

1 **Analysis of variance in soil research: let the analysis fit the design.**

2 R. WEBSTER^a & R. M. LARK^{b*}

3 ^a*Rothamsted Research, Harpenden AL5 2JQ, UK, and* ^b*British Geological Survey, Key-*
4 *worth, Nottingham NG12 5GG, UK*

5 Correspondence: R. Webster. E-mail: richard.webster@rothamsted.ac.uk

6 * Present address: School of Biosciences, The University of Nottingham, Sutton Bon-
7 ningtonCampus, Sutton Bonnington, Leicestershire LE12 5RD, UK

8 Running head:

9 *Analysis of variance in soil research*

10 **Summary**

11 Sound design for experiments on soil is based on two fundamental principles: repli-
12 cation and randomization. Replication enables investigators to detect and measure
13 contrasts between treatments against the back-drop of natural variation. Random al-
14 location of experimental treatments to units enables effects to be estimated without
15 bias and hypotheses to be tested. For inferential tests of effects to be valid an analysis
16 of variance (ANOVA) of the experimental data must match exactly the experimental
17 design. Completely randomized designs are usually inefficient. Blocking will usually
18 increase precision, and its role must be recognized as a unique entry in an ANOVA table.
19 Factorial designs enable questions on two or more factors and their interactions to be
20 answered simultaneously, and split-plot designs may enable investigators to combine
21 factors that require disparate amounts of land for each treatment. Each such design
22 has its unique correct ANOVA; no other ANOVA will do. One outcome of an ANOVA is

23 a test of significance. If it turns out to be positive then the investigator may examine
24 the contrasts between treatments to discover which themselves are significant. Those
25 contrasts should have been ones in which the investigator was interested at the outset
26 and which the experiment was designed to test. Post-hoc testing of all possible con-
27 trasts is deprecated as unsound, though the procedures may guide an investigator to
28 further experimentation. Examples of the designs with simulated data and programs
29 in GenStat and R for the analyses of variance are provided as supplementary material.

30 **Highlights**

- 31 • Replication and randomization are essential for sound experimentation on vari-
32 able soil.
- 33 • Analyses of variance of data from experiments must match the experimental
34 designs.
- 35 • Experiments should be designed to answer pre-planned questions and test hy-
36 potheses.
- 37 • Efficiency can be gained by blocking and factorial combinations of treatments.

38 **A little history**

39 In 1843 John Lawes, the then owner of the Rothamsted estate in Hertfordshire, Eng-
40 land, and his newly appointed scientist, Henry Gilbert, planned their experiment on
41 Broadbalk field to test and compare the responses of winter wheat to various combi-
42 nations of fertilizers. The experimental treatments were applied to long narrow strips
43 of land running the length of the field, which were divided in a perpendicular direction
44 into sections. Lawes and Gilbert weighed the yields, and they sampled both the crop
45 and the soil in every plot in every section so as to measure the off-take of nutrients and
46 the nutrient status of the soil. A few years later they laid down similar experiments
47 on spring barley (on Hoosfield, in 1852) and a meadow (Park Grass, in 1856), both
48 of which are still running. They also meticulously recorded the weather. Rotham-
49 sted Research (2006) has summarized the history and main findings of these long-term
50 experiments in its guide.

51 By the end of the First World War, during which Rothamsted began to receive
52 money from the British government for its research, a huge body of data had accrued
53 from these long-term experiments, and in 1919 R.A. Fisher was appointed to analyse
54 the data and make sense of them.

55 Fisher soon realized that without replication, which was the situation on Park Grass,
56 he could not discover how variable was the response to any one treatment. The treat-
57 ments on Broadbalk were replicated, but because the different plots for each treatment
58 lay in a single strip he could not separate the effects of the treatments from the soil's
59 natural variation as expressed in differences between the strips. This natural variation
60 and the treatment effects are said to be confounded. The treatments on the spring
61 barley experiment were replicated on plots that were separated from one another but
62 in a way that might be confounded with the natural variation in the field. So, again,
63 it was not possible to estimate the effects of the fertilizers alone.

64 Having recognized the serious shortcomings of those old trials, Fisher formalized and
65 systematized what had, hitherto, been inconsistently and erratically applied elements

66 of experimental design. One was replication, present in some of the experiments but
67 not all, and necessary to provide information on the variation in responses. The other
68 was randomization, necessary to avoid the bias which could arise if treatment effects
69 are confounded with sources of variation that are uncontrolled and might be unknown.
70 Fisher devised the analysis of variance (ANOVA) to separate the sources of variation in
71 data from such experiments, to estimate quantitatively the effects of different treat-
72 ments and to provide inferential tests to judge whether the observed differences could
73 have arisen by chance rather than as results of the imposed treatments. Fisher also
74 introduced blocking to remove effects such as trends across experiments. Trends of
75 this kind do not introduce bias if the experimental design is randomized, but block-
76 ing improves the sensitivity of the experiment to detect treatment effects against the
77 background variation represented by the trends.

78 Fisher's principles of experimental design and the concomitant analysis of variance
79 are as valid today as they were 90 years ago. They have been the foundation of
80 agronomic practice ever since, and statisticians collaborate with agronomists to ensure
81 that designs will produce data that can be analysed to answer the questions put at
82 the outset. Numerous text books are available to guide practitioners; two that we
83 can recommend unreservedly are the evergreen by Snedecor & Cochran (1989) and the
84 more recent book by Mead *et al.* (2003). Cochran & Cox (1957) remains a standard
85 text. You might like also to see the Statistical Checklists prepared by Jeffers (1978).

86 Sadly, many of today's soil scientists are working without the guidance or collabora-
87 tion of statisticians. One consequence is that they often plan experiments and surveys
88 that cannot or are unlikely to answer their questions; or having designed the experi-
89 ments soundly they vitiate the potential of the experiments to answer the questions
90 by improper sampling. Or they see opportunities to answer new questions that were
91 not envisaged when the original experiments were planned, either by themselves or
92 by other scientists, yet fail to appreciate the limitations inherent in the designs. A
93 further consequence is that despite having designed their experiments and surveys well
94 they analyse the data from them incorrectly. All too often they load their data into

95 a statistical package, press a few buttons on a menu without understanding, and copy
96 the output into their scripts.

97 We write in this critical vein from our experience as advisors to the journal's editors
98 in the last few years, and from the experience of the journal's statistical advisory panel.
99 It is no exaggeration to state that most of the papers on which the editors have sought
100 advice have embodied one or more of the above failings. In the first set of circumstances
101 we have felt obliged to judge the results of little worth and to advise the editors to reject
102 the papers. To paraphrase one of R.A. Fisher's remarks, it has been like conducting
103 post-mortems only to say what the experiments died of. In some instances we have
104 asked for further sampling. In the second we have seen that redemption is often possible
105 by fresh and correct analysis of the data.

106 In one short article we cannot describe all that investigators should do. Instead we
107 focus on the specific matter, namely analyses of variance that follow from the designs,
108 and in particular on the most frequent mismatches between design and analysis. At the
109 best such mismatches lead to loss of information and so to waste of the effort required
110 to do the experiment. At worst the inferences made from the analysis are unsafe and
111 lead to bad decisions. We have already remarked on this in an editorial (Webster *et*
112 *al.*, 2016). In the comic opera *The Mikado* by W.S. Gilbert and Arthur Sullivan the
113 Mikado himself demands that the punishment fit the crime. Here we demand that the
114 analysis fit the design.

115 **Designs**

116 We describe in detail below the commonest and most straightforward designs, starting
117 with the simplest, completely randomized schemes, introducing blocking, and progress-
118 ing to factorial and then split-plot designs. We have provided examples of these designs
119 with simulated data together with programs in GenStat and R for the correct anal-
120 yses of variance and the output from those analyses in the zip file **Supplementary**
121 **material.zip**.

122 *Completely randomized (CR) design*

123 We begin with the simplest design. Suppose that investigators wish to compare the
 124 effects of several manurial treatments on some property of the soil, say the microbial
 125 biomass, which we shall denote z . They replicate their treatments and assign them to
 126 the experimental plots in a completely randomized and independent way. Let there
 127 be n_1 treatments, each replicated n_2 times, so that there are $N = n_1 \times n_2$ plots, or
 128 units, of the design. Treatments are allocated to plots independently and at random.
 129 This means that the probability that the first plot in the experiment is allocated to the
 130 j th treatment is n_2/N , equivalently $1/n_1$. Subsequently when n_j replicates of the j th
 131 treatment remain to be assigned, the probability that any one of the N_u plots that have
 132 still to be assigned a treatment will ultimately receive treatment j is n_j/N_u . Figure 1
 133 shows one outcome of such assignment in which $n_1 = 4$ and $n_2 = 5$.

134 The files `exp1.*` in the Supplementary material contain data with this design and
 135 the programs for analysing them.

