

# MATH3014-6027 Design (and Analysis) of Experiments

Dave Woods

2022-02-16



# 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>25</b>
2.1 A unit-treatment linear model . . . . .	27
2.2 The partitioned linear model . . . . .	28
2.3 Reduced normal equations for the CRD . . . . .	29
2.4 Contrasts . . . . .	31
2.5 Treatment contrast estimators in the CRD . . . . .	31
2.6 Analysing CRDs in R . . . . .	33
2.7 Multiple comparisons . . . . .	35
2.8 Exercises . . . . .	38
<b>3 Blocking</b>	<b>43</b>
<b>4 Factorial experiments</b>	<b>45</b>
<b>5 Blocking in factorial designs</b>	<b>47</b>
<b>6 Fractional factorial designs</b>	<b>49</b>
<b>7 Response surface methodology</b>	<b>51</b>
<b>8 Optimal design of experiments</b>	<b>53</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.

Southampton prerequisites for this module are MATH2010 or MATH6174 and STAT6123 (or equivalent modules on linear modelling).



# 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. Collecting empirical evidence is key for scientific development, as described in terms of clinical trials by xked:

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 or MATH6174/STAT6123):

---

<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.



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

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 or MATH6174/STAT6123). 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

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

---

<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.

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

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

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

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

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

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

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]))
```

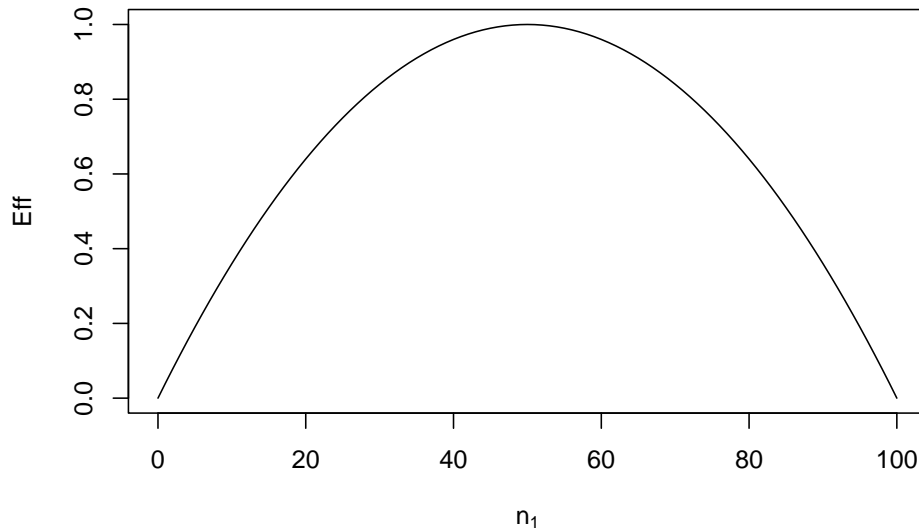


Figure 1.2: Efficiencies for designs for Example 1.1 with different numbers,  $n_1$ , of subjects assigned to treatment 1 when the total number of subjects is  $n = 100$ .

## 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	$x = +1$

- 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 or MATH6174 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?

Solution

- a. There are two treatments in the experiment - the two ingredients “milk first” and “tea first”.
- b. The experimental units are the “cups of tea”, made up from the tea and milk used and also the cup itself.
- c. The simplest method here might be to select four black playing cards and four red playing cards, assign one treatment to each colour, shuffle the cards, and then draw them in order. The colour drawn indicates the treatment that should be used to make the next cup of tea. This operation would give one possible randomisation.

We could of course also use R.

```
sample(rep(c("Milk first", "Tea first"), c(4, 4)), size = 8, replace = F)
```

```
## [1] "Tea first" "Tea first" "Tea first" "Milk first" "Tea first"
## [6] "Milk first" "Milk first" "Milk first"
```

---

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

- d. Type of cup could be considered as a blocking factor. One way of incorporating it would be to split the experiment into two (blocks), each with four cups (two milk first, two tea first). We would still wish to randomise allocation of treatments to units within blocks.

```
# block 1
sample(rep(c("Milk first", "Tea first"), c(2, 2)), size = 4, replace = F)

