sup>Check the Matrix Cookbook for matrix calculus, amongst much else.

Hence, in our example,

$$\begin{aligned}\text{Var}(\hat{\beta}) &= \frac{1}{n \sum x_j^2 - (\sum x_j)^2} \begin{bmatrix} \sum x_j^2 & -\sum x_j \\ -\sum x_j & n \end{bmatrix} \sigma^2 \\ &= \begin{bmatrix} \text{Var}(\hat{\beta}_0) & \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \text{Var}(\hat{\beta}_1) \end{bmatrix}.\end{aligned}$$

For estimating the difference between treatments, we are interested in

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \frac{n}{n \sum x_j^2 - (\sum x_j)^2} \sigma^2 \\ &= \frac{n}{n^2 - (\sum x_j)^2} \sigma^2,\end{aligned}$$

as  $x_j = \pm 1$ , therefore  $x_j^2 = 1$  for all  $j = 1, \dots, n$ , and hence  $\sum x_j^2 = n$ .

To achieve the most precise estimator, we need to minimise  $\text{Var}(\hat{\beta}_1)$  or equivalently minimise  $|\sum x_j|$ . This goal can be achieved through the choice of  $x_1, \dots, x_N$ :

- as each  $x_j$  can only take one of two values, -1 or +1, this is equivalent to choosing the numbers of subjects assigned to treatment 1 and treatment 2;
- call these  $n_1$  and  $n_2$  respectively, with  $n_1 + n_2 = N$

It is obvious that  $\sum x_j = 0$  if and only if  $n_1 = n_2$ . Therefore, assuming  $N$  is even, the **optimal design** has

- $n_1 = \frac{n}{2}$  subjects assigned to treatment 1 and
- $n_2 = \frac{n}{2}$  subjects assigned to treatment 2.

For  $N$  odd, we choose  $n_1 = \frac{n+1}{2}$ ,  $n_2 = \frac{n-1}{2}$ , or vice versa.

**Definition 1.3.** We can assess different designs using their **efficiency**:

$$\text{Eff} = \frac{\text{Var}(\hat{\beta}_1 | d^*)}{\text{Var}(\hat{\beta}_1 | d_1)} \quad (1.9)$$

where  $d_1$  is a design we want to assess and  $d^*$  is the optimal design with smallest variance. Note that  $0 \leq \text{Eff} \leq 1$ .

In Figure 1.2 below, we plot this efficiency for Example 1.1, using different choices of  $n_1$ . The total number of runs is fixed at  $n = 100$ , and the function **eff** calculates the efficiency from Definition 1.3 for a design with  $n_1$  subjects assigned to treatment 1. Clearly, efficiency of 1 is achieved when  $n_1 = n_2$  (equal allocation of treatments 1 and 2). If  $n_1 = 0$  or  $n_1 = 1$ , the efficiency is zero; we cannot estimate the difference between two treatments if we only allocate subjects to one of them.

```
n <- 100  
eff <- function(n1) 1 - ((2 * n1 - n) / n)^2  
curve(eff, from = 0, to = n, ylab = "Eff", xlab = expression(n[1]))
```