136 The analysis of variance for this design appears in Table 1. Note that this presenta-
 137 tion of the analysis of variance, and that for subsequent designs, hold for the balanced
 138 case in which the numbers of replicates of the treatments are equal. The texts to which
 139 we have referred provide further information on analysis in the unbalanced case, but
 140 the topic is beyond the scope of the paper. The total mean square is T :

$$T = \frac{1}{n_1 n_2 - 1} \sum_{j=1}^{n_1} \sum_{i=1}^{n_2} (z_{i,j} - \bar{z})^2, \quad (1)$$

141 where $z_{i,j}$ is the measured response of the i th replicate of the j th treatment and \bar{z} is
 142 the mean response over all $n_1 n_2$ plots. One can see that this quantity is a variance, the
 143 variance of the plot responses. The divisor of the sum of squares, $n_1 n_2 - 1$, is called
 144 the degrees of freedom in Table 1. It can be regarded as the number of independent
 145 pieces of information about the variation of the plot responses provided by the data.
 146 There are $n_1 n_2 - 1$ degrees of freedom rather than $n_1 n_2$ because each plot response is
 147 compared to the overall mean estimated from all the data. Because

$$\sum_{j=1}^{n_1} \sum_{i=1}^{n_2} (z_{i,j} - \bar{z}) = 0$$

148 it follows that, when we know the values of $n_1 n_2 - 1$ differences in the summation, the
 149 last one is fixed and so provides no new information.

150 The within-treatment mean square, W , is computed as

$$W = \frac{1}{n_1(n_2 - 1)} \sum_{j=1}^{n_1} \sum_{i=1}^{n_2} (z_{i,j} - \bar{z}_j)^2, \quad (2)$$

151 where \bar{z}_j is the average response of all plots in the j th treatment. The value estimated
 152 by W is the variance of plot responses within the treatments (i.e. the variance about
 153 the treatment means). This quantity is σ_W^2 in Table 1. It has $n_1(n_2 - 1)$ degrees of
 154 freedom in this simple balanced case because each of the n_1 treatments contributes
 155 $n_2 - 1$ degrees of freedom from the independent variations about the mean of its n_2
 156 replicates, from which the treatment mean is estimated.

157 The between-treatment mean square, called B in Table 1, is computed for this simple
 158 balanced case as

$$B = \frac{1}{n_1 - 1} \sum_{j=1}^{n_1} n_2 (\bar{z}_j - \bar{z})^2. \quad (3)$$

159 This is equivalent to the sum, over all plots, of the squared difference between the corre-
 160 sponding treatment mean and the overall mean, divided by the number of independent
 161 variations among the treatment means.

162 The residual mean square in an analysis of variance is a direct estimate of a variance
 163 component. In general, however, mean squares estimate combinations of more than
 164 one variance component. Table 1 shows that B estimates $\sigma_W^2 + n_2 \sigma_B^2$. The quantity
 165 σ_B^2 is the variance among the treatment means. If there were no differences between
 166 the treatments then this quantity would be zero, and, as can be seen in the table, B
 167 and W would both estimate σ_W^2 , and the ratio $F = B/W$ in the table would have
 168 an expected value of 1. We use the standard notation of the Roman letter s for an
 169 estimate of the underlying quantity σ , so by s_W^2 we denote the estimate of σ_W^2 provided
 170 by W in Table 1.

171 Apart from separating the sources of variation in the experiment and providing
 172 quantitative values of the variances attributed to those sources, the analysis enables us

173 to draw inferences. If the responses in z to the treatments differ from one another then
174 we should expect the ratio B/W to exceed 1. But B/W could exceed 1 purely through
175 random variation; so how can we tell that we have a real effect of the treatments?
176 We do so by putting forward the ‘null hypothesis’, often designated H_0 in statistics
177 textbooks. It is the hypothesis that there are no differences, and we consider the
178 strength of evidence against it. That evidence is the magnitude of B/W in relation to
179 the distribution of F if the null hypothesis were true. We can do so because, as a result
180 of our design, B and W would be independent estimates of σ_W^2 if the null hypothesis
181 were true. It follows from the independent random allocation of treatments to plots,
182 and it appears in the ANOVA table in the way that the $n_1n_2 - 1$ total degrees of freedom
183 are partitioned into the between-treatment and within-treatment (residual) degrees of
184 freedom.

185 In these circumstances the variance ratio has the F distribution under the null
186 hypothesis and the shape of the distribution that depends on the degrees of freedom for
187 the numerator and denominator of the ratio. One can therefore compute the probability
188 that an F ratio as large or larger than the value observed in the table would arise
189 under the null hypothesis through random variation. The smaller is this probability,
190 or P -value, the stronger is the experimental evidence that we should reject the null
191 hypothesis and say that the treatments have produced different responses. It is now a
192 short step to the common notion of statistical significance. It is conventional to take
193 $P = 0.05$ as a threshold. If P exceeds 0.05 investigators accept the null hypothesis.
194 Otherwise, with $P \leq 0.05$ they declare that the observed differences are ‘significant’—
195 and they decorate their tables of means with stars, which again we deprecate! One
196 may choose some other value of P depending largely on how serious it would be to
197 come to a false conclusion.

198 Inference from the analysis of an experiment like that above is based on assumptions
199 about the distribution of random quantities under the null hypothesis that are justified
200 by that design, the way it was laid out in the field, glasshouse or laboratory and on
201 the numbers of the degrees of freedom for the variance ratio. In this sense the analysis

202 (and ANOVA table) match the design.

203 *Randomized complete block (RCB) design*

204 Where investigators know of or suspect trends in fertility, drainage or pollutants that
205 might affect their results they typically replicate their treatments in blocks. In the
206 simplest case each treatment is replicated once and only once in each block. The allo-
207 cation of treatments within the blocks is done independently and at random. Figure 2
208 shows one realization of a RCB design for four treatments and five blocks, and so the
209 same total number of replicates as the completely randomized case in Figure 1. The
210 blocks are separated by the dotted lines; notice that in each block there is one plot
211 for each of the n_1 treatments. The blocks in this figure are laid out as rows across the
212 experimental layout and so would be suitable if a trend in soil properties was known
213 or suspected to occur from the top to the bottom of the site.

214 The files exp2.* in the Supplementary material contain data with this design and
215 the programs for analysing them.

216 The analysis of variance for this design, still with n_1 treatments each replicated
217 once in each of n_2 blocks, appears in Table 2. Here σ_W^2 and σ_B^2 are the underlying
218 variances for plots and treatments as before. There is an additional line in the table
219 for the between-block mean square with $n_2 - 1$ degrees of freedom; σ_A^2 is the variance
220 between blocks. The total degrees of freedom and the treatment degrees of freedom
221 are unchanged from Table 1, but there are $n_2 - 1$ fewer residual degrees of freedom.
222 This follows from simple arithmetic, but it also indicates that the random allocation
223 of treatments to plots is more constrained in the RCB design than in the CR design
224 (once one plot in block k has been assigned to the j th treatment we know that no other
225 plot in the block will receive it). For this reason there is somewhat less information in
226 the residual mean square than in the CR design with the same number of plots and
227 treatments.

228 Where does the between-block variance come from? It is natural variation in the
229 experimental environment which appears as between-block rather than within-block

230 variation. If blocking were not undertaken then this variation would be part of the
231 residual variance, σ_W^2 . This means that, if the between-block variance is large, then we
232 reduce the residual variance and so should increase the variance ratio B/W , making
233 the experiment and analysis more sensitive for comparing the differences between the
234 treatments. This is why blocking, appropriately planned, should be advantageous.
235 Snedecor & Cochran (1989) provide formulae for calculating the efficiency of blocking.
236 At its simplest they calculate it as the ratio of the residual variances:

$$\text{Efficiency} = s_{\text{CR}}^2 / s_{\text{RB}}^2, \quad (4)$$

237 where s_{CR}^2 is the residual variance on the assumption that the design was completely
238 randomized (CR) while s_{RB}^2 is the residual variance of the RCB design. You can find
239 further detail of the calculation on pages 263 and 264 of Snedecor & Cochran (1989).

240 An efficient blocking design is evidently one in which the differences between the
241 blocks are larger than the variation within the blocks. In practice one might achieve this
242 by keeping the blocks compact, although in a field where there is a strong trend in the
243 soil or environment in one direction rectangular blocks with the long side perpendicular
244 to the direction of the trend would be preferred. It is important to pay attention to the
245 structure of the blocks, because, as above, there is a small penalty for blocking from
246 the reduced residual degrees of freedom, and this will be worth paying only if there are
247 real differences between the blocks.

248 The variance ratio A/W appears in Table 2, and one could use it to test the null
249 hypothesis that the between-block variance, σ_A^2 , is zero. That would be of interest
250 only in that it shows whether the blocking is better than random assignment of plots
251 to blocks. Sometimes, however, the scientist, having found that the evidence for a
252 difference among the blocks is weak, ignores the blocking and reports an analysis of
253 variance appropriate for a CR design. Such an analysis does not fit the design. The
254 scientist might try to justify that analysis because the blocks have been shown not
255 to differ, but that misses the point. What the correct analysis shows us, and shows
256 explicitly in the ANOVA table, is how the actual allocation of treatments to plots was

257 undertaken; it shows that in the RCB case we have $(n_1 - 1) \times (n_2 - 1)$ degrees of
258 freedom, not $n_1(n_2 - 1)$. In short, the correct analysis reports the reduction, albeit
259 small, in information about the residual variance that follows from the constraints of
260 blocking. The extra n_2 residual degrees of freedom in the analysis as if the design were
261 completely randomized means that, other things being equal, a given variance ratio
262 appears to offer stronger evidence against the null hypothesis. This inference would be
263 unsafe, however, because the quoted degrees of freedom would not describe the actual
264 randomization. In practice this would mean that the variance ratio for a treatment
265 effect would be compared with the wrong distribution of the F statistic. The analysis
266 would not fit the design.

