



NTNU – Trondheim
Norwegian University of
Science and Technology

Analysis of Censored Data from Split-Plot Design

Marte Nevland Hansen

Master of Science in Physics and Mathematics

Submission date: February 2015

Supervisor: John Sølve Tyssedal, MATH

Norwegian University of Science and Technology
Department of Mathematical Sciences

Analysis of Censored Data from Split-Plot Design

Marte Nevland Hansen

TMA4905 - Industrial Mathematics

2015

Master's Thesis

Abstract

In reliability theory, there are often data missing due to censoring. Such incomplete datasets are usually difficult to analyse. The exact value of the censored data is not known, but some information exists. That is, the value is higher than the censoring limit if the data is right censored, or lower for left censoring. Statistical analysis methods assume complete data, thus the censored data needs to be estimated. The missing values are replaced with fictional values, found by different methods, making the dataset a fictional complete dataset. To get good results from the analysis, the estimated values of the missing data should be as close to the original data as possible. In this thesis, the goal has been to analyse censored data from split-plot design. A design performed in split-plot manner induces correlation among observations. Two censoring methods have been tested; the quick and dirty method and the maximum likelihood method combined with multiple imputation. In the latter, the variance of the different parts of split-plot design were estimated, and then used to estimate the effects of the factors. Some of the factors, the ones that seem to be of less importance, must be removed for maximum likelihood and multiple imputation to create the variances. If done carefully, the analysis gives information of the factors with most influence. The performances of the methods are evaluated through three examples, and two different types of censoring, right and left censoring. Numerical results are obtained from implementations in the programming language R.

The results in this thesis show that both methods give good estimates for the effects of the factors. However, the quick and dirty method is not a safe method if there are many censored observations or a big gap between the censoring limit and the true value of the censored observations. The outcome of this thesis indicate that multiple imputation using the maximum likelihood estimator is the most accurate and safe method.

Sammendrag

I pålitelighetsteori er det ofte at data mangler på grunn av sensurering. Slike ufullstendige datasett er vanligvis vanskelige å analysere. Den nøyaktige verdien av sensurerte data er ikke kjent, men noe informasjon finnes. Det vil si at verdien er høyere enn sensureringsgrensen dersom dataene er høyre sensurert, eller lavere for venstre sensurering. Statistiske analysemetoder antar komplette data, dermed må sensurerte data estimeres. Manglende verdier er erstattet med fiktive verdier, funnet ved ulike metoder, slik at datasettet blir et fiktivt komplett datasett. For å få gode resultater fra analysen, bør de estimerte verdiene av de manglende observasjonene være så nær den opprinnelige verdien som mulig. I denne oppgaven har målet vært å analysere sensurerte data fra splitt-plott design. Et design på splitt-plott form inducerer korrelasjon mellom observasjoner. To sensurerings metoder har blitt testet; quick og dirty metoden og sannsynlighetsmaksimeringsmetoden kombinert med multippel imputering. I sistnevnte ble variansen av de forskjellige delene av splitt-plott designet beregnet, og deretter brukt til å estimere effektene av faktorene. Noen av faktorene, de av mindre betydning, fjernes fra sannsynlighetsmaksimeringsmetoden og multippel imputering for å skape avvikene. Hvis det gjøres nøye, gir analysen informasjon om faktorene med mest innflytelse. Metodene evalueres gjennom tre eksempler, og to forskjellige typer sensurering, høyre og venstre sensurering. Numeriske resultater oppnås fra implementeringer i programmeringsspråket R.

Resultatene i denne avhandlingen viser at begge metodene gir gode estimater for effektene av faktorene. Imidlertid er quick og dirty metoden ikke en trygg metode dersom det er mange sensurerte observasjoner eller et stort gap mellom sensureringsgrensen og den sanne verdien av de sensurerte observasjonene. Utfallet av denne avhandlingen tyder på at multippel imputering med maksimal sannsynlighetsestimator er den mest nøyaktige og trygge metoden.

Preface

This report is the product of my master's thesis at NTNU, Department of Mathematical Sciences.

I am very grateful to my supervisor Professor John S. Tyssedal, who has helped me throughout my master's thesis. He helped me through the academics behind the master and to see the light at the end of the tunnel.

Magnus, I sincerely appreciate all your help and support. You have been by my side throughout my thesis, and always lifted me up when I was down. I would also like to thank my parents, for supporting me through my years at the university, and Lejla Begluk and Roger André Søråa for being great friends and helping me out whenever I needed it.

Trondheim, February, 2015

Marte Nevland Hansen

Contents

1	Introduction	1
2	Theory	3
2.1	Linear regression model	3
2.2	Censoring	3
2.3	The maximum likelihood	4
2.4	The maximum likelihood for right censored data	5
2.5	The maximum likelihood for left censored data	5
2.6	Multiple imputation	6
2.7	Truncation	7
3	Split-plot design and the multivariate normal distribution	9
3.1	Experimental design	9
3.2	Split-plot design	11
3.3	Split-plot design with mirror image pairs	12
3.3.1	SPMIP - Half factorial design	12
3.3.2	SPMIP - Full factorial design	14
3.4	Analysis of SPMIP designs	15
3.5	The multivariate normal distribution	16
3.5.1	Conditional distribution	16
3.6	Simulation of multivariate truncated Gaussian distribution	17
4	The examples	19
4.1	Example I	19
4.2	Example II	22
4.3	Example III	25
5	The methods	29
5.1	Previous work	29
5.2	R software	30
5.3	Censoring with the maximum likelihood and multiple imputation method	31
5.4	The quick and dirty method	33
6	Experiments and results	35
6.1	Example I - Right censoring	35
6.2	Example I - Left censoring	39
6.3	Example II - Right censoring	41

6.4	Example II - Left censoring	44
6.5	Example III - Right censoring	47
6.6	Example III - Left censoring	50
7	Discussion	53
7.1	Example I	53
7.2	Example II	54
7.3	Example III	54
7.4	Censoring with maximum likelihood and multiple imputation	54
7.5	Censoring with quick and dirty	55
8	Conclusion	57
	Appendix A	59
	Appendix B	63

Chapter 1

Introduction

All industrial experiments are split-plot experiments.

This provocative remark has been attributed to the famous industrial statistician Cuthbert Daniel, by Box et al. (2005)^[1] in their book on design of experiments. Split-plot experiments were introduced by Fisher (1925) and their importance in industrial experimentation is highly recognized.

Experimental design helps create a design that assures gaining desired information. The interest is often focused on the effects of the process in an experimental design, where the design is constructed to figure out these effects and their contribution to the experiment. When performing experiments, resources are rarely unlimited, nor in amount of time or money. An experimental design in split-plot manner saves resources, by means of limiting the amount of runs necessary to conduct the experiment.

Methods for dealing with censoring in experimental design have been tested by Sue-Chu^[2] and Støtvig^[3], among others. The conclusion is that multiple imputation with maximum likelihood gives the best estimations. The quick and dirty approach is concluded unsafe, although it may give a pointer to which effects that have the most influence on a product. In this thesis, both of these methods are tested. The desired result is whether or not the most significant effects in an experiment can be found if the dataset has some censored data. The datasets used in calculations are (1) an experiment about the uniformity in a single-wafer plasma etching process, (2) modification of the surface characteristics of a security paper with plasma treatment from Bisgaard et. al. and (3) the well known Box and Jones' optimal formulation of a cake mix. These were chosen since the analyses are known, which makes it easy to compare the results in this thesis to the original estimates. The experiments are not typical censoring experiments, that is, the limits are set arbitrarily. However, the datasets are not too large, thus the censored values are easily found manually. The censoring of the datasets produces artificial censored datasets, since there is no natural way of censoring when doing this experiment. For example, for the cake mix experiment, one could say the taste of the cake was "off-the-charts", since the right censoring limit is set to 6, when the scale is 1-7.

Conditional distributions are assumed appropriate for missing data where the failure time is not observed. Different variances concerning the censored split-plots are estimated and used in the scaled truncation which is combined with multiple imputation. This creates estimates for the censored data. The methods were im-

plemented in the programming language R, using both own code and embedded functions available in R. The R package *lm* was used to estimate all the effects of the factors for both the original and the quick and dirty censored dataset. Censoring of the datasets was done with the R package *censReg*, and this was used to estimate the effects of the factors for the maximum likelihood and multiple imputation method. For the truncation, the package *truncnorm* and the function *rtruncnorm* were used.

Chapter 2 concerns the basic theory used in the procedure in this thesis, i.e. the linear regression model, different types of censoring and the maximum likelihood method for the exponential distribution are defined. Then, multiple imputation and truncation are introduced. Chapter 3 is devoted to split-plot design and split-plot design with mirror image pairs. The multivariate normal distribution and simulating dependent values are also included in this chapter. The three examples are presented in Chapter 4. In Chapter 5 follows the description of the software, as well as the two methods used, after a quick introduction to the previous work. The results are described in Chapter 6, followed by a discussion in Chapter 7 and a conclusion in Chapter 8.

The estimated coefficient of every factor for each example is listed in Appendix A. That is, the original estimate, the result of the calculations for the maximum likelihood with multiple imputation and the result of the quick and dirty method. Appendix B contains the code for Example II. The code used for deriving the results in the other examples is very similar to this, thus they are omitted.

Chapter 2

Theory

2.1 Linear regression model

A regression model is a statistical technique for modelling the relationship between a response variable and one, or more, explanatory variables. The response variable depends on the explanatory variables, thus it is called the dependent variable. Regression analysis estimates the regression function, which describes how the response variable is related to the explanatory variables. The regression variable is called the independent variable.

The regression of a random variable y on the variables \mathbf{x} , is the expectation of y given the values of \mathbf{x} , that is $E(y|\mathbf{x})$. The linear regression model is expressed as follows

$$E(y|\mathbf{x}) = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_kx_k + \epsilon, \quad (2.1)$$

where y is the response variable, x_1, x_2, \dots, x_k are the explanatory variables, $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ are the regression coefficients and ϵ is the random error. Here, the errors are usually assumed uncorrelated and distributed by $N(0, \sigma^2)$. The regression coefficients determines to what extent each explanatory variable contributes to the response. Most commonly, the least square method is used for estimating the unknown β 's.

2.2 Censoring

To test how well a product works, or its lifetime, experiments are run on several units of the product. Such experiments can not run forever, that is, there must be some limiting conditions, called censoring. Limiting conditions can be time, economical reasons, loss of an object due to withdrawal from the study and so on. Thus, there are different types of censoring.

Type I censoring

A sample of n units are tested in the interval from time zero, t_0 , until the experiment is stopped at time t_k . Failures, or experiences, after time t_k are not observed. The experiencing of the event is random, but the total duration of the experiment is fixed.

Type II censoring

A sample of n units are observed until failure of the first r units. The r is predetermined, such that $r \leq n$. Since the experience of the event is random, the duration of the experiment is also random.

Right, left and interval censoring

Censoring is divided into three main categories. **Right censoring** occurs when there are still functioning units after the experiment is terminated. These are omitted from the analysis, i.e. censored. If some units have failed before the start of the experiment, they become **left censored**. When the censored data points lie between two values, i.e. the observed data lie outside this interval, the data is **interval censored**. In this thesis, right censoring and left censoring are considered.

2.3 The maximum likelihood

The maximum likelihood method^[2] is a method for estimating the parameters of a statistical model. It consists of maximizing the likelihood function. The likelihood function is the joint density of the independent random variables taken from a probability distribution. When estimating the likelihood, the log-likelihood is often maximized, finding the estimates for the parameters when the derivative of the log-likelihood function is set to zero.

As an example, let T_1, T_2, \dots, T_n be n independent random variables from the probability distribution $f(\mathbf{t}, \theta)$, where θ is a single parameter of the distribution. The likelihood function is given by

$$L(\theta) = \prod_{i=1}^n f(t_i, \theta).$$

If the T_i 's are exponentially distributed random variables with probability density function

$$f(t, \theta) = \frac{1}{\theta} e^{-\frac{t}{\theta}}.$$

With n observations, the likelihood function is given by

$$L(\theta) = \prod_{i=1}^n \frac{1}{\theta} e^{-\frac{t_i}{\theta}} = \frac{1}{\theta^n} \prod_{i=1}^n e^{-\frac{t_i}{\theta}}.$$

Taking the natural logarithm gives the log-likelihood,

$$l(\theta) = \ln L(\theta) = -n \ln(\theta) - \frac{\sum_{i=1}^n t_i}{\theta}.$$

Further, the derivative with respect to θ is set equal to zero,

$$\frac{\partial l(\theta)}{\partial \theta} = -\frac{n}{\theta} + \frac{\sum_{i=1}^n t_i}{\theta^2} = 0.$$

Then $\hat{\theta}$, the estimate of θ , is found,

$$\hat{\theta} = \frac{\sum_{i=1}^n t_i}{n}.$$

2.4 The maximum likelihood for right censored data

When dealing with censored data, the maximum likelihood has to be altered. Let $f(\mathbf{t}, \theta)$ denote the probability density function, $F(\mathbf{t}, \theta)$ the distribution function and $S(\mathbf{t}, \theta)$ the survival function. The probability that a unit survives the time interval $(0, t)$ is defined by

$$S(t_i, \theta) = P(T > t_i) = \int_{t_i}^{\infty} f(u, \theta) du = F(\infty, \theta) - F(t_i, \theta) = 1 - F(t_i, \theta).$$

Assume that n units are tested, and r units fail in the time interval. Let the lifetime and censoring be given as (Y_i, δ_i) , where

$$Y_i = \begin{cases} T_i, & \delta_i = 1 & \text{for uncensored data} \\ \min(T_i, C_i), & \delta_i = 0 & \text{for right censored data,} \end{cases}$$

where C_i is the censored time. If a unit fails at τ_i , the contribution to the likelihood is the density at the duration; $L_i = f(\tau_i, \theta)$. If a unit is still functioning, the lifetime exceeds τ_i ; $L_i = S(\tau_i)$. The likelihood can be written as follows,

$$L(\theta) = \prod_{i=1}^n L_i(\theta) = \prod_{\delta_i=1} f(\tau_i, \theta) \prod_{\delta_i=0} S(\tau_i, \theta) = \prod_{i=1}^r f(\tau_i)^{\delta_i} \prod_{i=r+1}^n S(\tau_i)^{1-\delta_i}.$$

2.5 The maximum likelihood for left censored data

The probability for the left censored observation is

$$P(T \leq t_i) = F(t_i, \theta) - F(-\infty, \theta) = F(t_i, \theta),$$

where the cumulative distribution function is

$$F(t_i, \theta) = 1 - S(t_i, \theta).$$

Say r units have failed, where some units started before the study began. The likelihood function is then defined as

$$L(\theta) = \prod_{\delta_i=1} f(t_i, \theta) \prod_{\delta_i=0} F(t_i, \theta),$$

where

$$\delta_i = \begin{cases} 1 & \text{for complete observations} \\ 0 & \text{for left censored observations,} \end{cases}$$

which is equivalent to

$$L(\theta) = \prod_{i=1}^r f(t_i, \theta) \prod_{i=r+1}^n F(t_i, \theta).$$

2.6 Multiple imputation

Many datasets are not complete. They miss some values, for example due to errors occurring while collecting them, or there was no value to observe at some points. This creates problems and limitations for analysis. Imputation is used to fill in missing data with credible data. Multiple imputation was proposed by Rubin^[4], where missing values are replaced by m imputed values to create a complete dataset. Each complete dataset is then analysed by standard procedures, and the results are combined to produce estimates. The m imputed values are drawn from a truncated distribution. It is a Monte Carlo technique and the missing values are replaced by $m > 1$ simulated values. A disadvantage of the multiple imputation is that it requires more work in both implementation and analysing the results. In this thesis, the number of imputations is set to 5.

The method for repeated-imputation inference has the following procedure:

A generic scalar quantity Q is to be estimated. The Q can, for example, represent the mean, correlation or odds ratio. Let Y denote the data. The data is split into two parts; the observed data, Y_{obs} , and missing data Y_{mis} . As if complete data were available, let $\hat{Q} = \hat{Q}(Y_{obs}, Y_{mis})$ denote the statistic to estimate Q . Also, let $\sigma^2 = \sigma^2(Y_{obs}, Y_{mis})$ denote the squared standard error. Thus, the normal approximation

$$\frac{\hat{Q} - Q}{\sqrt{\sigma^2}} \sim N(0, 1)$$

is appropriate when dealing with complete data. The Y_{mis} does not have any data. Suppose $m > 1$ independent simulated imputations $Y_{mis}^{(1)}, \dots, Y_{mis}^{(m)}$ are conducted. The imputed data estimates $\hat{Q}^{(k)} = \hat{Q}(Y_{obs}, Y_{mis}^{(k)})$ and their estimated variances $\sigma^{2(k)} = \sigma^2(Y_{obs}, Y_{mis}^{(k)})$, for $k = 1, \dots, m$ are calculated. The overall estimate of Q is then the average

$$\bar{Q} = m^{-1} \sum \hat{Q}^{(k)}.$$

The standard error for \bar{Q} can be found when calculating the between-imputation variance $V_b = (m - 1)^{-1} \sum (\hat{Q}^{(k)} - \bar{Q})^2$ and the within-imputation variance $V_i = m^{-1} \sum \sigma^{2(k)}$. The estimated total variance is

$$V_T = (1 + m^{-1})V_b + V_i,$$

where tests and confidence intervals are based on the t-approximation

$$\frac{\hat{Q} - Q}{\sqrt{V_T}} \sim t_\nu,$$

with degrees of freedom

$$\nu = (m - 1) \left[1 + \frac{V_i}{1 + m^{-1}V_b} \right]^2.$$

The V_T will reduce to V_i if Y_{mis} carries no information about Q , given that the imputed data estimated $\hat{Q}^{(k)}$ is identical. Thus, the relative increase in variance provoked by missing data is $r = (1 - m^{-1})V_b/V_i$. The rate of missing information in the system is $\lambda r/(1 + r)$, which combined with the equations above gives

$$\lambda = \frac{r + 2/(\nu + 3)}{1 + r}.$$

Multiple imputation is a simple and very general method that can be implemented to any data. The validity of this method is dependent on how the imputations $Y_{mis}^{(1)}, \dots, Y_{mis}^{(m)}$ are generated. If the imputations are created arbitrarily, it is not likely to obtain valid inferences in general. The imputations should on average give reasonable values for the missing data, and the variance should be within an appropriate degree of uncertainty.

Single imputation is not used in this thesis, since this method only considers one estimation. Doing multiple estimations and taking the mean of the results, seems more appropriate to get a more accurate estimation. This is based on the work of Sue-Chu^[2] and Støtvig^[3], where both concludes that multiple imputation with the maximum likelihood method produces the best estimates for censored data in experimental design.

2.7 Truncation

Truncation is described in Sue-Chu^[2]. In mathematics, truncation limits the number of digits in a number by discarding the least significant decimals. Statistical truncation refers to measurements that have been cut off at some value. When restricting the domain of the probability distribution, a truncated distribution is created. The cut of the domain creates a truncated sample. In this thesis, truncation is used to restrict the possible values for the estimates of the censored observations, created by multiple imputation.

As stated, when truncation is applied to a probability distribution, it leads to a new distribution. Let X be a random variable with distribution function $F(x)$, and let Y be a new random variable having the distribution of X truncated to the semi-open interval $(a, b]$. Thus Y has the distribution function

$$F_Y(y) = \begin{cases} 0, & : y \leq a, \\ \frac{F(y) - F(a)}{F(b) - F(a)} & : a < y \leq b, \\ 1 & : y > b. \end{cases}$$

Scaled truncation, combined with multiple imputation, can generate X for a potentially censored value. After restricting the domain of the probability function, the probability density of the random variable is needed. Let $y = (a, b]$ be the restricted domain. Then

$$f(y|a < Y \leq b) = \frac{g(y)}{F(b) - F(a)},$$

where

$$g(y) = \begin{cases} f(y) & : a < y \leq b, \\ 0 & : \text{otherwise.} \end{cases}$$

The truncated distribution with right censoring at a will then be

$$f(y|Y > a) = \frac{g(y)}{1 - F(y)},$$

where $g(y) = f(x)$ for $a < y$ and $g(x) = 0$ otherwise.

Chapter 3

Split-plot design and the multivariate normal distribution

3.1 Experimental design

Experimental design allows us to figure out how the response, or the output, responds when the settings of the input variables in a system are intentionally changed. Through an experiment, an investigator learns how the input variables affect the performance of a system, which provides a basis for choosing the optimal input settings. The motivation behind performing an experiment is often to identify significant factors. When performing an experiment, the factors are the input, i.e. the explanatory variables of a regression model. The response is the desired outcome. The levels describe the amount of magnitude of each factor in the different combinations, and for a two-level experiment, they are usually denoted as "high" and "low". An experiment considering f factors and l levels is expressed as a l^f factorial design.

