# Efficient Handling of Empirical Probability Distributions in RAMS Models

## Florian Marcel Carl Müller

# Preface

This master thesis was written during the spring semester 2017 at the Department of Mechanical and Industrial Engineering, Norwegian University of Science and Technology (Trondheim).

It is part of the two-year international master's program in RAMS (Reliability, Availability, Maintainability and Safety).

The topic is based on the specialization project done in the autumn semester 2016 titled "Comparison of Machine Learning Assisted Non-Parametrical Statistics with Parametrical Approaches". This was modified to focus on the efficient handling aspect of non-parametrical data.

The thesis is supervised by Professor Antoine Rauzy.

It is written for readers with basic understanding in the fields of reliability analysis and the application of failure models.

Trondheim, 2017-06-11

Florian Müller

# Acknowledgement

# Abstract

Performing reliability assessments always relies on utilizing data. Most often, this data is provided in the form of historic failure dates. To understand this data, models are used to derive reliability characteristics from it.

These models can be parametric, trying to describe the system by means of mathematical equations. They can also be empirical, letting the raw data describe the system without assuming a certain outcome.

Handling parametric models is convenient, as they are described by often just one value. Empirical probability distributions are built on all available data and hence requires them to be fully defined. Handling this amount of data is cumbersome.

Part of this thesis is proposing different methods to represent the empirical reliability estimator. These representations try to combine convenient usage while keeping accuracy.

Representing the empirical reliability graph by a reduced amount of linear segments is proposed and discussed. This is an efficient way to compress huge datasets to a low number of descriptive points to interpolate in.

Furthermore, the feasibility to use polynomial regression on the empirical probability distribution is evaluated.

The computational efficiency of all methods is compared. For all practical purposes, the time to retrieve a reliability estimate is negligible.

Parametric and empirical approaches are applied to various datasets and the results discussed. The empirical methods outperform the exponential estimator in all cases.

The given experiment hypothesis is validated on each of the four experiments: The empirical probability distributions do match sufficiently well the reference reliability and the computational efficiency is negligible for all practical purposes.

# Contents

# List of Abbreviations

| Abbreviation | Meaning |
| --- | --- |
| CSV | Character Separated Values |
| ISO | International Organization for Standardization |
| KM | Kaplan-Meier |
| MET | Maximum Event Time |
| MTTF | Mean Time To Failure |
| OREDA | Offshore and Onshore Reliability Data |
| PDS | Pålitelighet av Datamaskinbasert Sikkerhetssystemer |
| | (Norwegian acronym for "Reliability of Computer based Safety systems") |
| PFD | Probability of Failure on Demand |
| RAMS | Reliability, Availability, Maintainability & Safety |
| SIL | Safety Integrity Level |
| SIS | Safety Instrumented Systems |

# Chapter 1

# Introduction

## 1.1 Background

Since the early usage of machinery, stakeholders want to know how long their equipment will last. Systems eventually break down or their performance do no longer meet the requirements.

The process of failure estimation is always a balancing act between being too conservative and too risky. This can lead to either wasting potentially good equipment and money, or endangering people and the environment in case of a hazard.

The emerging field of statistics led to the widespread use of failure distribution models for this task. Parameter based models like the exponential one are commonly known and used. As soon as their parameters are found (given, assumed or estimated), the simple subsequent use is appealing to most engineers.

The selection of those parameters however is the critical step in the process. Its validity is defining the accuracy of all following calculations (Risk assessment, Definition of safety barriers, SIL allocation, etc.). The extent of consequences is often underestimated – a manufacturer stating a value for $\lambda$ in their data sheet does not necessarily know what the customer is using this value for.

When compressing all known information into a limited amount of parameters – mostly only one or two – information gets lost. There are however other methodologies to estimate the failure behavior of equipment. This loss can be avoided by using empirical methods. The information gained by applying this method is however not easy to represent.

With the emerging trend of "Big Data", Industry 4.0 and continuous condition monitoring, better failure prediction is highly sought after. This trend however is also making the matter of efficient data handling more important. The vast amount of generated and available data has to be processed in an intelligent way, often by means of truncation and compression.

## 1.2    Objectives

This thesis is highlighting the differences between parametric and empirical reliability estimation. It further proposes a method to represent the empirical distribution in efficient ways and applies them to numerous examples. Examining the required amount of failure data to give sufficient estimators is an additional objective. An efficiency and accuracy assessment is carried out to compare all approaches.

## 1.3    Limitations

Limitations are given by the lack of real raw data. Plenty of failure rates are available, but access to detailed listings of failure dates and fleet composition are rare. In order to increase the number of experiment cases, some data is self-generated based on mathematical expressions.

This thesis is not focusing on the statistical methodologies to evaluate the fit of various curves in a mathematical way. It uses simple, comprehensible tools to obtain a first evaluation of the proposed methods.

## 1.4    Approach

Different example data sets are presented with various origins and size. The two methods, parametric and empirical, are applied to each dataset and the results compared and discussed.

In order to implement the proposals and carry out the necessary calculations a Python-program is written as part of this thesis.

## 1.5    Structure of the Report

The thesis is opened with highlighting the current literature work around the topic of empirical reliability evaluation in chapter 2.

The following chapter 3 gives a motivating example including some background information on data sources and a case study illustrating the problem.

The conceptual developments are introduced in chapter 4. The proposed data representations are explained within this chapter as well a method to allow prediction beyond observed mission times. It also contains a section about the efficiency measures and results of the methods under evaluation.

The main chapter for the experiments carried out can be seen in chapter 5. Four datasets of varying origin are assessed for their reliability behavior.

A concluding statement and a proposal of future works can be seen at last in chapter 6.

The Python program developed is shown in the Appendix.

# Chapter 2

# Related Works

Parameter-based models and their benefits and disadvantages compared to empirical distribution models are widely discussed in various fields of sciences.

As soon as the unknown parameters are estimated, the parametric models are easy to use. If however the assumptions are false, then the resulting model can be misleading. Non-parametric methods perform better than poorly specified parametric models in almost all cases [21]. Perretti furthermore states "it is best to let the data tell us how the system works without imposing preconceived ideas on the outcome". This is also called the "true model myth" as stated in [8].

As reported by Bobrowski [3], empirical methods are recommended if there is no information available for a possible underlying distribution.

Similar conclusions are also drawn by Mokhtarian [16], stating that when both non-parametric and parametric methods are applicable to a problem, the parametric method is usually preferred because of its efficiency and simple use. However when the assumptions for those parameters are questionable, non-parametric methods are more suitable.

The oil- & gas industry, traditionally sensitive to safety and reliability matters, is following a dedicated standard for the acquisition of reliability data: ISO 14224 – "Petroleum, petrochemical and natural gas industries – Collection and exchange of reliability and maintenance data for equipment". [7] This standard however does not include methods for analysis and applying of reliability and maintenance data. It gives however principles on the calculation of basic reliability parameters in the appendix.

# Chapter 3

# Motivating Example

## 3.1   Why exponential distribution?

Many aspects of understanding the world we live in can be broken down to an essential question: When do events occur?

When will the remaining radiation be decreased to 50%?
When will next earthquake happen?
When will there be another fire in the forest?
$\vdots$

And lastly the question for all RAMS related tasks:
When will the next failure happen?

Describing – and predicting – the failure occurrence behavior of certain events is a major task of statisticians working in probability theory. Probability models are commonly used to achieve those two tasks.

Numerous probability distributions exist. In the RAMS aspect the most commonly [24] used ones are:

- Exponential distribution

- Weibull distribution

- Gamma distribution

- Homogeneous Poisson Processes

Of all available and suitable probability distributions, the exponential distribution is the prevalent one. Various reasons can be found for that:

1. **Ease of use**
   The exponential distribution is described by a single scalar parameter ("Hazard rate", "Failure rate", commonly used symbol: $\lambda$, commonly used unit: $\left[\frac{1}{h}\right]$).
   Having only one variable makes handling of data plain and simple. It can be entered by hand, looked up in tables and books and takes up almost no memory when stored digitally.

2. **Ease of comparison**
   When comparing different systems and components, the failure behavior is an important factor in the decision process. Comparing different failure rates yields instant results: A higher value represents a higher likelihood of a failure in a given time frame and hence a less preferable option.

   Additionally a confidence interval can be given in a generic $[\text{Value}]^{+[\text{Upper Margin}]}_{-[\text{Lower Margin}]}$ style.

3. **Suitability**
   A major criterion for using the exponential distribution is that many systems can be described good enough with it. It gives sufficient accuracy within a given confidence interval, as many setups seem to "follow" the exponential law.

4. **Constant failure rate**
   The nature of the exponential law is its constant failure rate. The probability of having a failure in a time interval of a given length is the same regardless of its position along the time axis. It has no memory. A full explanation of exponential distribution's properties can be found in [24]. This characteristic is required for many reliability tools, algorithms and classifications. This includes methods like Markov Chains, PFD calculations and SIL definition.

## 3.2 Current situation

When designing safety relevant systems, the RAMS engineer is relying on given information. This information is accessible in data sheets directly from the manufacturer, or gathered by collected and processed lifetime data. Latter is provided by their own company, if enough historic data is available, or by agencies specializing in aggregating relevant data.

### 3.2.1 Data sheets

The quality of information found in data sheets is depending on the kind of equipment.

Components for generic industrial or domestic usage often lack information about the failure behavior. Getting this knowledge is a costly and lengthy procedure. This is often not considered required information by the customer, either due to the low cost of the equipment or the short estimated usage time.

More advanced products often give information in the form of MTTF (Mean Time To Failure) values. While this value can be useful to compare competitors, it gives no information about the underlying failure probability distribution – a MTTF value can be calculated for any given

distribution [14]. It is merely the integral of the reliability distribution – see Equation 3.1.

$$MTTF = \int_0^\infty R(t) \tag{3.1}$$

Equipment intended for safety relevant applications give the most insightful information. Depending on the industry sector, some applicable standard might require the indication of certain information. This includes MTTF, failure rate values and naming the standard according to which this data is gathered.

### 3.2.2   Databases / Books

Various databases for failure data exist. Due to the sensitive nature of this data, the access to them is restricted. Access is usually granted by paying a subscription fee. The scope of this thesis is on data collections of the oil and gas industry.

#### 3.2.2.1   OREDA

The OREDA project is sponsored by eight oil and gas companies [17]. The fundamental task of this organization is to aggregate, process and exchange reliability related information. The full set of data is accessible to their partners in an online database and the reduced information level can be found in their printed OREDA handbooks ([18], [19]).