267 The Austrian philosopher Ludwig Wittgenstein was once impressed by an account
268 of a trial that took place following a car accident in Paris. During the trial, models
269 were used to represent the positions of the vehicles involved at the time of the collision
270 (Kenny, 2005). Inspired by this, he developed his picture theory by which a logical
271 proposition is equivalent to a picture of a state of affairs in the world. Such a proposition
272 may take different forms. It may, for example, be spoken, written or drawn. Let us
273 apply the idea in the present context to the design of field experiments.

274 Consider an experiment that has been done according to an RCB design. The design
275 could be illustrated with a diagram such as Figure 2. More often in scientific papers the
276 designs are described in words in *Methods* sections. The equivalent to Figure 2 would
277 be ‘The n_1 treatments were allocated independently and at random within each of n_2
278 blocks.’ Our contention is that the correct analysis of variance table for the experiment,
279 as shown in Table 2, is one more way in which we may express the same proposition.
280 The partition of the sum of squares between rows of the table represents the sources of
281 variation that the experimental design uniquely induces, and the numbers of degrees
282 of freedom show how many blocks and replicates were used as surely as does Figure 2
283 or the verbal statement.

284 That is one reason why this journal asks its authors to provide full ANOVA tables.
285 The request is sometimes misinterpreted as a request for a table of only a set of variance

286 ratios and corresponding P -values; but that is not what is required. The journal
287 requires a table like Tables 1 or 2 shown here, because such a table represents the
288 design definitively. When assessing an experiment both the reviewers and, ultimately,
289 readers must be able to see that the experiment as described in the methods section
290 accords with the ANOVA reported in the results.

291 *Factorial designs*

292 When an investigator is interested in the effects of several factors it is much more
293 efficient to include them in a single experiment than in a series of separate experiments,
294 one for each factor. This was recognized by Fisher (1926) who wrote:

295 *No aphorism is more frequently repeated in connection with field trials, than*
296 *that we must ask Nature few questions, or, ideally, one question, at a time.*
297 *The writer is convinced that this view is wholly mistaken. Nature, he sug-*
298 *gests, will best respond to a logical and carefully thought out questionnaire;*
299 *indeed, if we ask her a single question, she will often refuse to answer until*
300 *some other topic has been discussed.*

301 Yates (1937) set out the principles of factorial designs in his *Technical Communication*
302 *35*, which became the guiding text for fertilizer trials for many years. More recently
303 Carmer & Walker (1982) have urged investigators to take this course.

304 To illustrate the principles of the design and corresponding analysis we take a simple
305 example with three factors, the major plant nutrients, nitrogen (N), phosphorus (P)
306 and potassium (K). Factors are each applied at two or more ‘levels’; in this example
307 we assume that the nutrient is either applied or not (two levels). There are therefore
308 $2^3 = 8$ combinations of factor levels; these are our treatments. The treatments must
309 be replicated between units (plots in this case) according to a suitable design, and
310 analysed in accordance with that design. One might use CR or RCB designs as in the
311 examples already discussed.

312 Let us assume that there are, as before, n_2 replicates arranged in a CR design. We
313 could analyse the data as set out in Table 1 with $8 - 1 = 7$ degrees of freedom for the

314 treatments. This analysis would be quite correct, but it would not be very informative.
315 If we found that the treatments were significantly different then how should we interpret
316 this finding in terms of all our three factors? The factorial design allows us to do this.
317 We can partition the sum of squares due to differences among the treatments into what
318 are called main effects and interactions. There are three main effects in our example,
319 the differences between treatments with contrasting levels of N is one such, and the
320 other main effects are due to P and K. If these effects simply add to one another then
321 all of the treatment sum of squares will be accounted for by the sums of squares for
322 the three main effects. If, in contrast, the difference between plots that receive N and
323 those that receive none is not the same on plots that receive K and those that receive
324 no K then the factors K and N are said to interact. One can see that there are three
325 such interactions in our example: N.P, N.K and P.K. To complicate matters further,
326 if the N.K interaction differs between plots that receive P and those that receive none,
327 then there is a three-way interaction N.K.P. Note that we could express the same
328 three-way interaction in terms of an effect of, for example, the level of N on the P.K
329 interactions, so there is just one three-way interaction in a factorial experiment with
330 three factors. We use this ‘dot’ convention to indicate interactions as established by
331 Wilkinson & Rogers (1973).

332 Table 3 sets out the ANOVA for our example. Note that each main effect has a
333 single degree of freedom; this is because there are two levels of each factor, and so
334 the main effect consists of just the difference between the responses to these levels.
335 In general a factor with U_1 levels has $U_1 - 1$ degrees of freedom for its main effect.
336 Similarly the two-way interactions each have one degree of freedom, in general two
337 factors with U_1 and U_2 levels have an interaction with $(U_1 - 1) \times (U_2 - 1)$ degrees of
338 freedom. Equally the three-way interaction has 1 degree of freedom in our example.
339 In the general case where the third factor has U_3 levels, the three-way interaction has
340 $(U_1 - 1) \times (U_2 - 1) \times (U_3 - 1)$ degrees of freedom. The reader will note that in our
341 example the sum of the degrees of freedom for the main effects and interactions is 7,
342 the same as the treatment degrees of freedom. The treatment degrees of freedom are

343 partitioned between main effects and interactions as is the treatment sum of squares.

344 The quantity σ_W^2 in Table 3 is the underlying variance among the plots receiving the
345 same combination of treatments, and $\sigma_N^2, \sigma_P^2, \dots, \sigma_{NPK}^2$ are the variances attributed to
346 the nutrients and their combinations. The F ratio for any one entry is

$$F = \frac{\text{mean square for the treatments}}{\text{residual mean square}} . \quad (5)$$

347 The standard error of any of the treatment means is

$$SE_{\text{treatment}} = \sqrt{\text{residual mean square}/n_2} . \quad (6)$$

348 Where the investigator goes from there depends very much on the outcome of the
349 analysis. If it turns out that the interactions, especially the threefold interaction of
350 N, P and K, are non-significant and only the main effects of the three nutrients are
351 significant, the investigator may choose to focus on the main effects, i.e. on the means
352 of plots receiving each of the N, P and K averaged over all combinations that include
353 them. Their standard error is

$$SE_{\text{main effect}} = \sqrt{\text{residual mean square}/4n_2} . \quad (7)$$

354 The quantity 4 appears in the denominator because, in the example, n_2 replicates of
355 four treatments contribute to the estimate of the mean response for each level of one
356 of the factors.

357 We cannot consider here all the possible outcomes and their consequences; rather
358 we must leave readers to pursue them elsewhere. Again we recommend Snedecor &
359 Cochran (1989).

360 We include this account of factorial designs and analysis because all too often in
361 papers submitted to the journal the analysis does not match the design. Some authors,
362 having undertaken an experiment according to a factorial design, proceed to analyse
363 it in a series of one-way analyses for each of the main effects. This is bad practice for
364 two reasons. If all the data from the experiment are analysed in this way then the

365 influence of those main effects not considered in a particular analysis will inflate its
366 residual mean square. Further, when there is a substantial interaction between factors
367 the main effect may be small or negligible, even though the factor is an important
368 one. This is our interpretation of what Fisher mean by saying that nature ‘may refuse
369 to answer’ a particular question ‘until some other topic has been discussed.’ If the
370 design is factorial then the analysis should be so as well, otherwise it is very likely that
371 substantial information will be lost.

372 *Split plots*

373 Split-plot designs are common in agricultural experimentation. There are two general
374 circumstances in which they are used. The first is a factorial experiment in which one
375 of the factors can be replicated only between fairly large plots for logistical reasons.
376 A typical example is where one of the factors is an irrigation or drainage treatment.
377 Large plots are needed for these, but it would not be feasible to replicate such plots
378 in factorial combination with several fertilizer treatment as above. The experiment
379 would require too large an area to manage. The solution is to replicate the irrigation
380 factor between appropriate large plots (main plots in the jargon), and then to divide
381 each main plot into sub-plots, one sub-plot for each level or combination of levels of
382 the remaining factors which are allocated to sub-plots at random.

383 Let us suppose that the four manurial treatments of Figure 1 (M1, M2, M3, M4)
384 are to be combined in an experiment in which there are three irrigation treatments (I1,
385 I2, I3)—say no irrigation, irrigation when the soil has dried to half its available water
386 capacity, and irrigation at regular intervals regardless of the water deficit. Figure 3
387 shows a possible layout on the ground with the irrigation treatment replicated between
388 main plots in the blocks, and the manurial treatments replicated between sub-plots
389 within each main plot.

390 How would the data from this experiment be analysed? There are twelve treatments
391 (combinations of the four levels of the manure factor and the three levels of the irri-
392 gation factor). The treatments are replicated in four blocks. One might think that

393 Table 4 would partition the degrees of freedom for the ANOVA; the design is after all
394 a factorial one. An analysis with that structure would be wrong, however; the table
395 does not match the design. To see this reflect on the basic units of the experiments,
396 the sub-plots; there are twelve of them in each block. The ANOVA structure in Table 4
397 implies that there are no constraints on the randomization of the twelve treatments
398 between sub-plots within each block, but that is not the case. If we are told that a
399 plot in the top left corner of a block has treatment I3-M4 we can know, first, that all
400 plots in the same main plot receive level I3 of the irrigation factor, and, second, that
401 no other subplot in the main plot receives level M4 of the manure treatment. In short,
402 Table 4 fails to show that the levels of the irrigation factor were allocated to the main
403 plots while the levels of the manure factor were hen allocated to sub-plots within the
404 main plots.