## [1] "Tea first" "Milk first" "Milk first" "Tea first"

# block 2
sample(rep(c("Milk first", "Tea first"), c(2, 2)), size = 4, replace = F)

## [1] "Milk first" "Tea first" "Milk first" "Tea first"
```

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.

Solution

- a. The method of least squares minimises the sum of squared differences between the responses and the expected values, that is, minimises the expression

$$(y - X\beta)^T (y - X\beta) = y^T y - 2\beta^T X^T y + \beta^T X^T X \beta.$$

Differentiating with respect to the vector  $\beta$ , we obtain

$$\frac{\partial}{\partial \beta} = -2X^T y + 2X^T X \beta.$$

Set equal to 0 and solve:

$$X^T X \hat{\beta} = X^T y \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y.$$

The estimator  $\hat{\beta}$  is unbiased:

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

and has variance:

$$\begin{aligned} \text{Var}(\hat{\beta}) &= E \left\{ [\hat{\beta} - E(\hat{\beta})] [\hat{\beta} - E(\hat{\beta})]^T \right\} \\ &= E \left\{ [\hat{\beta} - \beta] [\hat{\beta} - \beta]^T \right\} \\ &= E \left\{ \hat{\beta} \hat{\beta}^T - 2\beta \hat{\beta}^T + \beta \beta^T \right\} \\ &= E \left\{ (X^T X)^{-1} X^T y y^T X (X^T X)^{-1} - 2\beta y^T X (X^T X)^{-1} + \beta \beta^T \right\} \\ &= (X^T X)^{-1} X^T E(y y^T) X (X^T X)^{-1} - 2\beta E(y^T) X (X^T X)^{-1} + \beta \beta^T \\ &= (X^T X)^{-1} X^T [\text{Var}(y) + E(y) E(y^T)] X (X^T X)^{-1} - 2\beta \beta^T X^T X (X^T X)^{-1} + \beta \beta^T \\ &= (X^T X)^{-1} X^T [I_N \sigma^2 + X \beta \beta^T X^T] X (X^T X)^{-1} - \beta \beta^T \\ &= (X^T X)^{-1} \sigma^2. \end{aligned}$$

b. As  $y \sim N(X\beta, I_N \sigma^2)$ , the likelihood is given by

$$L(\beta; y) = (2\pi\sigma^2)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) \right\}.$$

The log-likelihood is given by

$$l(\beta; y) = -\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) + \text{constant}.$$

Up to a constant, this expression is  $-1 \times$  the least squares equations; hence maximising the log-likelihood is equivalent to minimising the least squares equation.

3. Consider the two nested linear models

- (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 ( $\varepsilon \sim N(0, I_n \sigma^2)$ ).

- a. Construct an ANOVA table to compare model (ii) with the null model  $Y_j = \beta_0 + \varepsilon_j$ .

- b. 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).

- c. 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.

Solution

- a. From lectures

Source	Degrees of freedom	Sums of squares	Mean square
Regression	$p_1 + p_2$	$\hat{\beta}^T X^T X \hat{\beta} - n\bar{Y}^2$	$\left( \hat{\beta}^T X^T X \hat{\beta} - n\bar{Y}^2 \right) / (p_1 + p_2)$
Residual	$n - p_1 - p_2 - 1$	$(y - X\hat{\beta})^T (y - X\hat{\beta})$	$\text{RSS} / (n - p_1 - p_2 - 1)$
Total	$n - 1$	$y^T y - n\bar{Y}^2$	

where

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

b. To compare model (i) with the null model, we compute

$$\begin{aligned}\text{RSS}(\text{null}) - \text{RSS}(\text{i}) &= y^T y - N\bar{Y}^2 - (y - X_1 \hat{\beta}_1)^T (y - X_1 \hat{\beta}_1) \\ &= \hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1 - n\bar{Y}^2.\end{aligned}$$

To compare models (i) and (ii), we compare them both to the null model, and look at the difference between these comparisons:

$$\begin{aligned}[\text{RSS}(\text{null}) - \text{RSS}(\text{ii})] - [\text{RSS}(\text{null}) - \text{RSS}(\text{i})] &= \text{RSS}(\text{i}) - \text{RSS}(\text{ii}) \\ &= (y - X_1 \hat{\beta}_1)^T (y - X_1 \hat{\beta}_1) - (y - X \hat{\beta})^T (y - X \hat{\beta}) \\ &= \hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1.\end{aligned}$$

Source	Degrees of freedom	Sums of squares	Mean square
Regression	$p_1 + p_2$	$\hat{\beta}^T X^T X \hat{\beta} - n\bar{Y}^2$	$(\hat{\beta}^T X^T X \hat{\beta} - n\bar{Y}^2) / (p_1 + p_2)$
Model (i)	$p_1$	$\hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1 - n\bar{Y}^2$	$(\hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1 - n\bar{Y}^2) / p_1$
Extra due to Model (ii)	$p_2$	$\hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1$	$(\hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1) / p_2$
Residual	$n - p_1 - p_2 - 1$	$(y - X \hat{\beta})^T (y - X \hat{\beta})$	$\text{RSS} / (n - p_1 - p_2 - 1)$
Total	$n - 1$	$y^T y - n\bar{Y}^2$	

By definition, the Model (i) SS and the Extra SS for Model (ii) sum to the Regression SS.

a. The extra sum of squares for the terms in model (i) over and above those in model (ii) are obtained through comparison of the models

ia.  $y = X_2 \beta_2 + \varepsilon,$

ii.  $y = X_1 \beta_1 + X_2 \beta_2 + \varepsilon = X \beta + \varepsilon$

Extra sum of squares for model (iia):

$$\begin{aligned}
[\text{RSS}(\text{null}) - \text{RSS}(\text{iaa})] - [\text{RSS}(\text{null}) - \text{RSS}(\text{ia})] &= \text{RSS}(\text{ia}) - \text{RSS}(\text{iaa}) \\
&= (y - X_2 \hat{\beta}_2)^T (y - X_2 \hat{\beta}_2) - (y - X \hat{\beta})^T (y - X \hat{\beta}) \\
&= \hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_2^T X_2^T X_2 \hat{\beta}_2.
\end{aligned}$$

Hence, an ANOVA table for the extra sums of squares is given by

Source	Degrees of freedom	Sums of squares	Mean square
Regression	$p_1 + p_2$	$\hat{\beta} X^T X \hat{\beta} - n \bar{Y}^2$	$(\hat{\beta} X^T X \hat{\beta} - n \bar{Y}^2) / (p_1 + p_2)$
Extra Model (i)	$p_1$	$\hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_2^T X_2^T X_2 \hat{\beta}_2$	$(\hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_2^T X_2^T X_2 \hat{\beta}_2) / p_1$
Extra Model (ii)	$p_2$	$\hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1$	$(\hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}_1^T X_1^T X_1 \hat{\beta}_1) / p_2$
Residual	$n - p_1 - p_2 - 1$	$(y - X \hat{\beta})^T (y - X \hat{\beta})$	$\text{RSS} / (n - p_1 - p_2 - 1)$
Total	$n - 1$	$y^T y - n \bar{Y}^2$	

Note that for these *adjusted* sums of squares, in general the extra sum of squares for model (i) and (ii) do not sum to the regression sum of squares. This will only be the case if the columns of  $X_1$  and  $X_2$  are mutually orthogonal, i.e.  $X_1^T X_2 = 0$ .



# Chapter 2

# Completely randomised designs

The simplest form of experiment we will consider compares  $t$  different **unstructured** treatments. By unstructured, we mean the treatments form a discrete collection, not related through the settings of other experimental features (compare with factorial experiments in Chapter 4). We also make the assumption that there are no restrictions in the randomisation of treatments to experimental units (compare with Chapter 3 on blocking). A designs for such an experiment is therefore called a **completely randomised design** (CRD).

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

In a paper pulping mill, an experiment was run to examine differences between the reflectance (brightness; ratio of amount of light leaving a target to the amount of light striking the target) of sheets of pulp made by  $t = 4$  operators. The data are given in Table 2.1 below.