The OREDA organization is also involved in issuing the standard ISO 14224: "Petroleum, petrochemical and natural gas industries – Collection and exchange of reliability and maintenance data for equipment" [7].

The introduction chapter of the handbook describes the origin and processing of the data presented.

An important basic assumption is given early in this chapter:

"*All the failure rate estimates presented in this handbook are therefore based on the assumption that the failure rate function is* constant *and independent of time, in which case $z(t) = \lambda$, i.e. the failure rates are assumed to be exponentially distributed with parameter $\lambda$.*"

The method to calculate the maximum likelihood estimator $\lambda$ based on the collected information is given in Equation 3.2. This is also the method described in [24].

$$\hat{\lambda} = \frac{\text{Number of failures}}{\text{Aggregated time in service}} = \frac{n}{\tau} \tag{3.2}$$

The handbook further introduces the "OREDA estimator", a modified averaging estimator that considers data origins from different installations and samples [26].

### 3.2.2.2 ExproSoft Wellmaster

ExproSoft is a company founded in 2001 as an out-spring from the Sintef research society in Trondheim [6]. The original aim of the company was to improve the performance of Down-Hole Safety-Valves by use of statistical analysis of history reliability data. It expanded the scope towards complete installations. The company provides the tool "Wellmaster", an interactive page to access, analyze and display the available reliability data.

According to direct information requested by the company, they are using the same method to calculate an estimator for a constant failure rate $\lambda$ as seen before in Equation 3.2.

### 3.2.2.3 PDS Data Handbook

The PDS forum [20] is a collaboration of circa 25 participants. Those are representatives of operators, manufacturers and research societies with an interest in safety instrumented systems (SIS).

They regularly meet and release two matters: The "PDS Method Handbook" and the "PDS Data Handbook" [25].

The PDS Method Handbook gives an approach to implement and verify SIL requirements according to IEC 61508 / 61511 standards. The PDS Data Handbook is based on experience with operating SIS. The given data dossiers are based on multiple data sources. They range from using OREDA, direct vendor data, expert judgment and operational reviews.

## 3.3   Small case study

### 3.3.1   Purpose description

The reliability performance of a component or system is evaluated assuming an underlying exponential distribution. Doing so can give results, which are not as holistic compared to using empirical methods. This case study will demonstrate the differences on a small example system.

### 3.3.2   System description

Following scenario is used throughout the remaining section:

A reliability engineer wants to perform a reliability assessment of a system. The system is built up of six elements made out of three different types. A reliability block diagram illustrating the system can be seen in Figure 3.1.



Figure 3.1: Reliability block diagram illustrating the motivational case study

### 3.3.3   Input Data

All three components in this scenario (A, B, C) were used in previous similar installations. The maintenance department kept track of all the times a unit was taken out of service – either due to failure, or due to end of mission. This results in a right censored dataset for each component containing the event time and the information weather it was a failure or censoring.

The full set of used input data can be seen in Figure A.1 in the appendix.

### 3.3.4   Output

When assessing the reliability of a system, the probability of having the Top-Event is a critical indicator. There are several methods to calculate this value – some are exact, some are just approximations.

The methodology being used within this case study is described in subsection A.1.2.

In addition to that, the individual component's reliability is calculated and shown.

### 3.3.5   Parametric approach

The input data (failure times) are processed to give a constant failure rate for each component type (A, B, C). The results can be seen in Table 3.1 and the calculation in subsection A.1.1.

| Component Type | Failure rate $[\frac{1}{h}]$ |
|---|---|
| A | $4.11 \times 10^{-4}$ |
| B | $3.70 \times 10^{-4}$ |
| C | $1.35 \times 10^{-4}$ |

Table 3.1: Resulting failure rates based on exponential distribution

An illustration of the three corresponding reliability diagrams can be seen in Figure 3.2. They cover the time range up to one year (8760 hours).

Figure 3.2: Reliability diagrams for one year (8760 hours) for all three components (A, B, C) based on exponential distribution

Figure 3.3: Probability of having the Top-Event with using exponential distribution

### 3.3.6 Empirical approach

The traditional Kaplan-Meier estimator is used to gain information about the reliability characteristics of each component.

The Kaplan-Meier estimator for censored data used within this thesis is based on the original research paper [9].

If $T_{(1)} < T_{(2)} < \dots$ are the times with at least one failure, and $n_i$, $d_i$ are, respectively, the number at risk and the number of failures at T(i), then Equation 3.3 can be formulated.

$$\hat{R}(t) = \prod_{i:T_{(i)} \leq t} \frac{n_i - d_i}{n_i} \tag{3.3}$$

The results of those calculations are graphical depictions of failure probability over time, as seen in Figure 3.4.

It is worth mentioning that the Kaplan-Meier value is only available up to the highest observed event time. No information beyond that time is known, hence the abrupt ending of the graph. No assumptions are made for this time span. Furthermore, the maximum observed time is different for each component.

Figure 3.4: Reliability diagrams for all three components (A, B, C) based on empirical estimation

Similar to subsection 3.3.5, the Top-Event probability of the system is calculated using the empirical reliability information.



Figure 3.5: Probability of having the Top-Event with using Kaplan-Meier

### 3.3.7 Comparison

A visual comparison of the three individual reliability diagrams can be seen in Figure 3.6. Furthermore Figure 3.7 shows the probability of the Top-Event both based on empirical and exponential modelling.

Figure 3.6: Comparison of reliability diagrams (exponential and empirical)

Figure 3.7: Probability comparison of having the Top-Event (exponential and empirical)

### 3.3.8  Case study conclusion

The scope of this case study is to show the generic difference between using parametric and empirical reliability data in a simple reliability model consisting of only six elements.

Depending on the input data, a parametric approach can give approximations of varying quality. Just by visual comparison of the empirical Kaplan-Meier and the exponential distribution, it is apparent that the goodness of fit is prone to be insufficient for some applications.

Evaluating the probability comparison of having the Top-Event for exponential and empirical methods however reveals an issue with the latter: After a certain time, the empirical method does not return any estimator of the reliability. This can be critical in a system involving numerous components.

Not all components have sufficient data available to give empirical reliability estimates for the complete anticipated mission time. Even if just one of the components used is lacking the data, the overall system-reliability estimator is at stake.

# Chapter 4

# Conceptual Developments

## 4.1   Introduction to the problem

The case study as seen in section 3.3 illustrated the potential for using empirical distributions instead of exponential assumed distributions. The input data for the case-study is of reasonably small size, as can be seen in Figure A.1. For each of the three different components, the number of recorded events are between 30 and 35.

The amount of data for real systems can however be larger. Advancements like "Big Data", "Industry 4.0" and "Internet of Things" are creating vast amounts of data.

## 4.2   Characteristics, Assumptions and Limitations

The reliability of a component is commonly represented as a 2-dimensional function R(t). The resulting reliability distribution curve has certain inherent properties. They are caused by a few basic assumptions and limitations:

- At time t = 0, the beginning of the mission time, the equipment is assumed to be fully functioning. So R(0) = 1.

- Equipment is not repaired, the item can only degrade. This results in a graph which is monotonously decreasing.

- The data considered is right-censored. As soon one specimen is taken out of the observation, the reliability will never decrease down to zero, as there is always a lack of information when that component actually failed.

- If failure times data is self-generated, the dataset only contains failure times and no censoring times.

- The empirical probability distribution only gives data up to the maximum observed event time. After that time, there is no information. If one wants to predict beyond that, assumptions have to be made (see section 4.5).

## 4.3 Intoduction of modified Kaplan-Meier

The original Kaplan-Meier estimator results in a stepwise curve as seen in Figure 3.4. Degradation however is a continuous process. The chance of survival at $t + \Delta t$ is slightly lower than at $t$. Discrete steps are not representing this behavior appropriately.

Additionally, when comparing a step-curve to a continuous curve – as it is the case when for example comparing the Kaplan-Meier curve to the exponential curve – the resulting difference becomes unsteady. This is illustrated in Figure 4.1. The green curve shows the relative difference from the exponential value using the Kaplan-Meier value as the reference.



Figure 4.1: Illustration of unsteady differences between Kaplan-Meier curve and exponential curve

Different ways to make a smoother appearance are proposed and used by various authors ([30], [2], [10]).

These require varying mathematical implementation effort. The scope of this thesis is not to evaluate different smoothing algorithms, hence a simple, comprehensive method is proposed.

The method used within this thesis is carried out as follows:

1. The original Kaplan-Meier method is used to create a two-dimensional matrix containing the time, censoring indicator and reliability estimator. This can be seen in Table 4.1

| Time | Censoring | Kaplan-Meier estimator |
| --- | --- | --- |
| 0 | 1 | 1 |
| 16 | 1 | 0.965 |
| 35 | 0 | 0.965 |
| 63 | 1 | 0.824 |
| $\vdots$ | $\vdots$ | $\vdots$ |

Table 4.1: Example of Kaplan-Meier array

2. The first value is always a 100% reliability at mission start, so:
   $R_{Smooth}(0) = 1$

3. The last value is always the Kaplan-Meier value at the last known event time (Maximum Event Time, MET), so:
   $R_{Smooth}(Maximum\ Event\ Time) = R_{Kaplan\text{-}Meier}(Maximum\ Event\ Time)$

4. All censored values are taken out of the original Kaplan-Meier array. They do not contribute information to the definition of the curve.

5. At every step – the vertical path downwards – a new reliability estimator is generated for that time. The mean of the two Kaplan-Meier reliability values used for the step is the new value for this time.

6. The resulting points are connected by piecewise linear segments.

The resulting curve is demonstrated in Figure 4.2. It also shows the new coordinates used for the plot.

This smooth version of the Kaplan-Meier estimator is used instead of the original Kaplan-Meier values for all further comparisons.

When comparing the differences between the smooth Kaplan-Meier and the exponential curve, similar to Figure 4.1, the resulting graph shows no unsteady changes. There are still segments, but they are connected continuously and hence the differences are also continuous. The graphical representation of the differences over time can be seen in Figure 4.3.

Figure 4.2: Modified Kaplan-Meier estimator for smoother curve



Figure 4.3: Differences between smooth Kaplan-Meier curve and exponential curve

## 4.4 Different representations

### 4.4.1 Reduced Kaplan-Meier

#### 4.4.1.1 Description

Piecewise linear functions are extensively used to approximate functions [4]. The modified Kaplan-Meier, as introduced in section 4.3, is already a piecewise linear model.