When the number of factors increases, the number of runs in the experiment also increases. One way to reduce this number, is to choose a fraction of the total runs, to be used in the estimation. This selection is preferably chosen such that the main effects and the lower order interactions can be estimated, thus the higher order interactions are assumed negligible. This procedure is called a fractional factorial design. In the case of a two-level fractional factorial design, the notation becomes 2^{f-g} , where g is the number of generators. The fraction is denoted by $2^{-g} = \frac{1}{2^g}$. If the experimental design of levels in a factor equals the design of an interaction between other factors, the factor is said to be a generator of the design.

As an illustration on experimental design, say that an experiment with three important factors is investigated at two levels. Table 3.1 shows the design of the experiment, with factors A, B and C, and the response for each run of the experiment.

Table 3.1: A 2^3 experimental design.

Run no.	A	B	C	Value
1	-	-	-	y_1
2	+	-	-	y_2
3	-	+	-	y_3
4	-	-	+	y_4
5	+	+	-	y_5
6	+	-	+	y_6
7	-	+	+	y_7
8	+	+	+	y_8

An estimate for the main effect of A is found by taking the difference between the mean response at the high level and the mean response at the low level of the factor.

$$A = \frac{y_2 + y_5 + y_6 + y_8}{4} - \frac{y_1 + y_3 + y_4 + y_7}{4}.$$

When performing an experiment with more than one factor, interactions between two or more factors should be investigated. Interaction means that the effects of one factor may depend on the level of other factors. An estimate for the two factor interaction between A and B is found by adding the positive combinations of A and B and taking the average, before subtracting the average of the negative combinations of the two factors. From the design in Table 3.1, this interaction can be found by

$$AB = \frac{y_1 + y_4 + y_5 + y_8}{4} - \frac{y_2 + y_3 + y_6 + y_7}{4}.$$

Calculating the effect of the other factors and higher order interactions can be done by similar procedures.

Table 3.2 shows a 2^{4-1} experimental design. In this case there are four factors, but one of them is set to be a generator. Thus, the level of this factor, in each run, is decided by the levels of other factors. Here, $D = ABC$, that is, the level of factor D is based on the interaction of all the other factors.

Table 3.2: A 2^{4-1} experimental design.

Run no.	A	B	C	D
1	-	-	-	-
2	+	-	-	+
3	-	+	-	+
4	-	-	+	+
5	+	+	-	-
6	+	-	+	-
7	-	+	+	-
8	+	+	+	+

Furthermore, available degrees of freedom are used to estimate effects, i.e. the error cannot be estimated. Thus, normal probability plot is used to evaluate the effects.

3.2 Split-plot design

Split-plot designs are described in Box and Jones^[5], and are used for process and product optimization. Typically some of the factors are hard to change, that is, a complete randomization of the experiments is difficult, if not impossible, to conduct. Split-plotting can also occur (1) when two or more process steps are involved, (2) in robust product design experimentation and (3) when it is of interest to estimate some factors with higher precision than others, where the latter is then handled as subplot factors.

The subplot factors are easy to change, and are changed according to a design matrix called subplot design. The whole-plot factors are hard to change, thus they are changed less frequently according to a second design matrix called whole-plot design. For a randomly chosen level combination of the whole-plot factors, a design in the subplot factors is run in random order. That is, the total number of runs is the number of whole-plot level combinations times the number of runs in the design for the subplot factors.

When designing an experiment, it is crucial that the number of runs is at an acceptable level. In order to achieve this economy in the process, a two level experimental plan is often used.

The linear statistical model^[6] for the basic split-plot design in which observations are taken on s split-plots in the i 'th whole-plot can be written as

$$y_{ij} = \sum_{k=1}^m x_{ijk}\beta_k + u_{ij}, \quad j = 1, 2, \dots, s, \quad i = 1, 2, \dots, n,$$

i.e. an extension of eq. (2.1). Here, y_{ij} is the observed response value, the x_{ijk} are the m different control variables, the β_k are the m fixed unknown parameters, and u_{ij} is the unobservable random error. These errors consist of two components, a random element associated with the i 'th whole-plot, say ϵ_i^w , and a second independent random element associated with the j 'th subplot in the i 'th whole-plot, say ϵ_{ij}^s , i.e. $e_{ij} = \epsilon_i^w + \epsilon_{ij}^s$. The ϵ_i^w and ϵ_{ij}^s are assumed to be iid with zero mean and variances $\sigma_w^2 \geq 0$, $\sigma_s^2 > 0$ respectively. These assumptions imply

$$E_{e_{ij}e_{i'j'}} = \begin{cases} \sigma_w^2 + \sigma_s^2 & \text{if } i = i' \text{ and } j = j', \\ \sigma_w^2 & \text{if } i = i' \text{ and } j \neq j', \\ 0 & \text{if } i \neq i'. \end{cases}$$

The observations within each whole-plot is correlated. Thus, the analysis of split-plot design is generally based on the generalized least squares method.

The whole-plot effects contain a whole-plot error. An important characteristic of split-plot arrangements is that the subplot effects, and all their interactions with the whole-plot effects, are estimated with the same smaller subplot error. It follows that if the data from a split-plot arrangement are analysed graphically, two separated plots are needed.

If one is to investigate a two-level experiment, it is convenient to present "low" by a negative sign (-) and let a positive sign (+) represent "high". Then orthogonal factor columns are obtained and the coefficients are easily computed.

3.3 Split-plot design with mirror image pairs

Each split-plot dataset is divided into two parts, as described in Tyssedal and Kulachi^[7]. One property of the split-plot design with mirror image pairs (SPMIP) is that it divides the estimated effect into two orthogonal subspaces, separating subplot main effects and subplot by whole-plot interactions from the rest.

SPMIP designs have a design matrix that can be written as follows

$$\begin{bmatrix} \mathbf{W} & \mathbf{S} \\ \mathbf{W} & -\mathbf{S} \end{bmatrix},$$

where $\begin{bmatrix} \mathbf{W} \\ \mathbf{W} \end{bmatrix}$ contains the whole-plot factors and $\begin{bmatrix} \mathbf{S} \\ -\mathbf{S} \end{bmatrix}$ the subplot factors.

When using mirror image pairs, it is possible to construct half factorial design matrices; one with the mean, and the other with the difference between two response observations. The mean provides information about the whole-plot effects and possibly subplot by subplot interactions, and the difference provides information about the subplot effects and interactions between subplot and whole-plot effects. For the full factorial design, with four subplots per whole-plot, there will be three matrices, the whole-plot is the mean of all observations for each whole-plot combination, and the subplot effects are now separated in two matrices by different combinations within the whole-plot.

In this thesis, three types of examples are used. The first is a half factorial design, where only the combinations that give a high level are present, i.e. the interaction ABCDE is positive for all the experiments. There are two subplots per whole-plot in this design. The second is a 2^5 split-plot experiment with two subplot per whole-plot. For these first two cases, half factorial design method, Section 3.3.1, is used in the calculations. For the third example, there are four subplots per whole-plot, thus full factorial design method, Section 3.3.2, is appropriate.

3.3.1 SPMIP - Half factorial design

A cup-cake tray producer wants to find the best recipe to use in the cup-cake trays. The desired size of the cupcakes has already been found, and each cup in the tray is filled accordingly. The batter is a factory finished batter, where one just adds egg and water. The producer also wants to include baking cocoa. These three ingredients, [A, B, C], are the hard-to-change factors, while time and temperature, [D, E], are the easy-to-change factors. The batter will be made in large batches. The cupcakes are then given to a class in primary school, and rated by them, from 1 to 10, where 10 dictates the best cupcakes. The design of the experiment and the response are shown in Table 3.3.

To make a half factorial design, half of the data must be removed. In Table 3.3, the subplot factor E is set as a generator, $E=ABCD$, and only the runs that give a high level combination of the interaction between the five factors, $ABCDE=I$, are used in the half factorial design. There are two split-plot observations, y_{i1} and y_{i2} , for each whole-plot, i.e. $i = 1, 2, \dots, 8$. If the y_{i1} and y_{i2} are the response from a

Table 3.3: The half factorial design of the cup-cake example.

Recepie	A	B	C	E: + -	
				D: + -	
(1)	+	+	+	7	9
(2)	+	-	-	1	3
(3)	-	+	-	7	10
(4)	-	-	+	3	6
				E: + -	
				D: - +	
(5)	+	+	-	8	4
(6)	+	-	+	4	9
(7)	-	+	+	9	1
(8)	-	-	-	6	8

split-plot with mirror image pairs, and

$$Z_i^w = \frac{y_{i1} + y_{i2}}{2} \qquad Z_i^s = \frac{y_{i1} - y_{i2}}{2}, \qquad (3.1)$$

then y_{i1} is a function of the contributions from the whole-plot factors, the whole-plot noise, the main effect of the subplot factors and their interaction, the subplot noise and the interaction between subplot and whole-plot. The y_{i2} is the same function of the whole-plot factors and the whole-plot noise, but the subplot main effects and the interactions between the subplot and the whole-plot effects have the opposite sign and cancels out when they are added. The subplot interactions have the same sign in y_{i1} and y_{i2} and are therefore not cancelled out.

The two matrices are made from the complete dataset with the above equations, eq. (3.1), where the system of equations in Z_i^w gives the whole-plot matrix, \mathbf{W} , and Z_i^s gives the subplot matrix, \mathbf{S} . The level of each factor for the different runs will, just like the response, be put into the two equations. This tells which factors are whole-plot effects, and which factors are subplot effects, since they otherwise cancel out. The signs for each interaction column are derived by entry-wise multiplication of the signs of the constituent main effects. Linear models can be used to estimate the effect for both the whole-plot analysis and the subplot analysis, with subplot by whole-plot interactions included.

When the mirror image pairs, eq. (3.1), are applied to a dataset with two subplots per whole-plot, the number of rows in the matrices are halved. The whole-plot matrix is shown in Table 3.4, and Table 3.5 shows the subplot matrix. The matrices contain a system of linear equations which can be analysed to obtain information of the factors.

Table 3.4: The whole-plot matrix for the cup-cake example.

A	B	C	AB	AC	BC	DE	Response
+	+	+	+	+	+	+	8.0
+	-	-	-	-	+	+	2.0
-	+	-	-	+	-	+	8.5
-	-	+	+	-	-	+	4.5
+	+	-	+	-	-	-	6.0
+	-	+	-	+	-	-	6.5
-	+	+	-	-	+	-	5.0
-	-	-	+	+	+	-	7.0

Table 3.5: The subplot matrix for the cup-cake example.

D	E	AD	BD	CD	AE	BE	CE	Response
+	+	+	+	+	+	+	+	-1.0
+	+	+	-	-	+	-	-	-1.0
+	+	-	+	-	-	+	-	-1.5
+	+	-	-	+	-	-	+	-1.5
-	+	-	-	+	+	+	-	2.0
-	+	-	+	-	+	-	+	-2.5
-	+	+	-	-	-	+	+	4.0
-	+	+	+	+	-	-	-	-1.0

3.3.2 SPMIP - Full factorial design

Lets consider a 2^5 split-plot design with four subplots per whole-plot. Table 3.6 shows the setup of the subplots within one whole-plot. Here the level, denoted high or low, shows which combination that is considered high, ABCDE=I, and which is considered low, ABCDE=-I.

Table 3.6: Table.

D	E	Level
-	-	y_{i1} high
+	+	y_{i2} high
+	-	y_{i3} low
-	+	y_{i4} low

When applying the mirror image pairs in this case, the following equations are used,

$$\begin{aligned}
 Z_{i1}^w &= \frac{y_{i1} + y_{i2}}{2} & Z_{i2}^w &= \frac{y_{i3} + y_{i4}}{2} \\
 Z_{i1}^s &= \frac{y_{i1} - y_{i2}}{2} & Z_{i2}^s &= \frac{y_{i3} - y_{i4}}{2}.
 \end{aligned} \tag{3.2}$$

The whole-plot effects for the full factorial design are found by taking the mean of the four set-ups for each combination of the factors, that is, finding $\frac{Z_{i1}^w + Z_{i2}^w}{2}$ from eq. (3.2). The system of equations gained by this are combined to a whole-plot matrix, here denoted as \mathbf{W} . The subplot effects are divided in two matrices, \mathbf{S}^+ and \mathbf{S}^- . To find the effects in \mathbf{S}^+ , the system of equations from $\frac{Z_{i1}^w - Z_{i2}^w}{2}$ are used. Matrix \mathbf{S}^+ contains whole-plot by subplot interactions. The information in \mathbf{S}^- is gathered from $\frac{Z_{i1}^s + Z_{i2}^s}{2}$ and $\frac{Z_{i1}^s - Z_{i2}^s}{2}$, where both contain one subplot factor, interactions between whole-plot and subplot, and whole-plot interactions by subplot. Both of the subplot matrices contain subplot factors with the same error, which allows them to be plotted together in the same plot.

Full factorial experiments can be expensive and time-consuming, with f factors, l levels and R replications, the number of testes to be performed is Rl^f . In a fractional factorial experiment, some test combinations are eliminated. This means some information is lost, but if the experiment is planned well, only the effects that are believed to be unimportant are removed. Then a compromise between total information, experiment costs and experimental value is made. Higher order interactions are unlikely to have engineering meaning or to show statistical significance. Thus the full factorial experiment can give information that is not meaningful.

3.4 Analysis of SPMIP designs

In the following, w and s are used to denote whole-plot and subplot main effects respectively. Due to the way split-plot experiments are executed, there are two variance regimes. The form of the covariance matrix of the responses is given as

$$\mathbf{V} = \begin{bmatrix} \mathbf{C} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{C} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{C} \end{bmatrix},$$

here \mathbf{C} corresponds to each whole-plot and is a symmetric matrix on the form

$$\mathbf{C} = \begin{bmatrix} \sigma_w^2 + \sigma_s^2 & \sigma_w^2 & \dots & \sigma_w^2 \\ \sigma_w^2 & \sigma_w^2 + \sigma_s^2 & \dots & \sigma_w^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_w^2 & \sigma_w^2 & \dots & \sigma_w^2 + \sigma_s^2 \end{bmatrix},$$

where σ_w^2 and σ_s^2 are the between whole-plots and within whole-plot variances.

The two responses y_{i1} and y_{i2} have a common part that consists of a possible constant, whole-plot effects, whole-plot by whole-plot and subplot by subplot interactions and a whole-plot error ϵ_i^w . The part that differs consists of subplot main effects and whole-plot by subplot interactions and the subplot errors, ϵ_{i1}^s and ϵ_{i2}^s .

Consider an experiment that fits with the half factorial design method. Let the error part in Z_i^w be denoted by $u_i^w = \epsilon_i^w + \frac{\epsilon_{i1}^s + \epsilon_{i2}^s}{2}$, and the error part in Z_i^s denoted

by $u_i^s = \frac{\epsilon_{i1}^s - \epsilon_{i2}^s}{2}$. Then, Z_i^w , where $i = 1, 2, \dots, n/2$, are all independent with the same variance. This also applies to Z_i^s .

3.5 The multivariate normal distribution

Every observation in a split-plot experiment is dependent on one or more of the other observations. If each whole-plot contains 2^b subplots, where b is the number of subplot factors, then each observation is dependent on $2^b - 1$ observations. In the case of censoring, the censored value can be estimated by means of the observations it depends on. Thus, the multivariate normal distribution can be used to estimate the censored values. Multivariate analysis is described in Rencher and Christensen^[8].

If a random variable y , with mean μ and variance σ^2 , is normally distributed its density is given by

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\mu)^2/2\sigma^2}, \quad -\infty \leq y \leq \infty,$$

for the univariate normal distribution case. The density for the multivariate normal distribution case is similar. If \mathbf{y} has a multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, the density is given by

$$g(\mathbf{y}) = \frac{1}{(\sqrt{2\pi})^p |\boldsymbol{\Sigma}|^{1/2}} e^{-(\mathbf{y}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y}-\boldsymbol{\mu})/2}, \quad (3.3)$$

where p is the number of variables. When \mathbf{y} has density eq. (3.3), \mathbf{y} is distributed as $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

3.5.1 Conditional distribution

Let the observation vector be partitioned into two subvectors denoted by \mathbf{y} and \mathbf{x} , where \mathbf{y} is $n \times 1$ and \mathbf{x} is $m \times 1$. Then the expectation and covariance matrix become

$$E \begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_x \end{pmatrix}, \quad \text{cov} \begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix} = \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix},$$

i.e. $\begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix}$ is

$$N_{n+m} \left[\begin{pmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_x \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{pmatrix} \right].$$

If \mathbf{y} and \mathbf{x} are dependent, $\boldsymbol{\Sigma}_{yx} \neq \mathbf{0}$, the conditional distribution of \mathbf{y} given \mathbf{x} , $f(\mathbf{y}|\mathbf{x})$, is multivariate normal with

$$E(\mathbf{y}|\mathbf{x}) = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} (\mathbf{x} - \boldsymbol{\mu}_x), \quad (3.4)$$

$$\text{cov}(\mathbf{y}|\mathbf{x}) = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}.$$

Note that $E(\mathbf{y}|\mathbf{x})$ is a vector of linear functions of \mathbf{x} , and $\text{cov}(\mathbf{y}|\mathbf{x})$ is a matrix that is independent of \mathbf{x} .

In this thesis, \mathbf{y} are the censored observations and \mathbf{x} are the observed values that \mathbf{y} depends on. One special case of this distribution is included in this section.

Conditional distribution on the bivariate case

Let y_1 and y_2 be dependent. The conditional distribution of y_1 given by y_2 , $f(y_1|y_2)$, is then multivariate normal with

$$\begin{aligned} E(y_1|y_2) &= \mu_{y_1} + \Sigma_{y_1y_2}\Sigma_{y_2y_2}^{-1}(y_2 - \mu_{y_2}), \\ \text{cov}(y_1|y_2) &= \Sigma_{y_1y_1} - \Sigma_{y_1y_2}\Sigma_{y_2y_2}^{-1}\Sigma_{y_2y_1}. \end{aligned}$$

3.6 Simulation of multivariate truncated Gaussian distribution

When performing multiple imputation for censored data, there is a need to generate data from truncated distributions. A plausible scenario that can occur is that every observation within a whole-plot combination is censored. Chopin^[9] shows how to simulate such values.

Let $X = (X_1, \dots, X_d)$ be a d -dimensional Gaussian vector with mean μ and covariance matrix Σ , and let $[a_i, b_i]$ be d intervals, where b_i may be either a real number or ∞ . The distribution of X , conditional on the event that $X_i \in [a_i, b_i]$, $i = 1, \dots, d$, is usually called a truncated Gaussian distribution.

Truncation in the bi-dimensional case with semi-finite intervals

Consider the simulation of $X = (X_1, X_2) \sim N_2(\mu, \Sigma)$, subject to $X_1 \geq a_1$ and $X_2 \geq a_2$, i.e. for some truncation points a_1 and a_2 . Without loss of generality, set $\mu = (0, 0)^T$, $\Sigma = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$, and assume that $a_1 \geq a_2$; if necessary, swap components to impose the last condition. The joint density of the considered truncated density is, up to a constant:

$$p(x_1, x_2) \propto \exp\left\{-\frac{1}{2\nu^2}(x_1^2 + x_2^2 - 2\rho x_1 x_2)\right\} \times I(x_1 \geq a_1; x_2 \geq a_2), \quad (3.5)$$

where $\nu^2 = 1 - \rho^2$. The conditional distribution of $X_2|X_1 = x_1$ is a univariate Gaussian $N(\rho x_1, \nu^2)$ truncated to $X_2 \geq a_2$, which is denoted $TN_{[a_2, \infty)}(\rho x_1, \nu^2)$. The marginal density of X_1 is

$$\rho(x_1) \propto \varphi(x_1)\Phi\left(\frac{\rho x_1 - a_2}{\nu}\right)I(x_1 \geq a_1).$$

Here φ is the unit Gaussian probability density function, $\varphi(x) = \exp(-x^2/2)/\sqrt{2\pi}$, and Φ is the cumulative distribution function.

To derive a proposal distribution for eq. (3.5), $\Phi(\cdot)$ is derived with a simpler expression derived from the two following straightforward inequalities:

$$\frac{1}{2} \leq \Phi(x) \leq 1 \quad \text{for } x \geq 0,$$

$$\Phi(x) \leq c(x_0)\varphi(x) \quad \text{for } x \leq x_0 \leq 0,$$

where $c(x_0) = \min(\sqrt{\pi/2}, -1/x_0)$, for $x_0 < 0$, $c(0) = \sqrt{\pi/2}$. In split-plot experiments, the ρ will always be positive. There are then two relevant cases, S^+ and M^+ , for estimating the censored values. Here 'S' stands for 'Simple', and 'M' for 'Mixture'.

Case S^+

Let $\rho \geq 0$ and $\rho a_1 - a_2 \geq 0$. Simulate jointly (X_1, X_2) : sample $X_1 \sim TN_{[a_1, \infty)}(0, 1)$, $X_2|X_1 = x_1 \sim N(\rho x_1, \nu^2)$, and accept if $X_2 \geq a_2$; otherwise repeat.