```
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))  
  
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."  
)
```

The experiment has one factor (operator) with four levels (sometimes called a one-way layout). The CRD employed has equal replication of each treatment (operator).

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)
```

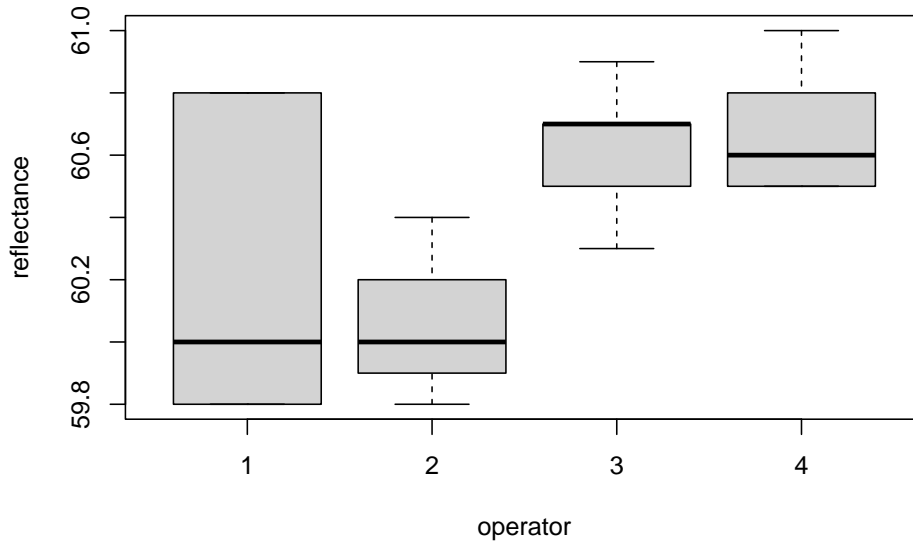


Figure 2.1: Pulp experiments: distributions of reflectance from the four operators.

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,

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

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 or STAT6123). 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 or MATH6174), 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  $E(y) = X_{2|1}\tau$ .

## 2.3 Reduced normal equations for the CRD

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

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

---

<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.

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 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 and STAT6123 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.