The number of linear segments in the the smoothed Kaplan-Meier Estimator is equal to the number of failures in the input dataset. This gives the most accurate level of results, but also requires the most data points and hence increases the data handling effort and reduces the efficiency of calculation.

The generic scope for this representation is to find the smallest number of segments which still succeed in approximating the smoothed Kaplan-Meier curve within a given threshold level.

#### 4.4.1.2 Mathematical background

Following algorithm is implemented to get the resulting reduced array

1. Create an array with two endpoints of a segment

    - The first value is always a 100% reliability at mission start, so:
      $R_{Reduced}(0) = 1$

    - The last value is always the Kaplan-Meier value at the last known event time, so:
      $R_{Reduced}(Maximum\ Event\ Time) = R_{Kaplan-Meier}(Maximum\ Event\ Time)$

2. Separate this segment $S_1$ in 10 equidistant times $t_{1..10}$

3. At each of those times: Evaluate the reliability value based on the smoothend Kaplan-Meier ($R_{KM}(t_{1..10})$) and also the linear interpolated values from that first segment ($R_{Reduced}(t_{1..10})$)

4. Calculate the differences between each of the 10 pairs, calculate the square of said value, summarize those squared values, built the mean (see Equation 4.1)

5. Compare the mean with a given threshold. If the value is below that value, the segment is considered to be sufficiently accurate. If not, the segment number will be saved in a list

6. Retrieve the list of all non-sufficient segments

7. For each of those segments, divide their time duration by two. Add another point to the array with this new time and a reliability estimate based on the smoothed Kaplan-Meier.

8. Repeat Step 2 to 5 until there are no more segments in need of treatment.

9. Unnecessary entries are removed afterwards. This is done by removing the second index (first one is $R_{Reduced}(0) = 1$) and then checked if all remaining segments are still considered to be good. If yes, then this entry is deleted from the resulted array. If not, it is kept.

Manual adjustments to the threshold parameter are necessary for each application.

An illustration of the development of the algorithm is seen in Figure 4.4.

$$\frac{\sum_{t=1}^{10}(R_{\text{KM}}(t_i) - R_{\text{Reduced}}(t_i))^2}{10} \tag{4.1}$$

$$\begin{bmatrix} 0 & 1 \\ MET & R_{\text{KM}}(MET) \end{bmatrix}$$

### 4.4.1.3 Data representation

Each segment of the smoothed Kaplan-Meier is defined by a data entry. They are essentially Cartesian coordinates for a two-dimensional curve.

An example can be seen in Listing 4.1

```
[[      0.              1.           ]
 [     93.             0.91170175]
 [    187.             0.90544117]
 [    375.             0.89292002]
 [    562.             0.88046547]
 [    750.             0.78191551]
 [   1501.             0.73017971]
 [   3002.             0.63500227]
 [   6004.             0.50253083]
 [  12008.             0.3523202 ]]
```

Listing 4.1: Data representation for piecewise linear approximation used in reduced Kaplan-Meier

### 4.4.1.4 Data access

The data as seen in Listing 4.1 can be saved as a CSV file and later loaded back into the program.

The Python package NumPy provides a dedicated routine to linearly interpolate values within a given array, named `numpy.interp` [27].

Running this methods requires 3 parameters:
The array giving the X-Coordinates (Time), the array giving the Y-Coordinates (Reliability) and the X-Coordinate (Time) to look up.

A discussion about the efficiency of this implementation is given in section 4.7.

(a) First step, the initial single segment

(b) Second step, 2 segments

(c) Third step, 4 segments

(d) Fourth step, 5 segments

(e) Fifth step, 7 segments

(f) Sixth step, 9 segments

(g) Seventh step, 10 segments

(h) Eights step, 7 segments after removal of unnecessary segments

Figure 4.4: Illustration of steps for the reduction algorithm

### 4.4.2 Polynomial regression

#### 4.4.2.1 Description

Polynomial regression is a widespread method to fit data to a mathematical equation.

This is done to by finding parameters of Equation 4.2

$$y = a_0 + a_1 x + a_2 x^2 + \cdots + a_n x^n \tag{4.2}$$

The amount of degrees $n$ is adjusted manually for every application.

The number of degrees and their corresponding values may inflict problems like overfitting or oscillation.

The polynomial regression does not incorporate known characteristics of the reliability curve like the monotonously decreasing shape and the starting condition at R(0) = 1.

#### 4.4.2.2 Mathematical background

`numpy.polyfit` uses a least square polynomial fit, mathematical details can be found in [29].

The principle scope is to minimize the squared error, as seen in Equation 4.3, with $p(x_j)$ being the polynomial value and $y_j$ the value to be fitted.

$$E = \sum_{j=0}^{k} |p(x_j) - y_j|^2 \tag{4.3}$$

#### 4.4.2.3 Data representation

Data is represented by the polynomial coefficients. The amount of coefficients is depending on the polynomial degree chosen.

These can be represented in a simple text file with one line per parameter, as seen in Listing 4.2.

```
1.234234234e1
4.123123123e2
5.2123123e−1
```

Listing 4.2: Data representation for a 3-degree polynomial equation

#### 4.4.2.4 Data access

To access the data, the coefficients have to be loaded into the program. A simple routine to read in the lines of the textfile (see Listing 4.2) is used.

Python/Numpy provides a special class for polynomial data: `poly1d` [28]. It is "a one-dimensional polynomial class" and values can be conveniently retrieved by a corresponding function.

## 4.5   Implementation of prediction

When using a parametric distribution, like the exponential distribution, the whole system is fully described by a mathematical equation and it gives results for times from $t = 0$ to $\infty$.

The empirical Kaplan-Meier estimator – as well as the smoothed version of it – is only defined for times up to the last observed event (which can be either a failure or censoring).

While polynomial equations also give results for all every point of time, they are only fitted for the known times. For any times beyond that, polynomial graphs tend to drop rapidly. Hence they cannot be used for prediction.

A major reason to use models is not only to understand the observed characteristics of a system, but also to predict the performance beyond the time of available data.

In order to implement this requirement, an addition to empirical models is given: For all times beyond the largest observed time, use data resulting from using a constant failure rate.

This results in an unsteady jump of reliability values at the highest observed event time. It can be in either direction and thus create a scenario in which the reliability is not monotonously decreasing over time. Hence a proposal is made:

At the maximum observed event time calculate both reliability values – based on constant failure rate and by empirical estimation.

If the exponential value is lower than the empirical, it is kept being used for prediction. An unsteady step to a lower value is conservative and thus accepted. An example of this behavior can be seen in Figure 4.5.

If the value is higher than the last known empirical value, a new constant failure rate is calculated. See Equation 4.4 to see the formula used. It's the inverse of the reliability at the maximum known time.

$$\lambda_{\text{Prediction}} = -\frac{\ln R_{\text{Smooth}}(\textit{Maximum Event Time})}{\text{Maximum Event Time}} \tag{4.4}$$

This new adapted $\lambda_{\text{Prediction}}$ is then used for predictions. It will give a continuously decreasing reliability. This is illustrated in Figure 4.6.

When utilizing this method for a reliability assessment, the user should be notified if, and for which times, the estimated reliability is based on an assumed constant failure rate.

Figure 4.5: Example of prediction based on unaltered constant failure rate



Figure 4.6: Example of prediction based on adapted constant failure rate

## 4.6    Accuracy of empirical probability distributions

The concept of reliability assessment itself builds on uncertainty and assumptions. There is no definite and known law behind the failure behavior of systems.

Evaluating the accuracy of any reliability model can be difficult. A "true" reference value is required to assess whether or not a result of a model is accurate or not.

If the data is self-generated by using a known generator with a defined mathematical model, this reference value is known and given. The estimated result from a model can be compared to the mathematically derived exact solution.

Real data however is not generated by an algorithm, but by a complex chain of events and consequences. There are no equations describing real failures. In this case, the smoothed Kaplan-Meier plot is used as a reference, as this estimator has no assumptions about any underlying failure distribution.

The essential task of evaluating the accuracy of the proposed empirical probability distribution is to compare two reliability plots with each other.

The Kolmogorov-Smirnov test is a non-parametric test to evaluate if two samples are drawn from the same distribution [22]. This test is used to give a numerical measure if the two curves under consideration are alike.

Following list describes how the Kolmogorov-Smirnov test value is implemented:

1. 30 uniformly random distributed times between 0 and the highest observed event time are selected

2. Reliability values of both models for each of those times (as illustrated with 10 time samples in Figure 4.7) are estimated

3. The values within each model for each discrete time step are accumulated

4. The value differences between those accumulation curves are calculated

5. The maximum difference found is the Kolmogorov-Smirnov test value (as indicated by red color in Figure 4.8)

6. Step 1. - 5. is repeated for 100 times and the mean of the resulting value calculated

This value is compared to the critical value for the two-sample Kolmogorov-Smirnov test for 30 compared samples. When using a significance value $\alpha = 0.05$ the value to compare against is 0.351 [5].

The simple example value (Figure 4.8) is 0.749. As this value is larger than the critical value for 10 samples (0.70), the two graphs shown are not based on the same distribution.

Figure 4.7: Illustration of the Kolmogorov-Smirnov test: First step



Figure 4.8: Illustration of the Kolmogorov-Smirnov test: Second step

## 4.7    Efficiency measures

Using a constant failure rate gives the benefit of calculating the estimated reliability at any given time by solving a simple equation ( Equation 4.5, [24]). Only one value ($\lambda$) has to be known beforehand to fully describe the model.

$$R(t) = e^{-\lambda \times t} \tag{4.5}$$

Retrieving reliability values from the proposed empirical estimators needs more steps. These are briefly mentioned in subsubsection 4.4.1.4 and subsubsection 4.4.2.4.

The calculation of the Kaplan-Meier estimator and the subsequent routine to reduce the data is applied as preparation. Fitting the polynomial equation to the failure data is the required preparation for the polynomial estimator. The computational efficiency of those required steps is out of scope for the reliability estimation as described within this work. Its efficiency is relevant when failure data is continuously monitored and processed online.

Four different functions are evaluated in terms of their computational efficiency:

- Parametric method: Calculate R(t) based on Equation 4.5

- Empirical method:

    1. Load data array into memory
    2. Interpolate values within this array to retrieve an estimator for R(t)
    3. Solve polynomial equation to get an estimator for R(t)

All steps are executed 1000000 times and the mean value is calculated. Python's provided `timeit` function is utilized [23].