Case M^+

Let $\rho \geq 0$ and $\rho a_1 - a_2 < 0$. If component 1 is selected, draw $X_1 \sim TN_{[a_2/\rho, \infty)}(0, 1)$, $X_2|X_1 = x_1 \sim N(\rho x_1, \nu^2)$, and accept simulated pair (x_1, x_2) if $x_2 \geq a_2$. Otherwise, draw $X_1 \sim TN_{[a_1, a_2/\rho]}(\theta, \nu^2)$, and accept with probability

$$\chi\left(\frac{a_2 - \rho x_1}{\nu}\right) / d\left(\frac{a_2 - \rho a_1}{\nu}\right).$$

Here $\theta = \rho(a_2 + \lambda\nu)$, $d(x_0) = \max(\sqrt{\pi/2}, \chi(-x_0))$ and $\chi(x) = e^{\lambda x} \Phi(-x) / \varphi(x)$, where λ is an optimal value, in terms of minimum acceptance rate. Chopin proposes to let λ equal 0.68. Upon acceptance, complete with

$$X_2|X_1 = x_1 \sim TN_{[a_2, \infty)}(\rho x_1, \nu^2).$$

Chapter 4

The examples

In this thesis, three datasets have been considered. These experiments were conducted to find out which factors were the most important for the outcome of the product. To find the effects of the factors for each matrix, the embedded function lm , linear models, in R is used. Every coefficient has 1 degree of freedom, and since all the degrees of freedom are used, the residuals are 0. The factors with the highest effects, in absolute value, are the most important for the experiment.

4.1 Example I

Example I^[10] is a 2^{5-1} split-plot experiment that considers the factors affecting uniformity in a single-wafer plasma etching process. There are three hard-to-change factors on the etching tool: A, the electrode gap, B, the gas flow and C, the pressure. The factors time and radio frequency power, denoted D and E respectively, are easy to change from run to run. The design generator is $E = ABCD$. Table 4.1 shows the design and the resulting uniformity data.

The experimental design is expanded to include all interactions between the factors. These are separated in whole-plot and subplot factors:

$$\mathbf{W}_I = [A, B, C, AB, AC, BC, DE], \quad \mathbf{S}_I = [D, E, AD, AE, BD, BE, CD, CE].$$

Table 4.1: The 2^{5-1} split-plot experiment for the plasma etching tool.

A	B	C	E:	-	+
			D:	-	+
+	+	+	70.31	81.03	
+	-	-	35.67	51.15	
-	+	-	41.80	37.01	
-	-	+	40.32	43.34	
			E:	-	+
			D:	+	-
+	+	-	48.67	91.09	
+	-	+	38.08	62.46	
-	+	+	41.03	31.99	
-	-	-	41.07	40.85	

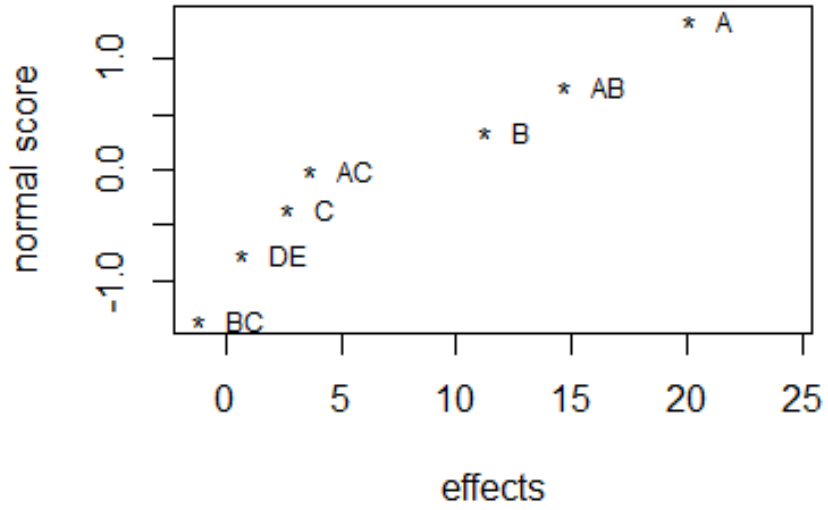
The estimated effects for all the factors and interactions are shown in Table 4.2. Here the first row shows the effects for all the factors. The following rows show the interactions effects, and the value of the intercept is placed last. The factors A, B, E and interactions AB and AE have the largest effects. The others seem to have no significant impact on the outcome of the product.

Table 4.2: The effects from the whole-plot and subplot analysis of Example I.

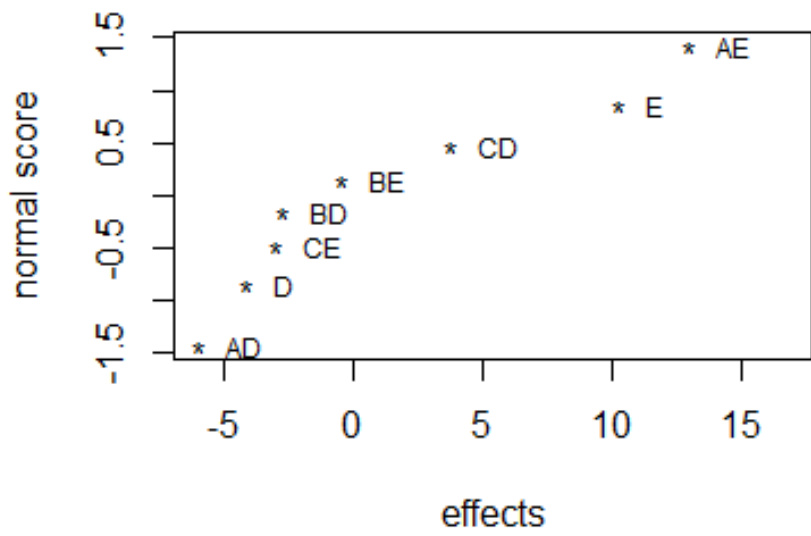
	A	B	C	D	E
	20.1312	11.2488	2.6562	-4.1388	10.2462
A		14.6862	3.6688	-6.0112	13.0038
B			-1.2088	-2.7238	-0.4188
C				3.7388	-2.9762
D					0.6738
Intercept:	49.7419				

Figure 4.1(a) shows the normal probability plot of the estimated effects for the whole-plot factors. Notice that factors A, B and the AB interaction have large effects compared to the others. Figure 4.1(b) shows the normal probability plot of the subplot effects. Only the main effect of E and the interaction AE are large.

From these plots it is easy to see that the factors that are most important for the outcome of the product are: the electrode gap, the gas flow, the interaction between these two, the radio frequency power and the interaction between the latter and the electrode gap.



(a) Whole-plot effects



(b) Subplot effects

Figure 4.1: Normal plot of the original effects of Example I.

4.2 Example II

Table 4.3 contains the design and response for a plasma-treated paper experiment from Bisgaard et. al.^[11]. There are four whole-plot factors A, B, C, D and one subplot factor, E. The factor A is pressure, B is the power, C is the gas flow rate and D is the gas type. Factor E, paper type, is the easy-to-change factor. The response is the "wettability" of the paper measured as the contact angle between the paper and a water droplet placed on the paper right after the plasma treatment.

Table 4.3: The 2^5 split-plot experiment for the plasma-treated paper.

A	B	C	D	E:	-	+
-	+	-	-		55.8	62.9
-	+	-	+		25.6	33.0
-	-	-	-		48.6	57.0
-	-	-	+		5.0	18.1
-	+	+	-		47.2	54.6
-	+	+	+		11.3	23.9
-	-	+	-		37.6	43.5
-	-	+	+		13.3	23.7
+	-	-	+		56.8	56.2
+	-	-	-		41.2	38.2
+	+	-	-		53.5	51.3
+	+	-	+		41.8	37.8
+	+	+	+		49.5	48.2
+	+	+	-		48.7	44.4
+	-	+	-		47.2	44.8
+	-	+	+		47.5	43.2

The expanded experiment includes the following whole-plot and subplot factors:

$$\mathbf{W}_{\text{II}} = [A, B, C, D, AB, AC, BC, AD, BD, CD, ABC, ABD, ACD, BCD, ABCD],$$

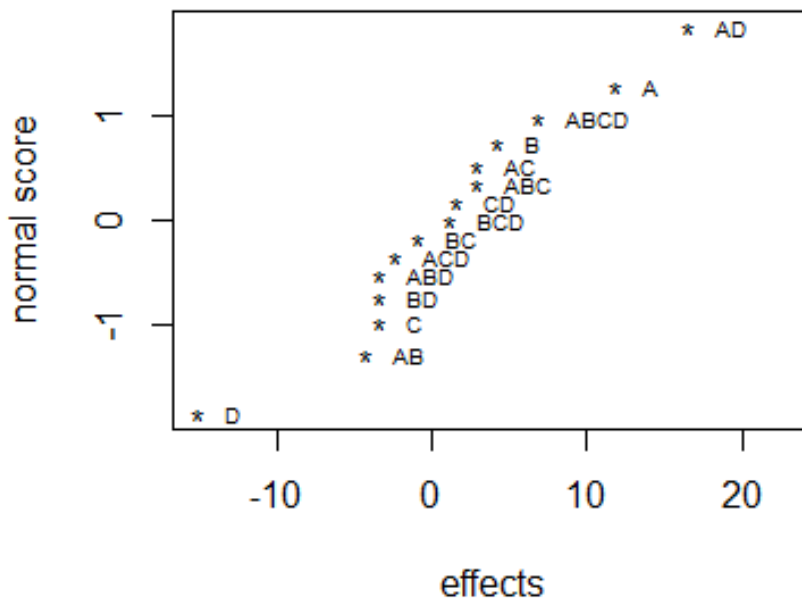
$$\mathbf{S}_{\text{II}} = [E, AE, BE, CE, DE, ABE, ACE, ADE, BCE, BDE, CDE, ABCE, ABDE, ACDE, BCDE, ABCDE].$$

The results of the analysis of Example II are shown in Table 4.4. All interactions with factor E, except the interaction between A and E are very small. Thus, these effects are negligible. The factors A, D and the interaction AD are much higher than the other effects. Therefore, they will influence most on the result of the experiment.

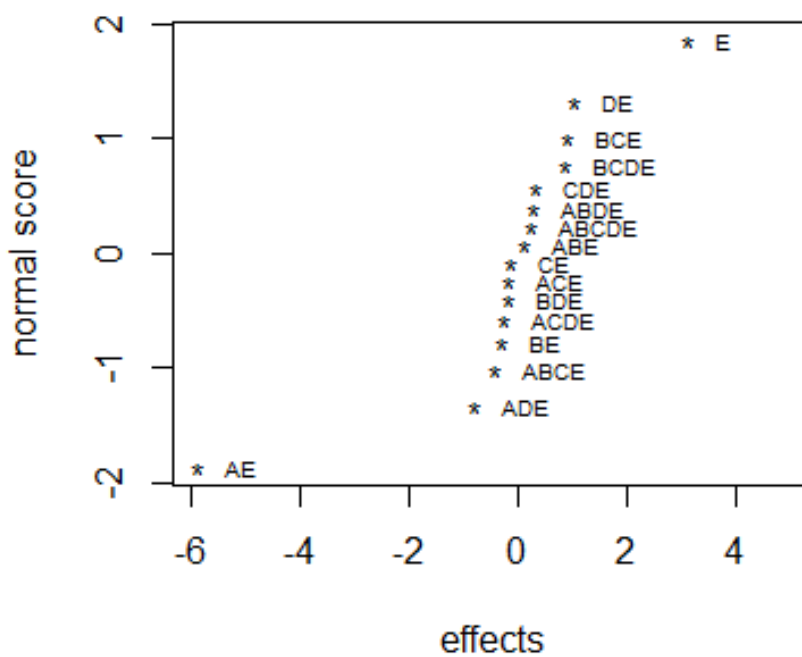
Table 4.4: The effects from the whole-plot and subplot analysis of Example II.

		D	E	DE
		-15.1000	3.1375	1.0250
A	11.8250	16.5626	-5.9000	-0.8125
B	4.2250	-3.3124	-0.3000	-0.1875
C	-3.3876	1.6750	-0.1375	0.3250
AB	-4.2126	-3.3000	0.1125	0.2750
AC	2.9750	-2.3126	-0.1750	-0.2625
BC	-0.8500	1.2374	0.9000	0.8875
ABC	2.8624	6.8500	-0.4375	0.2500
Intercept:	40.9813			

Figure 4.2(a) shows the normal probability plot of the estimated effects for the whole-plot factors. The factors D, A and the interaction AD are the largest in absolute values and do not line up with the rest of the effects. Therefore they are the most important of the whole-plot factors for the result. Figure 4.2(b) shows the estimated effects of the subplot factors. The interaction AE and factor E stand out as the most influential to the result of the experiment. Thus, the outcome of the wettability depends on the pressure, the gas type, the interaction between these two, the paper type and the interaction between the latter and the pressure.



(a) Whole-plot effects



(b) Subplot effects

Figure 4.2: Normal plot of the original effects of Example II.

4.3 Example III

Example III is taken from a report by Box and Jones^[5]. A package-foods manufacturer wished to develop an optimal formulation of a cake mix. These cake mixes are made in large batches, therefore the ingredient factors are hard-to-change. Here A is the amount of flour, B is the amount of shortening, and C the amount of egg powder in the mixture. Many packages are produced from one batch, and the individual packages of cake mix can be baked using different baking times and temperatures. That is, the subplot factors are temperature and time, denoted D and E respectively, since these are easy-to-change. In the experiment, there were 32 runs. The responses from this experiment was obtained from a taste panel, measuring how good the cake tasted, on a scale from 1 - 7.

Table 4.5 shows the original design of Example III and the response. There are four experiments for each combination of the eight different cake mixtures. The dataset is complete with no censored observations.

Table 4.5: The original design of Example III.

Recipe	E:			D:				Average
	A	B	C	-	+	-	+	
(1)	-	-	-	1.1	1.4	1.0	2.9	1.6
(2)	+	-	-	1.8	5.1	2.8	6.1	3.95
(3)	-	+	-	1.7	1.6	1.9	2.1	1.825
(4)	+	+	-	3.9	3.7	4.0	4.4	4
(5)	-	-	+	1.9	3.8	2.6	4.7	3.25
(6)	+	-	+	4.4	6.4	6.2	6.6	5.9
(7)	-	+	+	1.6	2.1	2.3	1.9	1.975
(8)	+	+	+	4.9	5.5	5.2	5.7	5.325

For the fully expanded experiment, there are in this case two parts of subplot effects. Thus, eq. (3.2) is used. The following shows which factors that belong in each part of the split-plot, whole-plot and two parts of subplot respectively;

$$\mathbf{W}_{\text{III}} = [A, B, C, AB, AC, BC, ABC],$$

$$\mathbf{S}_{\text{III}}^- = [D, E, AD, AE, BD, BE, CD, CE, ABD, ABE, ACD, ACE, BCD, BCE, ABCD, ABCE],$$

$$\mathbf{S}_{\text{III}}^+ = [DE, ADE, BDE, CDE, ABDE, ACDE, BCDE, ABCDE].$$

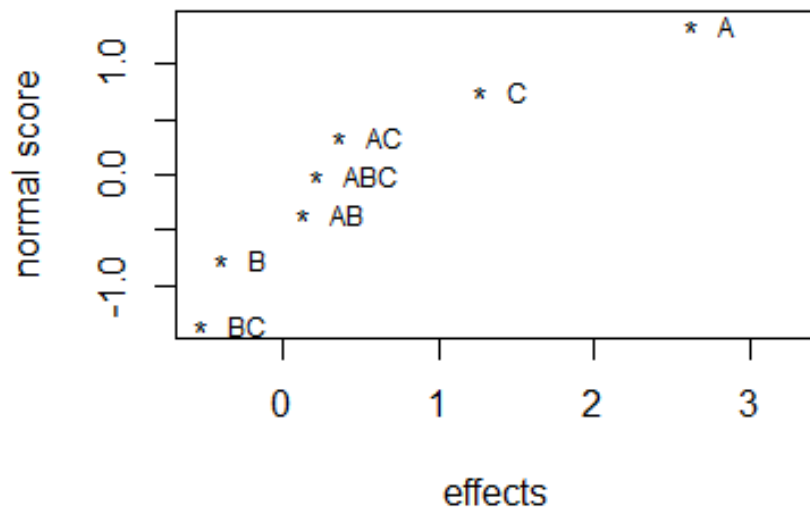
The effects found from the analysis of Example III are shown in Table 4.6. Factor A has clearly the largest effect. The factors C and D are large enough to be significant. Higher order interactions seem to have negligible impact on the outcome of the product.

Table 4.6: The effects from the whole-plot and subplot analysis of Example III.

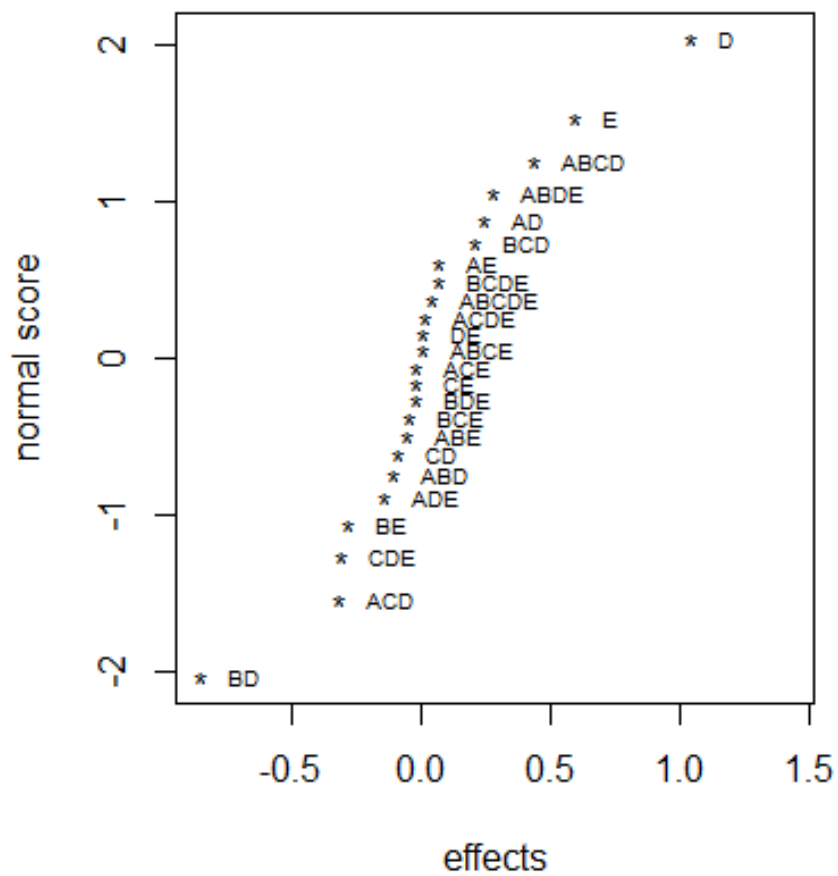
		D	E	DE
		1.0438	0.5938	0.0063
A	2.6313	0.2438	0.0688	-0.1438
B	-0.3938	-0.8563	-0.2813	-0.0188
C	1.2688	-0.0938	-0.0188	-0.3063
AB	0.1313	-0.1063	-0.0563	0.2813
AC	0.3688	-0.3188	-0.0188	0.0188
BC	-0.5313	0.2063	-0.0438	0.0688
ABC	0.2188	0.4313	0.0063	0.0438
Intercept:	3.4781			

The effect of the whole-plot factors are shown in Figure 4.3(a). The factors A and C seem to be of most importance for the outcome of the cake. Figure 4.3(b) shows the subplot effects and the factor D stands out as the most significant of these effects. The factor E and interaction BD might influence the result.

The effects in $\mathbf{S}_{\text{III}}^+$ are small, and insignificant, since they are mainly higher order interactions.



(a) Whole-plot effects



(b) Subplot effects

Figure 4.3: Normal plot of the original effects of Example III.

Chapter 5

The methods

First, the previous work, leading up to this thesis, is discussed, then the software used in the calculations is introduced. Thereafter, the layout of the two methods, maximum likelihood with multiple imputation and the quick and dirty method, are presented.

5.1 Previous work

The specialising project leading up to this thesis concerns analysing right censored data from split-plot design with mirror image pairs^[12], using the maximum likelihood and multiple imputation method. If censored datasets are analysed using SPMIP, the runtime of analysis may be much reduced. The size of the new datasets will be a fraction of the originals, thus the analysis will have a considerably lower runtime. Also the problem with correlation between observations will be solved. This was the motivation for the project.