Let's study equation (2.5) in more detail. We have

$$\begin{aligned} X_{2|1}^T X_{2|1} &= X_2^T (I_n - H_1) X_2 \\ &= X_2^T X_2 - X_2^T H_1 X_2 \\ &= \text{diag}(n) - \frac{1}{n} X_2^T J_n X_2 \\ &= \text{diag}(n) - \frac{1}{n} n n^T, \end{aligned}$$

where  $n^T = (n_1, \dots, n_t)$ . Hence, the reduced normal equations become

$$\left[ \text{diag}(n) - \frac{1}{n} n n^T \right] \hat{\tau} = X_2^T y - \frac{1}{n} X_2^T J_n y \quad (2.6)$$

$$= X_2^T y - n \bar{y}_{..}, \quad (2.7)$$

where  $\bar{y}_{..} = \frac{1}{n} \sum_{i=1}^t \sum_{j=1}^{n_i} y_{ij}$ , i.e. the overall average of the observations from the experiment.

From (2.7) we obtain a system of  $t$  equations, each having the form

$$\hat{\tau}_i - \hat{\tau}_w = \bar{y}_{i.} - \bar{y}_{..}, \quad (2.8)$$

where  $\hat{\tau}_w = \frac{1}{n} \sum_{i=1}^t n_i \hat{\tau}_i$  and  $\bar{y}_{i.} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$  ( $i = 1, \dots, t$ ).

These  $t$  equations are not independent; when multiplied by the  $n_i$ , they sum to zero due to the linear dependency between the columns of  $X_{2|1}$ . Hence, there is no unique solution to  $\hat{\tau}$  from equation (2.7). However, we can estimate certain linear combinations of the  $\tau_i$ , called *contrasts*.

---

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

<sup>3</sup>If we recalled the material on “dummy” variables from MATH2010 or STAT6123, we would already have realised this.

## 2.4 Contrasts

**Definition 2.1.** A treatment **contrast** is a linear combination  $c^T\tau$  with coefficient vector  $c^T = (c_1, \dots, c_t)$  such that  $c^T\mathbf{1} = 0$ ; that is,  $\sum_{i=1}^t c_i = 0$ .

For example, assume we have  $t = 3$  treatments, then the following vectors  $c$  all define contrasts:

1.  $c^T = (1, -1, 0)$ ,
2.  $c^T = (1, 0, -1)$ ,
3.  $c^T = (0, 1, -1)$ .

In fact, they define all  $\binom{3}{2} = 3$  pairwise comparisons between treatments. The following are also contrasts:

4.  $c^T = (2, -1, -1)$ ,
5.  $c^T = (0.5, -1, 0.5)$ ,

each comparing the sum, or average, of expected responses from two treatments to the expected response from the remaining treatment.

The following are not contrasts, as  $c^T\mathbf{1} \neq 0$ :

6.  $c^T = (2, -1, 0)$ ,
7.  $c^T = (1, 0, 0)$ ,

with the final example once again demonstrating that we cannot estimate the individual  $\tau_i$ .

## 2.5 Treatment contrast estimators in the CRD

We estimate a treatment contrast  $c^T\tau$  in the CRD via linear combinations of equations (2.8):

$$\begin{aligned} \sum_{i=1}^t c_i \hat{\tau}_i - \sum_{i=1}^t c_i \hat{\tau}_w &= \sum_{i=1}^t c_i \bar{y}_{i.} - \sum_{i=1}^t c_i \bar{y}_{..} \\ \Rightarrow \sum_{i=1}^t c_i \hat{\tau}_i &= \sum_{i=1}^t c_i \bar{y}_{i.} , \end{aligned}$$

as  $\sum_{i=1}^t c_i \hat{\tau}_w = \sum_{i=1}^t c_i \bar{y}_{..} = 0$ , as  $\sum_{i=1}^t c_i = 0$ . Hence, the unique estimator of the contrast  $c^T\tau$  has the form

$$\widehat{c^T\tau} = \sum_{i=1}^t c_i \bar{y}_{i.} .$$

That is, we estimate the contrast in the treatment effects by the corresponding contrast in the treatment means.

The variance of this estimator is straightforward to obtain:

$$\begin{aligned}\text{var}\left(\widehat{c^T\tau}\right) &= \sum_{i=1}^t c_i^2 \text{var}(\bar{y}_{i.}) \\ &= \sigma^2 \sum_{i=1}^t c_i^2 / n_i,\end{aligned}$$

as, under our model assumptions, each  $\bar{y}_{i.}$  is an average of independent observations with variance  $\sigma^2$ . Similarly, from model (2.1) we can derive the distribution of  $\widehat{c^T\tau}$  as

$$\widehat{c^T\tau} \sim N\left(c^T\tau, \sigma^2 \sum_{i=1}^t c_i^2 / n_i\right).$$

Confidence intervals and hypothesis tests for  $c^T\tau$  can be constructed/conducted using this distribution, e.g.