### 4.7.1    Efficiency of parametric estimation

In order to assess the time needed to calculate a reliability value based on Equation 4.5, it is implemented in a simple Python function as seen in Listing 4.3.

```
def GetExponentialValue(Time, EstimatedLambda):
    return math.e ** (-(EstimatedLambda*Time))
```
Listing 4.3: GetExponentialValue function used for efficience measurement

A uniformly random distributed number between 0 and 10000 is used for each iteration as a time of interest, a uniformly distributed number between 0.001 and 0.00001 is chosen as the failure rate. The resulting time computational time per instance is $3{,}73 \times 10^{-7}$ seconds.

### 4.7.2    Efficiency of empirical estimation

#### Loading data array into memory

Loading the data into memory is a necessary first step for further processing. The data is stored in a CSV text file and loaded into the program with the NumPy `loadtxt` function. The function

implementing this routine used for measurement is seen in Listing 4.4.

This step is necessary for both the reduced Kaplan-Meier and the polynomial estimator. The polynomial estimator however has very little amount of data, as the number of reasonably used degrees is rather low.

```
def LoadArray():
    KM_Array = np.loadtxt("failure_dates/E_1000000_Efficiency_Array.txt")
```
Listing 4.4: LoadArray function used for efficiency measurement

### Interpolating values

To evaluate the computational efficiency of the interpolation routine, a function as seen in Listing 4.5 is used.

```
def GetKaplanMeierValue(Time):
    return np.interp(Time, KM_Array[:,0], KM_Array[:,-1])
```
Listing 4.5: GetKaplanMeierValue function used for efficiency measurement

### Polynomial values

Retrieving reliability values from the polynomial estimation is done by utilizing the polynomial functions from Numpy. The testing function used to evaluate the efficiency can be seen in Listing 4.6. A polynomial array with 6 degrees is used for assessment. Each parameter is a uniformly distributed random number between -1 and 1.

```
def GetPolyValue(Time, Polyobject):
    return Polyobject(Time)
```
Listing 4.6: GetPolyValue function used for efficiency measurement

### Results

Unlike the parametric estimator, the empirical estimation steps require a dataset. The size of it has influence on the necessary time to run execute of both functions.

Four different datasets are used, with 6, 50, 1000 and 1000000 entries.

The results can be seen in Table 4.2.

| | Size of Array | | | |
| --- | --- | --- | --- | --- |
| | 6 | 50 | 1000 | 1000000 |
| Time to load array into memory [s] | $6,62 \times 10^{-4}$ | $6,76 \times 10^{-4}$ | $1,17 \times 10^{-2}$ | $12,1$ |
| Time to estimate reliability value [s] | $4,89 \times 10^{-6}$ | $4,73 \times 10^{-6}$ | $6,58 \times 10^{-6}$ | $6,24 \times 10^{-3}$ |
| Time to retrieve polynomial value [s] | $1,61 \times 10^{-5}$ | - | - | - |

Table 4.2: Empirical efficiency measurement results

### 4.7.3   Discussion and background information

The time needed to calculate a reliability value based on the exponential law is negligible.

For the empirical approach, the first step (loading of the data) has the biggest influence on overall execution time. This sub-step is also influenced the most by the size of the input data.

For this assessment, the data is saved and loaded to/from plaintext CSV files. NumPy offers a binary format optimized for storing and reading numerical data. Using an universal plaintext format has the benefit of being readable by any application and humans.

After applying the reduction algorithm for the Kaplan-Meier estimator, the amount of data is unlikely to exceed values above 50. In that range, the combined time of loading the array and estimating the reliability is in a magnitude of $10^{-4}$ seconds.

Even for the extreme case of 1 million entries, the loading time of the array is in the order of several seconds, and as soon as the array is loaded, the actual estimation time is still less than 0,1 seconds.

All calculations are performed on a machine with following brief specifications:

- Windows 7 Enterprise

- Intel Core i7-3770, 3.40 GHz

- 16 GB Memory

- Python 3.6.0

- NumPy 1.11.3

# Chapter 5

# Experiments

## 5.1  Experiment hypothesis

The hypothesis to be tested within this chapter is described as follows:

With enough data points provided, empirical probability distributions match sufficiently well the true probability distributions.

Additionally: It is possible to calculate reliability values, based on an empirical model, with sufficiently good computational efficiency.

Consequently, these empirical models can be used as a generic input of RAMS tools.

## 5.2  Experiment description

In the first step the data used is described. Its origin, any pre-existing knowledge and basic information like size and amount of failures.

If feasible (depending of the size) the whole data set is given, otherwise an excerpt of the data structure with a reference to the full data.

As a next step the data is processed. This is done in a parametric way and in both described empirical ways.

The parametric way is an estimation of a constant failure rate $\lambda$. The empirical methods as described in section 4.4 are applied to the data.

The accuracy of all results are compared. The reasoning and method behind the comparison is given in section 4.6. Both a qualitative, visual comparison and a quantitative evaluation are given.

## 5.3    First Data Set: Wellmaster

### 5.3.1    Data origin

The data for the first data set is extracted from the Wellmaster database by ExproSoft (See subsubsection 3.2.2.2). For the purpose of this thesis only limited access to the database is given – only the component "TRSCSSV" ("Tubing retrievable, surface controlled, subsurface safety valve", commonly known as "Downhole safety valve") is populated with failure data.

### 5.3.2    Data characteristics

The Wellmaster interface is showing basic characteristics about the failure data of the selected component. Figure 5.1 shows the information as they are presented. This includes the failure rate, confidence intervals and also service time and number of failures.

Additionally a graphical illustration of the survival probability can be accessed directly within Wellmaster – see Figure 5.2. It also shows the plot based on the constant failure rate as per Figure 5.1.

The Wellmaster software provides access to a CSV file with the recorded service times. Table 5.1 shows the content of it. Column `ServiceTimeInDays` gives the service time in days and column `Failures` is an indicator whether or not it is a failure or censoring (`1` corresponds to a failure, `0` to censoring). Only those two columns where considered for further processing. Column `LengthKm` had the value `0` for all entries, column `ServiceTimeLength` is completely empty and the `FailureId` is not taken into account.

In total the data consists of 6072 events – 1427 failures and 4645 censorings.

| Average failure rate | Mean time to failure | Weibull parameters | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| **Average failure rate** | | | | | | | |
| TRSCSSV | Per year | Per 10⁶ h | LCL | UCL | Operating intervals | Service time (y) | Failures in operation |
| Filter 1 | **0.041** | **4.698** | 4.498 | 4.911 | 6,072 | 34,654 | 1,427 |

Figure 5.1: Average failure rate according to Wellmaster



Figure 5.2: Survival probability plot according to Wellmaster

| LengthKm | ServiceTimeInDays | ServiceTimeLength | Failures | FailureId |
|----------|-------------------|-------------------|----------|-----------|
| 0 | 3336 | | 0 | |
| 0 | 2265 | | 0 | |
| 0 | 1726 | | 0 | |
| 0 | 4150 | | 0 | |
| 0 | 1628 | | 0 | |
| 0 | 8603 | | 0 | |
| 0 | 3036 | | 0 | |
| 0 | 1228 | | 0 | |
| 0 | 511 | | 0 | |
| 0 | 2552 | | 0 | |
| 0 | 950 | | 0 | |
| 0 | 5073 | | 0 | |
| 0 | 4390 | | 0 | |
| 0 | 177 | | 0 | |
| 0 | 252 | | 0 | |
| 0 | 1368 | | 1 | 1704 |
| 0 | 1605 | | 0 | |
| 0 | 3112 | | 0 | |
| 0 | 2199 | | 0 | |
| 0 | 5526 | | 1 | 2410 |
| 0 | 2390 | | 0 | |
| 0 | 123 | | 0 | |
| 0 | 1628 | | 0 | |
| 0 | 1552 | | 0 | |
| 0 | 1370 | | 1 | 335 |
| 0 | 405 | | 1 | 336 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

Table 5.1: Raw service times data provided by Wellmaster

### 5.3.3   Data processing

#### 5.3.3.1   Parametric

The estimation of a constant failure rate is carried out according to Equation 5.1.

$$\lambda_1 = \frac{\text{Number of failures}}{\text{Aggregated time in service}} = \frac{1427}{303565392} = 4.70 \times 10^{-6} \frac{1}{h} \qquad (5.1)$$

This value corresponds with the estimator given by Wellmaster as seen in Figure 5.1. The small deviation is most likely due to rounding differences for the service time. A figure showing the survival probability over time in comparison with the smoothed Kaplan-Meier is seen in Figure 5.3.



Figure 5.3: Survival probability plot for the first dataset based on exponential failure distribution

### 5.3.3.2   Empirical (Reduced Kaplan-Meier)

When applying the smoothed Kaplan-Meier, as introduced in section 4.3, a reliability plot is generated. The shape is similar to the survival probability plot given by Wellmaster (Figure 5.2).

The algorithm to reduce the amount of segments, as introduced in subsection 4.4.1, is applied to this dataset. The threshold parameter is set to 0.0003. An illustration of the resulting segments together with the smoothed Kaplan-Meier plot can be seen in Figure 5.4.



Figure 5.4: Survival probability plot for the first dataset based on reduced Kaplan-Meier

### 5.3.3.3 Empirical (Polynomial)

Applying the polynomial fit to the given data set gives Equation 5.2. The resulting graph is seen in Figure 5.5 in comparison with the smoothed Kaplan-Meier.

$$R(t) = 0.9679 - 6.792 \times 10^{-6}t + 7.937 \times 10^{-11}t^2 - 4.933 \times 10^{-16}t^3 + 9.964 \times 10^{-22}t^4 \quad (5.2)$$



Figure 5.5: First Dataset: Comparison of smoothed Kaplan-Meier and polynomial fit

### 5.3.4    Accuracy comparison

As this dataset is not generated by mathematical equations, the smoothed Kaplan-Meier is used as the reference to compare accuracy.

A visual comparison of all estimators can be seen in Figure 5.6. Based on the shape of the curves it can be seen that the exponential model is too optimistic in the first 50.000 hours of operation. In times beyond that, the estimation is too pessimistic and potentially sufficient equipment might not be utilized completely.