Through the project, some problems occurred. There is a huge risk that uncensored data will be censored and lost. Especially information about the subplot effects. For each censored observation, a censoring limit had to be calculated. As an example, consider an experiment with four subplots per whole-plot. For the whole-plot, the censoring limit for whole-plot combination i is found by

$$\frac{\sum_{j=1}^4 y_{ij}}{4} \geq \frac{nc + \sum_{j=2}^{4-n} o_j}{4} = c_i, \quad (5.1)$$

where c is the original censoring limit, o_j is the value of observations less than c , y_{ij} are the original response values within whole-plot i , and c_i is the new censoring limit for whole-plot combination i . If this inequality holds for any i , the observation becomes censored, and c_i is set as the censoring limit. This test is done for all i 's, thus this method requires a lot of calculations that can produce errors. The censoring limit depends on all the observations within each whole-plot. If all of these are censored, the limit is set to be the original censoring limit. Otherwise the limit is calculated by means of the observed values.

Now, consider right censoring of an experiment with two subplots per whole-plot. Let y_1 and y_2 be two observations within the same whole-plot. If $y_1 \geq c$ and $y_2 \geq c$, the whole-plot combination is censored with limit c . If $y_1 \geq c$ and $y_2 = o$,

the whole-plot combination is censored by $c^* = \frac{c+o}{2}$. The censoring limits for the subplot effects require more work. Let $Z_1 = \frac{y_1+y_2}{2}$ and $Z_2 = \frac{y_1-y_2}{2}$, such that $Z_1 + Z_2 = y_1$. The observation becomes right censored if

$$y_1 \geq cZ_2 = y_1 - Z_1 \geq c - Z_1$$

$$\text{or } y_1 \geq oZ_2 = y_1 - Z_1 \geq o - Z_1.$$

The calculations of the subplot effects are likely to produce even larger errors. This can happen if the difference between two dependent observations, both below the censoring limit, is higher than the difference between two censored observations. The same apply in the case where one observation is censored and the other is not. This results in that the censoring limit should be set high enough for all the censored observations to become censored, without losing any observed values. Some censored values will not be high enough for this condition, and therefore, their value must be set equal to a fictional value.

The computation of the censored data can only handle one censoring limit. Thus, the limit becomes the lowest of the calculated c_i 's. This limit might be set too low, i.e. some observed values are censored, thereby lost.

The motivation for this thesis is to find a better method for analysing split-plot data with censored observations, by means of the variances in the split-plot design.

5.2 R software

R^[13] is an open-source statistical programming language which is widely used for data analysis and statistical computing. The software provides statistical and graphical techniques, with classical statistical tests. The R code for Example II can be found in Appendix B. It produces helpful graphs, such as showing which distribution fits the data best, and whether or not the factors included are significant (Danielplot). A Danielplot is a normal plot of effects from a two-level factorial experiment. Effects that show a linear trend are viewed as insignificant. If one or more effects fall out of this linear trend, they are significant.

There are many embedded functions in R. In this thesis, the most used functions are *censReg* and *lm*. The *censReg*-function takes in a dataset with the corresponding response values, and sets a censoring limit for the response. It assumes that the data are from a normal distribution, which sometimes can limit the function. In this function, the maximum likelihood is calculated by the Newton-Raphson method. Linear models (*lm*) can be used to perform regression analysis, where it returns, among other, the coefficients of the specified model. In this thesis, it is used to find the effects from the uncensored dataset, to compare with our results, and to estimate the effects for the censored cases. The output of the estimates of the factors in *censReg* and *lm* is the coefficient of the factors. Thus, to find the effects, the coefficient of each term has to be doubled.

When working with the truncated normal distribution, the package *truncnorm* with the function *rtruncnorm* is used. This function generates n random deviates in a defined interval (a, b) from a mean and standard deviation.

5.3 Censoring with the maximum likelihood and multiple imputation method

When censored data is analysed with maximum likelihood and multiple imputation, each part of the split-plot is analysed separately, since they have different error terms. The upper, or lower, limit should be set where it is natural, i.e. where there is a big leap between the observed values, or to get some censored values. When using R to create estimates of the two variances, some factors must be omitted from *censReg*. The factors with the smallest coefficient should be the ones that are left out. Here, this is done by exhaustive search. Each combination of factors is tested in *censReg*, and the factor with the lowest coefficient is removed. This is done until every factor is placed after significance. These factors are then brought back if they influence the size of $\hat{\sigma}_w^2$ and $\hat{\sigma}_s^2$, i.e. makes them smaller. It should be noted that neither of these variances are allowed to be negative. The following method is based on the theory from Section 3.5 and 3.6.

Consider an experiment with two subplots per whole-plot. Using multiple imputation and drawing from truncated distributions, let y_{i1} be a censored observation, and y_{i2} be observed. The two are dependent, with $\text{cov}(y_{i1}, y_{i2}) = \sigma_w^2$ and correlation $\rho = \frac{\sigma_w^2}{\sigma_w^2 + \sigma_s^2}$. After subtracting estimated expected value, there are, for each i , two noise terms left. Let these be $e_{i1} = \epsilon_i^w + \epsilon_{i1}^s$ and $e_{i2} = \epsilon_i^w + \epsilon_{i2}^s$, where only the latter can be estimated. Then, e_{i1} and e_{i2} are bi-normal distributed, i.e. $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and $\boldsymbol{\Sigma} = (\sigma_w^2 + \sigma_s^2) \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$. When finding the distribution of e_{i1} given that $e_{i1} > c$ and $e_{i2} = a_{i2}$ i.e.

$$F(a_{i1}) = P(e_{i1} \leq a_{i1} | e_{i1} > c \cap e_{i2} = a_{i2}) = \frac{P(c < e_{i1} \leq a_{i1} | e_{i2} = a_{i2})}{P(e_{i1} > c | e_{i2} = a_{i2})},$$

which is a truncated normal distribution with mean ρa_{i2} and variance $(\sigma_w^2 + \sigma_s^2)(1 - \rho^2)$.

In other words, the coefficients are multiplied with their associated levels and added up with the intercept, to find the estimated expected response for each row of the matrix which is censored. Multiple imputation using normal distribution is then implemented five times for each censored value. When the five values for each censored value are found, these values are substituted into the original response in place of their associated censored value. Now, there are five new sets of responses and linear models are used to estimate the effects of each factor for each of the sets. The mean of the effects of the factors is then calculated to give the estimated effects found by the maximum likelihood and the multiple imputation method. In the following, let $K_s = (n/2) - k_s$, where n is the number of runs in the experiment and k_s is the number of censored observations in the subplot matrix. The subplot error can be found by this equation,

$$\frac{\hat{\sigma}_s^2}{2} = K_1 = \frac{\sum_{i=1}^{K_s} (Z_i^s - E(y_{ij}))^2}{K_s - d_s}, \quad j = 1, 2, \quad (5.2)$$

where d_s is the number of factors used in *censReg* and $K - d > 0$. The estimate for $K_2 = \hat{\sigma}_w^2 + \frac{\hat{\sigma}_s^2}{2}$ is found similarly.

Next follows the proposed procedure for right censoring of an experiment with two subplots per whole-plot.

1. Identify the most important subplot main effects and whole-plot by subplot interaction by using *censReg*.
 - Let d_s denote the number of subplot factors used in *censReg*.
 - Calculate the estimate $E_s(y_i)$ with subplot coefficients and their associated level for every i .
 - If the two dependent observations are not censored, find $Z_i^s = \frac{y_{i1} - y_{i2}}{2}$.
 - Calculate K_1 by eq. (5.2).
2. Use *censReg* to identify the most important whole-plot main effects, the whole-plot by whole-plot and subplot by subplot interactions.
 - Let d_w denote the number of whole-plot factors used in *censReg*.
 - Calculate the estimate $E_w(y_i)$ with whole-plot coefficients and their associated level for every i .
 - If the two dependent observations are not censored, find $Z_i^w = \frac{y_{i1} + y_{i2}}{2}$.
 - Calculate K_2 similar to K_1 .
3. If $K_2 \geq K_1$ and $K_2, K_1 \geq 0$, and both K_1 and K_2 are fairly small, stop. Else, recover one factor/interaction in *censReg* and compute again.
4. Find $\hat{\sigma}_s^2, \hat{\sigma}_w^2$ and ρ and use both the whole-plot and subplot factors/interactions to estimate $E(y_{ij}), j = 1, 2$, for all set-ups.
 - Calculate $B_{ij} = Z_{ij} - E(y_{ij})$ and $A_{ij} = \rho B_{ij}$, for those who are not censored.
 - If $y_{i1} \geq c, y_{i2} = o$, estimate values for y_{i1} by conditional distribution: let $\mu = A_{i2}, \sigma = \sqrt{(\hat{\sigma}_w^2 + \hat{\sigma}_s^2)(1 - \rho^2)}$ and $\alpha = c - E(y_{i1})$, and draw from the truncated distribution with mean μ and standard deviation σ on the interval $[\alpha, \infty)$.
 - If $y_{i1} \geq c$ and $y_{i2} \geq c$, decide a_1, a_2 and x , from $c - E(y_{ij})$, and check if M^+ or S^+ is appropriate. Estimate from truncated distributions explained in Section 3.6.
5. Create five new datasets, and exchange the censored observations with the estimated values, added to $E(y_{ij})$.
 - Apply mirror image pairs on the new datasets, find the new estimated coefficients and take the mean of each from the five datasets. Multiply the mean coefficients with 2 to get the effects.

5.4 The quick and dirty method

The quick and dirty (QnD) method is a censoring method that sets an upper limit, for right censoring, or a lower limit, for left censoring. Every value greater than this limit, or lower for the left censoring case, will be put equal to the limit. Then, the dataset will be complete, and analysed by standard methods for complete data. If the limit is set such that about half of the values in a complete dataset are censored, there will be poor estimates for the effects of the factors. Furthermore, in the right censored case, if the highest values in the complete dataset are much higher than the limit, the estimation of the effects will most likely not be the same, or close, to the original dataset. If the highest values in the complete dataset are close to the chosen limit, the estimates will be very close to the original estimates.

The method is known to be unstable, but it is included here such that it can be compared with the other method. It is also an easy method to implement, and is therefore useful to test in a quick analysis. QnD can also be used to check estimates found by other methods, by comparing the analyses. If they produce estimates that are extremely different, one should consider to run the analysis again to ensure good results.

Chapter 6

Experiments and results

In the following sections R and L will denote right and left censoring, respectively. The sets of factors will be indexed after which example they belong to. To reduce the amount of text, the short term MI, multiple imputation, is used in tables throughout the thesis.

6.1 Example I - Right censoring

The right censoring limit for Example I is set to 70. There are three observations that are higher than this limit, i.e. three values are censored. After testing all possible combinations of factors, these are the ones that acquired the lowest pair of $\hat{\sigma}_w^2$ and $\hat{\sigma}_s^2$: $\mathbf{W}_I^R = [A, B, AB]$, $\mathbf{S}_I^R = [E, AE, CD]$. The calculated values for the variances and the correlation are shown in Table 6.1.

Table 6.1: Estimated values for the variances and the correlation for the right censored data in Example I, for the maximum likelihood with multiple imputation method.

Whole-plot variance:	$\hat{\sigma}_w^2$	4.8758
Subplot variance:	$\hat{\sigma}_s^2$	8.9809
Correlation:	ρ	0.3519

The expected value for the censored observations and the standard deviation calculated from the values in Table 6.1, are used in the truncation. The estimated values for the censored observations are shown in Table 6.2. Conditional distribution, eq. (3.4), is used on the first censored value, since observation 10 is dependent on observation 9 which is not censored. While the theory from Section 3.6 is used on the two other censored values, given that observation 1 and 2 are dependent and both censored.

The estimated value of observation 10 differs a lot from run to run. This can influence the result for each estimation, but taking the mean of the estimated effects from multiple runs will create a more reliable result. The estimations for run 1 are very close to the original value. The values for the other runs are not as good as run 1, but they seem reasonable. After filling in the values from Table 6.2 to create five complete fictional datasets, the estimated effects are calculated. Figure 6.1 shows

Table 6.2: Estimated values from multiple imputation for the right censored data in Example I.

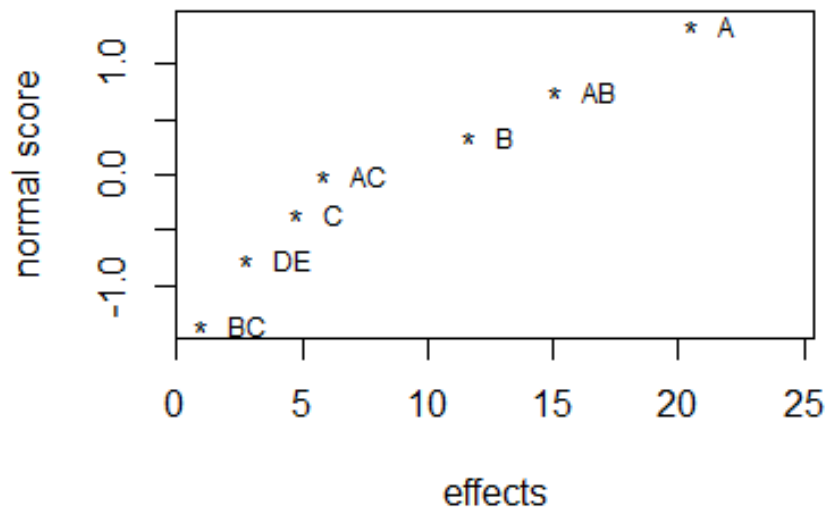
Observation	Run nr:					Original
	1	2	3	4	5	
1	70.0332	70.0534	70.0126	70.2023	70.0955	70.31
2	91.4760	92.3104	91.9745	91.7790	92.4852	81.03
10	84.1512	74.5250	85.6074	78.3539	86.6941	91.09

the normal plot of the right censored effects, where (a) shows the whole-plot effects and (b) the corresponding subplot effects. The largest whole-plot effects are factor A, B, and the interaction between them, AB. These are also the effects that fall off the linear trend. The largest subplot effects are the interaction AE and factor E. When looking at the plot of the subplot effects, the interaction CD, along with E and AE, does not lie on the linear trend with the other effects. Although this interaction has a much lower effect than the largest, it may influence the outcome of the product.

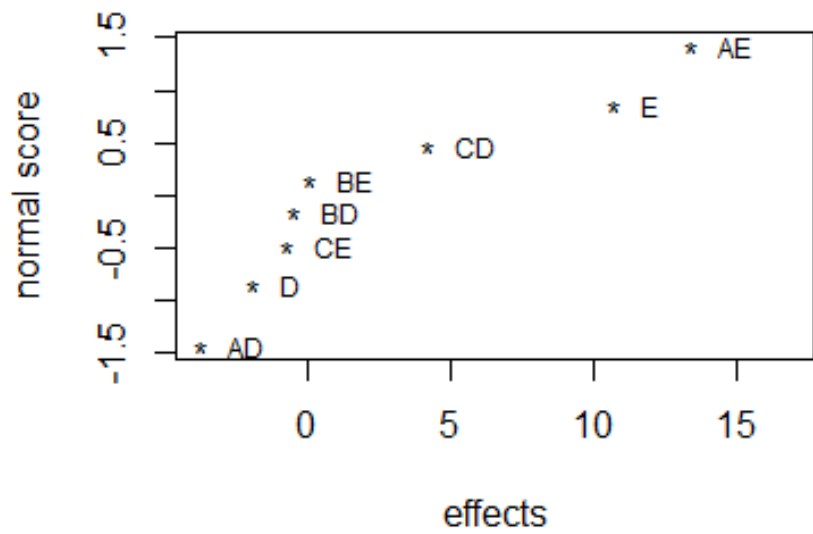
In the QnD case, the calculations are easy to compute. Figure 6.2 shows the effects after censoring by the QnD method. From the whole-plot effects in Figure 6.2(a) it is easy to see that factor A stands out as the most important factor, followed by the interaction AB. The effect of factor B is still large, but in this case, it falls onto the linear trend, which makes it harder to decide if it is significant.

The whole-plot effects in the QnD case seem to have a more linear trend than the whole-plot effects in Figure 6.1(a). For the subplot effects in Figure 6.2(b), the interaction AE and factor E are the largest, and the interaction AD might be important for the outcome of the experiment, since it seems to fall slightly off the linear trend.

The numerical results for both methods are shown in Table 6.3. The estimates for the QnD method are consequently lower than those for the estimates found by maximum likelihood with multiple imputation. Compared to the original estimates in Table 4.2, the maximum likelihood with multiple imputation gives the best results.

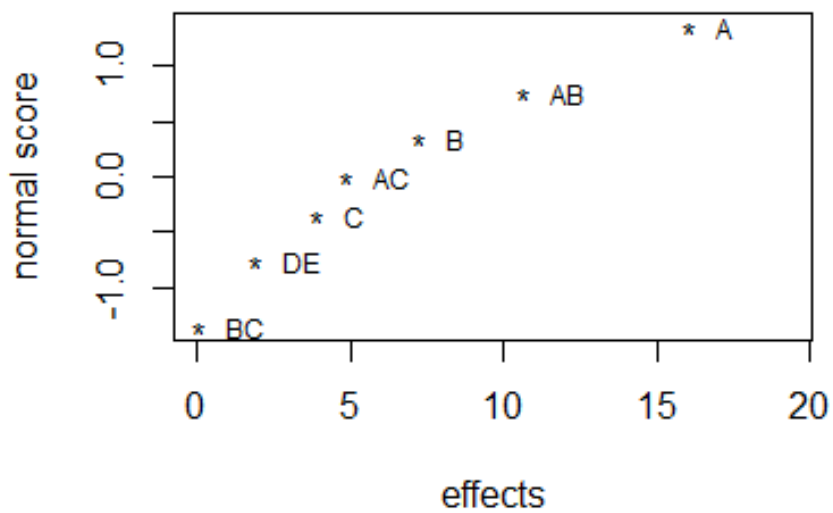


(a) Whole-plot effects

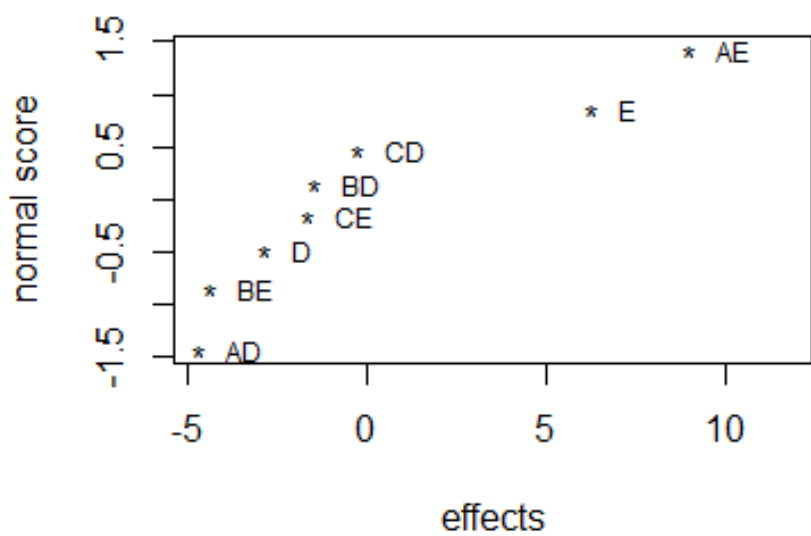


(b) Subplot effects

Figure 6.1: Normal plot of the right censored estimated effects of Example I using multiple imputation.



(a) Whole-plot effects



(b) Subplot effects

Figure 6.2: Normal plot of the effects with right censoring by QnD from Example I.

Table 6.3: Estimated effects of the factors from the two methods for right censoring of Example I.