- a  $100(1 - \frac{\alpha}{2})\%$  confidence interval:

$$c^T\tau \in \sum_{i=1}^t c_i \bar{y}_{i.} \pm t_{n-t, 1-\frac{\alpha}{2}} s \sqrt{\sum_{i=1}^t c_i^2 / n_i},$$

where  $t_{n-t, 1-\frac{\alpha}{2}}$  is the  $1 - \frac{\alpha}{2}$  quantile of a  $t$ -distribution with  $n - t$  degrees of freedom and

$$s^2 = \frac{1}{n-t} \sum_{i=1}^t \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.})^2 \quad (2.9)$$

is the estimate of  $\sigma^2$ .

- the hypothesis  $H_0 : c^T\tau = 0$  against the two-sided alternative  $H_0 : c^T\tau \neq 0$  is rejected using a test of with confidence level  $1 - \alpha/2$  if

$$\frac{|\sum_{i=1}^t c_i \bar{y}_{i.}|}{s \sqrt{\sum_{i=1}^t c_i^2 / n_i}} > t_{n-t, 1-\frac{\alpha}{2}}.$$



Table 2.2: 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

## 2.6 Analysing CRDs in R

Let's return to Example 2.1.

```
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."
)
```

Clearly, we could directly calculate, and then compare, mean responses for each operator. However, there are (at least) two other ways we can proceed which use the fact we are fitting a linear model. These will be useful when we consider more complex models.

1. Using `pairwise.t.test`.

```
with(pulp,
  pairwise.t.test(reflectance, operator, p.adjust.method = 'none'))

##
## Pairwise comparisons using t tests with pooled SD
##
## data: reflectance and operator
##
##    1      2      3
## 2 0.396 -      -
## 3 0.084 0.015 -
## 4 0.049 0.008 0.775
##
## P value adjustment method: none
```

This function performs hypothesis tests for all pairwise treatment comparisons (with a default confidence level of 0.95). Here we can see that operators 1 and 4, 2 and 3, and 2 and 4 have statistically significant differences.

2. Using `lm` and the `emmeans` package.

```
pulp.lm <- lm(reflectance ~ operator, data = pulp)
anova(pulp.lm)

## Analysis of Variance Table
##
## Response: reflectance
##           Df Sum Sq Mean Sq F value Pr(>F)
## operator   3    1.34   0.447    4.2  0.023 *
## Residuals 16    1.70   0.106
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

pulp.emm <- emmeans::emmeans(pulp.lm, ~ operator)
pairs(pulp.emm, adjust = 'none')

## contrast estimate      SE df t.ratio p.value
## 1 - 2          0.18 0.206 16   0.873  0.3960
## 1 - 3         -0.38 0.206 16  -1.843  0.0840
## 1 - 4         -0.44 0.206 16  -2.134  0.0490
## 2 - 3         -0.56 0.206 16  -2.716  0.0150
## 2 - 4         -0.62 0.206 16  -3.007  0.0080
## 3 - 4         -0.06 0.206 16  -0.291  0.7750
```

Here, we have first fitted the linear model object. The `lm` function, by default, will have set up dummy variables with the first treatment (operator) as a baseline (see MATH2010 or STAT6123). We then take the intermediate step of calculating the ANOVA table for this experiment, and use an F-test to compare the model accounting for operator differences to the null model; there are differences between operators at the 5% significance level,

The choice of dummy variables in the linear model is unimportant; any set could be used<sup>4</sup>, as in the next line we use the `emmeans` function (from the package of the same name) to specify that we are interested in estimating contrasts in the factor `operator` (which specifies our treatments in this experiment). Finally, the `pairs` command performs hypothesis tests for all pairwise comparisons between the four treatments. The results are the same as those obtained from using `pairwise.t.test`.

Our preferred approach is using method 2 (`lm` and `emmeans`), for four main reasons:

- a. The function `contrasts` in the `emmeans` package can be used to estimate arbitrary treatment contrasts (see `help("contrast-methods")`).

---

<sup>4</sup>Although  $\mu$  and  $\tau$  are not uniquely estimable, fitted values  $\hat{y}_i = \hat{\mu} + \hat{\tau}_i$  are, and hence it does not matter which dummy variables we use in `lm`.

```
# same as `pairs` above
emmmeans::contrast(pulp.emm, "pairwise", adjust = "none")

## contrast estimate SE df t.ratio p.value
## 1 - 2          0.18 0.206 16  0.873  0.3960
## 1 - 3         -0.38 0.206 16 -1.843  0.0840
## 1 - 4         -0.44 0.206 16 -2.134  0.0490
## 2 - 3         -0.56 0.206 16 -2.716  0.0150
## 2 - 4         -0.62 0.206 16 -3.007  0.0080
## 3 - 4         -0.06 0.206 16 -0.291  0.7750

# estimating single contrast c = (1, -.5, -.5)
# comparing operator 1 with operators 2 and 3
contrast1v23.emmc <- function(levs)
  data.frame('t1 v avg t2 t3' = c(1, -.5, -.5, 0))
emmmeans::contrast(pulp.emm, 'contrast1v23')
```

```
## contrast estimate SE df t.ratio p.value
## t1.v.avg.t2.t3    -0.1 0.178 16 -0.560  0.5830
```

- b. It more easily generalises to the more complicated models we will see in Chapter 3.
- c. It explicitly acknowledges that we have fitted a linear model, and so encourages us to check the model assumptions (see Exercise 3).
- d. It is straightforward to apply adjustments for multiple comparisons.