A quantitative accuracy comparison is given in Table 5.2. Based on those Kolmogorov-Smirnov test results, the exponential distribution is not matching the reference reliability values (as 0,932 > 0,351). The accuracy of the reduced Kaplan-Meier is sufficient (as 0,071 < 0,351). Using a polynomial regression is giving comparable good results (with 0,098 < 0,351).

|                                                        | Value |
| ------------------------------------------------------ | ----- |
| Reduced Kaplan-Meier vs. Smoothed Kaplan-Meier         | 0,071 |
| Polynomial regression vs. Smoothed Kaplan-Meier        | 0,098 |
| Exponential distribution vs. Smoothed Kaplan-Meier     | 0,932 |
| Critical Kolmogorov-Smirnov test value                 | 0,351 |

Table 5.2: Comparison of Kolmogorov-Smirnov test results for the first data set



Figure 5.6: First Dataset: Comparison of smoothed Kaplan-Meier, reduced Kaplan-Meier, polynomial regression and exponential distribution

## 5.4   Second Data Set: Hard drives

### 5.4.1   Data origin

The company "Backblaze", based in California / USA, is operating a data hosting center for cloud based storage solutions. Within their server farm they operate more than 60.000 hard drives. Since 2013 the company is releasing raw statistical data on all their hard drives in use [1].

This data is freely available for everyone interested.

### 5.4.2   Data characteristics

The data is presented in CSV files. For each day of operation, one file is provided. An excerpt of one file can be seen in Table 5.3. There are 90 additional columns provided which give information retrieved by the on-board diagnostics. Those are omitted.

| date | serial_number | model | capacity_bytes | failure |
|------|---------------|-------|----------------|---------|
| 2016-10-01 | MJ0351YNG9Z0XA | Hitachi HDS5C3030ALA630 | 3000592982016 | 0 |
| 2016-10-01 | MJ0351YNG9WJSA | Hitachi HDS5C3030ALA630 | 3000592982016 | 0 |
| 2016-10-01 | PL1321LAG34XWH | Hitachi HDS5C4040ALE630 | 4000787030016 | 0 |
| 2016-10-01 | MJ0351YNGABYAA | Hitachi HDS5C3030ALA630 | 3000592982016 | 0 |
| 2016-10-01 | Z305B2QN | ST4000DM000 | 4000787030016 | 0 |
| 2016-10-01 | PL2331LAGN2YTJ | HGST HMS5C4040BLE640 | 4000787030016 | 0 |
| 2016-10-01 | WD-WMC4N2899475 | WDC WD30EFRX | 3000592982016 | 0 |
| 2016-10-01 | Z302A0YH | ST4000DM000 | 4000787030016 | 0 |
| 2016-10-01 | Z305BT0W | ST4000DM000 | 4000787030016 | 0 |
| 2016-10-01 | MJ0351YNG9Z7LA | Hitachi HDS5C3030ALA630 | 3000592982016 | 0 |
| 2016-10-01 | Z302A0YE | ST4000DM000 | 4000787030016 | 0 |
| 2016-10-01 | Z302PGH8 | ST4000DM000 | 4000787030016 | 0 |
| 2016-10-01 | Z3023VGH | ST4000DM000 | 4000787030016 | 0 |
| 2016-10-01 | PL1311LAG2205A | Hitachi HDS5C4040ALE630 | 4000787030016 | 0 |

Table 5.3: Raw data as provided by Backblaze

This data representation is complete, but inconvenient in its usage for the used tool. Ross Lazarus ([12], [13]) has written a script to process this data into a different format. It keeps track of the appearance and disappearance of each individual drive (based on the unique serial_number field). The output is a summary file with one line per drive stating its observation time and reason of removal (failure or censoring). This approach neglects the operating hours of the drives before the first available daily report.

Data from the 2013 - 2016 period is used.

In total, the data consists of 92328 events – 5762 failures and 86566 censorings.

### 5.4.3   Data processing

#### 5.4.3.1   Parametric

The estimator for a constant failure rate for the data can be seen in Equation 5.3. A plot showing the reliability based on this value is given in Figure 5.7. The scaling of the reliability ordinate is changed to show the range from 0,8 to 1,0 for better visibility.

$$\lambda_1 = \frac{\text{Number of failures}}{\text{Aggregated time in service}} = \frac{5762}{1486091064} = 3.88 \times 10^{-6} \frac{1}{\text{h}} \tag{5.3}$$



Figure 5.7: Second Dataset: Exponentially assumed failure distribution

#### 5.4.3.2 Empirical (Reduced Kaplan-Meier)

Applying the reduced Kaplan-Meier algorithm to the data with a chosen threshold value of 0.00001 gives a representation with 5 segments. This can be seen in in comparison with the smoothed Kaplan-Meier in Figure 5.8.



Figure 5.8: Second Dataset: Reduced Kaplan-Meier

### 5.4.3.3  Empirical (Polynomial)

Applying the polynomial fit to the second data set gives Equation 5.4. The resulting graph is seen in Figure 5.9 in comparison with the smoothed Kaplan-Meier.

$$R(t) = 0.9957 + 3.825 \times 10^{-7}t - 5.5 \times 10^{-10}t^2 + 2.24 \times 10^{-14}t^3 - 2.762 \times 10^{-19}t^4 \quad (5.4)$$



Figure 5.9: Second Dataset: Comparison of smoothed Kaplan-Meier and polynomial fit

### 5.4.4 Accuracy comparison

This data set is based on real failure data.

A visual comparison of all estimators can be seen in Figure 5.10.

All estimators are visually matching the trend of the data well, especially considering the adjusted scaling of the ordinate.

A quantitative comparison of accuracy is given in Table 5.4. Both the parametric and empirical distributions are meeting the given requirement and represent the reference sufficiently.

Therefore this dataset is used to compare the data requirements in section 5.7.

|  | Value |
| --- | --- |
| Reduced Kaplan-Meier vs. Smoothed Kaplan-Meier | 0,027 |
| Polynomial regression vs. Smoothed Kaplan-Meier | 0,011 |
| Exponential distribution vs. Smoothed Kaplan-Meier | 0,051 |
| Critical Kolmogorov-Smirnov test value | 0,351 |

Table 5.4: Comparison of Kolmogorov-Smirnov test results for the second data set



Figure 5.10: Second Dataset: Comparison of smoothed Kaplan-Meier, reduced Kaplan-Meier, polynomial regression and exponential distribution

## 5.5  Third Data Set: Modified Bathtub

### 5.5.1  Data origin

The traditional bath-tub curve, as illustrated in Figure 5.11, is described in almost every standard reliability textbook [11]. This curve is used as the basis to introduce a new distribution.



Figure 5.11: Traditional bathtub curve

The modification is based on an artificial idea that even during the "useful period" (flat part after infant mortality) a constant failure rate is unlikely. Fluctuations of the failure rate are still possible within this period. Sections of increased or decreased failure rate could occur. A visualization of this behavior can be seen in figure Figure 5.12.

The implementation of this idea used in this experiment is using a discreet interpretation of this curve. Figure 5.13 shows the failure rate for a set of parameters. The defining parameters can be seen in Table 5.5. This distribution is also used in the motivational case study (see section 3.3). All failure times for this experiment were self-generated based on this introduced distribution.

Figure 5.12: Modified bathtub curve

| Point No. | Time [h] | Rate [$h^{-1}$] |
|-----------|----------|-----------------|
| #1 | 0 | 0.002 |
| #2 | 100 | 0.00009 |
| #3 | 600 | 0.00009 |
| #4 | 620 | 0.0008 |
| #5 | 850 | 0.0008 |
| #6 | 900 | 0.0002 |
| #7 | 6000 | 0.0002 |
| #8 | 6050 | 0.000035 |
| #9 | 7500 | 0.00035 |
| #10 | 7550 | 0.0002 |
| #11 | 9000 | 0.0002 |
| #12 | 10000 | 0.003 |

Table 5.5: Used parameters in selfmade function

Figure 5.13: Selfmade distribution failure rate

## 5.5.2    Data characteristics

One thousand failure times are drawn based on this proposed distribution. All of them are listed as failures, no censoring. As the data is self-generated, the output is adjusted according to the required input of the processing tools. This is a CSV text-file with one line per event. The first column represents the event time and the second column the type (failure or censoring).

An excerpt of the file can be seen in Table 5.6.

| Event Time | Type |
|---|---|
| 9.150000000000000000e+02 | 1 |
| 9.370000000000000000e+02 | 1 |
| 9.510000000000000000e+02 | 1 |
| 9.580000000000000000e+02 | 1 |
| 9.660000000000000000e+02 | 1 |
| 9.810000000000000000e+02 | 1 |
| 1.002000000000000000e+03 | 1 |
| 1.003000000000000000e+03 | 1 |

Table 5.6: Excerpt of structure of the input data used for experiment 3

### 5.5.3 Data processing

#### 5.5.3.1 Parametric

The estimator for a constant failure rate for the data can be seen in Equation 5.5. A plot showing the reliability based on this value in comparison with the true reliability is given in Figure 5.14.

$$\lambda_1 = \frac{\text{Number of failures}}{\text{Aggregated time in service}} = \frac{1000}{3612389} = 2.77 \times 10^{-4} \frac{1}{h} \tag{5.5}$$
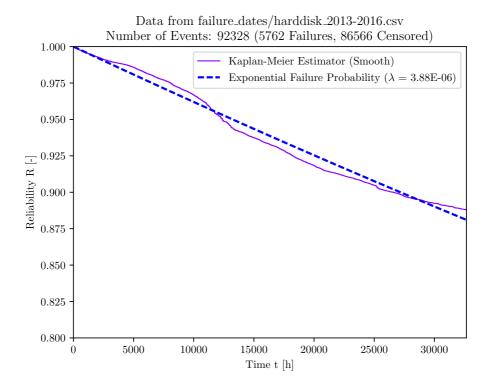


Figure 5.14: Third Dataset: Exponentially assumed failure distribution in comparison to the true reliability

### 5.5.3.2  Empirical (Reduced Kaplan-Meier)

Applying the reduced Kaplan-Meier algorithm to the data with a chosen threshold value of 0.0005 gives a representation with 11 segments. The resulting model can be seen in Figure 5.15.



Figure 5.15: Third Dataset: Reduced Kaplan-Meier in comparison to the true reliability

### 5.5.3.3   Empirical (Polynomial)

Applying the polynomial regression to this data set gives Equation 5.6. The resulting plot is seen in Figure 5.16 in comparison with the true reliability distribution.

$$R(t) = 0.9556 - 2.703 \times 10^{-4}t + 3.714 \times 10^{-8}t^2 - 1.669 \times 10^{-12}t^3 - 3.016 \times 10^{-17}t^4 \quad (5.6)$$



Figure 5.16: Third Dataset: Polynomial regression in comparison to the true reliability

### 5.5.4 Accuracy comparison

An illustration of all estimator shapes compared to the true reliability is found in Figure 5.17. It is visible that the exponential distribution is not able to match the unique characteristics of the underlying failure distribution.