405 Table 5 sets out the correct analysis for this experiment with the three levels of the
406 irrigation factor randomly allocated between main plots in each of four blocks, and the
407 four levels of the manure factor randomly allocated to the sub-plots within each main
408 plot.

409 The files exp3.* in the Supplementary material contain data with this design and
410 the programs for analysing them.

411 Notice how the F ratios are calculated in Table 5. The denominator for the irrigation
412 F ratio is the main-plot error mean square. That for the manures and the interaction
413 between the irrigation and manures is the sub-plot error mean square. In such a
414 design the sub-plot error variance is smaller than the main-plot error variance. These
415 variances follow through to different standard errors for the means. In this example the
416 manurial treatments are compared more sensitively than the irrigation treatments. If
417 the data from this experiment were mistakenly analysed as in Table 4 then one would
418 underestimate the main-plot error variance and overestimate the sub-plot variance.

419 In an experiment like the one above the treatments, say, manurial and irrigation,
420 are laid out in split-plot designs from the start. While such experiments are not
421 always correctly analysed in papers submitted to the journal, problems more often arise

422 when split-plots are introduced into experiments later on. Consider an original RCB
423 experiment with four treatments like that above. Let us suppose that the treatments
424 are four different kinds of manure and that the investigator planned to compare rates of
425 respiration in the soil between these treatments. Having seen the results he or she then
426 introduces a second factor, the soil water potential. Two soil cores are taken from each
427 plot of the original experiment and equilibrated at one of two soil water potentials, and
428 then the respiration rate of each is measured. The plots in such an experiment are not
429 physically split, and authors are sometimes puzzled when we tell them that they have
430 split-plot designs. They need to recognize that in such a situation the experiment has
431 a split-plot design with manures replicated between main plots and the cores extracted
432 from each main plot serve as sub-plots between which the levels of the water-potential
433 factor are randomized. This should be reflected in an ANOVA table like Table 5. Too
434 often we receive papers in which such experiments are analysed as if they had simple
435 RCB factorial designs.

436 *Sampling within experimental plots*

437 One can rarely measure soil properties of whole plots; almost always the most one can
438 do is to sample the soil and measure the properties of interest on the samples. If one
439 were to take one sample, whether as a single core or a bulked sample from several cores,
440 one would analyse the measurements as above according to the design; i.e. completely
441 randomized or blocked.

442 However, one might well measure the property on each of several cores from each
443 plot. This would provide information on the variation within the plots, and one could
444 elaborate the analysis of variance accordingly. Suppose that one takes n_3 cores of soil
445 from each and every plot, as illustrated in Figure 4 in which there are $n_1 = 4$ treatments
446 replicated $n_2 = 5$ times in a completely randomized arrangement, and $n_3 = 3$ cores
447 per plot. The correct analysis of variance for this design is set out in Table 6. The
448 quantities σ_W^2 and σ_B^2 are the underlying variances between plots within treatments
449 and between treatment means respectively, and σ_C^2 is the variance among cores within

450 plots. This table is comparable to one for a split-plot design with cores as the sub-
 451 plots. The difference is that no factor is replicated randomly at the core level. The
 452 replication is simply to improve estimates of the plot means. Nonetheless, the between-
 453 treatment mean square must be compared with the correct residual, the between-plots
 454 within-treatments mean square, because the treatments are randomized at the plot
 455 level.

456 The standard error of a plot mean is $SE_{\text{plot}} = \sqrt{C/n_3}$, where C is the variance
 457 between cores within plots. If we denote the estimated variance between plots within
 458 treatments by s_{W}^2 we obtain the standard error per treatment mean as

$$SE_{\text{treatment}} = \sqrt{\frac{C}{n_3 n_2} + \frac{s_{\text{W}}^2}{n_2}}. \quad (8)$$

459 If the replicates were arranged in blocks then there would be a corresponding addi-
 460 tional entry for blocks in the analysis.

461 *Pseudo replication*

462 In the previous example, with the ANOVA as in Table 6, the experimenter recognizes
 463 that treatments are replicated and randomized at the plot level, even though measure-
 464 ments are made on n_3 cores in each plot. If, incorrectly, the experimenter treated this
 465 design as one with $n_3 \times n_2$ independent replicates of each treatment, it would be a case
 466 of what statisticians call ‘pseudo replication’. We introduce the topic of pseudo repli-
 467 cation here because many authors of the papers we see commit it either inadvertently
 468 or knowingly without appreciating its inferential consequences. We distinguish three
 469 situations.

- 470 1. The investigator misguidedly regards all $n_2 \times n_3$ observations on each treatment
 471 as the units of the design and for a CR design analyses the data as in Table 1.
 472 He or she then tests the treatment mean against a residual mean square with
 473 $n_1 \times n_2 \times n_3 - n_1$ degrees of freedom. This comprises a form of pseudo replication
 474 because the replicates within plots are not true replicates of the experimental
 475 treatments. Fortunately no serious damage is done; once alerted to the mistake

476 the investigator can re-analyse the data correctly according to Table 6.

477 2. A similar situation arises when a scientist takes either a single core from each plot
478 or bulks multiple cores from each and then splits them into several sub-samples
479 for measurement in the laboratory. These replicate measurements cannot be
480 regarded as independent units in the design. They are pseudo replicates. They
481 may be averaged and analysed as in Table 1, or they may be analysed as individual
482 values as in Table 6. In latter case the variance σ_C^2 represents the variance due
483 to sub-sampling of a single core or composite sample, rather than within-plot
484 variance.

485 3. Most serious of all is when an investigator takes multiple cores of soil from an
486 experiment which itself has few replicates, perhaps only one, and believes that
487 treating the numerous cores as units will compensate for lack of replication of the
488 main plots and analyses the data according to Table 1. The correct analysis is that
489 exemplified in Table 6. With few true replicates of the treatments, however, the
490 experiment is unlikely to be sufficiently sensitive to reveal any but the biggest and
491 most obvious differences. Here the shortcoming is in the design; the experiment
492 should have been planned with more replication in the field and more resources
493 allocated to its execution.

494 The situation arises more often in surveys where investigators want to know
495 how the soil differs from one cultural practice or environment to another. The
496 main difficulty here is in finding sufficient replicates of each kind of practice or
497 environment, especially if access and travel between them are time-consuming and
498 expensive. What usually happens is that the investigator replicates observations
499 at the few sites that can be reached, often only one of each kind.

500 Mean values for the sites actually sampled might be estimated precisely, but
501 differences between practices or environments would not be. If the latter are not
502 replicated, perhaps because replication was impossible, then the investigator can
503 say at the end only by how much the sites themselves differ from one another; any

504 inference about the populations they represent cannot be based on the statistics.

505 *Repeated measurements*

506 The last couple of decades have seen increasing interest in the behaviour of soil over
507 time. Soil scientists have monitored the soil and planned experiments with installa-
508 tions such as static chambers in which to collect gaseous emissions—see, for example,
509 González-Méndez *et al.* (2015) and their repeated measurements of the associated
510 redox potentials from electrodes buried in the soil (González-Méndez *et al.*, 2017),
511 lysimeters in which to monitor leachates passing through the soil, laboratory reactors
512 in which to organic matter is mineralized (e.g. Coban *et al.*, 2016) and microcosms
513 in which to measure the responses of bacteria to imposed treatments over time. The
514 scientists quite properly design their experiments by assigning their treatments to the
515 units, whether chambers, electrodes, lysimeters, reactors or microcosms, with replica-
516 tion and randomization. Then at intervals they make their measurements on every
517 unit. This is especially easy when the measurement is non-invasive, for example by
518 spectrometers. It is also feasible to do so by repeated sub-sampling soil from micro-
519 cosms or field plots. (The soil in long-term experimental plots at Rothamsted has been
520 sampled at intervals over the years since they were first established.)

521 If measurements are made on only two occasions then an appropriate analysis of the
522 data depends on the specific objectives of the experiment. If the variable of interest
523 is the difference between the two observations (e.g. the change in a soil property
524 between the start of a growing season and the end) then the difference may be computed
525 directly for each experimental unit and, being replicated at the level of these units,
526 may be analysed in a straightforward way. If the two observations on each unit are to
527 be analysed together then we have a split-plot design with the chambers, electrodes,
528 lysimeters or microcosms as replicated main plots and the two occasions as sub-plots
529 within the main plots. One can analyse the data quite correctly as set out in Table 5.

530 In situations when observations are repeated on the same units, and they are made
531 on more than two occasions, one must take into account possible correlations between

532 the repeated measurements on any one unit. These correlations might depend on the
533 interval in time between the observations, which the simple split-plot analysis can-
534 not accommodate. The successive measurements on any one installation cannot be
535 regarded as independent. For the purpose of the statistical analysis the chambers,
536 electrodes, lysimeters or microcosms are the units. The data comprise repeated mea-
537 surements on those units, and special techniques that take into account the possible
538 correlations, are required to analyse them. The techniques often go under name of
539 ‘longitudinal analysis’.