## 2.7 Multiple comparisons

When we perform hypothesis testing, we choose the critical region (i.e. the rule that decides if we reject  $H_0$ ) to control the probability of a type I error; that is, we control the probability of incorrectly rejecting  $H_0$ . If we need to test multiple hypotheses, e.g. to test all pairwise differences, we need to consider the overall probability of incorrectly rejecting **one or more** null hypothesis. This is called the **experiment-wise** or **family-wise** error rate.

For Example 2.1, there are  $\binom{4}{2} = 6$  pairwise comparisons. Under the assumption that all tests are independent<sup>5</sup>, assuming each individual test has type I error 0.05, the experiment-wise type I error rate is given by:

```
alpha <- 0.05
1 - (1 - alpha)^6
```

```
## [1] 0.265
```

---

<sup>5</sup>They aren't, but it simplifies the maths!

An experiment-wise error rate of 0.265 is substantially greater than 0.05. Hence, we would expect to make many more type I errors than may be desirable. `xkcd` has a fun example:

```
alpha <- 0.05
1 - (1 - alpha)^20

## [1] 0.642
```

Therefore it is usually desirable to maintain some control of the experiment-wise type I error rate. We will consider two methods.

1. The **Bonferroni method**. An upper bound on the experiment-wise type I error rate for testing  $k$  hypotheses can be shown to be

$$\begin{aligned} P(\text{wrongly reject at least one of } H_0^1, \dots, H_0^k) &= P\left(\bigcup_{i=1}^k \{\text{wrongly reject } H_0^i\}\right) \\ &\leq \sum_{i=1}^k \underbrace{P(\text{wrongly reject } H_0^i)}_{\leq \alpha} \\ &\leq k\alpha. \end{aligned}$$

Hence a *conservative*<sup>6</sup> adjustment for multiple comparisons is to test each hypothesis at size  $\alpha/k$ , i.e. for the CRD compare to the quantile  $t_{n-t, 1-\frac{\alpha}{2k}}$  (or multiply each individual p-value by  $k$ ).

For Example 2.1, we can test all pairwise comparisons, each at size  $\alpha/k$  using the `adjustment` argument in `pairs`.

```
pairs(pulp.emm, adjust = 'bonferroni')

## contrast estimate      SE df t.ratio p.value
## 1 - 2          0.18 0.206 16   0.873  1.0000
## 1 - 3         -0.38 0.206 16  -1.843  0.5030
## 1 - 4         -0.44 0.206 16  -2.134  0.2920
## 2 - 3         -0.56 0.206 16  -2.716  0.0920
## 2 - 4         -0.62 0.206 16  -3.007  0.0500
## 3 - 4         -0.06 0.206 16  -0.291  1.0000
##
## P value adjustment: bonferroni method for 6 tests
```

Now, only one comparison is significant at an experiment-wise type I error rate of  $\alpha = 0.05$  (operators 2 and 4).

2. **Tukey's method**. An alternative approach that gives an exact experiment-wise error rate of  $100\alpha\%$  compares the  $t$  statistic to a critical

---

<sup>6</sup>So the experiment-wise type I error will actually be less than the prescribed  $\alpha$

value from the studentised range distribution<sup>7</sup>, given by  $\frac{1}{\sqrt{2}}q_{t,n-t,1-\frac{\alpha}{2}}$  with  $q_{t,n-t,1-\frac{\alpha}{2}}$  the  $1 - \frac{\alpha}{2}$  quantile from the studentised range distribution (available in R as `qtukey`).

For Example 2.1:

```
pairs(pulp.emm)
```

```
## contrast estimate SE df t.ratio p.value
## 1 - 2      0.18 0.206 16  0.873 0.8190
## 1 - 3     -0.38 0.206 16 -1.843 0.2900
## 1 - 4     -0.44 0.206 16 -2.134 0.1840
## 2 - 3     -0.56 0.206 16 -2.716 0.0660
## 2 - 4     -0.62 0.206 16 -3.007 0.0380
## 3 - 4     -0.06 0.206 16 -0.291 0.9910
##
```