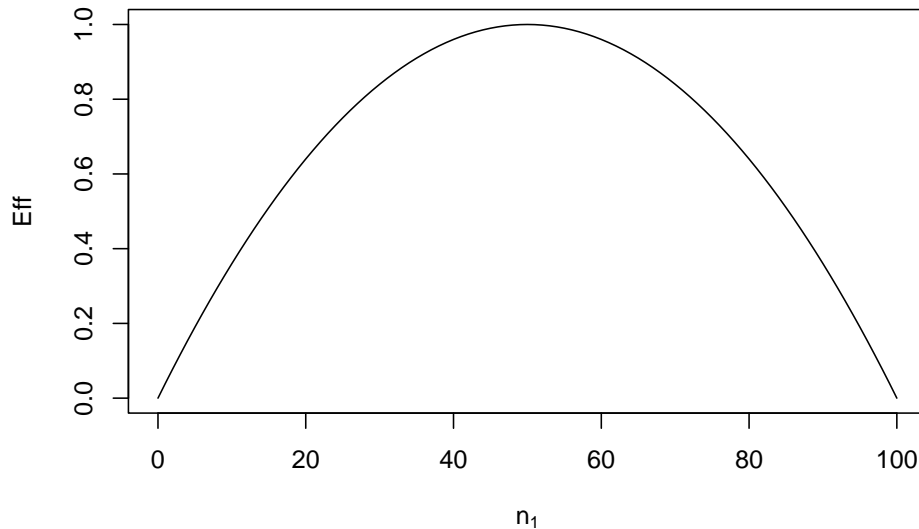


Figure 1.2: Efficiencies for designs for Example 1.1 with different numbers,  $n_1$ , of subjects assigned to treatment 1 when the total number of subjects is  $n = 100$ .

## 1.2 Aims of experimentation and some examples

Some reasons experiments are performed:

1. Treatment comparison (Chapters 2 and 3)
  - compare several treatments (and choose the best)
  - e.g. clinical trial, agricultural field trial
2. Factor screening (Chapters 4, 5 and 6)
  - many complex systems may involve a large number of (discrete) factors (controllable features)
  - which of these factors have a substantive impact?
  - (relatively) small experiments
  - e.g. industrial experiments on manufacturing processes
3. Response surface exploration (Chapter 7)
  - detailed description of relationship between important (continuous) variables and response

- typically second order polynomial regression models
  - larger experiments, often built up sequentially
  - e.g. alcohol yields in a pharmaceutical experiments
4. Optimisation (Chapter 7)
- finding settings of variables that lead to maximum or minimum response
  - typically use response surface methods and sequential “hill climbing” strategy

### 1.3 Some definitions

**Definition 1.4.** The **response**  $Y$  is the outcome measured in an experiment; e.g. yield from a chemical process. The response from the  $n$  observations are denoted  $Y_1, \dots, Y_n$ .

**Definition 1.5. Factors** (discrete) or **variables** (continuous) are features which can be set or controlled in an experiment;  $m$  denotes the number of factors or variables under investigation. For discrete factors, we call the possible settings of the factor its **levels**. We denote by  $x_{ij}$  the value taken by factor or variable  $i$  in the  $j$ th run of the experiment ( $i = 1, \dots, m; j = 1, \dots, n$ ).

**Definition 1.6.** The **treatments** or **support points** are the *distinct* combinations of factor or variable values in the experiment.

**Definition 1.7.** An experimental **unit** is the basic element (material, animal, person, time unit, ...) to which a treatment can be applied to produce a response.

In Example 1.1 (comparing two treatments):

- Response  $Y$ : Measured outcome, e.g. protein level or pain score in clinical trial, yield in an agricultural field trial.
- Factor  $x$ : “treatment” applied
- Levels
 

treatment 1	$x = -1$
treatment 2	$x = +1$
- Design point: factor level applied to  $j$ th subject;  $x_j = \pm 1$
- Treatment or support point: Two treatments or support points
- Experimental unit: Subject (person, animal, plot of land, ...).

### 1.4 Principles of experimentation

Three fundamental principles that need to be considered when designing an experiment are:

- replication
- randomisation
- stratification (blocking)

### 1.4.1 Replication

Each treatment is applied to a number of experimental units, with the  $j$ th treatment replicated  $r_j$  times. This enables the estimation of the variances of treatment effect estimators; increasing the number of replications, or replicates, decreases the variance of estimators of treatment effects. (Note: proper replication involves independent application of the treatment to different experimental units, not just taking several measurements from the same unit).

### 1.4.2 Randomisation

Randomisation should be applied to the allocation of treatments to units. Randomisation protects against **bias**; the effect of variables that are unknown and potentially uncontrolled or subjectivity in applying treatments. It also provides a formal basis for inference and statistical testing.

For example, in a clinical trial to compare a new drug and a control random allocation protects against

- “unmeasured and uncontrollable” features (e.g. age, sex, health)
- bias resulting from the clinician giving new drug to patients who are sicker.

Clinical trials are usually also *double-blinded*, i.e. neither the healthcare professional nor the patient knows which treatment the patient is receiving.