540 There is no single correct way of analysing repeated measurements, and we cannot
541 delve into the detail of any of them. Webster & Payne (2002), in this journal, reviewed
542 several options. They described in detail one in which the order of correlations were
543 estimated first by an antedependence analysis, as devised by Kenward (1987), and the
544 results of which were then incorporated into an analysis of differences between treat-
545 ments by residual maximum likelihood (REML). Other options in which the variations
546 in time are modelled as autoregressive processes are available—see again Coban *et al.*
547 (2016).

548 In whatever way data of repeated measurements are analysed that way must honour
549 the design. If you wish to investigate processes in the soil over time with fixed instal-
550 lations such as static chambers or lysimeters or in the laboratory with microcosms
551 then plan your experiments in consultation with a professional statistician and know
552 in advance how you will analyse the data. Of course, you should always know how
553 you will analyse data from any experiment you plan, and for the more straightforward
554 cases you can find recipes in textbooks.

555 **Inferences and comparisons**

556 *Orthogonal contrasts*

557 Obtaining a statistically significant result from an ANOVA, one say for which $P < 0.05$,
558 is never the end of an investigation. On its own it is of limited interest. Far more
559 important are the differences between the means: which of the differences contributed

560 to the result? And are they the ones about which the investigator wanted to know
561 when the experiment was designed?

562 Consider an experiment in which a scientist wants to compare the effects of organic
563 additions to the soil on the respiration rate. The materials to be added are barley
564 straw, wheat straw, cattle slurry and pig slurry. In addition to these four treatments
565 there is a fifth treatment, a control where nothing is added. When this experiment is
566 complete the ANOVA table will include a treatment mean square with four degrees of
567 freedom. This mean square may be compared with the residual mean square to test
568 the null hypothesis that there are no differences in response to the different treatments.
569 Let us suppose that the P -value is so small that the null hypothesis is rejected. Now,
570 which differences contributed to the result? Did the respiration caused by the addition
571 of straw differ from that caused by the addition of slurry? Did the kind of straw affect
572 the result? How did the additions of these organic materials affect the respiration rate
573 in relation to the control? These are the pre-planned questions that the scientist might
574 reasonably have had in mind when the experiment was designed, and the design should
575 have been such as to answer those questions and test the hypotheses underlying them
576 by the appropriate analysis.

577 Why pre-planned questions? With five different treatments there are ten different
578 comparisons that can be made between pairs of treatments, and there are more com-
579 parisons between combinations of treatments. One might test a comparison between
580 the means of two treatments with a t test. The standard error for the difference be-
581 tween two treatment means is $\sqrt{2W/n_2}$, so the test is easy to do. Indeed, for the simple
582 balanced case with n_2 replicates per treatment one may compute the least significant
583 difference for comparison between any pair: $LSD = t\sqrt{2W/n_2}$. With so many possi-
584 ble comparisons it is likely that some will appear ‘significant’ purely through random
585 variation, and with the human eye and brain well-adapted to pick out large differences
586 in tables of means, any inference out of these multiple comparisons is unlikely to be
587 safe. Lark (2017) and Webster (2007) have discussed this matter in greater depth. The
588 meaning of the P -value for a null hypothesis holds when the comparison is planned at

589 the outset; it does not hold for examination of differences after one has inspected a
590 table of means and noted ones that look interesting.

591 Pre-planned questions can be expressed conveniently as a set of orthogonal contrasts.
592 A contrast is a comparison between two treatments, or two groups of treatments. In
593 the example above one contrast might be between soils receiving cattle manure and
594 those receiving pig manure. If we consider the treatments in order:

595 Control; Pig Manure; Cattle Manure; Barley Straw; Wheat Straw,

596 then the contrast mentioned can be expressed by a vector of coefficients

$$\mathbf{c}_1 = [0, -1, 1, 0, 0] .$$

597 This contrast is a comparison between the two manures. There are zero entries that
598 correspond to treatments not in the contrast, and the difference in sign expresses the
599 fact that we are interested in the difference between the two manure treatments.

600 Another contrast one could consider is between the control and all the treatments
601 with additions to the soil. This would be expressed by the coefficients

$$\mathbf{c}_2 = [4, -1, -1, -1, -1] .$$

602 Note that the mean for the control has a coefficient of 4, balancing the -1 entry for
603 each of the treatments with an organic amendment, and the coefficients therefore sum
604 to zero, as in the previous example.

605 We have yet to explain what we mean by an orthogonal contrast. Consider the
606 two examples given. Neither of these contrasts contributes in any way to the other.
607 That is because the second contrast is between the control and all the treatments with
608 an amendment, whereas the first is a contrast between two treatments in the latter
609 group. If I know that the first contrast is large it tells us nothing about the second.
610 Mathematically this is expressed by the fact that the inner product of the two contrast
611 vectors, the sum of the products of their corresponding elements, is zero

$$\mathbf{c}_1 \cdot \mathbf{c}_2 = 0 ,$$

612 as can easily be verified.

613 We can specify two more contrasts, \mathbf{c}_3 and \mathbf{c}_4 , such that the full set are mutually
614 orthogonal. These are

$$\mathbf{c}_3 = [0, 0, 0, -1, 1] ,$$

615 and

$$\mathbf{c}_4 = [0, -1, -1, 1, 1] .$$

616 The contrast \mathbf{c}_3 is between wheat straw and barley straw, and the contrast \mathbf{c}_4 is between
617 straw and manure. The reader can check that any pair of contrasts drawn from the set
618 $\{\mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_3, \mathbf{c}_4\}$ is orthogonal.

619 Note that there are four orthogonal contrasts in this set, which is complete: no
620 additional contrast could be found that is orthogonal to all in this set of four. The
621 number of orthogonal contrasts among a set of treatments is equal to the treatment
622 degrees of freedom. In fact, the orthogonal contrasts can be put into the ANOVA table,
623 one line each, in place of the treatment effects. The treatment sum of squares is
624 partitioned between the contrasts exactly, and each has one degree of freedom. Each
625 contrast can be tested by the ratio of its mean square to the appropriate residual mean
626 square in the design. Note also that orthogonal contrasts can be used in the analysis
627 of a factorial experiment, in which case contrasts can be examined between groups of
628 levels of each factor, and the interaction sum of squares may also be partitioned into
629 corresponding components, each with one degree of freedom.

630 The use of orthogonal contrasts is much to be commended. It requires experimenters
631 to think in advance about their hypotheses, to express them in terms of contrasts and
632 so to embed them in the experimental design. By pre-specifying the orthogonal sets of
633 contrasts experimenters ensure that the P -values they use to test their hypotheses can
634 be interpreted validly.

635 Often investigators notice, at the end of an experiment, contrasts of interest that
636 they had not expected and for which their design did not cater. Should they apply
637 tests for them? The short answer is 'no'; the only safe way to test the hypothesis

638 implied by such a contrast is to design a new experiment for the purpose.

639 Several methods have been proposed to test all comparisons post-hoc. They include
640 Scheffé's critical difference, the Newman–Keuls test, Tukey's 'honest significant differ-
641 ence' and Duncan's multiple range test. The idea underlying them is that by setting
642 the critical limit of P according to the total number of possible comparisons one can
643 identify which specific contrasts can be regarded as significant. Numerous papers sub-
644 mitted to the journal contain results of these methods to test all comparisons between
645 treatment means, and authors then express the results by littering bar charts or ta-
646 bles of treatment means with letters such that all means with the letter 'a' appended
647 cannot be regarded as significantly different, and so on. This is poor practice. It is
648 of the essence of experimental science to advance hypotheses and to test them; that is
649 the scientist's responsibility. It cannot be delegated to an algorithm. Furthermore, the
650 practice wastes the statistical power of a well-designed experiment which is only fully
651 exploited by the proper analysis of a set of orthogonal preplanned contrasts. That is
652 why, with the backing of two of the most experienced statistical analysts of the last
653 century—Nelder (1971) and Finney (1988)—and the allegorical exposition by Carmer
654 & Walker (1982), this journal eschews routine multiple comparisons from tests.

655 Nevertheless, these tests can have merit if they are used in what we might call
656 the 'wash-up' phase of the experimental analysis after the primary hypotheses have
657 been tested. They may be used legitimately to 'screen' differences and help investi-
658 gators to decide whether further research is warranted and to design new experiments
659 accordingly.

660 In summary, good scientific practice identifies a set of hypotheses that can be ex-
661 pressed as particular pre-planned contrasts between the mean responses of treatments
662 or groups of treatments. This is part of the experimental design. The analysis fits
663 the design when the ANOVA table includes the specific orthogonal contrasts as single
664 lines, with one degree of freedom for each mean square, to be tested against the correct
665 residual mean square given constraints on randomization of the treatments between
666 units. If other contrasts catch the experimenter's eye then some of the 'post-hoc' tests

667 listed above might be invoked to screen them.

668 **Some thoughts on sampling**

669 In this paper we have focused on the designs of experiments and the analyses of variance
670 for inference from data obtained according to those designs. Similar considerations
671 apply to sampling to estimate, for example, the mean values of soil properties within
672 regions of interest. We have described suitable designs elsewhere (Webster & Lark,
673 2013), and we cannot go into detail here. Readers can find the general principles in the
674 classic text by Cochran (1977) and their application to spatial sampling in de Gruijter
675 *et al.* (2006).