Figure 5.17: Third Dataset: Comparison of exponential, reduced Kaplan-Meier, polynomial regression and the true reliability diagram

A quantitative comparison of accuracy is given in Table 5.7. The shape of the true reliability characteristics of the underlying distribution cannot be matched with a constant failure rate $(1,106 > 0,351)$. The reduced Kaplan-Meier however is sufficient for doing so (as $0,164 < 0,351$). Utilizing a polynomial regression gives an estimator which is less accurate, but still sufficient $(0,288 < 0,351)$.

|  | Value |
| --- | --- |
| Reduced Kaplan-Meier vs. True reliability | 0,164 |
| Polynomial regression vs. True reliability | 0,288 |
| Exponential distribution vs. True reliability | 1,106 |
| Critical Kolmogorov-Smirnov test value | 0,351 |

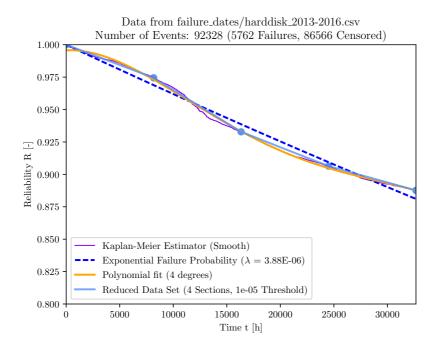Table 5.7: Comparison of Kolmogorov-Smirnov test results for the third data set

## 5.6    Fourth data set: Dynamic reliability analysis

### 5.6.1    Data origin

Based on the work of Manno et. al. [15], a failure distribution according to a dynamic reliability analysis is proposed. In their example a cooling unit is more likely to fail when used. It still has a small probability of failure due to aging even when not used. During the year, the change of weather is also changing the demand to use this cooling unit.

A simplified adaption of this behavior is implemented in a new distribution model. Piecewise linear decreasing reliability is defined for sections of usage and non-usage. Different rates are used for each.

The reliability distribution of three years of operation can be seen in Figure 5.18.

During the first half of a year (4380 hours), the reliability drops by 0.3 (Time in which the unit is used). In the following second half of the year, the reliability further drops by 0.033 (Time in which the unit is not used). At the end of a three-year cycle, the reliability reaches 0.



Figure 5.18: Dynamic reliability model

### 5.6.2 Data characteristics

As this dataset is self-generated, the data characteristic is similar to the $3^{rd}$ experiment (see subsection 5.5.2). Only 100 failure dates are drawn from the distribution.

### 5.6.3 Data processing

#### 5.6.3.1 Parametric

The estimator for a constant failure rate for the data can be seen in Equation 5.7. A plot showing the reliability based on this value is given in Figure 5.19.

$$\lambda_1 = \frac{\text{Number of failures}}{\text{Aggregated time in service}} = \frac{100}{1126397} = 8.88 \times 10^{-5} \frac{1}{h} \qquad (5.7)$$



Figure 5.19: Fourth Dataset: Exponentially assumed failure distribution in comparison to the true reliability

### 5.6.3.2 Empirical (Reduced Kaplan-Meier)

Applying the reduced Kaplan-Meier algorithm to the data with a chosen threshold value of 0.0005 gives a representation with 7 segments. The resulting model can also be seen in Figure 5.20.
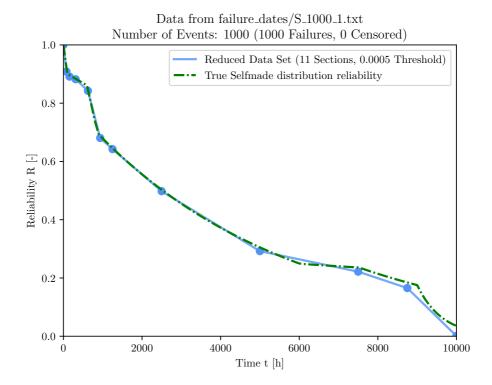


Figure 5.20: Fourth Dataset: Reduced Kaplan-Meier in comparison to the true reliability

### 5.6.3.3 Empirical (Polynomial)

Applying the polynomial regression to this data set returns Equation 5.8. The resulting reliability estimator is seen in Figure 5.21 in comparison with the true reliability distribution.

Having such a high number of resulting polynomial terms reveals problems. This estimator does not start at R(0) = 1 and at the very end a slight increase of the values can be seen, thus conflicting with the continuously decreasing nature of reliability.

$$
\begin{aligned}
R(t) =\ & 1.129 - 4.142 \times 10^{-4}t + 3.127 \times 10^{-7}t^2 \\
& - 1.363 \times 10^{-10}t^3 + 3.306 \times 10^{-14}t^4 - 4.725 \times 10^{-18}t^5 \\
& + 4.136 \times 10^{-22}t^6 - 2.242 \times 10^{-26}t^7 + 7.331 \times 10^{-31}t^8 \\
& - 1.325 \times 10^{-35}t^9 + 1.017 \times 10^{-40}t^{10}
\end{aligned}
\tag{5.8}
$$



Figure 5.21: Fourth Dataset: Polynomial regression

### 5.6.4   Accuracy comparison

A visual comparison of the true dynamic reliability model, the reduced Kaplan-Meier, the polynomial estimator and the exponential estimator can be seen in Figure 5.22.

The constant failure rate is not able to reveal the changes throughout the observed three year period.
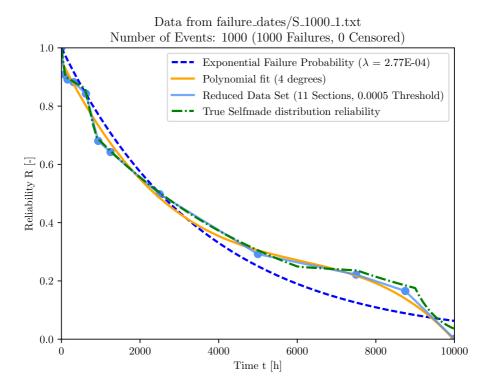


Figure 5.22: Fourth Dataset: Comparison of exponential estimator, reduced Kaplan-Meier, polynomial regression empirical and the true reliability diagram

A quantitative comparison of accuracy is given in Table 5.8. The shape of the true reliability characteristics of the underlying distribution cannot be matched with a constant failure rate (1,451 > 0,351). The reduced Kaplan-Meier however is sufficient for doing so (as 0,125 < 0,351). Less accurate but still sufficient is the estimator based on polynomial regression (0,239 < 0,351).

|                                                  | Value |
| ------------------------------------------------ | ----- |
| Reduced Kaplan-Meier vs. True reliability        | 0,125 |
| Polynomial regression vs. True reliability       | 0,239 |
| Exponential distribution vs. True reliability    | 1,451 |
| Critical Kolmogorov-Smirnov test value           | 0,351 |

Table 5.8: Comparison of Kolmogorov-Smirnov test results for the fourth data set

## 5.7   Data requirements

The second dataset, as discussed in section 5.4, is the only experiment which yielded sufficient accurate reliability results for the parametric exponential estimation and both proposed empirical estimators.

The accuracy of both parametric and empirical reliability models is depending on the amount of input data. In case of a self-generated dataset, the number of failure dates is unlimited. Real reliability assessments however require the availability of real data. In order to evaluate the required amount of failure dates, a simulation is carried out.

A number $N$ of failure dates is evaluated. In case of the harddisk dataset, those $N$ failure dates are randomly drawn from the full dataset.

Based on those failure dates, the reliability estimators are calculated. Those may be the smoothend Kaplan-Meier, the polynomial fit and the constant failure rate. The reduced Kaplan-Meier is not considered, as its computational implementation made it not feasible to be included in the Monte-Carlo simulation.

The resulting reliability is then compared to the reference as described in section 4.6.

The result of the Kolmogorov-Smirnov test is stored together with the number of failure dates ($N$).

This procedure is repeated for $N + 1$ failure dates. The maximum amount $N_{max}$ of used failure dates is set to 6000 Each set of simulations from $N$ to $N_{max}$ is carried out 100 times. Afterwards the mean Kolmogorov-Smirnov test value is calculated for each $N$.

This gives an information how many samples in average have to be considered to give sufficient results.

In average minimum 380 failure dates are required to produce a smoothed Kaplan-Meier estimator, which is sufficiently accurate. To retrieve an accurate estimator based on a constant failure rate, in average minimum 333 failure dates have to be assessed. Using the polynomial regression, 355 failure dates have to be processed to give a satisfying result.

A graphical illustration of the Monte-Carlo simulation forcan be seen in Figure 5.23. It shows the how the Kolmogorov-Smirnov test value is decreasing with the increasing number of $N$. The green horizontal line indicates the critical value of the Kolmogorov-Smirnov test (0,351).

Figure 5.23: Comparison of data requirements for all three estimators

# Chapter 6

# Conclusion & Future Works

Based on the results for each of the four experiments, several conclusions can be drawn. With a given number failure dates, empirical distributions reveal more precise insights into the failure behavior. This is valid for all of the four presented experiments.

The exponential distribution is only sufficient enough for one of those ($2^{nd}$ experiment, see section 5.4).

All of the calculated empirical estimators are superior in their accuracy, both from a qualitative visual evaluation and supported by a quantitative comparison.

According to the computational efficiency evaluation (see section 4.7), the processing time to retrieve reliability values for those estimators is negligible short.

Both parts of the experiment hypothesis could be validated on each of the four experiments: The empirical probability distributions do match sufficiently well the reference reliability and the computational efficiency is sufficiently good.

Implementing the proposed estimators as generic input of RAMS tools is possible and would be beneficial for reliability assessments.

Assessing the amount of required data to achieve the demonstrated accuracy levels revealed that, for the given dataset, each method requires a similar amount of data to be sufficiently accurate.

During the work on this thesis, several challenges and unsolved issues got revealed. These are to be addressed in future works.

- **Access to OREDA/ExproSoft database**
  Both OREDA and ExproSoft provide online databases for reliability data. All collected reliability data is available in those with rich detail level. Their policy however does not allow academic institutes to gain access to the full data. The underlying raw data is a useful resource to test and optimize empirical algorithms and compare them with the provided $\lambda$ estimates.