### 1.4.3 Stratification (or blocking)

We would like to use a wide variety of experimental units (e.g. people or plots of land) to ensure **coverage** of our results, i.e. validity of our conclusions across the population of interest. However, if the sample of units from the population is too heterogeneous, then this will induce too much random variability, i.e. increase  $\sigma^2$  in  $\varepsilon_j \sim N(0, \sigma^2)$ , and hence increase the variance of our parameter estimators.

We can reduce this extraneous variation by splitting our units into homogenous sets, or **blocks**, and including a blocking term in the model. The simplest blocked experiment is a **randomised complete block design**, where each block contains enough units for all treatments to be applied. Comparisons can then be made *within* each block.

Basic principle: block what you can, randomise what you cannot.

Later we will look at blocking in more detail, and the principle of **incomplete blocks**.

## 1.5 Revision on the linear model

Recall:  $Y = X\beta + \varepsilon$ , with  $\varepsilon \sim N(0, I_n\sigma^2)$ . Let the  $j$ th row of  $X$  be denoted  $x_j^T$ , which holds the values of the predictors, or explanatory variables, for the  $j$ th observation. Then

$$Y_j = x_j^T \beta + \varepsilon_j, \quad j = 1, \dots, n.$$

For example, quite commonly, for continuous variables

$$x_j = (1, x_{1j}, x_{2j}, \dots, x_{mj})^T,$$

and so

$$x_j^T \beta = \beta_0 + \beta_1 x_{1j} + \dots + \beta_m x_{mj}.$$

The least squares estimators are given by

$$\hat{\beta} = (X^T X)^{-1} X^T Y,$$

with

$$\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2.$$

### 1.5.1 Variance of a Prediction/Fitted Value

A prediction of the mean response at point  $x_0$  (which may or may not be in the design) is

$$\hat{Y}_0 = x_0^T \hat{\beta},$$

with

$$\begin{aligned} \text{Var}(\hat{Y}_0) &= \text{Var}(x_0^T \hat{\beta}) \\ &= x_0^T \text{Var}(\hat{\beta}) x_0 \\ &= x_0^T (X^T X)^{-1} x_0 \sigma^2. \end{aligned}$$

For a linear model, this variance depends only on the assumed regression model and the design (through  $X$ ), the point at which prediction is to be made ( $x_0$ ) and the value of  $\sigma^2$ ; it does not depend on data  $Y$  or parameters  $\beta$ .

Similarly, we can find the variance-covariance matrix of the fitted values:

$$\text{Var}(\hat{Y}) = \text{Var}(X\hat{\beta}) = X(X^T X)^{-1} X^T \sigma^2.$$

### 1.5.2 Analysis of Variance and $R^2$ as Model Comparison

To assess the goodness-of-fit of a model, we can use the residual sum of squares

$$\begin{aligned} \text{RSS} &= (Y - X\hat{\beta})^T(Y - X\hat{\beta}) \\ &= \sum_{j=1}^n \{Y_j - x_j^T \hat{\beta}\}^2 \\ &= \sum_{j=1}^n r_j^2, \end{aligned}$$

where

$$r_j = Y_j - x_j^T \hat{\beta}.$$

Often, a comparison is made to the null model

$$Y_j = \beta_0 + \varepsilon_j,$$

i.e.  $Y_i \sim N(\beta_0, \sigma^2)$ . The residual sum of squares for the null model is given by

$$\text{RSS}(\text{null}) = Y^T Y - m\bar{Y}^2,$$

as

$$\hat{\beta}_0 = \bar{Y} = \frac{1}{n} \sum_{j=1}^n Y_j.$$

How do we compare these models?