676 In sampling, as with experiments, the principle that the analysis should fit the
677 design still holds good. In the context of sampling our objective is estimation, and
678 an estimate should be accompanied by a confidence interval to indicate its precision.
679 There are standard methods to compute such confidence intervals, but the method that
680 is used must accord with the sampling design if it is to be safe. For example, most soil
681 scientists would recognize the procedure of computing the sample variance, s^2 , from a
682 set of N observations and then calculating the standard error of the sample mean as

$$\frac{s}{\sqrt{N}}. \quad (9)$$

683 One can compute the confidence interval for the sample mean by multiplying the
684 standard error by the value of Student's t for which the distribution function with
685 $n - 1$ degrees of freedom takes an appropriate value (e.g. 0.975 for the 95% confidence
686 interval). This simple analysis is appropriate, however, only when the N samples have
687 been collected independently and completely at random (also known as simple random
688 sampling). Without the independence, which independent random sampling ensures,
689 the computation of the standard error in Equation (9) is wrong.

690 Too often the journal receives papers in which the analysis of sample data does not fit
691 the design. Most commonly that is because the authors use Equation (9) to compute
692 the standard error of a sample mean based on N samples which were not collected

693 independently and at random, either because the sampling was not randomized (sample
694 sites may have been selected purposively to cover a range of soil variation) or because
695 the samples were collected according to a systematic design (a grid or transect). In the
696 latter, once the positions of one or two sampling sites have been chosen the positions
697 of all the others in the designs are determined by the interval of the grid or transect.
698 One may compute a correct standard error for an estimated mean where sampling
699 has been done systematically on several transects provided the starting points of the
700 transects are chosen at random (de Gruijter *et al.*, 2006) and the analysis fits the
701 design appropriately. Alternatively, model-based estimation may be used (Lark &
702 Cullis, 2004).

703 Other sampling designs may be appropriate. Stratified random sampling is directly
704 analogous to the RCB experimental design discussed above. The domain of interest
705 is divided into strata, which one hopes are less variable internally than the domain as
706 a whole. The estimates are likely to be more precise than those from simple random
707 sampling because the estimation variances are based on the variances within the strata
708 rather than on that of the whole domain. Each stratum is sampled independently
709 and at random, the stratum sample means are combined to obtain an estimate of the
710 domain mean, and the stratum variances are similarly combined to obtain a variance
711 of the estimated mean. If stratification has been used in the sampling design then it
712 must be accounted for in the analysis.

713 **Departures from assumptions**

714 We have stressed throughout that the correct analysis of variance fits the design; no
715 other will do. The conclusions that you may draw from such analyses, however, are
716 based on the assumption that the effects of the various factors (treatments and blocks
717 and their combinations) are additive, that the residuals are normally and independently
718 distributed, and that the variances are homogeneous. Small departures from these ideal
719 conditions are unlikely to affect your conclusions—the analysis of variance is robust in
720 this respect. Large ones, on the other hand, might. Testing for serious departures and

721 the transformations required to make data conform to the assumptions are substantial
722 subjects in their own right, and we cannot deal with them here. Instead we refer you
723 to Chapter 15, pages 273–296, in Snedecor & Cochran (1989), and Chapter 8, pages
724 159–181, in Mead *et al.* (2003).

725 **Epilogue**

726 This paper is not a comprehensive account of the design and analysis of experiments;
727 it was never our intention that it should be. Rather, we have wanted to stress the
728 importance of sound experimental designs, of doing experiments according to those
729 designs and then subsequently analysing the data that accrue likewise. Readers can
730 find details of the designs we mention in the texts we have cited; those texts should
731 cover their requirements.

732 Sound inferences about the effects of treatments on the soil demand that treatments
733 are replicated and assigned to experimental units at random. The natural variability
734 of the soil is substantial, and many replicates might be needed to reveal the effects
735 of the treatments against this back-drop of natural variation. One can often reduce
736 the amount of replication, and increase the efficiency of an investigation, by blocking.
737 Whether a completely randomized design is used, or a randomized complete block
738 design, the design must be accounted for in the analysis, and it should be made explicit
739 by the full ANOVA table. If your paper does not contain such a table then readers cannot
740 be sure that you have analysed your data in a way that fits the design and is valid
741 therefore.

742 More complex experimental designs might be needed for practical reasons. We have
743 given the example of split plots, but others include designs with incomplete blocks
744 and designs in which certain interactions are deliberately confounded and so cannot
745 be estimated. In all cases the experimental design constrains the analysis, and the
746 degrees of freedom in the ANOVA table, and the residual mean square against which an
747 effect is tested, must accord with the design as described. The same holds for repeated
748 measures on the same experimental units, and for experiments when replicated samples

749 from within the basic experimental units are analysed separately.

750 Finally, we have stressed that scientists have the responsibility to propose hypotheses
751 and to design experiments accordingly. By pre-planning particular comparisons scien-
752 tists embed their hypotheses in those designs. Their analyses partition the treatment
753 sums of squares into components corresponding to the orthogonal contrasts.

754 Soil scientists nowadays use some of the most advanced techniques from nuclear
755 magnetic resonance to shallow geophysics, and we like to think that they take advice
756 from specialists beforehand. They should do the same when they apply statistical
757 methods. Modern software provides a wide range of readily available tools for statistical
758 analysis. But when misused by investigators who lack proper understanding they lead
759 to flawed inferences, and those can have damaging consequences if they lead in turn to
760 bad decisions by farmers, environmental managers, statutory authorities and agencies
761 responsible for public health.

762 We encourage soil scientists to think hard about how they design their experiments
763 and then analyse the data. We encourage educators in soil science to ensure that statis-
764 tics, taught by specialists, has an essential place in curricula at both undergraduate
765 and postgraduate level. Finally, we urge soil scientists to consult statisticians when
766 they plan their experiments, and not go along to them at the end and ask them how
767 to analyse their data. Neither you nor we want Fisher to look down and pronounce yet
768 another post-mortem on your experiment.

769 **Supplementary material**

770 As mentioned above, we have provided examples of CR, RCB and split-plot designs
771 with simulated data together with programs in GenStat and R for the correct anal-
772 yses of variance and the output from those analyses in the zip file **Supplementary**
773 **material.zip**. This file can be down-loaded for immediate use. Alternatively, you
774 may obtain it from us directly.

775 **References**

- 776 Carmer, S.G. & Walker, W.M. 1982. Baby Bear's dilemma: a statistical tale. *Agronomy Journal*, **74**, 122-124.
777
- 778 Coban, H., Miltner, A., Centler, F. & Kästner, M. 2016. Effects of compost, biochar
779 and manure on carbon mineralization of biogas residues applied to soil. *European
780 Journal of Soil Science*, **67**, 217–225.
- 781 Cochran, W.G. & Cox, G.M. 1957. *Experimental Designs*, 2nd edition. John Wiley
782 & Sons, New York.
- 783 Cochran, W.G, 1977. *Sampling Techniques*, 3rd edition. John Wiley & Sons, New
784 York.
- 785 De Gruijter, J.J., Brus, D.J., Bierkens, M.F.P. & Knotters, M. 2006. *Sampling for
786 Natural Resources Monitoring*. Springer-Verlag, Berlin.
- 787 Finney, D.J. 1988. Was this in your statistics textbook? III. Design and analysis.
788 *Experimental Agriculture*, **24**, 421–432.
- 789 Fisher, R.A. 1926. The arrangement of field experiments. *Journal of the Ministry of
790 Agriculture of Great Britain*, **33**, 503–513.
- 791 González-Méndez, B., Webster, R., Fiedler, S., Loza-Reyes, E., Hernández, J.M.,
792 Ruíz-Suárez, L.G. & Siebe, C. 2015. Short-term emissions of CO₂ and N₂O in
793 response to periodic flood irrigation with waste water in the Mezquital Valley of
794 Mexico. *Atmospheric Environment*, **101**, 116–124.
- 795 González-Méndez, B., Webster, R., Fiedler, S. & Siebe, C. 2017. Changes in soil
796 redox potential in response to flood irrigation with waste water in central Mexico.
797 *European Journal of Soil Science*, **68**, in press. doi: 10.1111/ejss.12484.
- 798 Jeffers, J.N.R. 1978. *Design of Experiments. Statistical Checklist 1*. NERC Institute
799 of Terrestrial Ecology, Grange-over-Sands, U.K.
800 Available at <http://nora.nerc.ac.uk/5271/> [Accessed 21st July 2017]