- **Extend to more applications**
  The comparison of empirical and parametric methods is a generic issue. The work in this thesis is focusing on simple reliability estimation. In the given RAMS context, this discussion can also be extended to maintenance models or risk quantification.

- **Improve the algorithm for data reduction**
  The implementation for data-reduction done for this thesis is rather simple. It is not efficient and does give the best possible results. Setting the threshold value is cumbersome and the removal of unnecessary points is frequently omitting redundant points. Developing an approximation algorithm is not scope of this work. Implementing a well defined and tested model will be beneficial.

- **Improve polynomial fit**
  The polynomial regression as used throughout this thesis is a very basic mathematical implementation. It does not respect the limitations and assumptions of a reliability distribution. Implementing these (Always start at R(0) = 1, monotonous decreasing) would greatly improve the applicability.

- **Implementation into RAMS tools**
  RAMS engineers do not rely on academic examples and hand-calculations. Numerous software tools for RAMS related assessments are available. Further review is necessary to get an overview over the different tools and their current functionality in terms of empirical probability distributions. An implementation of the proposed method would allow to conduct comparing experiments.

- **Implementation of uncertainty**
  Having an understanding of uncertainty and confidence intervals is a critical criteria in assessing and comparing reliability information. There are measures to identify the confidence interval for the Kaplan-Meier estimator. These should be tested and expanded for the reduced Kaplan-Meier estimator and similar for the polynomial estimator. Some mathematical metrics exist to evaluate the quality of a polynomial regression, and these should be adapted to the field of reliability

# Bibliography

[1] Backblaze [2017], 'Hard drive test data - determining failure rates and more'. Accessed: 2017-05-25.
**URL:** *https://www.backblaze.com/b2/hard-drive-test-data.html*

[2] Bae, W., Choi, H., Park, B. U. and Kim, C. [2005], 'Smoothing techniques for the bivariate kaplan–meier estimator', *Communications in Statistics - Theory and Methods* **34**(7), 1659–1674.

[3] Bobrowski, S., Chen, H., Döring, M., Jensen, U. and Schinköthe, W. [2015], 'Estimation of the lifetime distribution of mechatronic systems in the presence of a covariate: A comparison among parametric, semiparametric and nonparametric models', *Reliability Engineering & System Safety* **139**, 105–112.

[4] Camponogara, E. and Nazari, L. F. [2015], 'Models and algorithms for optimal piecewise-linear function approximation', *Mathematical Problems in Engineering* **2015**, 1–9.
**URL:** *https://doi.org/10.1155%2F2015%2F876862*

[5] *Critical Values for the Two-sample Kolmogorov-Smirnov test (2-sided)* [n.d.]. Accessed: 2017-06-02.
**URL:** *https://www.webdepot.umontreal.ca/Usagers/angers/ MonDepotPublic/STT3500H10/Critical_KS.pdf*

[6] ExproSoft [2017], 'Exprosoft - about us'. Accessed: 2017-05-07.
**URL:** *http://www.exprosoft.com/about-us/*

[7] ISO 14224:2016 [2016], Petroleum, petrochemical and natural gas industries – Collection and exchange of reliability and maintenance data for equipment, Standard, International Organization for Standardization, Geneva.

[8] Jabot, F. [2015], 'Why preferring parametric forecasting to nonparametric methods?', *Journal of Theoretical Biology* **372**, 205–210.

[9] Kaplan, E. L. and Meier, P. [1958], 'Nonparametric estimation from incomplete observations', *Journal of the American Statistical Association* **53**(282), 457–481.
**URL:** *http://dx.doi.org/10.1080/01621459.1958.10501452*

[10] Kim, C., Park, B., Kim, W. and Lim, C. [2003], 'Bezier curve smoothing of the kaplan-meier estimator', *Annals of the Institute of Statistical Mathematics* **55**(2), 359–367.

[11] Klutke, G., Kiessler, P. and Wortman, M. [2003], 'A critical look at the bathtub curve', *IEEE Transactions on Reliability* **52**(1), 125–129.
**URL:** *http://dx.doi.org/10.1109/TR.2002.804492*

[12] Lazarus, R. [2016*a*], 'Hard drive test data - determining failure rates and more'. Accessed: 2017-05-25.
**URL:** `https://github.com/fubar2/backblazeKM`

[13] Lazarus, R. [2016*b*], 'Survival analysis of hard disk drive failure data'. Accessed: 2017-05-25.
**URL:** `http://bioinformare.blogspot.no/2016/02/survival-analysis-of-hard-disk-drive.html`

[14] Lewis, E. E. [1995], *Introduction to Reliability Engineering*, Wiley.

[15] Manno, G., Zymaris, A., Kakalis, N. M. P., Chiacchio, F., Cipollone, F. E., Compagno, L., Urso, D. D. and Trapani, N. [2013], *Dynamic reliability analysis of three nonlinear aging components with different failure modes characteristics*, CRC Press, pp. 3047–3055. doi:10.1201/b15938-457.
**URL:** *http://dx.doi.org/10.1201/b15938-457*

[16] Mokhtarian, P., Namzi-Rad, M. R., Ho, T. K. and Suesse, T. [2013], Bayesian nonparametric reliability analysis for a railway system at component level, IEEE Computer Society, pp. 197–202.

[17] OREDA [2017], 'Oreda - about us'. Accessed: 2017-05-07.
**URL:** `https://www.oreda.com/about-us/`

[18] OREDA (project), SINTEF, NTNU. [2015*a*], *OREDA : offshore and onshore reliability data handbook : Vol. 1 : Topside equipment*, OREDA Participants, Norway.

[19] OREDA (project), SINTEF, NTNU. [2015*b*], *OREDA : offshore and onshore reliability data handbook : Vol. 2 : Subsea equipment*, OREDA Participants, Norway.

[20] PDS Forum [2017], 'PDS - Reliability of Safety Instrumented Systems'. Accessed: 2017-05-22.
**URL:** `http://www.sintef.no/pds`

[21] Perretti, C. T. and Munch, S. B. [2015], 'On estimating the reliability of ecological forecasts'.

[22] Pons, O. [2013], *Statistical Tests of Nonparametric Hypotheses : Asymptotic Theory*, STATISTICAL TESTS OF NONPARAMETRIC HYPOTHESES, World Scientific Publishing Company, Singapore.

[23] Python Software Foundation [2017], 'timeit — measure execution time of small code snippets'. Accessed: 2017-06-09.
**URL:** `https://docs.python.org/3/library/timeit.html`

[24] Rausand, M. [2004], *System reliability theory : Models, statistical methods, and applications*, Wiley-Interscience, Hoboken, NJ.

[25] Solfrid Håbrekke, Stein Hauge, T. O. [2013], Reliability data for safety instrumented systems : PDS data handbook, Sintef rapport, SINTEF Technology and Society, Safety Research, Trondheim.

[26] Spjøtvoll, E. [1985], 'Estimation of failure rate from reliability data bases'.

[27] The SciPy community [2017*a*], 'numpy.interp'. Accessed: 2017-05-15.
    **URL:** *https://docs.scipy.org/doc/numpy/reference/generated/numpy.interp.html*

[28] The SciPy community [2017*b*], 'numpy.poly1d'. Accessed: 2017-06-05.
    **URL:** *https://docs.scipy.org/doc/numpy/reference/generated/numpy.poly1d.html*

[29] The SciPy community [2017*c*], 'numpy.polyfit'. Accessed: 2017-06-05.
    **URL:** *https://docs.scipy.org/doc/numpy/reference/generated/numpy.polyfit.html*

[30] Youndjé, E. [2016], 'Some heuristics about bandwidth selection for the smooth kaplan–meier estimator', *Journal of the Korean Statistical Society* **45**(4), 568 – 580.
    **URL:** *http://www.sciencedirect.com/science/article/pii/S1226319215301344*

# Appendix A

# Appendix

## A.1    Additional information for the case study

### A.1.1    Failure rate and times

Example calculation for $\lambda_A$:

$$\lambda_A = \frac{\text{Number of failures for component A}}{\text{Aggregated time in service for component A}} = 4.11 \times 10^{-4} [\frac{1}{h}] \qquad \text{(A.1)}$$

The failure data for component "A" is drawn from a Weibull distribution with shape-parameter $\alpha = 1.3$ and scale-paramter $\lambda = 0.0005$ based on equation Equation A.2. The censoring is done by a uniform randomization.

$$R(t) = e^{-(\lambda t)^{\alpha}} \qquad \text{(A.2)}$$

The failure data for component "B" is drawn from an exponential distribution with failure rate $\lambda = 0.0005$ based on equation Equation A.3 The censoring is done by a uniform randomization.

$$R(t) = e^{-(\lambda t)} \qquad \text{(A.3)}$$

The failure data for component "C" is drawn from a self created distribution based on a modified bathtub-shaped failure rate. More information on this distribution is given in section 5.5. The censoring is done by a uniform randomization.

A table of all drawn values can be found in Figure A.1.

```
[[  139.     1.]        [[   92.     1.]        [[   16.     1.]
 [  255.     1.]         [  147.     1.]         [   33.     0.]
 [  339.     0.]         [  165.     0.]         [   35.     1.]
 [  445.     1.]         [  166.     1.]         [   62.     0.]
 [  464.     0.]         [  313.     1.]         [   64.     0.]
 [  464.     1.]         [  371.     1.]         [   66.     1.]
 [  621.     1.]         [  507.     0.]         [  623.     0.]
 [  621.     0.]         [  575.     1.]         [  638.     1.]
 [  701.     0.]         [  576.     1.]         [  640.     1.]
 [  701.     1.]         [  612.     1.]         [  651.     0.]
 [  776.     1.]         [  660.     1.]         [  662.     0.]
 [  875.     0.]         [  689.     0.]         [  716.     1.]
 [ 1094.     1.]         [  908.     1.]         [  747.     0.]
 [ 1260.     1.]         [  965.     1.]         [  821.     0.]
 [ 1321.     1.]         [ 1045.     1.]         [  830.     1.]
 [ 1430.     0.]         [ 1110.     1.]         [ 2099.     0.]
 [ 1486.     0.]         [ 1177.     0.]         [ 2117.     0.]
 [ 1645.     1.]         [ 1341.     1.]         [ 2607.     1.]
 [ 1989.     1.]         [ 1607.     1.]         [ 2625.     0.]
 [ 2172.     1.]         [ 1665.     0.]         [ 2838.     1.]
 [ 2283.     0.]         [ 1717.     1.]         [ 5111.     1.]
 [ 2309.     1.]         [ 1955.     1.]         [ 5135.     0.]
 [ 2504.     1.]         [ 2330.     1.]         [ 5543.     0.]
 [ 2566.     1.]         [ 2391.     1.]         [ 5769.     1.]
 [ 2709.     0.]         [ 2391.     0.]         [ 6004.     0.]
 [ 3136.     1.]         [ 2458.     0.]         [ 6796.     0.]
 [ 3186.     1.]         [ 2556.     1.]         [ 9376.     1.]
 [ 3259.     1.]         [ 2769.     1.]         [10256.     0.]
 [ 3291.     0.]         [ 3595.     1.]         [11649.     0.]
 [ 4639.     1.]]        [ 4151.     0.]         [12008.     1.]]
                         [ 5302.     1.]
                         [ 5307.     1.]
                         [ 5315.     1.]
                         [ 6248.     0.]
                         [ 7030.     1.]]
```