1. Ratio of residual sum of squares:

$$\begin{aligned} R^2 &= 1 - \frac{\text{RSS}}{\text{RSS}(\text{null})} \\ &= 1 - \frac{(Y - X\hat{\beta})^T(Y - X\hat{\beta})}{Y^T Y - n\bar{Y}^2}. \end{aligned}$$

The quantity  $0 \leq R^2 \leq 1$  is sometimes called the **coefficient of multiple determination**:

- high  $R^2$  implies that the model describes much of the variation in the data;



- **but** note that  $R^2$  will always increase as  $p$  (the number of explanatory variables) increases, with  $R^2 = 1$  when  $p = n$ ;
- some software packages will report the adjusted  $R^2$ .

$$\begin{aligned} R_a^2 &= 1 - \frac{\text{RSS}/(n-p)}{\text{RSS}(\text{null})/(n-1)} \\ &= 1 - \frac{(Y - X\hat{\beta})^T(Y - X\hat{\beta})/(n-p)}{(Y^T Y - n\bar{Y}^2)/(n-1)}; \end{aligned}$$

- $R_a^2$  does not necessarily increase with  $p$  (as we divide by degrees of freedom to adjust for complexity of the model).
2. Analysis of variance (ANOVA): An ANOVA table is compact way of presenting the results of (sequential) comparisons of nested models. You should be familiar with an ANOVA table of the following form.

Table 1.1: A standard ANOVA table.

Source	Degress of Freedom	(Sequential) Sum of Squares	Mean Square
Regression	$p - 1$	By subtraction; see (1.12)	Reg SS/ $(p - 1)$
Residual	$N - p$	$(Y - X\hat{\beta})^T(Y - X\hat{\beta})^6$	RSS/ $(N - p)$
Total	$N - 1$	$Y^T Y - N\bar{Y}^2$ <sup>7</sup>	

In row 1 of Table 1.1 above,

$$\text{Regression SS} = \text{Total SS} - \text{RSS} = Y^T Y - n\bar{Y}^2 - (Y - X\hat{\beta})^T(Y - X\hat{\beta}) \quad (1.10)$$

$$= -n\bar{Y}^2 - \hat{\beta}^T (X^T X) \hat{\beta} + 2\hat{\beta}^T X^T Y \quad (1.11)$$

$$= \hat{\beta}^T (X^T X) \hat{\beta} - n\bar{Y}^2, \quad (1.12)$$

with the last line following from

$$\begin{aligned} \hat{\beta}^T X^T Y &= \hat{\beta}^T (X^T X)(X^T X)^{-1} X^T Y \\ &= \hat{\beta}^T (X^T X) \hat{\beta} \end{aligned}$$

<sup>6</sup>Residual sum of squares for the full regression model.

<sup>7</sup>Residual sum of squares for the null model.

This idea can be generalised to the comparison of a *sequence* of nested models - see Problem Sheet 1.

Hypothesis testing is performed using the mean square:

$$\frac{\text{Regression SS}}{p-1} = \frac{\hat{\beta}^T (X^T X) \hat{\beta} - n\bar{Y}^2}{p-1}.$$

Under  $H_0 : \beta_1 = \dots = \beta_{p-1} = 0$

$$\begin{aligned} \frac{\text{Regression SS}/(p-1)}{\text{RSS}/(N-p)} &= \frac{(\hat{\beta}^T (X^T X) \hat{\beta} - n\bar{Y}^2)/(p-1)}{(Y - X\hat{\beta})^T (Y - X\hat{\beta})/(n-p)} \\ &\sim F_{p-1, n-p}, \end{aligned}$$

an  $F$  distribution with  $p-1$  and  $n-p$  degrees of freedom; defined via the ratio of two independent  $\chi^2$  distributions.

Also,

$$\frac{\text{RSS}}{n-p} = \frac{(Y - X\hat{\beta})^T (Y - X\hat{\beta})}{n-p} = \hat{\sigma}^2$$

is an unbiased estimator for  $\sigma^2$ , and

$$\frac{(n-p)}{\sigma^2} \hat{\sigma}^2 \sim \chi_{n-p}^2.$$

This is a Chi-squared distribution with  $N-p$  degrees of freedom (see MATH2010 notes).

## 1.6 Exercises

1. (Adapted from Morris, 2011) A classic and famous example of a simple hypothetical experiment was described by Fisher (1935):

A lady declares that by tasting a cup of tea made with milk she can discriminate whether the milk or the tea infusion was added first to the cup. We will consider the problem of designing an experiment by means of which this assertion can be tested. For this purpose let us first lay down a simple form of experiment with a view to studying its limitations and its characteristics, both those that are essential to the experimental

method, when well developed, and those that are not essential but auxiliary.