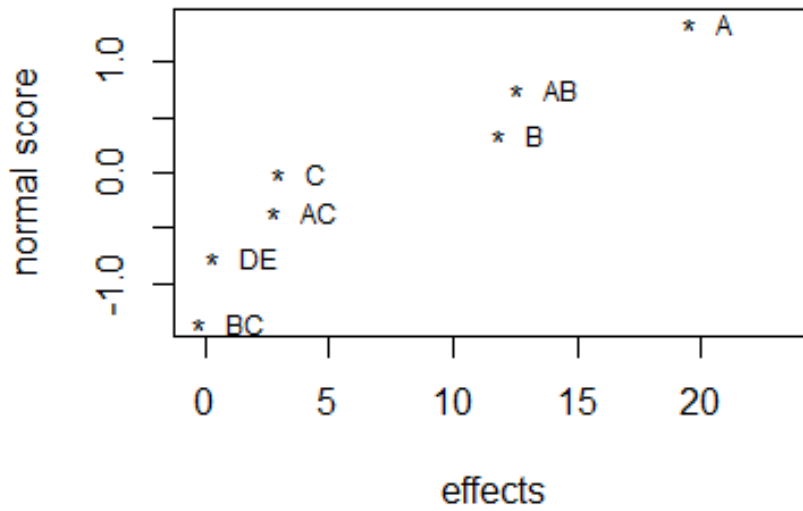
	MI:	QnD:
(I)	49.8369	47.7150
A	20.3214	16.0775
B	11.4388	7.1950
C	5.1522	3.8750
AB	14.8764	10.6325
AC	6.1648	4.8875
BC	1.2872	0.0100
DE	3.1698	1.8925
D	-1.5850	-2.8425
E	10.4940	6.2700
AD	-3.4576	-4.7150
AE	13.2514	9.0275
BD	-0.1700	-1.4275
BE	-0.1710	-4.3950
CD	3.9864	-0.2375
CE	-0.4226	-1.6800

6.2 Example I - Left censoring

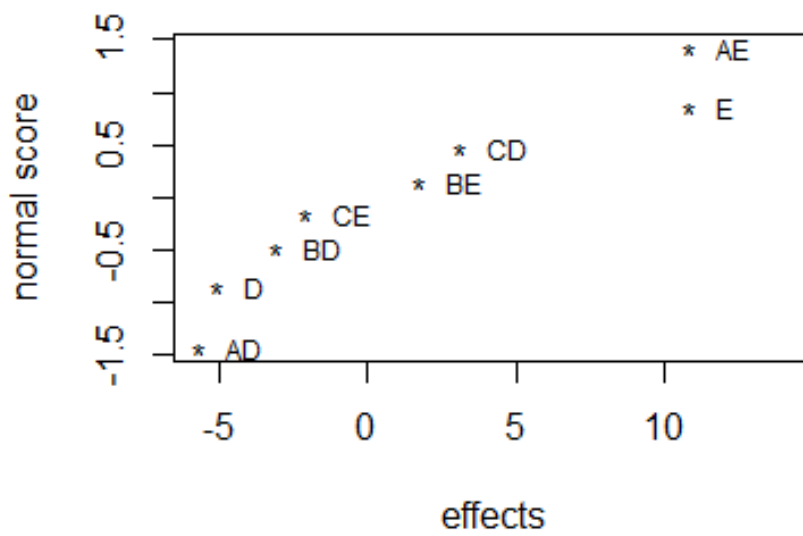
For the left censoring of Example I, there is no solution where $K_2 \geq K_1$. Different limits have been tested, without success. Thus, it is not possible to estimate the values of the factors effects in this case. To get around this, $\hat{\sigma}_w^2$ can be set to a static value; usually zero. This is omitted in this thesis, since the tested procedure failed to execute.

In the QnD case, there are no limiting restrictions, thus it can be performed no matter what. The censoring limit is set to 40, that is, four values are censored. Figure 6.3 shows the normal plot of the effects when left censored by the QnD method. Figure 6.3(a) shows the whole-plot effects and (b) shows the subplot effects. The whole-plot factors A, B and the interaction AB are by far the largest. Thus these effects are the most important. The subplot factor E and the interaction AE stand out as the most significant subplot factors for the outcome of the product of Example I.

The values for the estimated coefficients can be found in Table A.1, Appendix A. The QnD method for the left censoring case gives the closest estimates to the original values overall, compared to the values from the right censoring case. A reason for this is that the left censoring limit is closer to the censored observations original values, than that of the right censoring case.



(a) Whole-plot effects



(b) Subplot effects

Figure 6.3: Normal plot of the effects with left censoring by QnD from Example I.

6.3 Example II - Right censoring

The right censoring limit for Example II is set to 55. There are five observations censored by this limit. After testing all possible combinations of factors, these are the ones that acquired the lowest pair of $\hat{\sigma}_w^2$ and $\hat{\sigma}_s^2$:

$\mathbf{W}_{\text{II}}^R = [\text{A, B, D, AC, AD, ABD, ABCD}]$, $\mathbf{S}_{\text{II}}^R = [\text{E, AE, DE, ACE, ADE, BCDE}]$.

The different variances and the correlation are shown in Table 6.4. These values are used to create an estimate for the uncensored observations.

Table 6.4: Estimated values for the variances and the correlation for the right censored data in Example II, for the maximum likelihood with multiple imputation method.

Whole-plot variance:	$\hat{\sigma}_w^2$	9.7756
Subplot variance:	$\hat{\sigma}_s^2$	1.5774
Correlation:	ρ	0.8611

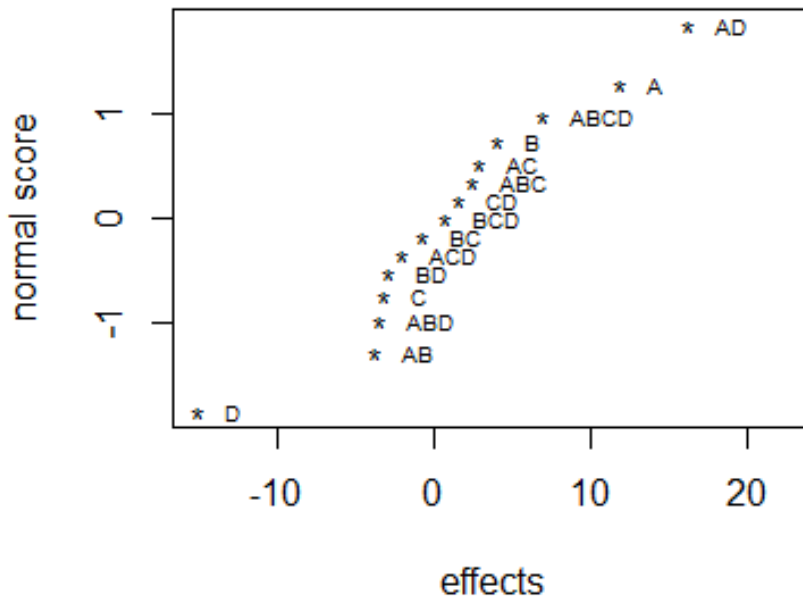
The calculations of the multiple imputation give the values for each observation shown in Table 6.5. These values are set into the censored dataset to create five fictional complete datasets. Two pairs of observations are estimated based on Section 3.6 and truncation. For each run, the values obtained through multiple imputation are very similar. The last censored observation is estimated by means of conditional distribution, eq. (3.4). All the estimated values are close to the original values. This is due to that the censoring limit is close to every censored value.

Table 6.5: Estimated values from multiple imputation for the right censored data in Example II.

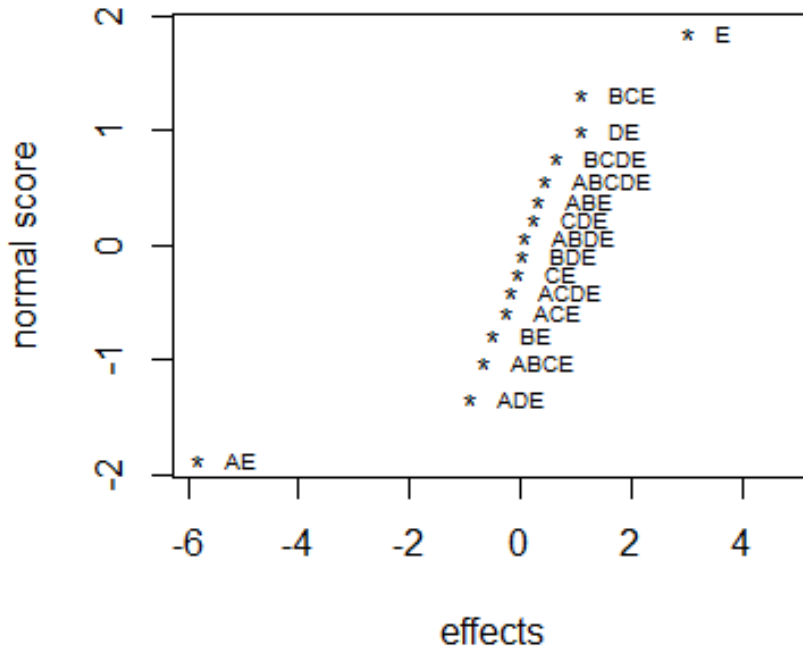
Observation	Run nr:					Original
	1	2	3	4	5	
1	55.2387	55.0019	55.0057	55.0572	55.1124	55.8
2	60.0856	59.8699	59.7799	59.7815	59.8631	62.9
6	55.9403	55.0221	55.1994	55.4620	58.1196	57.0
17	56.0013	55.7628	56.0069	55.9434	55.8964	56.8
18	55.0859	55.0025	55.1183	55.1436	55.0923	56.2

Regression is used on each of the datasets, and the result for the whole-plot analysis is shown in Figure 6.4(a), and Figure 6.4(b) shows the corresponding subplot results. The whole-plot factors A, D and the interaction AD are the largest of the whole-plot effects, and thus the most significant of these. For the subplot factors, the interaction AE and the factor E stand out as the most significant for the outcome of the result.

The result of the QnD method is shown in Figure 6.5, with the whole-plot effects in (a), and the subplot effects in (b). Just like in Figure 6.4(a), the factors A, D and the interaction AD stand out as the most important whole-plot factors. For the subplot effects, the interaction AE is the most significant. Factor E is still the second largest, but here it lies on the linear trend. It is thus not possible to tell if this factor is significant.

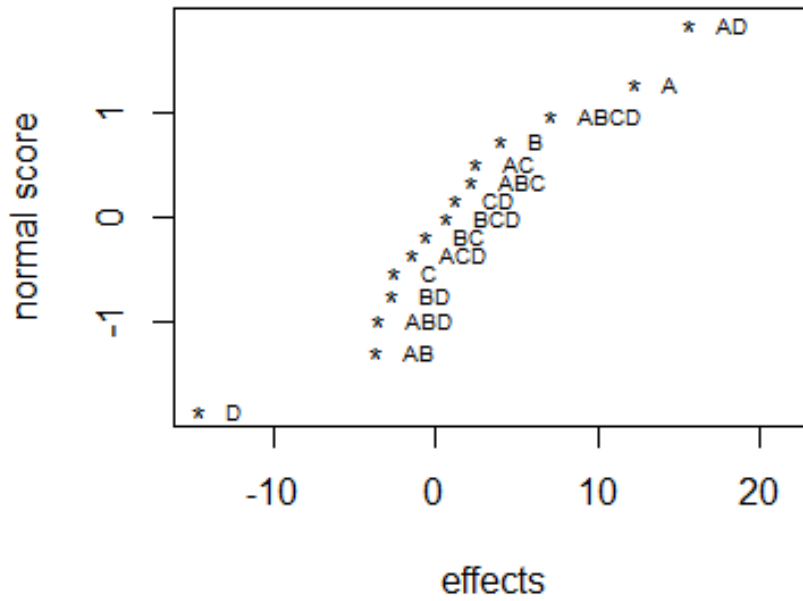


(a) Whole-plot effects

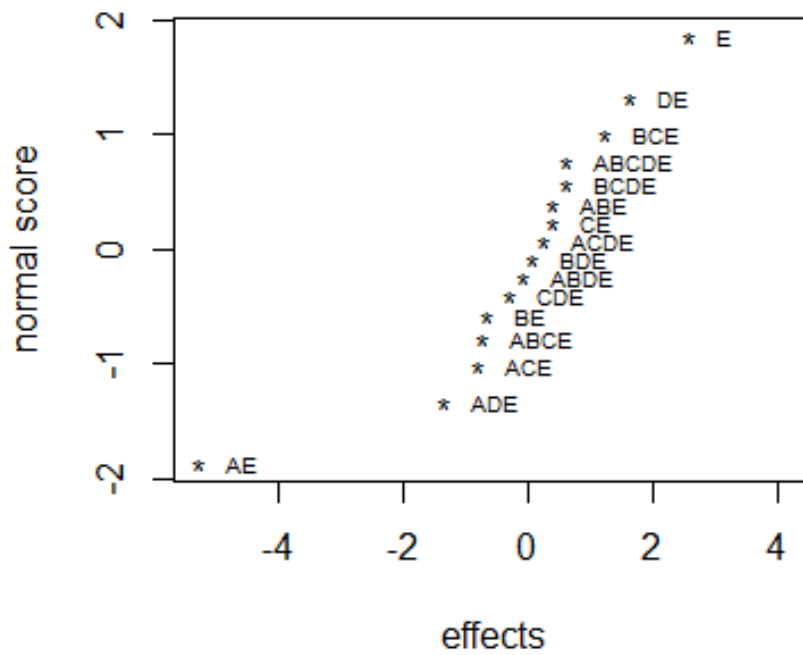


(b) Subplot effects

Figure 6.4: Normal plot of the right censored effects of Example II using multiple imputation.



(a) Whole-plot effects



(b) Subplot effects

Figure 6.5: Normal plot of the effects with right censoring by QnD from Example II.

The estimated coefficients for both methods can be found in Table A.2, Appendix A. Compared to the original effects, multiple imputation with maximum likelihood gives the best estimation.

6.4 Example II - Left censoring

For the left censoring of Example II, the censoring limit is set to 25, that is, six observations are censored. The factors that creates the pair of $\hat{\sigma}_w^2$ and $\hat{\sigma}_s^2$ with the lowest acceptable values are:

$$\mathbf{W}_{\text{II}}^L = [A, C, D, AC, AD, BD, ABCD], \quad \mathbf{S}_{\text{II}}^L = [AE, DE, ABE, ADE, ABCDE].$$

Table 6.6 shows numbers for the different variances and the correlation.

Table 6.6: Estimated values for the variances and the correlation for the left censored data in Example II, for the maximum likelihood with multiple imputation method.

Whole-plot variance:	$\hat{\sigma}_w^2$	7.0372
Subplot variance:	$\hat{\sigma}_s^2$	2.3343
Correlation:	ρ	0.7509

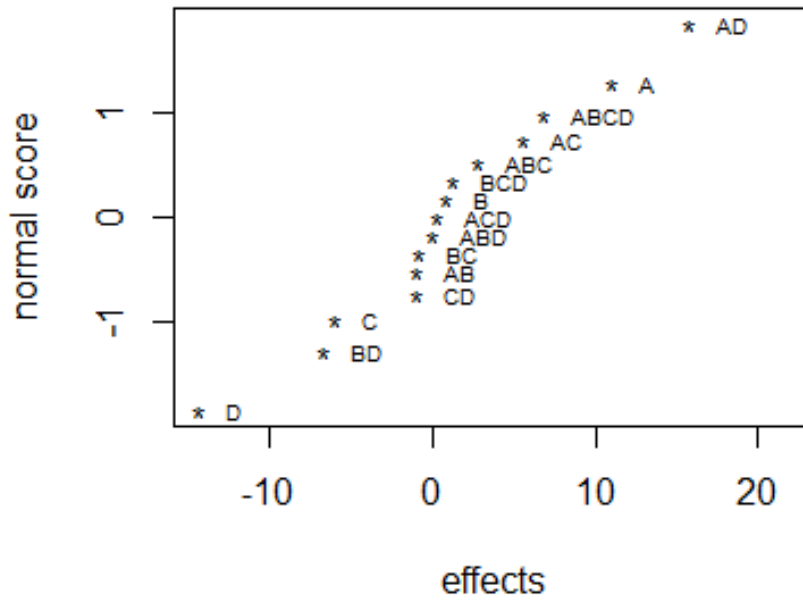
There are three whole-plot combinations that are censored in this case. Thus, the method where both dependent observations are censored, Section 3.6, is used on each combination. Multiple imputation creates estimates for each censored value, shown in Table 6.7.

Table 6.7: Estimated values from multiple imputation for the left censored data in Example II.

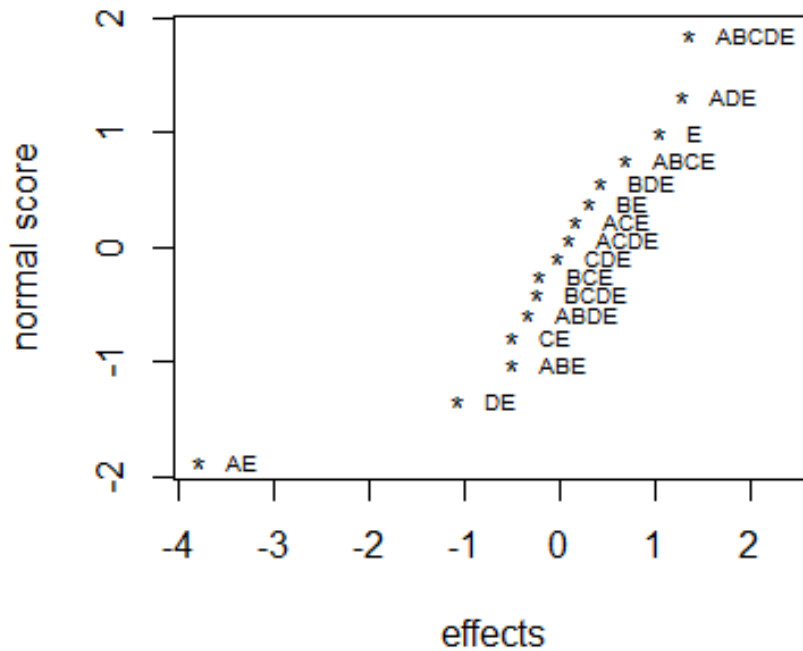
Observation	Run nr:					Original
	1	2	3	4	5	
7	25.2207	25.5315	24.7438	24.8516	25.1883	5.0
8	24.4113	24.9805	24.2776	24.3963	24.5432	18.1
11	6.8982	8.0627	7.0220	8.9404	8.3585	11.3
12	7.6602	8.9095	8.0828	9.3735	8.80711	23.9
15	20.0818	19.7337	22.1194	21.1430	19.7989	13.3
16	22.7173	21.9018	24.1054	23.3482	22.1078	23.7

Notice from Table 6.7 that only the estimates for run 16 are close to the original value. What is common for all the censored values is that one of the dependent values are fairly close to the censoring limit, while the other is far away from it. This influences the estimates.

Regression analysis of the artificial completed datasets are shown as normal plots in Figure 6.6. The whole-plot effects in Figure 6.6(a) shows that the factor D and the interaction AD are the most significant. Factor A is the third largest, but it seems to lie on a linear trend, which makes it uncertain whether or not it is significant. For the subplot effects in Figure 6.6(b), only one subplot effect stands out as significant; the interaction AE.



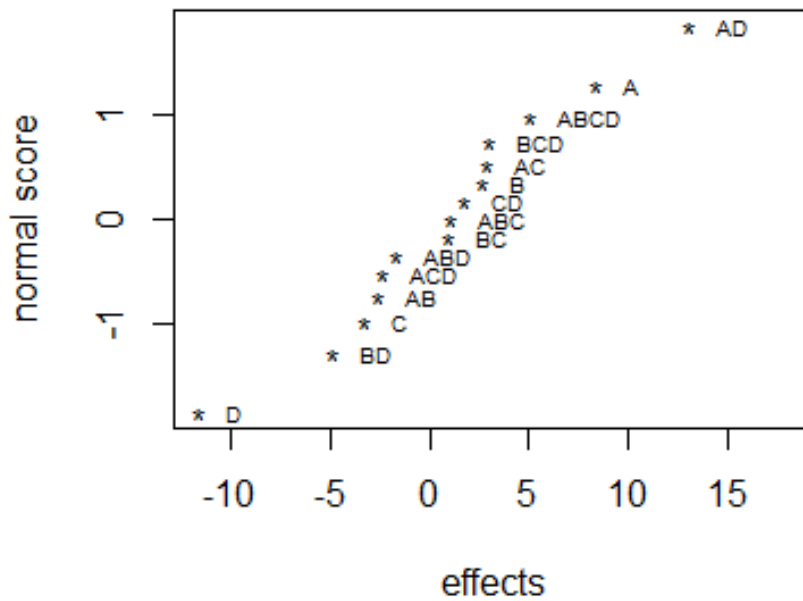
(a) Whole-plot effects



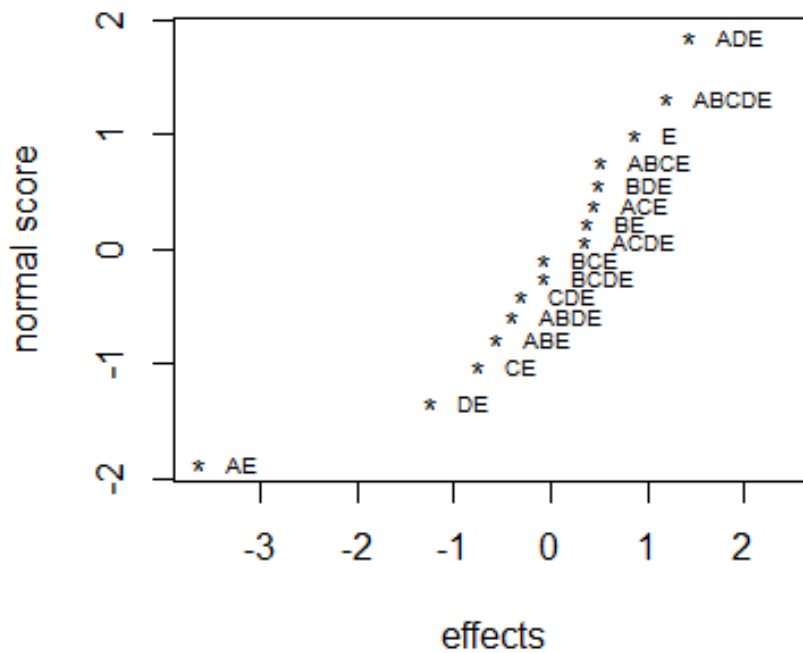
(b) Subplot effects

Figure 6.6: Normal plot of the left censored estimated effects of Example II using multiple imputation.