- 801 Kenny, A.J. 2005. *Wittgenstein*. Wiley–Blackwell, Oxford.
- 802 Kenward, M.G. 1987. A method for comparing profiles of repeated measurements.
803 *Applied Statistics*, **36**, 296–308.
- 804 Lark, R.M. 2017. Controlling the marginal false discovery rate in inferences from a
805 soil data set with α -investment. *European Journal of Soil Science*, **68**, 221–234.
- 806 Lark, R.M. & Cullis, B.R. 2004. Model-based analysis using REML for inference
807 from systematically sampled data on soil. *European Journal of Soil Science*. **55**,
808 799–813.
- 809 Mead, R., Curnow, R.N. & Hasted, A.M. 2003. *Statistical Methods in Agriculture and*
810 *Experimental Biology*. Chapman and Hall/CRC, Boca Raton, Florida.
- 811 Nelder, J.A. 1971. Discussion on papers by Wynn, Bloomfield, O’Neill and Wetherill
812 (1971). *Journal of the Royal Statistical Society, B*, **33**, 244–246.
- 813 Rothamsted Research 2006. *Guide to the Classical and other Long-term Experiments,*
814 *Datasets and Sample Archive*. Lawes Agricultural Trust, Harpenden, UK.
- 815 Snedecor, G.W & Cochran, W.G. 1989. *Statistical Methods*, 8th edition. Iowa State
816 University Press, Ames, Iowa.
- 817 Webster, R. 2007. Analysis of variance, inference, multiple comparisons and sampling
818 effects in soil research. *European Journal of Soil Science*, **58**, 74–82.
- 819 Webster, R. & Lark, R.M. 2013. *Field Sampling for Environmental Science and*
820 *Management*. Routledge, London.
- 821 Webster, R., Oliver, M.A. & Lark, R.M. 2016. Editorial: statistics in the journal.
822 *European Journal of Soil Science*, **67**, 133–134.
- 823 Webster, R. & Payne, R.W. 2002. Analysing repeated measurements in soil monitor-
824 ing and experimentation. *European Journal of Soil Science*, **53**, 1–13.

- 825 Wilkinson, G.N. & Rogers, C.E. 1973. Symbolic description of factorial models for
826 analysis of variance. *Applied Statistics*, **22**, 392–399.
- 827 Yates, F. 1937. *The Design and Analysis of Factorial Experiments*. Technical Com-
828 munication 35. Commonwealth Bureau of Soil Science, Harpenden, UK.

829 **Table 1** Analysis of variance for n_1 treatments replicated n_2 times in a completely
 830 randomized (CR) design

Source	Degrees of freedom	Mean squares	Parameters estimated	F ratio
Between treatments	$n_1 - 1$	B	$\sigma_W^2 + n_2\sigma_B^2$	B/W
Within treatments (residual)	$n_1(n_2 - 1)$	W	σ_W^2	
Total	$n_1n_2 - 1$	T		

832 **Table 2** Analysis of variance for n_1 treatments replicated n_2 times in a randomized
 833 complete block (RCB) design

834

Source	Degrees of freedom	Mean squares	Parameters estimated	F ratio
Blocks	$n_2 - 1$	A	$\sigma_W^2 + n_2\sigma_A^2$	A/W
835 Between treatments	$n_1 - 1$	B	$\sigma_W^2 + n_2\sigma_B^2$	B/W
Within treatments (residual)	$(n_1 - 1) \times (n_2 - 1)$	W	σ_W^2	
Total	$n_1n_2 - 1$	T		

836 **Table 3** Three-way analysis of variance for three factors, N, P and K, each at two
 837 levels replicated n_2 times in a CR design

Source	Degrees of freedom	Parameters estimated by mean squares	F ratio
Between treatments	7	$\sigma_W^2 + n_2\sigma_B^2$	
N	1	$\sigma_W^2 + n_2\sigma_N^2$	
P	1	$\sigma_W^2 + n_2\sigma_P^2$	
K	1	$\sigma_W^2 + n_2\sigma_K^2$	
N.P	1	$\sigma_W^2 + n_2\sigma_{NP}^2$	
N.K	1	$\sigma_W^2 + n_2\sigma_{NK}^2$	
P.K	1	$\sigma_W^2 + n_2\sigma_{PK}^2$	
N.P.K	1	$\sigma_W^2 + n_2\sigma_{NPK}^2$	
Within treatments (residual)	$8 \times (n_2 - 1)$	σ_W^2	
Total	$8 \times n_2 - 1$	σ_T^2	

838

839 **Table 4** Incorrect partial analysis of variance table for the factorial experiment with
840 manure and irrigation factors illustrated in Figure 3.

Source	Degrees of freedom
Between blocks	3
Between treatments	11
Manure	3
Irrigation	2
Manure×Irrigation	6
Residual	33
Total	47

841

842 **Table 5** Analysis of variance for the split plot experiment with three levels of the
 843 irrigation factor replicated between main plots within blocks, and four levels of the
 844 manure factor replicated between sub-plots within each main plot.

Source	Degrees of freedom	Mean squares	F ratio
Main plots			
Block	3	B_B	B_B/W_{MP}
Irrigation	2	B_I	B_I/W_{MP}
Main plot error	6	W_{MP}	
Sub-plots			
Manures	3	B_M	B_M/W_{SP}
Irrigation \times manures	6	B_{IM}	B_{IM}/W_{SP}
Sub-plot error	27	W_{SP}	
Total	47	T	

846 The subscripts are B for block, I for irrigation, M for manures, MP for main plot, SP
 847 for sub-plot, and MPE and SPE denote the main-plot and sub-plot errors.

848 **Table 6** Analysis of variance for n_1 treatments replicated n_2 times on plots in a com-
 849 plete randomized block design with n_3 measurements per plot

850

Source	Degrees of freedom	Mean squares	Paramaters estimated	F ratio
Between treatments	$n_1 - 1$	B	$\sigma_C^2 + n_3\sigma_W^2 + n_2n_3\sigma_B^2$	B/W
851 Between plots within treatments	$n_1(n_2 - 1)$	W	$\sigma_C^2 + n_3\sigma_W^2$	
Between cores within plots	$n_1n_2(n_3 - 1)$	C	σ_C^2	
Total	$n_1n_2n_3 - 1$	T		

852 **Figure captions**

- 853 1. An example lay-out of a completely randomized balanced experimental design in
854 which five replicates of each of four manurial treatments, M1, M2, M3 and M4,
855 are independently and randomly allocated to plots.
- 856 2. An example lay-out of a randomized blocked experimental design in which the
857 plots are grouped in blocks of four (separated by the dotted lines) and one repli-
858 cate of each of four manurial treatments, M1, M2, M3 and M4, is independently
859 and randomly allocated to a plot within each block. There are five blocks in
860 total, separated by dotted lines in the Figure.
- 861 3. An example layout of a split plot design with blocks. Three main plots are in
862 each block, and one replicate of each of three levels of an irrigation factor, I1,
863 I2 and I3, is independently and randomly allocated to a main plot within each
864 block. The three levels of the irrigation factor are distinguished in this figure by
865 dark grey, light grey or white shading. Within each main plot are four sub plots
866 and one replicate of each of four manurial treatments, M1, M2, M3 and M4, is
867 independently and randomly allocated to a sub plot within each main plot.
- 868 4. An example lay-out of the same completely randomized balanced experimental
869 design exemplified in Figure 1 with sites for collection of three soil cores (black
870 discs) independently and randomly located within each plot.