Our experiment consists in mixing eight cups of tea, four in one way and four in the other, and presenting them to the subject for judgement in a random order. The subject has been told in advance of what the test will consist, namely that she will be asked to taste eight cups, that these shall be four of each kind, and that they shall be presented to her in a random order, that is an order not determined arbitrarily by human choice, but by the actual manipulation of the physical apparatus used in games of chance, cards, dice, roulettes, etc., or, more expeditiously, from a published collection of random sampling numbers purporting to give the actual results of such manipulation<sup>8</sup>. Her task is to divide the 8 cups into two sets of 4, agreeing, if possible, with the treatments received.

- a. Define the treatments in this experiment.
- b. Identify the units in this experiment.
- c. How might a “physical apparatus” from a “game of chance” be used to perform the randomisation. Explain one example.
- d. Suppose eight tea cups are available for this experiment but they are not identical. Instead they come from two sets. Four are made from heavy, thick porcelain; four from much lighter china. If each cup can only be used once, how might this fact be incorporated into the design of the experiment?

2. Consider the linear model

$$y = X\beta + \varepsilon,$$

with  $y$  an  $n \times 1$  vector of responses,  $X$  a  $n \times p$  model matrix and  $\varepsilon$  a  $n \times 1$  vector of independent and identically distributed random variables with constant variance  $\sigma^2$ .

- a. Derive the least squares estimator  $\hat{\beta}$  for this multiple linear regression model, and show that this estimator is unbiased. Using the definition of (co)variance, show that

$$\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2.$$

- b. If  $\varepsilon \sim N(0, I_n \sigma^2)$ , with  $I_n$  being the  $n \times n$  identity matrix, show that the maximum likelihood estimators for  $\beta$  coincide with the least squares estimators.

3. Consider the two nested linear models

---

<sup>8</sup>Now, we would use routines such as `sample` in R.

- (i)  $Y_j = \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2j} + \dots + \beta_{p_1} x_{p_1 j} + \varepsilon_j$ , or  $y = X_1 \beta_1 + \varepsilon$ ,  
(ii)  $Y_j = \beta_0 + \beta_1 x_{1j} + \beta_2 x_{2j} + \dots + \beta_{p_1} x_{p_1 j} + \beta_{p_1+1} x_{(p_1+1)j} + \dots + \beta_{p_2} x_{p_2 j} + \varepsilon_j$ ,  
or  $y = X_1 \beta_1 + X_2 \beta_2 + \varepsilon$

with  $\varepsilon_j \sim N(0, \sigma^2)$ , and  $\varepsilon_j, \varepsilon_k$  independent ( $\varepsilon \sim N(0, I_n \sigma^2)$ ).

- Construct an ANOVA table to compare model (ii) with the null model  $Y_j = \beta_0 + \varepsilon_j$ .
- Extend this ANOVA table to compare models (i) and (ii) by further decomposing the regression sum of squares for model (ii).

**Hint:** which residual sum of squares are you interested in to compare models (i) and (ii)?

You should end up with an ANOVA table of the form

Source	Degrees of freedom	Sums of squares	Mean square
Model (i)	$p_1$	?	?
Model (ii)	$p_2$	?	?
Residual	$n - p_1 - p_2 - 1$	?	?
Total	$n - 1$	$y^T y - n\bar{Y}^2$	

The second row of the table gives the **extra sums of squares** for the additional terms in fitting model (ii), over and above those in model (i).

- Calculate the extra sum of squares for fitting the terms in model (i), over and above those terms only in model (ii), i.e. those held in  $X_2 \beta_2$ . Construct an ANOVA table containing both the extra sum of squares for the terms only in model (i) and the extra sum of squares for the terms only in model (ii). Comment on the table.