The results of the QnD method are shown in Figure 6.7(a), with the whole-plot effects, and in Figure 6.7(b), with the subplot effects. The most important whole-plot factor is the interaction AD, followed by factor D. Factor A also might be of importance for the outcome of the experiment. For the subplot effects, only the interaction AE seems to be significant.



(a) Whole-plot effects



(b) Subplot effects

Figure 6.7: Normal plot of the estimated effects with left censoring by QnD from Example II.

The estimated coefficients are listed in Appendix A, Table A.2. Both methods produces good estimates, but for the most significant factors, QnD gives the poorest results.

6.5 Example III - Right censoring

The right censoring limit for Example III is set to 6. Four observations are censored due to this limit. The following factors create the optimal pair of $\hat{\sigma}_w^2$ and $\hat{\sigma}_s^2$: $\mathbf{W}_{\text{III}}^R = [A, B, C, BC]$, $\mathbf{S}_{\text{III}}^{-R} = [D, BD]$. All of the factors in $\mathbf{S}_{\text{III}}^{+R}$ are left out during the calculations, since none of them seem to be important for the result of the experiment.

The estimated variances and the correlation for the right censoring of Example III are shown in Table 6.8. Note that the variances are generally smaller than in the other examples. This is natural since the original data is restricted to a much smaller interval.

Table 6.8: Estimated values for the variances and the correlation for the right censored data in Example III, for the maximum likelihood with multiple imputation method.

Whole-plot variance:	$\hat{\sigma}_w^2$	0.0734
Subplot variance:	$\hat{\sigma}_s^2$	0.2475
Correlation:	ρ	0.2288

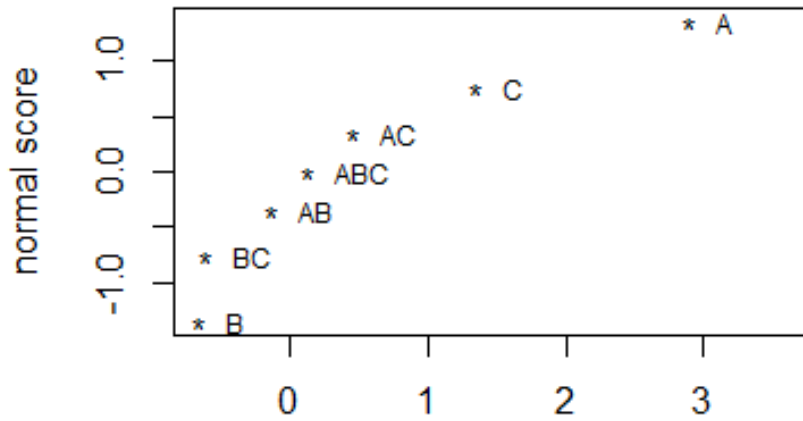
There are one observation within one of the whole-plots, and three observations within another, that are censored. The methods of calculation are straight forward using conditional distribution, eq. (3.4), first with three dependent observations for the single censored, then one for the set with three censored observations. The results of the multiple imputation are shown in Table 6.9.

Table 6.9: Estimated values from multiple imputation for the right censored data in Example III.

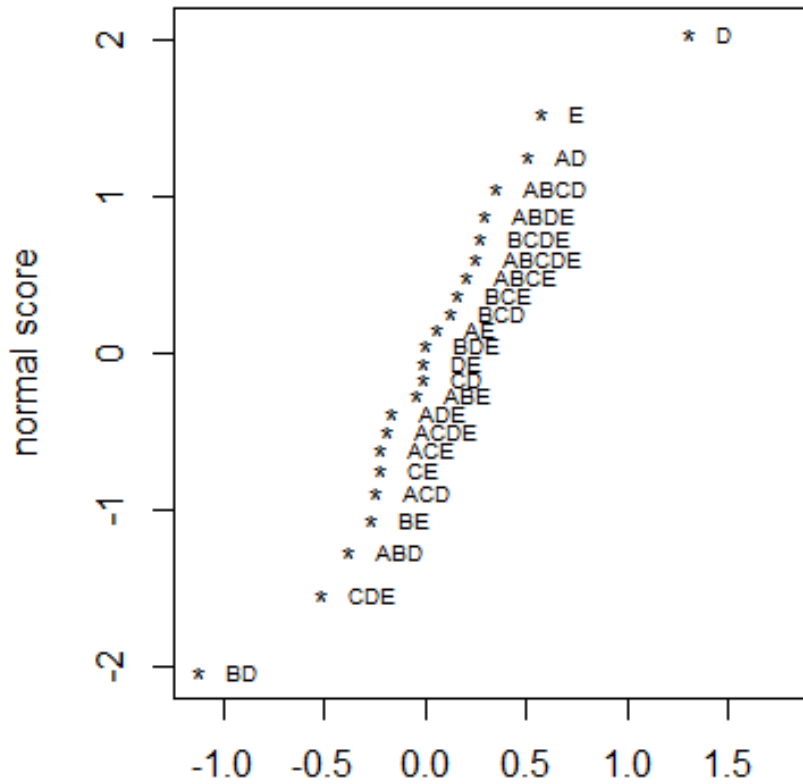
Observation	Run nr:					Original
	1	2	3	4	5	
8	7.5878	6.4237	7.0643	6.3209	7.1783	6.1
22	7.1397	7.0891	7.7956	6.7407	6.5224	6.4
23	8.6784	7.1870	6.6056	6.9442	7.2784	6.2
24	6.2343	6.0161	6.3656	6.1624	6.3194	6.6

Some of the estimated values exceed 7, which was the highest score a combination could gain in the original experiment. As stated, this is not a typical censoring experiment, thus all the values are acceptable.

The estimated whole-plot effects are shown in Figure 6.8(a). Factor A is clearly significant. Factor C might be important, but it lies closer to the linear trend, thus it is more difficult to decide. Figure 6.8(b) shows the subplot effects. Factor D and the interaction BD stand out as the most important of the effects. As expected, the $\mathbf{S}_{\text{III}}^{+R}$ effects are not significant for the outcome of the product.



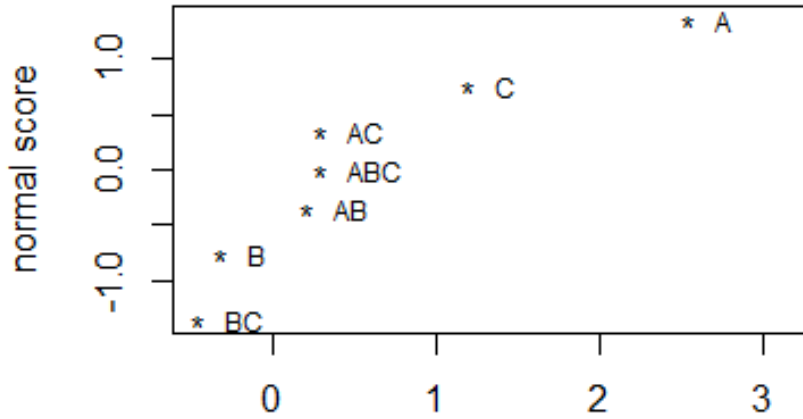
(a) Whole-plot effects



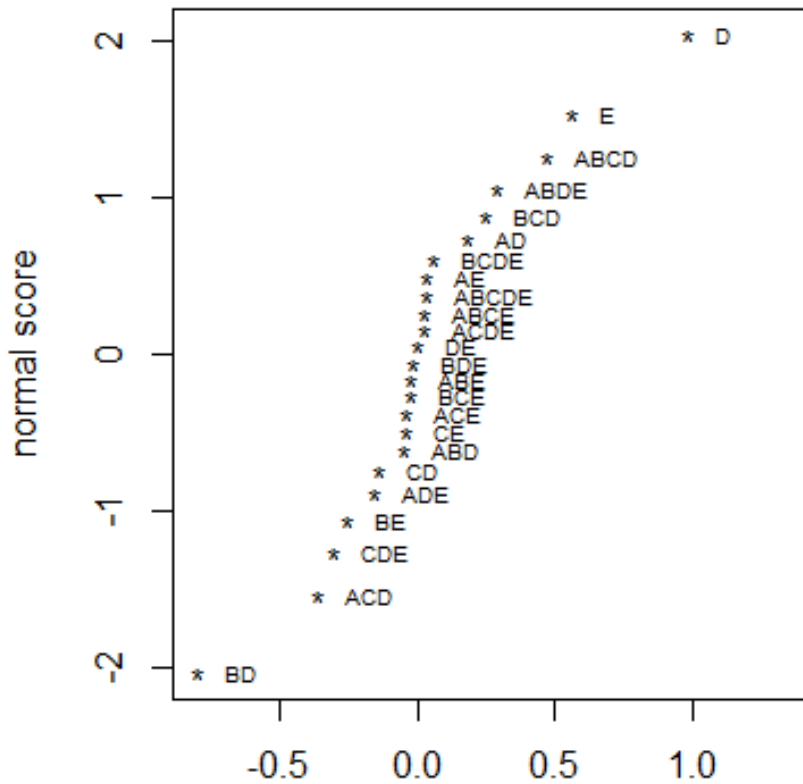
(b) Subplot effects

Figure 6.8: Normal plot of the right censored effects of Example III using multiple imputation.

Figure 6.9 shows the results of the QnD method for right censoring of Example III. As in Figure 6.8, the whole-plot effects from the QnD analysis, Figure 6.9(a), shows that the factors A and C are the most significant. The same applies to the subplot effects, Figure 6.9(b), where factor D and the interaction BD stand out as the most important.



(a) Whole-plot effects



(b) Subplot effects

Figure 6.9: Normal plot of the estimated effects with right censoring by QnD from Example III.

The whole-plot gained from the QnD analysis is very similar to the original estimates. The reason for this is that the censoring limit is close to the original values. The difference between the estimated values in Table 6.9 and the original values are generally higher than that between the latter and the censoring limit. In this case, the QnD method provides the best estimates. The numerical results can be found in Appendix A, Table A.3.

6.6 Example III - Left censoring

For the left censoring of Example III, the censoring limit is set to 1.6. Five observations are censored by this limit. The factors used to estimate the optimal pair of $\hat{\sigma}_w^2$ and $\hat{\sigma}_s^2$ are: $\mathbf{W}_{\text{III}}^L = [A, C, BC]$, $\mathbf{S}_{\text{III}}^{-L} = [D, E, BD, BE, ABCD]$.

As for the right censoring case, the factors for $\mathbf{S}_{\text{III}}^{+L}$ are omitted due to low significance. The values for estimated variances and the correlation in the left censoring case are shown in Table 6.10.

Table 6.10: Estimated values for the variances and the correlation for the left censored data in Example III, for the maximum likelihood with MI method.

Whole-plot variance:	$\hat{\sigma}_w^2$	0.0275
Subplot variance:	$\hat{\sigma}_s^2$	0.2349
Correlation:	ρ	0.1048

All the censored observations are within three whole-plot combinations, with three observations in one combination, and singles in the others. Eq. (3.4) can therefore be used straight forward. The estimated values for the censored observations are listed in Table 6.11. Just like in Section 6.5, some of the estimated values are outside the scale. As previously stated, this is allowed, since the experiment is no ordinary censoring experiment.

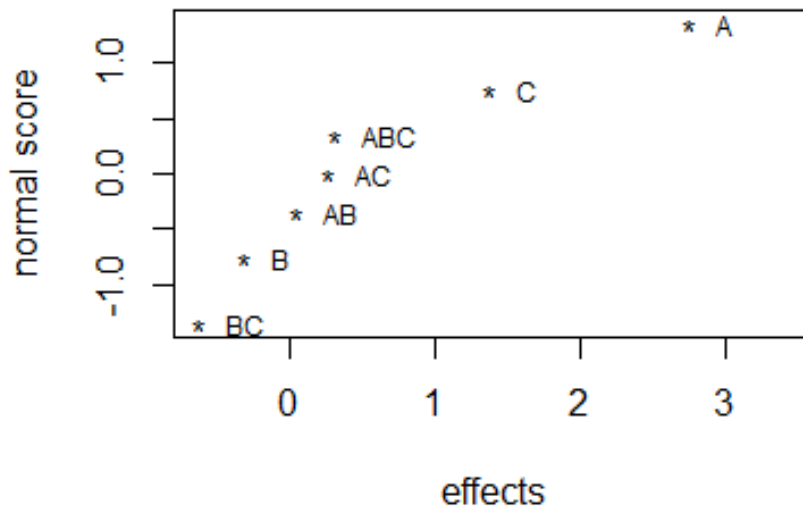
Table 6.11: Estimated values from multiple imputation for the left censored data in Example III.

Observation	Run nr:					Original
	1	2	3	4	5	
1	0.6179	-0.2167	-1.3438	-0.5575	-0.0242	1.1
2	1.2539	1.2329	1.1162	1.0510	0.7223	1.4
3	0.0387	0.9129	0.8304	0.4496	0.5598	1.0
10	1.4158	1.3008	1.5372	1.5630	1.5179	1.6
25	1.4466	1.5751	1.4308	1.5525	1.1524	1.6

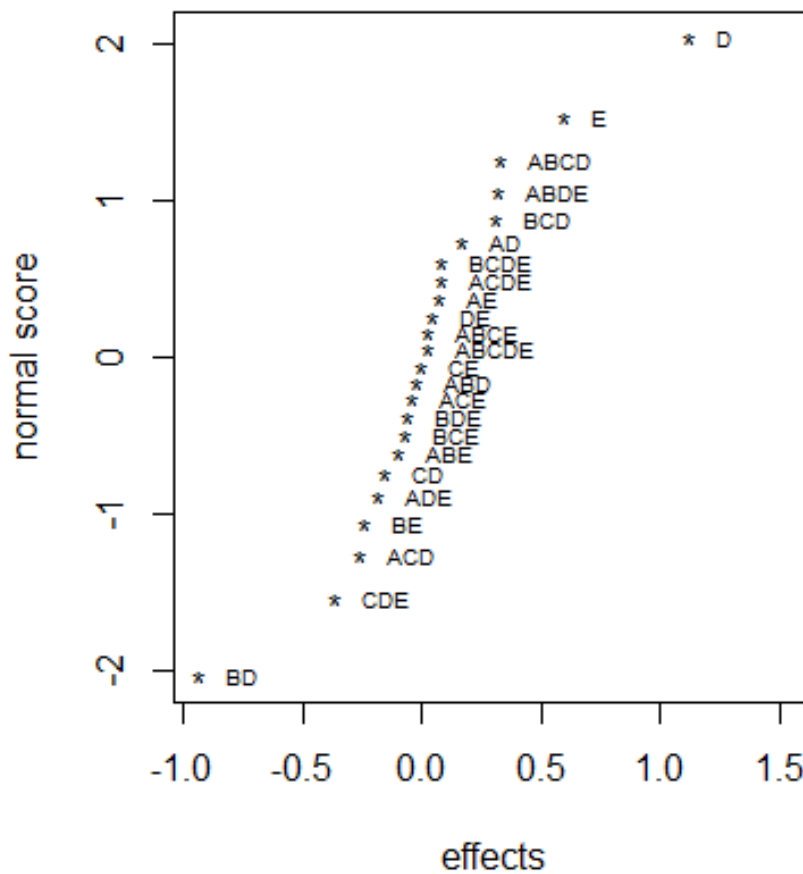
The resulting whole-plot effects are shown in Figure 6.10(a), and the subplot effects in Figure 6.10(b). These are very similar to the ones from the right censoring case, Figure 6.8. The only difference is the order of the factors with the linear trend. Thus, the factors A, C, D and BD are the most important for the outcome of the product.

The results of the QnD analysis for the left censoring of Example III are shown in Figure 6.11. The plot of the whole-plot effects, Figure 6.11(a), shows that the factors A and C are the most significant. The corresponding plot of the subplot effects shows that the factor D and interaction BD are the most important.

For both methods, the estimated subplot effects are good compared to the originals. The whole-plots are good as well, but the linear trend differs in both cases. In Figure 6.10(a) factor C stands out from the linear trend, therefore it is significant for the result. Factor C in Figure 6.11(a) lies closer to the linear trend. In



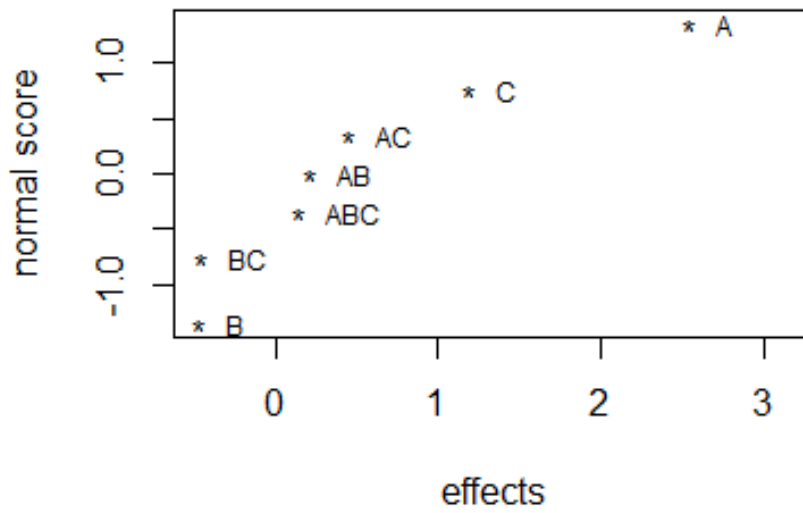
(a) Whole-plot effects



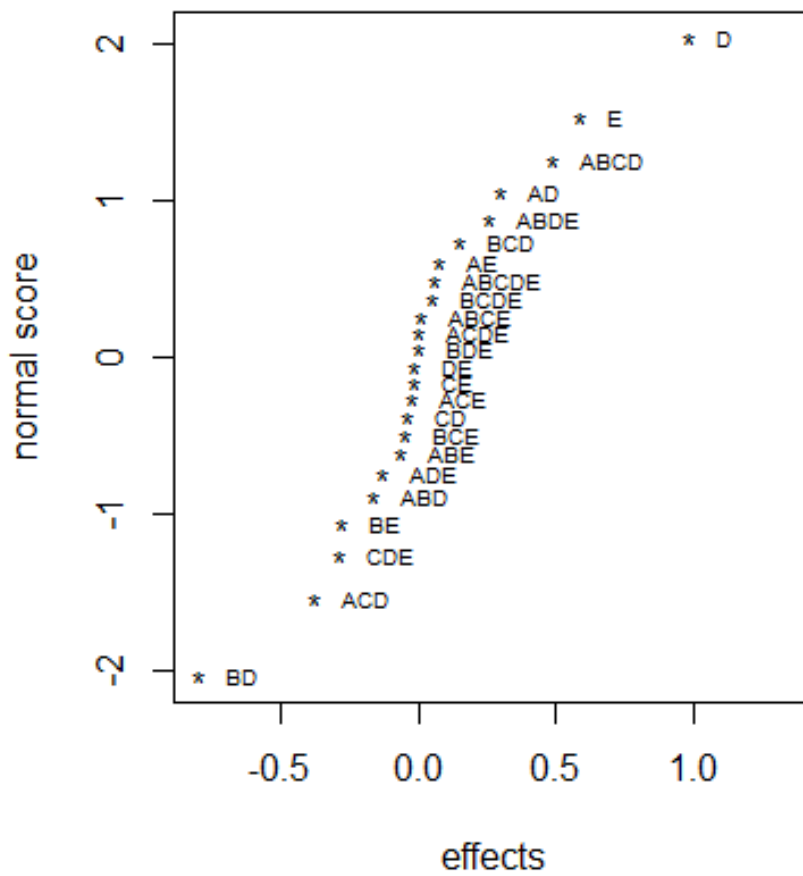
(b) Subplot effects

Figure 6.10: Normal plot of the left censored estimated effects of Example III using multiple imputation.

both plots, C is found to be significant. The numerical results are listed in Table A.3 in Appendix A.



(a) Whole-plot effects



(b) Subplot effects

Figure 6.11: Normal plot of the estimated effects with left censoring by QnD from Example III.

Chapter 7

Discussion

In this chapter, the performances of both methods will be summarized for each of the examples and then evaluated separately.

7.1 Example I

The right limit for the censoring of Example I is set where there is a gap between the response values. This results in a gap between the censoring limit and most of the values that become censored. Since this is an experiment with 16 runs, only three values were censored. Run number 1 and 2 are both censored, and the estimated values for the first, found through multiple imputation, is very close to the original values. The values for run 2 are much higher than the original, thus not a very good fit. The estimated values for run number 10 differs from one simulation to the next, and none of them gives as good estimates compared to the two other censored observations. This is most likely due to that the dependent observations are far below the censoring limit.