## P value adjustment: tukey method for comparing a family of 4 estimates

The default adjustment in the `pairs` function is the Tukey method. Comparing the p-values for each comparison using unadjusted t-tests, the Bonferroni and Tukey methods:

```
pairs.u <- pairs(pulp.emm, adjust = 'none')
pairs.b <- pairs(pulp.emm, adjust = 'bonferroni')
pairs.t <- pairs(pulp.emm)
data.frame(transform(pairs.b)[, 1:5], Bonf.p.value = transform(pairs.b)[, 6], Tukey.p.value = tra
```

```
## contrast estimate SE df t.ratio Bonf.p.value Tukey.p.value
## 1 1 - 2      0.18 0.206 16  0.873      1.0000      0.8185
## 2 1 - 3     -0.38 0.206 16 -1.843      0.5034      0.2903
## 3 1 - 4     -0.44 0.206 16 -2.134      0.2918      0.1845
## 4 2 - 3     -0.56 0.206 16 -2.716      0.0915      0.0658
## 5 2 - 4     -0.62 0.206 16 -3.007      0.0501      0.0377
## 6 3 - 4     -0.06 0.206 16 -0.291      1.0000      0.9911
## unadjust.p.value
## 1      0.39551
## 2      0.08389
## 3      0.04864
## 4      0.01525
## 5      0.00835
## 6      0.77476
```

Although the decision on which hypotheses to reject (comparison 2 - 4) is the same here for both methods, the p-values from the Bonferroni method are all

<sup>7</sup> Given two independent samples  $u_1, \dots, u_l$  and  $v_1, \dots, v_m$  from the same distribution, the studentised range distribution is the distribution of  $\frac{R}{\sqrt{2}S}$ , where  $R = u_{\max} - u_{\min}$  is the range of the first sample, and  $S = \sqrt{\frac{1}{m-1} \sum_{i=1}^m (v_i - \bar{v})^2}$  be the sample standard deviation of the second sample.

larger, reflecting its more conservative nature.

## 2.8 Exercises

1. a. For Example 2.1, calculate the mean response for each operator and show that the treatment differences and results from hypothesis tests using the results in Section 2.5 are the same as those found in Section 2.6 using `pairwise.t.test`, and `emmeans`.
- b. Also check the results in Section 2.7 by (i) adjusting individual p-values (for Bonferroni) and (ii) using the `qtukey` command.

Solution

As a reminder, the data from the experiment is as follows.

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

The mean response from each treatment is given by

operator	n_i	mean	variance
1	5	60.2	0.268
2	5	60.1	0.058
3	5	60.6	0.052
4	5	60.7	0.047

The sample variance,  $s^2 = 0.106$ , from (2.9). As  $\sum_{i=1}^t c_i^2/n_i = \frac{2}{5}$  for contrast vectors  $c$  corresponding to pairwise differences, the standard error of each pairwise difference is given by  $\sqrt{\frac{2s^2}{5}} = 0.206$ . Hence, we can create a table of pairwise differences, standard errors and test statistics.

contrast	estimate	SE	df	t.ratio	unadjust.p.value	Bonferroni	Tukey
1 - 2	0.18	0.206	16	0.873	0.396	1.000	0.819
1 - 3	-0.38	0.206	16	-1.843	0.084	0.503	0.290
1 - 4	-0.44	0.206	16	-2.134	0.049	0.292	0.184
2 - 3	-0.56	0.206	16	-2.716	0.015	0.092	0.066
2 - 4	-0.62	0.206	16	-3.007	0.008	0.050	0.038
3 - 4	-0.06	0.206	16	-0.291	0.775	1.000	0.991

Unadjusted p-values are obtained from the t-distribution, as twice the tail probabilities ( $2 * (1 - \text{pt}(\text{abs}(\text{t.ratio}), 16))$ ). For Bonferroni, we simply multiply these p-values by  $\binom{t}{2} = 6$ , and then take the minimum of this value and

1. For the Tukey method, we use `ptukey(abs(t.ratio) * sqrt(2), 4, 16)` (see `?ptukey`).

Alternatively, to test each hypothesis at the 5% level, we can compare each `t.ratio` to (i) `qt(0.975, 16) = 2.12` (unadjusted); (ii) `qt(1 - 0.025/6, 16) = 3.008` (Bonferroni); or (iii) `qtukey(0.95, 4, 16) / sqrt(2) = 2.861`.

2. (Adapted from Wu and Hamada, 2009) The bioactivity of four different drugs  $A$ ,  $B$ ,  $C$  and  $D$  for treating a particular illness was compared in a study and the following ANOVA table was given for the data:

Source	Degrees of freedom	Sums of squares	Mean square
Treatment	3	64.42	21.47
Residual	26	62.12	2.39
Total	29	126.54	

- What considerations should be made when assigning drugs to patients, and why?
- Use an  $F$ -test to test at the 0.01 level the null hypothesis that the four drugs have the same bioactivity.
- The average response from each treatment is as follows:  $\bar{y}_A = 66.10$  ( $n_A = 7$  patients),  $\bar{y}_B = 65.75$  ( $n_B = 8$ ),  $\bar{y}_C = 62.63$  ( $n_C = 9$ ), and  $\bar{y}_D = 63.85$  ( $n_D = 6$ ). Conduct hypothesis tests for all pairwise comparisons using the Bonferroni and Tukey methods for an experiment-wise error rate of 0.05.
- In fact,  $A$  and  $B$  are brand-name drugs and  $C$  and  $D$  are generic drugs. Test the null hypothesis that brand-name and generic drugs have the same bioactivity.