# Completely randomised designs

**Example 2.1.** Pulp experiment (Wu and Hamada, 2009, ch. 2)

```
pulp <- data.frame(operator = rep(factor(1:4), 5),  
                  repetition = rep(1:5, rep(4, 5)),  
                  reflectance = c(59.8, 59.8, 60.7, 61.0, 60.0, 60.2, 60.7, 60.8,  
                                60.8, 60.4, 60.5, 60.6, 60.8, 59.9, 60.9, 60.5, 59.8, 60.0, 60.5, 60.0, 60.5, 60.5, 60.5, 60.5))  
  
knitr::kable(  
  tidyr::pivot_wider(pulp, names_from = operator, values_from = reflectance)[, -1],  
  col.names = paste("Operator", 1:4),  
  caption = "Pulp experiment: reflectance values (unitless) from four different operators."  
)
```

21

Table 2.1: Pulp experiment: reflectance values (unitless) from four different operators.

Operator 1	Operator 2	Operator 3	Operator 4
59.8	59.8	60.7	61.0
60.0	60.2	60.7	60.8
60.8	60.4	60.5	60.6
60.8	59.9	60.9	60.5
59.8	60.0	60.3	60.5

We can informally compare the responses from these four treatments graphically.

```
boxplot(reflectance ~ operator, data = pulp)
```

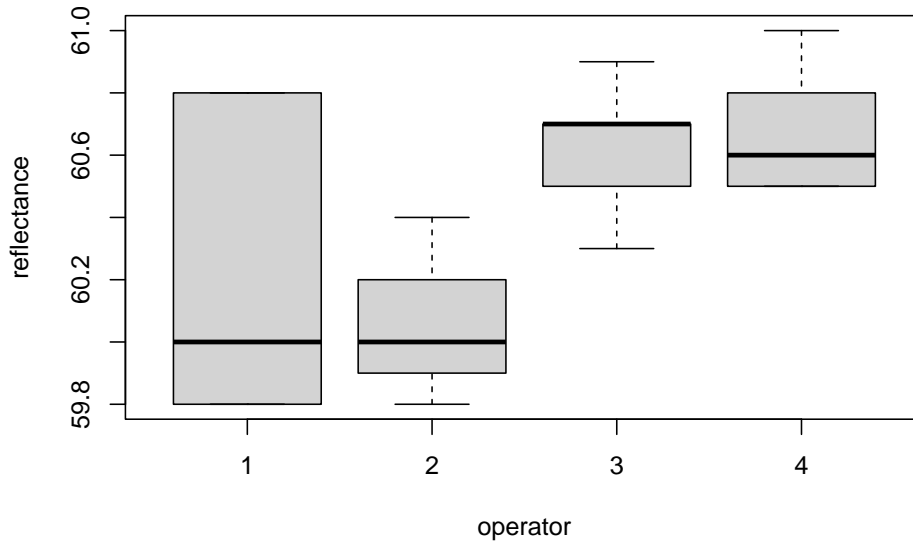


Figure 2.1: Pulp experiments: distributions of reflectance from the four operators.

Figure 2.1 shows that, relative to the variation, there may be a difference in the mean response between treatments 1 and 2, and 3 and 4. In this chapter, we will see how to make this comparison formally using linear models, and to assess how the choice of design impacts our results.

Throughout this chapter we will assume the  $i$ th treatment is applied to  $n_i$  experimental unit, with total number of runs  $n = \sum_{i=1}^t n_i$  in the experiment.

## 2.1 A unit-treatment linear model

An appropriate, and common, model to describe data from such experiments when the response is continuous is given by

$$y_{ij} = \mu + \tau_i + \varepsilon_{ij}, \quad i = 1, \dots, t; j = 1, \dots, n_i, \quad (2.1)$$

where  $y_{ij}$  is the response from the  $j$ th application of treatment  $i$ ,  $\mu$  is a constant parameter,  $\tau_i$  is the effect of the  $i$ th treatment, and  $\varepsilon_{ij}$  is the random individual effect from each experimental unit with  $E(\varepsilon_{ij})$  and  $\text{Var}(\varepsilon_{ij}) = \sigma^2$ . All random errors are assumed independent and here we also assume  $\varepsilon_{ij} \sim N(0, \sigma^2)$ .

Model (2.1) assumes that each treatment can be randomly allocated to one of the  $n$  experimental units, and that this allocation is independent of the allocation of all the other treatments.