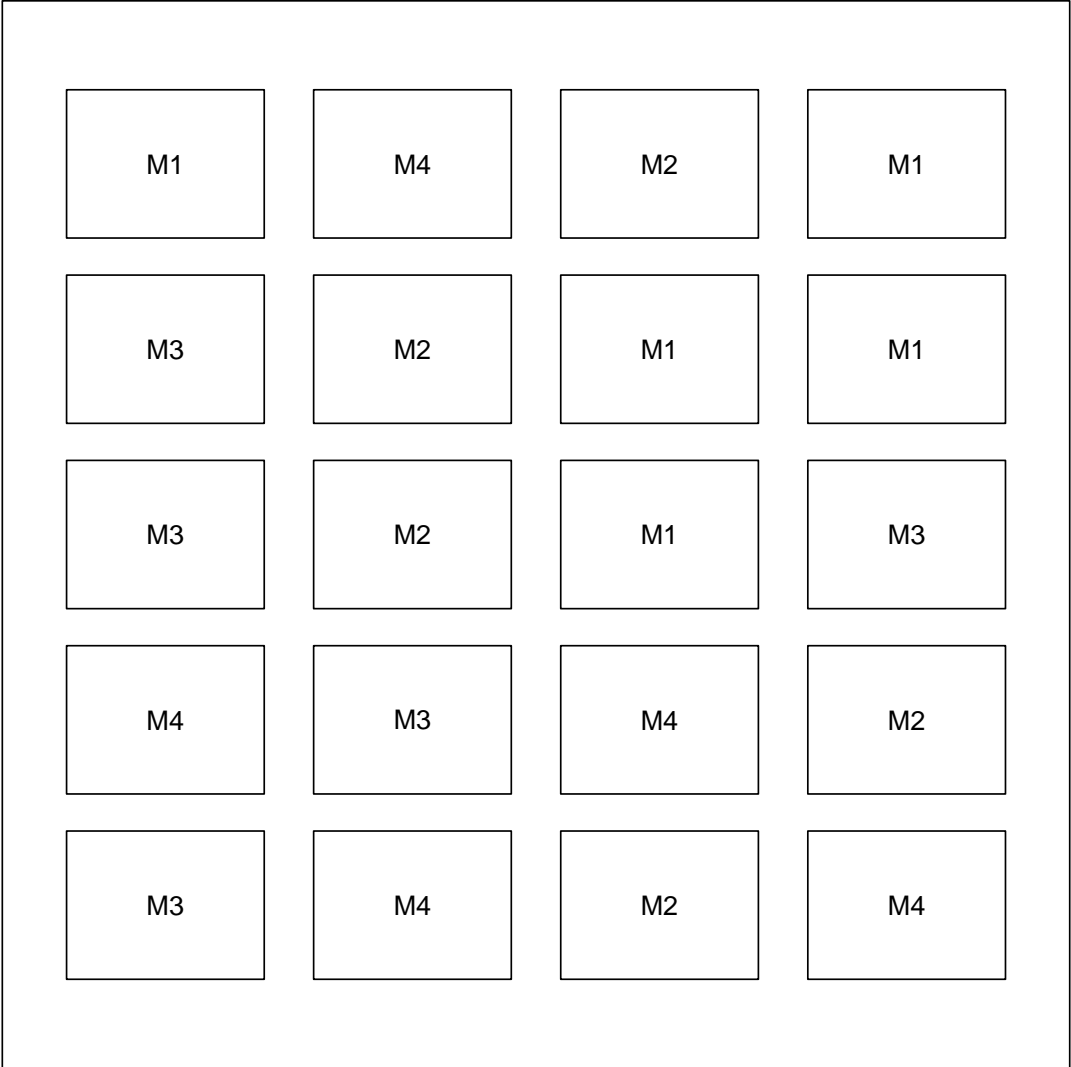


Figure 1: Fig 1

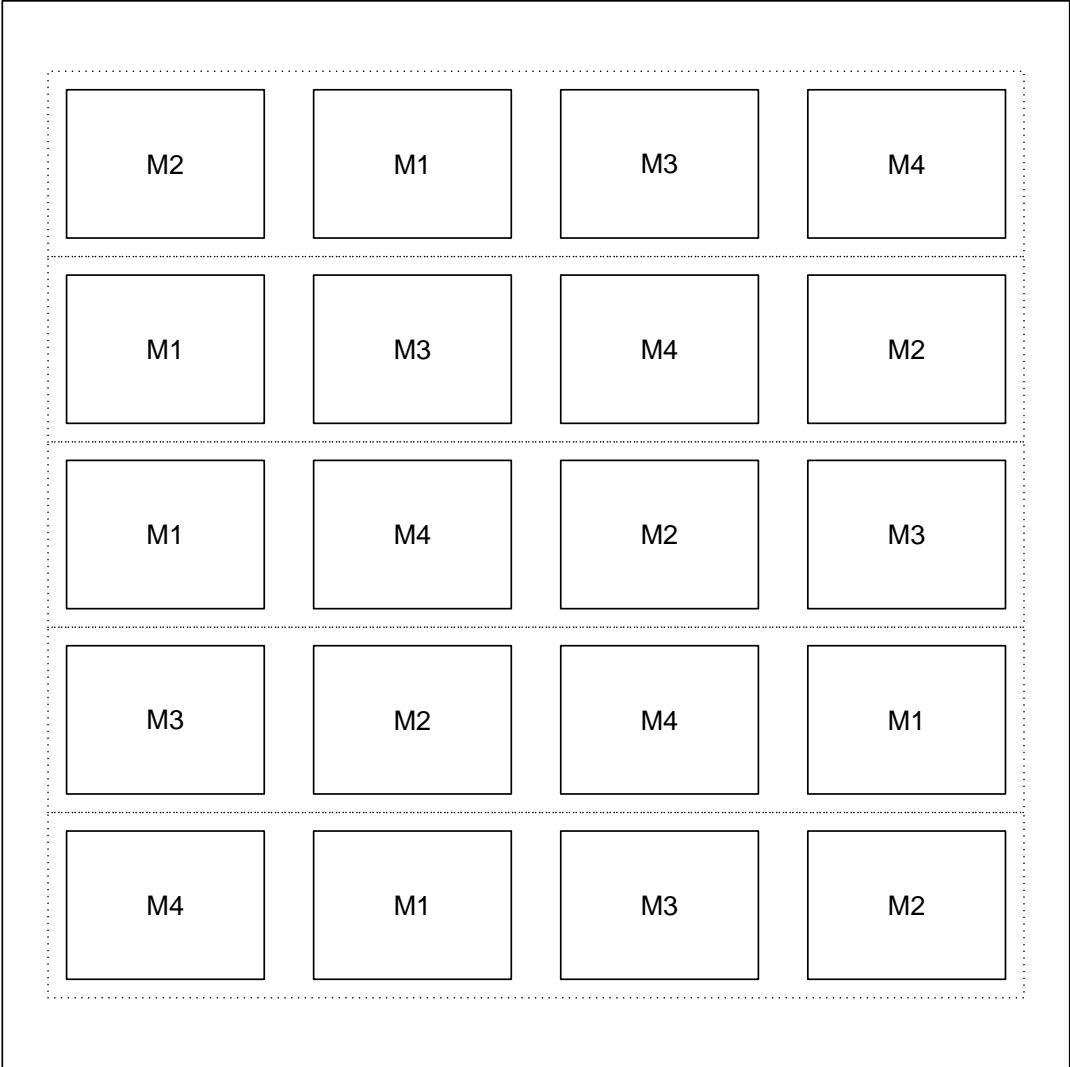


Figure 2: Fig 2

M3	M4	M3	M4	M2	M4
M1	M1	M1	M1	M4	M3
M4	M3	M4	M2	M3	M1
M2	M2	M2	M3	M1	M2
M4	M2	M1	M4	M3	M2
M1	M3	M2	M3	M4	M4
M2	M4	M4	M2	M1	M1
M3	M1	M3	M1	M2	M3

Figure 3: Fig 3

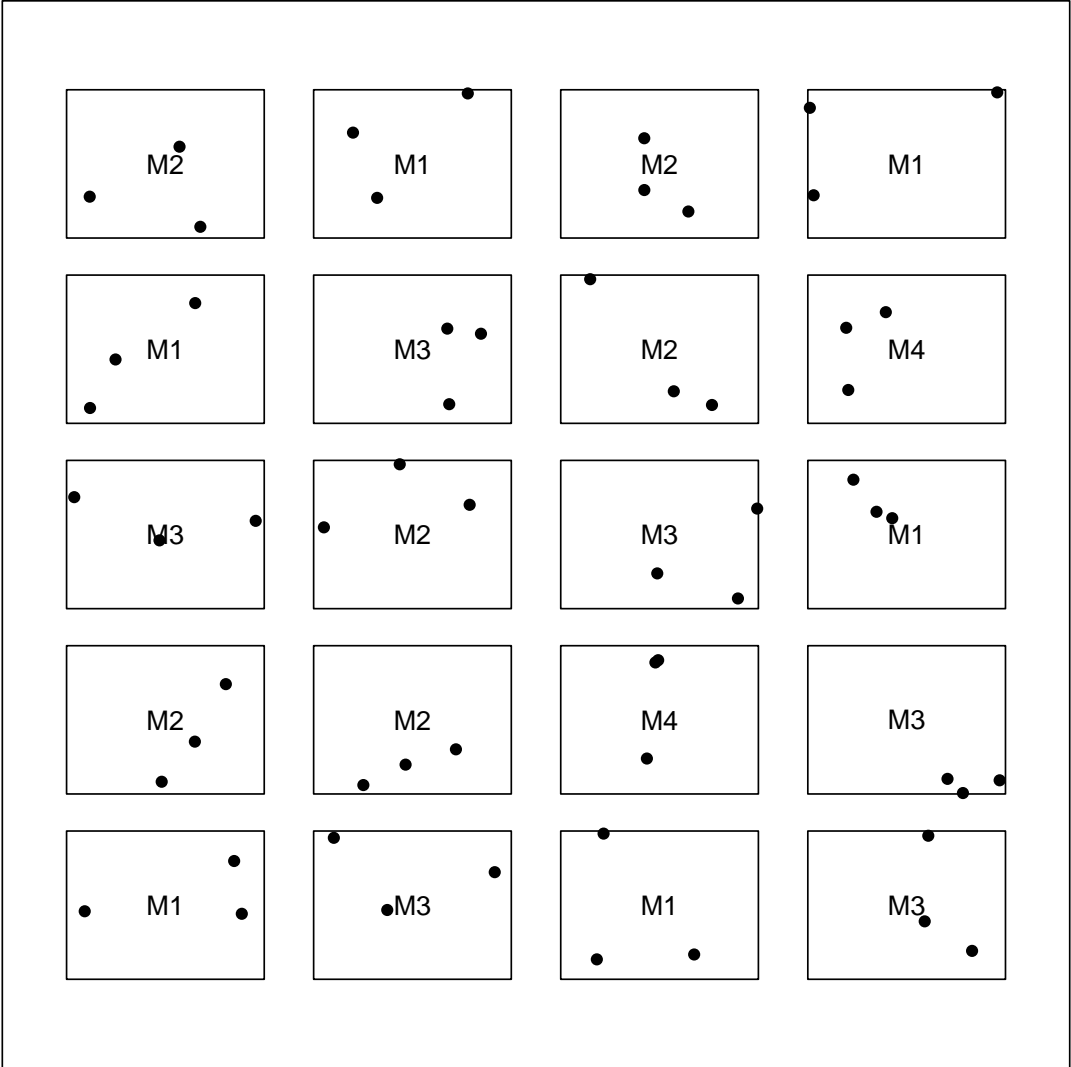


Figure 4: Fig 4