#### Solution

- Each patient should be randomly allocated to one of the drugs. This is to protect against possible bias from lurking variables, e.g. demographic variables or subjective bias from the study administrator (blinding the study can also help to protect against this).
- Test statistic = (Treatment mean square)/(Residual mean square) =  $21.47/2.39 = 8.98$ . Under  $H_0$ : no difference in bioactivity between the drugs, the test statistic follows an  $F_{3,26}$  distribution, which has a 1% critical value of `qf(0.99, 3, 26) = 4.637`. Hence, we can reject  $H_0$ .
- For each difference, the test statistic has the form

$$\frac{|\bar{y}_i - \bar{y}_j|}{s \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}},$$

for  $i, j = A, B, C, D; i \neq j$ . The treatment means and repetitions are given in the question (note that not all  $n_i$  are equal). From the ANOVA table, we get  $s^2 = 62.12/26 = 2.389$ . The following table summarises the differences between drugs:

	$A - B$	$A - C$	$A - D$	$B - C$	$B - D$	$C - D$
Abs. difference	0.35	3.47	2.25	3.12	1.9	1.22
Test statistic	0.44	4.45	2.62	4.15	2.28	1.50

The Bonferroni critical value is  $t_{26, 0.01/12} = 3.507$ . The Tukey critical value is  $q_{4, 26}/\sqrt{2} = 0.304$  (available R as `qtukey(0.01, 4, 26) / sqrt(2)`). Hence under both methods, bioactivity of drugs  $A$  and  $C$ , and  $B$  and  $C$ , are significantly different.

- iv. A suitable contrast has  $c = (0.5, 0.5, -0.5, -0.5)$ , with  $c^T \tau = (\tau_A + \tau_B)/2 - (\tau_C + \tau_D)/2$  (the difference in average treatment effects).

An estimate for this contrast is given by  $(\bar{y}_A. + \bar{y}_B.)/2 - \bar{y}_C. + \bar{y}_D.)/2$ , with variance

$$\text{Var} \left( \frac{1}{2}(\bar{y}_A. + \bar{y}_B.) - \frac{1}{2}(\bar{y}_C. + \bar{y}_D.) \right) = \frac{\sigma^2}{4} \left( \frac{1}{n_A} + \frac{1}{n_B} + \frac{1}{n_C} + \frac{1}{n_D} \right).$$

Hence, a test statistic for  $H_0 : \frac{1}{2}(\tau_A + \tau_B) - \frac{1}{2}(\tau_C + \tau_D) = 0$  is given by

$$\frac{\frac{1}{2}(\bar{y}_A. + \bar{y}_B.) - \frac{1}{2}(\bar{y}_C. + \bar{y}_D.)}{\sqrt{\frac{s^2}{4} \left( \frac{1}{n_A} + \frac{1}{n_B} + \frac{1}{n_C} + \frac{1}{n_D} \right)}} = \frac{2.685}{\frac{\sqrt{2.389}}{2} \sqrt{\frac{1}{7} + \frac{1}{8} + \frac{1}{9} + \frac{1}{6}}} = 4.70.$$

The critical value is  $t_{26, 0.01/2} = 2.78$ . Hence, we can reject  $H_0$  and conclude there is a difference between brand-name and generic drugs.

3. The below table gives data from a completely randomised design to compare six different batches of hydrochloric acid on the yield of a dye (naphthalene black 12B).

```
napblack <- data.frame(batch = rep(factor(1:6), rep(5, 6)),
  repetition = rep(1:5, 6),
  yield = c(145, 40, 40, 120, 180, 140, 155, 90, 160, 95,
            195, 150, 205, 110, 160, 45, 40, 195, 65, 145,
            195, 230, 115, 235, 225, 120, 55, 50, 80, 45)
)
knitr::kable(
tidyr::pivot_wider(napblack, names_from = batch, values_from = yield)[, -1],
```



Table 2.5: Naphthalene black experiment: yields (grams of standard colour) from six different batches of hydrochloric acid.

Batch 1	Batch 2	Batch 3	Batch 4	Batch 5	Batch 6
145	140	195	45	195	120
40	155	150	40	230	55
40	90	205	195	115	50
120	160	110	65	235	80
180	95	160	145	225	45

```
col.names = paste("Batch", 1:6),  
caption = "Naphthalene black experiment: yields (grams of standard colour) from six different batches of hydrochloric acid."  
)
```

Conduct a full analysis of this experiment, including

- exploratory data analysis;
- fitting a linear model, and conducting an F-test to compare to a model that explains variation using the six batches to the null model;
- usual linear model diagnostics;
- multiple comparisons of all pairwise differences between treatments.



## 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.