Why is this model appropriate and commonly used? The expected response from the application of the  $i$ th treatment is

$$E(y_{ij}) = \mu + \tau_i.$$

The parameter  $\mu$  can be thought of as representing the impact of many different features particular to **this** experiment, and  $\tau_i$  is the deviation due to applying treatment  $i$ . From the application of two different hypothetical experiments, A and B, the expected response from treatment  $i$  may be different due to a different overall mean. From experiment A:

$$E(y_{ij}) = \mu_A + \tau_i.$$

From experiment B:

$$E(y_{ij}) = \mu_B + \tau_i.$$

But the **difference** between treatments  $k$  and  $l$  ( $k, l = 1, \dots, t$ )

$$\begin{aligned} E(y_{kj}) - E(y_{lj}) &= \mu_A + \tau_k - \mu_A - \tau_l \\ &= \tau_k - \tau_l, \end{aligned}$$

is constant across different experiments. This concept of **comparison** underpins most design of experiments, and will be applied throughout this module.

## 2.2 The partitioned linear model

In matrix form, we can write model (2.1) as

$$y = X_1\mu + X_2\tau + \varepsilon,$$

where  $X_1 = 1_n$ , the  $n$ -vector with every entry equal to one,

$$X_2 = \begin{bmatrix} 1_{n_1} & 0_{n_1} & \cdots & 0_{n_1} \\ 0_{n_2} & 1_{n_2} & \cdots & 0_{n_2} \\ \vdots & & \ddots & \vdots \\ 0_{n_t} & 0_{n_t} & \cdots & 1_{n_t} \end{bmatrix},$$

with  $0_{n_i}$  is the  $n_i$ -vector with every entry equal to zero,  $\tau = [\tau_1, \dots, \tau_t]^T$  and  $\varepsilon = [\varepsilon_{11}, \dots, \varepsilon_{tn_t}]^T$ .

Why are we partitioning the model? Going back to our discussion of the role of  $\mu$  and  $\tau_i$ , it is clear that we are not interested in estimating  $\mu$ , which represents an experiment-specific contribution to the expected mean. Our only interest is in estimating the (differences between the)  $\tau_i$ . Hence, we can treat  $\mu$  as a nuisance parameter.

If we define  $X = [X_1 | X_2]$  and  $\beta^T = [\mu | \tau^T]$ , we can write the usual least squares equations

$$X^T X \hat{\beta} = X^T y \tag{2.2}$$

as a system of two matrix equations

$$\begin{aligned} X_1^T X_1 \hat{\mu} + X_1^T X_2 \hat{\tau} &= X_1^T y \\ X_2^T X_1 \hat{\mu} + X_2^T X_2 \hat{\tau} &= X_2^T y. \end{aligned}$$

Assuming  $(X_1^T X_1)^{-1}$  exists, which it does in this case, we can pre-multiply the first of these equations by  $X_2^T X_1 (X_1^T X_1)^{-1}$  and subtract it from the second equation to obtain

$$\begin{aligned} X_2^T [I_n - X_1 (X_1^T X_1)^{-1} X_1^T] X_1 \hat{\mu} + X_2^T [I_n - X_1 (X_1^T X_1)^{-1} X_1^T] X_2 \hat{\tau} \\ = X_2^T [I_n - X_1 (X_1^T X_1)^{-1} X_1^T] y. \end{aligned}$$

Writing  $H_1 = X_1 (X_1^T X_1)^{-1} X_1^T$ , we obtain



$$X_2^T[I_n - H_1]X_1\hat{\mu} + X_2^T[I_n - H_1]X_2\hat{\tau} = X_2^T[I_n - H_1]y. \quad (2.3)$$

The matrix  $H_1$  is a “hat” matrix for a linear model containing only the term  $\mu$ , and hence  $H_1X_1 = X_1$  (see MATH2010). Hence the first term in (2.3) is zero, and we obtain the **reduced normal equations** for  $\tau$ :

$$X_2^T[I_n - H_1]X_2\hat{\tau} = X_2^T[I_n - H_1]y. \quad (2.4)$$

Note that the solutions from (2.4) are not different from the solution to  $\hat{\tau}$  that would be obtained from solving (2.2); equation (2.4) is simply a re-expression, where we have eliminated the nuisance parameter  $\mu$ . This fact means that we rarely need to solve (2.4) explicitly.