The estimates produced by the maximum likelihood with multiple imputation method gives the best results for the effects. The QnD method creates poorer estimates, but it still points out the decidedly largest effects.

Left censoring on Example I did not work well with the multiple imputation method. No combination of factors satisfied the condition $K_2 \geq K_1$. The reason for this is that dependent observations, in one or more cases, had a big difference between one another. The estimated K_1 becomes much larger than K_2 for each tried censoring limit, and the procedure fails to execute. This shows that there are pitfalls with the proposed procedure. To get around this, a dataset should be tested before it's analysed. If there are two observations within each whole-plot, half the difference between them should not exceed the mean of the two. There is no guaranty that the procedure will fail if this condition is violated, but it should be viewed as a warning.

The QnD method has no restrictions. There is no natural gap to place the left limit, that is, it was placed to get some censored values. Four values were censored, and the created estimates are good compared to the original effects. This is expected, since the censoring limit is close to most of the original values.

7.2 Example II

Example II is an experiment with 32 runs. The right limit could not be set at a natural value, thus it was placed to get some censored observations. The limit is close to all the censored observations original values. It is then obvious that the QnD method performs well. Maximum likelihood with multiple imputation shows again that the best estimates are those who are made for whole-plot combinations where all observations are censored. This method provides good estimates close to the original values, and the plot of the subplot effects are more similar to the original than the plot found with the QnD method.

In the left censoring case, the limit is placed close to some of the censored values, and far from the others. Both analysing methods manage to find the most significant effects, but none of them give satisfactory results. The linear trends are harder to separate from the significant values, and especially factor E disappears into the linear trend and seems to be absolutely not important for the outcome of the product. The main cause lies probably in that the original values of the dependent observations have very different sizes.

7.3 Example III

The dataset used in Example III has 32 runs, and the response is restricted from 1 to 7. The censoring limits will therefore lie close to the original values no matter where they are placed. It is thereby expected that both methods will perform well.

Conditional distribution is used on all the censored values in the right and left censoring cases. For the maximum likelihood with multiple imputation method this results in estimated values that vary much from run to run. Yet the mean of the estimated effects is somewhat similar to the original effects. The biggest difference lies in the linear trend, which is more conspicuous for the estimated effects. As stated, the censoring limit is close to the original values, thus the QnD method provides good estimates.

7.4 Censoring with maximum likelihood and multiple imputation

For all the examples, maximum likelihood with multiple imputation gives good estimates. This is not surprising, considered that the method imputes m values drawn from a truncated distribution and then uses the average value in computations. In almost all cases it provided the most significant effects that influence the outcome of the experiment. The estimates found by multiple imputation show that if dependent observations have a big gap between the responses, the results of the analyses are influenced and often results in poorer estimates for the effects. On the other hand, if all observations within a whole-plot combination are censored, the simulated values vary little from run to run and are usually closer to the original values.

To acquire the lowest set of $\hat{\sigma}_w^2$ and $\hat{\sigma}_s^2$, in some cases there were used higher order interactions. Normally, these are considered to be zero, but when testing all

possible combinations, some of them reduced the size of either $\hat{\sigma}_w^2$ or $\hat{\sigma}_s^2$.

Even though this procedure seems to work well, there are some pitfalls. If two censored values are dependent, but lie far apart, the error term is drastically affected. In some cases, like left censoring of Example I, the $\hat{\sigma}_w^2$ becomes negative. Thus, the procedure in this thesis will fail to estimate values for the factors effects. A condition is proposed in Section 7.1, to check if a dataset should not be analysed with this method. Also, if the true failure times within a whole-plot lie far apart, estimates of censored values should not be blindly trusted. Of course, there is no way to know this beforehand. On the other hand, the analyst should be aware of cases where one value is far from the censoring limit, while the other is censored.

7.5 Censoring with quick and dirty

The different examples have a different number of censored observations. The fewer observations censored and the closer the censoring limit is to the original values, the better the QnD method works.

The QnD method performs well for all the examples. In all the cases, the censoring limit is either close to the original values or there is a small number of censored observations. If this is not the case, the QnD method should not be trusted. This is supported by the left censoring of Example I, Section 6.2, and earlier research^[12]. The method's greatest aspect is that it is fast and easy to implement. The disadvantage is that it is not possible to know in advance if the results obtained are satisfactory or not, i.e. there is always a possibility of obtaining poor results. When this is said, the QnD method is an appropriate method when the number of resources is limited, and the expected level of accuracy is not strict.

Chapter 8

Conclusion

In this thesis, it has been investigated how to manage censored data from a split-plot experiment. Three examples have been tested, all with both left and right censoring, and two methods have been implemented; the quick and dirty method and the maximum likelihood with multiple imputation method. The results through these test show that both methods give reasonable results on every example. Even though the maximum likelihood method with multiple imputation does not produce estimates that are very close to the original values, this method is considered the best of the two tested in this thesis. The number of runs in each experiment and the number of censored data are comprehensible, thus the computations could be done manually. If an experiment requires a larger number of runs, and the number of censored values is kept low, the methods used in this thesis could easily be transferred with acceptable results.

If a dataset in split-plot manner is censored and a positive variance can be estimated, then the procedure in this thesis might give a good estimate for the censored observations. It is not possible to draw a final conclusion for the general case beyond this, since the number of datasets, models and censoring limits should be incredibly higher to verify these procedures. However, the analysis gives a good indication of what effects that are most important for the outcome of an experiment.

Bibliography

- [1] Jones, B. and Nachtsheim, C.J., *Split-plot Designs: What, why and how*, Journal of Quality Technology, Vol. 41, No. 4, October 2009.
- [2] Sue-Chu, Arja M Stout, *Multiple Imputation for Censored Data using Experimental Design*, NTNU, Autumn 2012.
- [3] Støtvig, J. Gunnekleiv, *Censored Weibull Distributed Data in Experimental Design*, NTNU, February 2014.
- [4] Rubin, D.B., *Statistical Analysis with Missing Data*, New York, Wiley, 1987.
- [5] Box, G. and Jones, S., *Split-Plots for Robust Product and Process Experimentation*, Report No. 178, *Center for Quality and Productivity Improvement*, University of Wisconsin-Madison, 2000.
- [6] Landois, Luis Leon, *Generalized Least Squares Analysis of the Split-Plot Model using an Estimated Variance-Covariance Matrix*, Inst. of Statistics Mimeo Series #1352, 1981.
- [7] Tyssedal, J. and Kulahci, M., *Analysis of Split-Plot Designs with Mirror Image Pairs as Sub-Plots*, 2005.
- [8] Rencher, Alvin C., Christensen, William F., *Methods of Multivariate Analysis*, Wiley, 2012.
- [9] Chopin, Nicolas, *Fast simulation of truncated Gaussian distribution*, Springer Science & Business Media, 2010.
- [10] Montgomery, Douglas C., *Design and Analysis of Experiments*, Wiley, 2012
- [11] Mee, Robert, *A Comprehensive Guide to Factorial Two-Level Experimentation*, Springer Science & Business Media, 2009.
- [12] Hansen, M. Nevland, *Analysis of censored data from split-plot design with mirror image pairs*, NTNU, Spring 2014.
- [13] *The R Project for Statistical Computing*, www.r-project.org, www.rstudio.com.

Appendix A

Numerical results

Table A.1: Coefficients of the factors for the different estimations of Example I.

		Original:	MI: Right	QnD: Left	Right
W:	(I)	49.7419	49.8369	50.8200	47.7150
	A	10.0656	10.1607	9.7688	8.0388
	B	5.6244	5.7194	5.9213	3.5975
	C	1.3281	2.5761	1.4913	1.9375
	AB	7.3431	7.4382	6.2650	5.3163
	AC	1.8344	3.0824	1.3700	2.4438
	BC	-0.6044	0.6436	-0.1400	0.0050
	DE	0.3369	1.5849	0.1738	0.9463
S:	D	-2.0694	-0.7925	-2.5338	-1.4213
	E	5.1231	5.2470	5.4200	3.1350
	AD	-3.0056	-1.7288	-2.8425	-2.3575
	AE	6.5019	6.6257	5.4238	4.5138
	BD	-1.3619	-0.0850	-1.5250	-0.7138
	BE	-0.2094	-0.0855	0.8688	-2.1975
	CD	1.8694	1.9932	1.5725	-0.1188
	CE	-1.4881	-0.2113	-1.0238	-0.8400

Table A.2: Coefficients of the factors for the different estimations of Example II.

		Original:	MI:		QnD:	
			Right	Left	Right	Left
W:	(I)	40.9813	40.7693	41.4239	40.5531	42.6906
	A	5.9125	6.0001	5.4699	6.1531	4.2031
	B	2.1125	2.0906	0.4963	1.9969	1.3281
	C	-1.6938	-1.4818	-2.9092	-1.2656	-1.6656
	D	-7.5500	-7.4624	-7.1074	-7.3094	-5.8406
	AB	-2.1063	-1.9600	-0.4901	-1.8031	-1.3219
	AC	1.4875	1.3999	2.7030	1.2469	1.4594
	BC	-0.4250	-0.4031	-0.3831	-0.3094	0.4719
	AD	8.2813	8.0693	7.8386	7.8531	6.5719
	BD	-1.6562	-1.5100	-3.2724	-1.3531	-2.4406
	CD	0.8375	0.7499	-0.3780	0.5969	0.8656
	ABC	1.4312	1.2850	1.3894	1.1281	0.5344
	ABD	-1.6500	-1.6719	-0.0338	-1.7656	-0.8656
	ACD	-1.1563	-0.9443	0.0592	-0.7281	-1.1844
	BCD	0.6187	0.4725	0.6606	0.3156	1.5156
ABCD	3.4250	3.4469	3.3831	3.5406	2.5281	
S:	E	1.5688	1.4565	0.5152	1.3031	0.4406
	AE	-2.9500	-2.8523	-1.8964	-2.6469	-1.8219
	BE	-0.1500	-0.1819	0.1605	-0.3281	0.1906
	CE	-0.0688	0.0435	-0.2670	0.1969	-0.3781
	DE	0.5125	0.6102	-0.5411	0.8156	-0.6156
	ABE	0.0563	0.1028	-0.2542	0.1969	-0.2844
	ACE	-0.0875	-0.1852	0.1107	-0.3906	0.2219
	ADE	-0.4063	-0.5185	0.6473	-0.6719	0.7219
	BCE	0.4500	0.4819	-0.0949	0.6281	-0.0281
	BDE	-0.0938	-0.0472	0.2167	0.0469	0.2469
	CDE	0.1625	0.0648	-0.0357	-0.1406	-0.1469
	ABDE	0.1375	0.1056	-0.1730	-0.0406	-0.2031
	ABCE	-0.2188	-0.2653	0.3261	-0.3594	0.2594
	ACDE	-0.1313	-0.0190	0.0670	0.1344	0.1781
	BCDE	0.4438	0.3972	-0.1011	0.3031	-0.0344
ABCDE	0.1250	0.1569	0.6699	0.3031	0.6031	

Table A.3: Coefficients of the factors for the different estimations of Example III.

		Original:	MI:		QnD:	
			Right	Left	Right	Left
W:	(I)	3.4781	3.5478	3.4008	3.4375	3.5188
	A	1.3156	1.3853	1.3929	1.2750	1.2750
	B	-0.1969	-0.2666	-0.1384	-0.1563	-0.2375
	C	0.6344	0.6531	0.7011	0.6000	0.5938
	AB	0.0656	-0.0041	0.0072	0.1063	0.1063
	AC	0.1844	0.2031	0.1176	0.1500	0.2250
	BC	-0.2656	-0.2844	-0.3346	-0.2313	-0.2250
	ABC	0.1094	0.0906	0.1783	0.1438	0.0688
S⁻:	D	0.5219	0.5904	0.5705	0.4938	0.4938
	E	0.2969	0.3079	0.3465	0.2813	0.2938
	AD	0.1219	0.1904	0.0732	0.0938	0.1500
	AE	0.0344	0.0454	-0.0153	0.0188	0.0375
	BD	-0.4281	-0.4966	-0.4746	-0.4000	-0.4000
	BE	-0.1406	-0.1517	-0.1714	-0.1250	-0.1375
	CD	-0.0469	-0.0293	-0.0850	-0.0688	-0.0188
	CE	-0.0094	-0.0493	-0.0485	-0.0188	-0.0063
	ABD	-0.0531	-0.1216	-0.0067	-0.0250	-0.0813
	ABE	-0.0281	-0.0392	0.0027	-0.0125	-0.0313
	ACD	-0.1594	-0.1418	-0.1212	-0.1813	-0.1875
	ACE	-0.0094	-0.0493	0.0298	-0.0188	-0.0125
	BCD	0.1031	0.0856	0.1601	0.1250	0.0750
	BCE	-0.0219	0.0180	0.0195	-0.0125	-0.0250
	ABCD	0.2156	0.1981	0.1586	0.2375	0.2438
ABCE	0.0031	0.0430	-0.0382	0.0125	0.0063	
S⁺:	DE	0.0031	0.0129	-0.0179	3.331e-16	-0.0063
	ADE	-0.0719	-0.0621	-0.0508	-0.0750	-0.0625
	BDE	-0.0094	-0.0192	0.0095	-0.0063	4.398e-16
	CDE	-0.1531	-0.1942	-0.1426	-0.1500	-0.1438
	ABDE	0.1406	0.1308	0.1218	0.1438	0.1313
	ACDE	0.0094	-0.0317	-0.0012	0.0125	-1.124e-16
	BCDE	0.0344	0.0755	0.0050	0.0313	0.0250
	ABCDE	0.0219	0.0630	0.0512	0.0188	0.0313

Appendix B

Code - Example II

The computations for the three examples are very similar. Therefore, only the R code for Example II is included in this thesis.

```
library(censReg)
library(FrF2)
library(plyr)
library(truncnorm)
#####

# Set up the original experiment

A = c(-1,-1,-1,-1,-1,-1,-1,-1,-1,-1,-1,-1,-1,-1,-1,-1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1)
B = c( 1, 1, 1, 1,-1,-1,-1,-1, 1, 1, 1, 1,-1,-1,-1,-1,
-1,-1,-1,-1, 1, 1, 1, 1, 1, 1, 1, 1,-1,-1,-1,-1)
C = c(-1,-1,-1,-1,-1,-1,-1,-1, 1, 1, 1, 1, 1, 1, 1, 1,
-1,-1,-1,-1,-1,-1,-1,-1, 1, 1, 1, 1, 1, 1, 1, 1)
D = c(-1, 1,-1, 1,-1, 1,-1, 1,-1, 1,-1, 1,-1, 1,-1, 1,
 1,-1, 1,-1,-1, 1,-1, 1, 1,-1, 1,-1,-1, 1,-1, 1)
E = c(-1,-1, 1, 1,-1,-1, 1, 1,-1,-1, 1, 1,-1,-1, 1, 1,
-1,-1, 1, 1,-1,-1, 1, 1,-1,-1, 1, 1,-1,-1, 1, 1)
Y = c(55.8, 25.6, 62.9, 33.0, 48.6, 5.0, 57.0, 18.1, 47.2, 11.3,
54.6, 23.9, 37.6, 13.3, 43.5, 23.7, 56.8, 41.2, 56.2, 38.2, 53.5,
41.8, 51.3, 37.8, 49.5, 48.7, 48.2, 44.4, 47.2, 47.5, 44.8, 43.2)

Ex2 <- as.data.frame(matrix(0,32,6))
colnames(Ex2) = c("A", "B", "C", "D", "E", "Y")
Ex2$A = A
Ex2$B = B
Ex2$C = C
Ex2$D = D
Ex2$E = E
Ex2$Y = Y
```

```

Left = 25 # Left censoring limit
Right = 55 # Right censoring limit

Ex2 <- arrange(Ex2, A, B, C, D)

#####
# Big dataset - Expand Ex2
#####

BigEx2<- as.data.frame(matrix(0,nrow(Ex2),32))
colnames(BigEx2) = c("A", "B", "C", "D", "E", "AB", "AC", "BC", "AD", "AE",
"BD", "BE", "CD", "CE", "DE", "ABC", "ABD", "ABE", "ACD", "ACE", "ADE", "BCD",
"BCE", "BDE", "CDE", "ABCD", "ABDE", "ABCE", "ACDE", "BCDE", "ABCDE", "Value")
for ( i in 1:(nrow(Ex2)))
{
  A = Ex2$A[i]
  B = Ex2$B[i]
  C = Ex2$C[i]
  D = Ex2$D[i]
  E = Ex2$E[i]

  BigEx2$A[i] = A
  BigEx2$B[i] = B
  BigEx2$C[i] = C
  BigEx2$D[i] = D
  BigEx2$E[i] = E
  BigEx2$Value[i] = Ex2$Y[i]

  BigEx2$AB[i] = A*B
  BigEx2$AC[i] = A*C
  BigEx2$BC[i] = B*C
  BigEx2$AD[i] = A*D
  BigEx2$AE[i] = A*E
  BigEx2$BD[i] = B*D
  BigEx2$BE[i] = B*E
  BigEx2$CD[i] = C*D
  BigEx2$CE[i] = C*E
  BigEx2$DE[i] = D*E

  BigEx2$ABC[i] = A*B*C
  BigEx2$ABD[i] = A*B*D
  BigEx2$ABE[i] = A*B*E
  BigEx2$BCD[i] = B*C*D
  BigEx2$BCE[i] = B*C*E
  BigEx2$ACD[i] = A*C*D
  BigEx2$ACE[i] = A*C*E
  BigEx2$ADE[i] = A*D*E

```

```

BigEx2$BDE[i] = B*D*E
BigEx2$CDE[i] = C*D*E

BigEx2$ABCD[i] = A*B*C*D
BigEx2$ABDE[i] = A*B*D*E
BigEx2$ABCE[i] = A*B*C*E
BigEx2$ACDE[i] = A*C*D*E
BigEx2$BCDE[i] = B*C*D*E

BigEx2$ABCDE[i] = A*B*C*D*E
}
E2 <- BigEx2
#####

# Create whole and sub matrices

Whole <- as.data.frame(matrix(0,(nrow(E2)/2),ncol(E2)))
colnames(Whole) = c("A", "B", "C", "D", "AB", "AC", "AD", "BC", "BD", "CD", "ABC",
"ABD", "ACD", "BCD", "ABCD")
Sub <- as.data.frame(matrix(0,(nrow(E2)/2),ncol(E2)))
colnames(Sub) = c("E", "AE", "BE", "CE", "DE", "ABE", "ACE", "ADE", "BCE", "BDE",
"CDE", "ABCE", "ABDE", "ACDE", "BCDE", "ABCDE")

for(i in 1:(nrow(E2)/2))
{
  Whole$A[i] = (E2$A[2*i-1] + E2$A[2*i])/2
  Whole$B[i] = (E2$B[2*i-1] + E2$B[2*i])/2
  Whole$C[i] = (E2$C[2*i-1] + E2$C[2*i])/2
  Whole$D[i] = (E2$D[2*i-1] + E2$D[2*i])/2

  Whole$AB[i] = (E2$AB[2*i-1] + E2$AB[2*i])/2
  Whole$AC[i] = (E2$AC[2*i-1] + E2$AC[2*i])/2
  Whole$AD[i] = (E2$AD[2*i-1] + E2$AD[2*i])/2
  Whole$BC[i] = (E2$BC[2*i-1] + E2$BC[2*i])/2
  Whole$BD[i] = (E2$BD[2*i-1] + E2$BD[2*i])/2
  Whole$CD[i] = (E2$CD[2*i-1] + E2$CD[2*i])/2

  Whole$ABC[i] = (E2$ABC[2*i-1] + E2$ABC[2*i])/2
  Whole$ABD[i] = (E2$ABD[2*i-1] + E2$ABD[2*i])/2
  Whole$ACD[i] = (E2$ACD[2*i-1] + E2$ACD[2*i])/2
  Whole$BCD[i] = (E2$BCD[2*i-1] + E2$BCD[2*i])/2
  Whole$ABCD[i] = (E2$ABCD[2*i-1] + E2$ABCD[2*i])/2
  Whole$Value[i] = (E2$Value[2*i-1] + E2$Value[2*i])/2

  Sub$E[i] = (E2$E[2*i-1] - E2$E[2*i])/2
  Sub$AE[i] = (E2$AE[2*i-1] - E2$AE[2*i])/2
  Sub$BE[i] = (E2$BE[2*i-1] - E2$BE[2*i])/2