Recalling that for a hat matrix,  $I_n - H_1$  is idempotent and symmetric (see MATH2010), if we define

$$X_{2|1} = (I_n - H_1)X_2,$$

then we can rewrite equation (2.4) as

$$X_{2|1}^T X_{2|1} \hat{\tau} = X_{2|1}^T y, \quad (2.5)$$

which are the normal equations for a linear model with expectation  $E(y) = X_{2|1}\tau$ .

For the CRD discussed in this chapter,  $X_1^T X_1 = n$ , the total number of runs in the experiment<sup>1</sup>. Hence  $(X_1^T X_1)^{-1} = 1/n$  and  $H_1 = \frac{1}{n}J_n$ , with  $J_n$  the  $n \times n$  matrix with all entries equal to 1.

The adjusted model matrix then has the form

$$\begin{aligned} X_{2|1} &= (I_n - H_1)X_2 \\ &= X_2 - \frac{1}{n}J_n X_2 \\ &= X_2 - \frac{1}{n}[n_1 1_n | \cdots | n_t 1_n]. \end{aligned}$$

That is, every column of  $X_2$  has been adjusted by the subtraction of the column mean from each entry<sup>2</sup>. Also notice that each row of  $X_{2|1}$  has a row-sum equal

<sup>1</sup>In later chapters we will see examples where  $X_1$  has  $> 1$  columns, and hence  $X_1^T X_1$  is a matrix.

<sup>2</sup>Often called “column centred”

to zero ( $= 1 - \sum_{i=1}^t n_i/n$ ). Hence,  $X_{2|1}$  is not of full column rank, and so the reduced normal equations do not have a unique solution<sup>3</sup>.

In MATH2010 we fitted models with categorical variables by defining a set of dummy variables and estimating a reduced model. Here, we will take a slightly different approach and study which combinations of parameters from (2.1) are estimable, and in particular which linear combinations of the treatment parameters  $\tau_i$  we can estimate. We start with a short digression to discuss the concept of a generalised inverse.

### 2.3 Generalised inverses

A generalised inverse of a square, symmetric matrix  $A$  is any matrix  $A^-$  that satisfies

$$AA^-A = A.$$

If matrix  $A$  is of full column rank, then there is a single generalised inverse  $A^- = A^{-1}$ , the usual inverse matrix. In other cases, there may be infinitely many choices for  $A^-$ .

### 2.4 Solutions to the reduced normal equations

We can use a generalised inverse to generate solutions to (2.5) of the form:

$$\hat{\tau} = (X_{2|1}^T X_{2|1})^- X_{2|1}^T y.$$

Substituting this expression into (2.5), we obtain

$$\begin{aligned} X_{2|1}^T X_{2|1} (X_{2|1}^T X_{2|1})^- X_{2|1}^T y &= X_{2|1}^T H_{2|1} y \\ &= X_{2|1}^T y, \end{aligned}$$

where  $H_{2|1} = X_{2|1} (X_{2|1}^T X_{2|1})^- X_{2|1}^T$  and  $X_{2|1}^T H_{2|1} = X_{2|1}^T$ .

---

<sup>3</sup>If we recalled the material on “dummy” variables from MATH2010 chapter 5, we would already have realised this.

## Chapter 3

# Blocking



## Chapter 4

# Factorial experiments



## Chapter 5

# Blocking in factorial designs





## Chapter 6

# Fractional factorial designs



## Chapter 7

# Response surface methodology



## Chapter 8

# Optimal design of experiments



# Bibliography

- Fisher, R. A. (1935). *The Design of Experiments*. Oliver and Boyd, Edinburgh.
- Luca, M. and Bazerman, M. H. (2020). *The Power of Experiments: Decision Making in a Data-Driven World*. MIT Press, Cambridge.
- Morris, M. D. (2011). *Design of Experiments: An Introduction based on Linear Models*. Chapman and Hall/CRC Press, Boca Raton.
- Wu, C. F. J. and Hamada, M. (2009). *Experiments: Planning, Analysis, and Parameter Design Optimization*. Wiley, New York, 2nd edition.