```

```

Sub$CE[i] = (E2$CE[2*i-1] - E2$CE[2*i])/2
Sub$DE[i] = (E2$DE[2*i-1] - E2$DE[2*i])/2

Sub$ABE[i] = (E2$ABE[2*i-1] - E2$ABE[2*i])/2
Sub$ACE[i] = (E2$ACE[2*i-1] - E2$ACE[2*i])/2
Sub$ADE[i] = (E2$ADE[2*i-1] - E2$ADE[2*i])/2
Sub$BCE[i] = (E2$BCE[2*i-1] - E2$BCE[2*i])/2
Sub$BDE[i] = (E2$BDE[2*i-1] - E2$BDE[2*i])/2
Sub$CDE[i] = (E2$CDE[2*i-1] - E2$CDE[2*i])/2

Sub$ABCE[i] = (E2$ABCE[2*i-1] - E2$ABCE[2*i])/2
Sub$ABDE[i] = (E2$ABDE[2*i-1] - E2$ABDE[2*i])/2
Sub$ACDE[i] = (E2$ACDE[2*i-1] - E2$ACDE[2*i])/2
Sub$BCDE[i] = (E2$BCDE[2*i-1] - E2$BCDE[2*i])/2
Sub$ABCDE[i] = (E2$ABCDE[2*i-1] - E2$ABCDE[2*i])/2
Sub$Value[i] = (E2$Value[2*i-1] - E2$Value[2*i])/2
}

Wholelm = lm(Value ~ A+B+C+D+AB+AC+BC+AD+BD+CD+ABC+ABD+ACD+BCD+
ABCD, data=Whole) #original whole-plot estimates
Sublm = lm(Value ~ E+AE+BE+CE+DE+ABE+ACE+ADE+BCE+BDE+CDE+ABDE+
ABCE+ACDE+BCDE+ABCDE-1, data=Sub) #original subplot estimates

#####
# Quick and Dirty - Right censoring
#####

QnDRight<- E2
for(i in 1:nrow(QnDRight)) {
  if(QnDRight$Value[i] >= Right)
  {QnDRight$Value[i] = Right}
  next
}

QnDRSub <- Sub
QnDRWhole <- Whole
for(i in 1:(nrow(QnDRight)/2))
{
  QnDRSub$Value[i] = (QnDRight$Value[2*i-1] - QnDRight$Value[2*i])/2
  QnDRWhole$Value[i] = (QnDRight$Value[2*i-1] + QnDRight$Value[2*i])/2
}

QnDRrightlm_S = lm(Value ~ E+AE+BE+CE+DE+ABE+ACE+ADE+BCE+BDE+CDE+ABDE+
ABCE+ACDE+BCDE+ABCDE-1, data=QnDRSub)
QnDRrightlm_W = lm(Value ~ A+B+C+D+AB+AC+BC+AD+BD+CD+ABC+ABD+ACD+BCD+
ABCD, data=QnDRWhole)

```

```
#####
# Quick and Dirty - Left censoring
#####

QnDLeft<- E2

for(i in 1:nrow(QnDLeft)) {
  if(QnDLeft$Value[i] <= Left)
  {QnDLeft$Value[i] = Left}
  next
}

QnDLSub <- Sub
QnDLWhole <- Whole
for(i in 1:(nrow(QnDLeft)/2))
{
  QnDLSub$Value[i] = (QnDLeft$Value[2*i-1] - QnDLeft$Value[2*i])/2
  QnDLWhole$Value[i] = (QnDLeft$Value[2*i-1] + QnDLeft$Value[2*i])/2
}

QnDLeftlm_W = lm(Value ~ A+B+C+D+AB+AC+BC+AD+BD+CD+ABC+ABD+ACD+BCD+
ABCD, data=QnDLWhole)
QnDLeftlm_S = lm(Value ~ E+AE+BE+CE+DE+ABE+ACE+ADE+BCE+BDE+CDE+ABDE+
ABCE+ACDE+BCDE+ABCDE-1, data=QnDLSub)

#####
# Estimating variance
#####

censRegRight <- censReg(Value ~ BCDE+AE+A+D+AD+ABCD+ABD+B+AC+DE+ADE+
ACE+E, data = E2, left = -Inf, right = Right) # Right censoring
E2R = coef(censRegRight)
NoWhole = 7 # number of whole factors
NoSub = 6 # number of sub factors

# Sub-plot variance
VarS <- as.data.frame(matrix(0,16, 5))
colnames(VarS) = c("Sub", "No1", "No2", "Diff", "X")
for( i in 1:nrow(VarS))
{
  VarS$No1[i] = E2$Value[2*i-1]
  VarS$No2[i] = E2$Value[2*i]
  SubR = E2R[2]*E2$BCDE[2*i-1] + E2R[3]*E2$AE[2*i-1] + E2R[11]*E2$DE[2*i-1] +
E2R[12]*E2$ADE[2*i-1] + E2R[13]*E2$ACE[2*i-1] + E2R[14]*E2$E[2*i-1]

  VarS$Sub[i] = SubR
  Diff = (E2$Value[2*i-1] - E2$Value[2*i])/2
}
```

```

VarS$Diff[i] = Diff
VarS$X[i] = Diff-SubR

VarS$Diff[i] = signif(as.numeric(VarS$Diff[i]), digits=4)
VarS$X[i] = signif(as.numeric(VarS$X[i]), digits=4)

if(VarS$No1[i] >=Right)
{VarS$Diff[i] = "NULL"
  VarS$No1[i] = "NULL"
  VarS$X[i] = "Cen"

  if(VarS$No2[i] >= Right)
  {VarS$No2[i] = "NULL"
    next}
}
if(VarS$No2[i] >= Right)
{VarS$Diff[i] = "NULL"
  VarS$No2[i] = "NULL"
  VarS$X[i] = "Cen"
  next}
}

# Whole-plot variance
VarW <- as.data.frame(matrix(0,16, 5))
colnames(VarW) = c("Add", "No1", "No2", "Mean", "X")
for( i in 1:nrow(VarW))
{
  VarW$No1[i] = E2$Value[2*i-1]
  VarW$No2[i] = E2$Value[2*i]
  Add = E2R[1] + E2R[4]*E2$A[2*i-1] + E2R[5]*E2$D[2*i-1] + E2R[6]*E2$AD[2*i-1] +
  E2R[7]*E2$ABCD[2*i-1] + E2R[8]*E2$ABD[2*i-1] + E2R[9]*E2$B[2*i-1] +
  E2R[10]*E2$AC[2*i-1]
  VarW$Add[i] = Add
  Mean = (E2$Value[2*i-1] + E2$Value[2*i])/2
  VarW$Mean[i] = Mean
  VarW$X[i] = Mean-Add

  VarW$Mean[i] = signif(as.numeric(VarW$Mean[i]), digits=4)
  VarW$X[i] = signif(as.numeric(VarW$X[i]), digits=4)

  if(VarW$No1[i] >=Right)
  {VarW$Mean[i] = "NULL"
    VarW$No1[i] = "NULL"
    VarW$X[i] = "Cen"

    if(VarW$No2[i] >= Right)
    {VarW$No2[i] = "NULL"

```



```

    next}
  }

  if(VarW$No2[i] >= Right)
  {VarW$Mean[i] = "NULL"
  VarW$No2[i] = "NULL"
  VarW$X[i] = "Cen"
  next}
}

censRegLeft <- censReg(Value ~ AD+D+AE+A+BD+ABCD+C+AC+ADE+ABCDE+DE+
ABE, data = E2, left = Left, right = Inf) # Left censoring
E2L = coef(censRegLeft)
NoWholeLeft = 7 # number of whole factors
NoSubLeft = 5 # number of sub factors

# Sub-plot
VarSLeft <- as.data.frame(matrix(0,16, 5))
colnames(VarSLeft) = c("Sub", "No1", "No2", "Diff", "X")
for(i in 1:nrow(VarSLeft))
{
  VarSLeft$No1[i] = E2$Value[2*i-1]
  VarSLeft$No2[i] = E2$Value[2*i]
  SubL = E2L[4]*E2$AE[2*i-1] + E2L[10]*E2$ADE[2*i-1] +
E2L[11]*E2$ABCDE[2*i-1] + E2L[12]*E2$DE[2*i-1] + E2L[13]*E2$ABE[2*i-1]

  VarSLeft$Sub[i] = SubL
  Diff = (E2$Value[2*i-1] - E2$Value[2*i])/2
  VarSLeft$Diff[i] = Diff
  VarSLeft$X[i] = Diff-SubL

  VarSLeft$Diff[i] = signif(as.numeric(VarSLeft$Diff[i]), digits=4)
  VarSLeft$X[i] = signif(as.numeric(VarSLeft$X[i]), digits=4)

  if(VarSLeft$No1[i] <= Left)
  {VarSLeft$Diff[i] = "NULL"
  VarSLeft$No1[i] = "NULL"
  VarSLeft$X[i] = "Cen"

  if(VarSLeft$No2[i] <= Left)
  {VarSLeft$No2[i] = "NULL"
  next}
}
if(VarSLeft$No2[i] <= Left)
{VarSLeft$Diff[i] = "NULL"
  VarSLeft$No2[i] = "NULL"

```

```

    VarSLeft$X[i] = "Cen"
  next}
}

# Whole-plot
VarWLeft <- as.data.frame(matrix(0,16, 5))
colnames(VarWLeft) = c("Add", "No1", "No2", "Mean", "X")
for( i in 1:nrow(VarWLeft))
{
  VarWLeft$No1[i] = E2$Value[2*i-1]
  VarWLeft$No2[i] = E2$Value[2*i]
  Add = E2L[1] + E2L[2]*E2$AD[2*i-1] + E2L[3]*E2$D[2*i-1] + E2L[5]*E2$A[2*i-1] +
  E2L[6]*E2$BD[2*i-1] + E2L[7]*E2$ABCD[2*i-1] + E2L[8]*E2$C[2*i-1] +
  E2L[9]*E2$AC[2*i-1]

  VarWLeft$Add[i] = Add
  Mean = (E2$Value[2*i-1] + E2$Value[2*i])/2
  VarWLeft$Mean[i] = Mean
  VarWLeft$X[i] = Mean-Add

  VarWLeft$Mean[i] = signif(as.numeric(VarWLeft$Mean[i]), digits=4)
  VarWLeft$X[i] = signif(as.numeric(VarWLeft$X[i]), digits=4)

  if(VarWLeft$No1[i] <=Left)
  {VarWLeft$Mean[i] = "NULL"
  VarWLeft$No1[i] = "NULL"
  VarWLeft$X[i] = "Cen"

  if(VarWLeft$No2[i] <= Left)
  {VarWLeft$No2[i] = "NULL"
  next}
}
if(VarWLeft$No2[i] <= Left)
{VarWLeft$Mean[i] = "NULL"
  VarWLeft$No2[i] = "NULL"
  VarWLeft$X[i] = "Cen"
  next}
}

Df <- as.data.frame(matrix(0,1,4)) # Degrees of freedom
Df[1,1] = NoWhole
Df[1,2] = NoSub
Df[1,3] = NoWholeLeft
Df[1,4] = NoSubLeft

Variance <- as.data.frame(matrix(0,16,4))
colnames(Variance) = c("WholeR", "SubR", "WholeL", "SubL")

```

```

Variance$WholeR = VarW$X
Variance$SubR = VarS$X
Variance$WholeL = VarWLeft$X
Variance$SubL = VarSLeft$X

K <- as.data.frame(matrix(0,1,4))
colnames(K) = c("KWR", "KSR", "KWL", "KSL")
for(m in 1:ncol(Variance))
{
  EK = 0
  k = 0
  for(i in 1:nrow(Variance))
  {
    if(Variance[i,m] == "Cen")
    {k = k + 1
     next}
    EK = (as.numeric(Variance[i,m]))^2 + EK
  }
  K[1,m] = EK/(nrow(Variance)-k-Df[1,m]) #Estimates K1 and K2
}

```

```

#####
# RIGHT CENSORING
#####

SigSRight = 2*K$KSR[1] #sigma sqared, sub
SigWRight = K$KWR[1] - (K$KSR[1]) #sigma squared, whole
RhoRight = SigWRight/(SigWRight+SigSRight) #rho
VAR = (SigWRight + SigSRight)*(1-RhoRight^2) #variance
SD = sqrt(VAR) #standard deviation
n = 5 #number of runs

Row <- as.data.frame(matrix(0,nrow(E2), 3))
colnames(Row) = c("Row","Adjust", "Expect")
for(i in 1:nrow(E2))
{
  row = E2R[1] + E2R[4]*E2$A[i] + E2R[5]*E2$D[i] + E2R[6]*E2$AD[i] +
  E2R[7]*E2$ABCD[i] + E2R[8]*E2$ABD[i] + E2R[9]*E2$B[i] + E2R[10]*E2$AC[i] +
  E2R[2]*E2$BCDE[i] + E2R[3]*E2$AE[i] + E2R[11]*E2$DE[i] + E2R[12]*E2$ADE[i] +
  E2R[13]*E2$ACE[i] + E2R[14]*E2$E[i]
  Row[i,1] = row
  adjust = E2$Value[i] - row
  Row[i,2] = adjust
  Row[i,3] = adjust*RhoRight #Expected value, mean

  if(E2$Value[i] >= Right)
  {
    Row[i,2] = "Cen"
    Row[i,3] = "Cen"
    next
  }
}

#####
# Multiple imputation - Conditional distribution
#####

# Cen row 2
a = Right - Row[2,1]
mEAn = as.numeric(Row[1,3])

Y1 = rtruncnorm(n, mean = mEAn, sd = SD, a = a, b = Inf)

#####
# TWO ROWS CENSORED
#####

```

```

rho = RhoRight
nu = 1-rho^2

swap <- function(matrix, row1, row2){
  row3 <- matrix[row1,]
  matrix[row1,] = matrix[row2,]
  matrix[row2,] = row3
  return (matrix)
}

t = 4
cens = matrix(0,t,7)
cens[,1] = c(9,10,19,20)

for(i in 1:nrow(cens)){
  cens[i,2] = Right - Row[cens[i,1],1]
}

for(i in 1:(nrow(cens)/2)){
  if(cens[2*i-1,2] > cens[2*i,2]){
    cens <- swap(cens, 2*i-1, 2*i)}
}

for(i in 1:(nrow(cens)/2)){
  a2 = cens[2*i-1,2] # Small
  a1 = cens[2*i,2] # Large
  x = a1*rho - a2

  # Case M+
  if(x <= 0){
    cens[2*i,3:7] = rtruncnorm(n, mean = 0, sd = 1, a = (a2/rho), b = Inf)
    for(k in 3:(n+2)){
      cens[2*i-1,k] = rtruncnorm(1, mean = rho*cens[2*i,k], sd = nu^2, a = a2, b =
    }
  }

  # Case S+
  if(x >= 0){
    cens[2*i,3:7] = rtruncnorm(n, mean = 0, sd = 1, a = a1, b = Inf)
    for(k in 3:(n+2)){
      cens[2*i-1,k] = rnorm(1, mean = rho*cens[2*i,k], sd = nu^2)
      if(cens[2*i-1,k] < a2){
        print("Error: New value required")
        print(cens[2*i-1,1])
      }
    }
  }
}

```

```

}
cens <- as.data.frame(cens)
cens <- arrange(cens,cens[,1])
#####
v = E2$Value
NEW = matrix(v,32,5) # Make a matrix with 5 columns, and change censored values
for(m in 1:n){
  NEW[2,m] = Y1[m] + Row[2,1]
  for (j in 1:nrow(cens)) {
    NEW[cens[j,1],m] = cens[j,m+2] + Row[cens[j,1],1]
  }
}

#####
# Estimating effects
#####

MatWR = as.data.frame(matrix(0, (nrow(E2)/2), n))
MatSR = as.data.frame(matrix(0, (nrow(E2)/2), n))
EffectsWR = 0
EffectsSR = 0

for(m in 1:n){
  for(i in 1:(nrow(E2)/2)){
    MatWR[i,m] = (NEW[2*i -1,m] + NEW[2*i,m])/2
    MatSR[i,m] = (NEW[2*i -1,m] - NEW[2*i,m])/2
  }
  SubR = Sub
  SubR$Value = MatSR[,m]
  lmS = lm(Value ~ E+AE+BE+CE+DE+ABE+ACE+ADE+BCE+BDE+CDE+ABDE+ABCE+ACDE+
            BCDE+ABCDE-1, data = SubR)

  WholeR = Whole
  WholeR$Value = MatWR[,m]
  lmW = lm(Value ~ A+B+C+D+AB+AC+BC+AD+BD+CD+ABC+ABD+ACD+BCD+
            ABCD, data = WholeR)

  EffectsWR = EffectsWR + (coef(lmW)/5)
  EffectsSR = EffectsSR + (coef(lmS)/5)
}

```

```

#####
# Left CENSORING
#####

SigSLeft = 2*K$KSL[1] #sigma squared, sub
SigWLeft = K$KWL[1] - (K$KSL[1]) #sigma squared, whole
RhoLeft = SigWLeft/(SigWLeft+SigSLeft) #rho
VARLeft = (SigWLeft + SigSLeft)*(1-RhoLeft^2) #variance
SD = sqrt(VARLeft) #standard deviation
n = 5 # number of runs

RowLeft <- as.data.frame(matrix(0,nrow(E2), 3))
colnames(RowLeft) = c("Row","Adjust", "Expect")
for(i in 1:nrow(E2)){
  row = E2L[4]*E2$AE[i] + E2L[10]*E2$ADE[i] + E2L[11]*E2$ABCDE[i] +
  E2L[12]*E2$DE[i] + E2L[13]*E2$ABE[i]+ E2L[1] + E2L[2]*E2$AD[i] +
  E2L[3]*E2$D[i] + E2L[5]*E2$A[i] + E2L[6]*E2$BD[i] + E2L[7]*E2$ABCD[i] +
  E2L[8]*E2$C[i] + E2L[9]*E2$AC[i]

  RowLeft[i,1] = row
  adjust = E2$Value[i] - row
  RowLeft[i,2] = adjust
  RowLeft[i,3] = adjust*RhoLeft #Expected value, mean

  if(E2$Value[i] <= Left){
    RowLeft[i,2] = "Cen"
    RowLeft[i,3] = "Cen"
    next
  }
}

#####
# TWO ROWS CENSORED
#####

rho = RhoLeft
nu = 1-rho^2

swap <- function(matrix, row1, row2){
  row3 <- matrix[row1,]
  matrix[row1,] = matrix[row2,]
  matrix[row2,] = row3
  return (matrix)
}

t = 6
censL = matrix(0,t,7)

```

```

censL[,1] = c(3,4,7,8,15,16)

for(i in 1:nrow(censL)){
  censL[i,2] = Left - RowLeft[censL[i,1],1]
}

for(i in 1:(nrow(censL)/2)){
  if(censL[2*i-1,2] > censL[2*i,2]){
    censL <- swap(censL, 2*i-1, 2*i)}
}

for(i in 1:(nrow(censL)/2)){
  a2 = censL[2*i-1,2] # Small
  a1 = censL[2*i,2] # Large
  x = a1*rho - a2

  # Case M+
  if(x <= 0){
    censL[2*i,3:7] = rtruncnorm(n, mean = 0, sd = 1, a = (a2/rho), b = Inf)
    for(k in 3:(n+2)){
      censL[2*i-1,k] = rtruncnorm(1, mean = rho*censL[2*i,k], sd = nu^2, a = a2, b
    )
    }
  }

  # Case S+
  if(x >= 0){
    censL[2*i,3:7] = rtruncnorm(n, mean = 0, sd = 1, a = a1, b = Inf)
    for(k in 3:(n+2)){
      censL[2*i-1,k] = rnorm(1, mean = rho*censL[2*i,k], sd = nu^2)
      if(censL[2*i-1,k] < a2){
        print("Error: New value required")
        print(censL[2*i-1,1])
      }
    }
  }
}

censL <- as.data.frame(censL)
censL <- arrange(censL,censL[,1])
#####
v = E2$Value
NEWLeft = matrix(v,32,5) # Make a matrix with 5 colums, and change censored values
for(m in 1:n){
  for (j in 1:nrow(censL)) {
    NEWLeft[cens[j,1],m] = censL[j,m+2] + RowLeft[censL[j,1],1]
  }
}

```



```

#####
# Estimating effects
#####

MatWL = as.data.frame(matrix(0, (nrow(E2)/2), n))
MatSL = as.data.frame(matrix(0, (nrow(E2)/2), n))
EffectsWL = 0
EffectsSL = 0

for(m in 1:n){
  for(i in 1:(nrow(E2)/2)){
    MatWL[i,m] = (NEWLeft[2*i -1,m] + NEWLeft[2*i,m])/2
    MatSL[i,m] = (NEWLeft[2*i -1,m] - NEWLeft[2*i,m])/2
  }
  SubL = Sub
  SubL$Value = MatSL[,m]
  lmSL = lm(Value ~ E+AE+BE+CE+DE+ABE+ACE+ADE+BCE+BDE+CDE+ABDE+ABCE+ACDE+
             BCDE+ABCDE-1, data = SubL)

  WholeL = Whole
  WholeL$Value = MatWL[,m]
  lmWL = lm(Value ~ A+B+C+D+AB+AC+BC+AD+BD+CD+ABC+ABD+ACD+BCD+
             ABCD, data = WholeL)

  EffectsWL = EffectsWL + (coef(lmWL)/5)
  EffectsSL = EffectsSL + (coef(lmSL)/5)
}

```