(a) Component A             (b) Component B             (c) Component C

Figure A.1: Used failure dates for all three component types in the case study

## A.1.2   Calculation of Top-Event probability

The TOP Event probability is calculated using standard textbook [24] methodology, not using the approximation. Listing A.1 shows the Python code used to calculate the probability.

```python
A_System_Probability = A_Probability * A_Probability

B_System_Probability = B_Probability

C_SubSystem_Probability = C_Probability * C_Probability

C_System_Probability = 1 - ((1 - C_SubSystem_Probability) * (1 -
    C_SubSystem_Probability) * (1 - C_SubSystem_Probability))

TOP_Probability = 1 - ((1 - A_System_Probability) * (1 - B_System_Probability)
    * (1 - C_System_Probability))
```

Listing A.1: Calculation of probability of TOP event

## A.2  Python program source code

```
 1  # Import numerical python module with short handle np
 2  import numpy as np
 3  import sys
 4  # Import math module
 5  import math
 6  import random
 7  # Import Plotting
 8  import matplotlib.pyplot as plt
 9  from matplotlib import rc
10  from matplotlib.offsetbox import AnchoredText
11  from matplotlib.pyplot import cm
12  import matplotlib.patches as patches
13  import matplotlib.lines as lines
14  # Import OS related module
15  import os
16  import glob
17  import generator
18  from scipy import optimize
19  from scipy import stats
20  from decimal import Decimal
21  from statistics import mean
22  import dynamic_reliability
23
24  np.set_printoptions(suppress=True, linewidth=1000)
25
26  SelfmadeData = generator.Florian(100.0,        # t1
27  600.0,                       # t2
28  620.0,                       # t3
29  850.0,                       # t4
30  900.0,                       # t5
31  6000.0,                      # t6
32  6050.0,                      # t7
33  7500.0,                      # t8
34  7550.0,                      # t9
35  9000.0,                      # t10
36  10000.0,                     # t11
37  0.002,                       # z1
38  0.00009,                     # z2
39  0.0008,                      # z3
40  0.0002,                      # z4
41  0.000035,                    # z5
42  0.0002,                      # z6
43  0.003)                       # z7
44
45  DynamicParameterArray = np.array([  [0,0],
46                                      [4380,0.300],
47                                      [8760,0.333],
48                                      [13140,0.633],
49                                      [17520,0.666],
50                                      [21900,0.966],
51                                      [26280,1.0]
52                                   ])
53
54  TrueDynamic = dynamic_reliability.Dynamic(DynamicParameterArray)
55
56  class Dataset:
57
58      def __init__(self, filename, failure_indicator, censor_indicator, timebase, NumberOfReducedSections, ExpansionFactor,
                ReductionThreshold, PolyFitDegrees, IsNumpyArray, ArrayName):
59          # Filename of the dataset
60          self.filename = filename
61          self.IsNumpyArray = bool(IsNumpyArray)
62          self.ArrayName = ArrayName
63          # What are the Censor and Failure indicators?
64          self.censor_indicator = censor_indicator
65          self.failure_indicator = failure_indicator
66
67          # What is the timebase?
68          # May be days or years
69          # Standard is "hours"
70          self.timebase = timebase
71
72          # Prepare the variable for the maximum/highest event time based on first column of KaplanMeierArray
73          self.MaximumEventTime = 0
74
75          # Number Of Sections for the reduced section method
76          self.NumberOfReducedSections = NumberOfReducedSections
77
78          # Degrees for PlotPolyFit
79          self.PolyFitDegrees = PolyFitDegrees
80
81          # Factor for prediction time line Expansion
82          self.ExpansionFactor = ExpansionFactor
83
84          # Threshold for section reduction
85          self.ReductionThreshold = ReductionThreshold
86          self.AmountOfAutomaticallyFoundSegments = 0
87
88          # Load "filename" into numpy array of name RawDataArray
89
90
```

```
 91            if self.IsNumpyArray == False:
 92                try:
 93                    self.RawDataArray = np.loadtxt(str(self.filename), delimiter=";")
 94                except:
 95                    sys.stderr.write('Unable to open file "%s"\n' % filename)
 96                    sys.stderr.flush()
 97                    exit()
 98                    return
 99            if self.IsNumpyArray == True:
100                self.RawDataArray = self.ArrayName
101
102            # If timebase is hours, do nothing
103
104            if self.timebase == "days":
105                self.RawDataArray[:,0] *= 24
106
107            if self.timebase == "years":
108                self.RawDataArray[:,0] *= 24*365.25
109
110            # Highest observed event time
111            self.MaximumEventTime = np.max(self.RawDataArray[:,0])
112
113            # Execute KaplanMeier() and place resulting array in GlobalKaplanMeierArray variable for future common use
114            self.GlobalKaplanMeierArray = self.KaplanMeier()
115            self.GlobalKaplanMeierArrayWithoutCensored = self.KaplanMeierArrayWithoutCensored()
116            self.GlobalReducedArray = self.ReducedDataIntelligent(verbose=0)
117
118            self.EstimatedLambda = self.EstimateLambda()
119
120        def GetNumberOfFailures(self):
121
122            # Read RawDataArray, only select the parts of the array where the second column (censoring value) is equal to Zero
123            # This equals a Failure. Then give out the shape of the array and from that information take the first value. This is
            the amount of rows.
124
125            NumberOfFailures = self.RawDataArray[self.RawDataArray[:, 1] == self.failure_indicator].shape[0]
126            return NumberOfFailures
127
128        def GetNumberOfCensored(self):
129
130            # Read RawDataArray, only select the parts of the array where the second column (censoring value) is equal to One
131            # This equals a Censoring. Then give out the shape of the array and from that information take the first value. This is
            the amount of rows.
132
133            NumberOfCensored = self.RawDataArray[self.RawDataArray[:, 1] == self.censor_indicator].shape[0]
134            return NumberOfCensored
135
136        def GetNumberOfEvents(self):
137
138            # Read RawDataArray. Then give out the shape of the array and from that information take the first value. This is the
            amount of rows.
139
140            NumberOfFailures = self.RawDataArray.shape[0]
141            return NumberOfFailures
142
143        def AccumulatedServiceTime(self):
144
145            # Sum up all the event times.
146            AccumulatedServiceTime = np.sum(self.RawDataArray, axis=0)[0]
147            return AccumulatedServiceTime
148
149        def EstimateLambda(self):
150
151            # Estimate Lambda based on the approach 'Number of Failures' / 'Accumulated Service Time'
152
153            NumberOfFailures = self.GetNumberOfFailures()
154            AccumulatedServiceTime = self.AccumulatedServiceTime()
155            EstimatedLambda = NumberOfFailures / AccumulatedServiceTime
156
157            return EstimatedLambda
158
159        def GetUnworthySegments(self, ReducedArray, verbose=0):
160            self.ReducedArray = ReducedArray
161            self.verbose = 0
162
163            CheckerNumber = 10
164
165            UnworthySegments = []
166            NumberOfSegments = int(ReducedArray.shape[0]-1)
167
168            if self.verbose == 1:
169                print("Running GetUnworthySegments")
170                print("Number of Segments in Reduced Array: " + str(NumberOfSegments))
171                print("This is the ReducedArray to start with:")
172                print(self.ReducedArray)
173                print("")
174
175            for i in range (1,NumberOfSegments+1):
176
177                StartingTime = self.ReducedArray[i-1,0]
178                EndTime = self.ReducedArray[i,0]
179                TimeDifference = EndTime - StartingTime
180
181                # Create CheckerArray.
182                # n (CheckerNumber) values evenly distributed between starting and ending time of each Segments
183
```

```
184                    DifferenceSquareList = []
185                    CheckerArray = np.linspace(StartingTime,EndTime,CheckerNumber)
186
187                    if self.verbose == 1:
188                        print("Checking Segment number " + str(i))
189                        print("\t\StartingTime: " + str(StartingTime))
190                        print("\tEndTime: " + str(EndTime))
191                        print("\tTimeDifference: " + str(TimeDifference))
192                        print("\tCheckerArray: " + str(CheckerArray))
193
194                    for c in np.nditer(CheckerArray):
195                        KM_SmoothValue = self.GetKaplanMeierValueSmooth(c)
196                        KM_ReducedValue = np.interp(c, self.ReducedArray[:,0], self.ReducedArray[:,-1])
197                        DifferenceSquare = (KM_SmoothValue - KM_ReducedValue)**2
198                        DifferenceSquareList.append(DifferenceSquare)
199                        if self.verbose == 1:
200                            print("\t\tChecking time: " + str(c))
201                            print("\t\tKM_SmoothValue: " + str(KM_SmoothValue))
202                            print("\t\tKM_ReducedValue: " + str(KM_ReducedValue))
203                            print("\t\tDifferenceSquare: " + str(DifferenceSquare))
204                            print("\n")
205
206                    MeanOfSquares = mean(DifferenceSquareList)
207                    if self.verbose == 1:
208                        print("Mean of Squares for this segment:")
209                        print(MeanOfSquares)
210
211                    if MeanOfSquares > self.ReductionThreshold:
212                        UnworthySegments.append(i)
213                        if self.verbose == 1:
214                            print("\t--> Not Worthy! (Larger than " + str(self.ReductionThreshold) + ")")
215                            print("\tAdding " + str(i) + " to UnworthySegments List")
216                    else:
217                        if self.verbose == 1:
218                            print("\t--> Worthy! (Smaller than " + str(self.ReductionThreshold) + ")")
219                        pass
220
221            if self.verbose == 1:
222                print("")
223                print("List of unworthy segments: " + str(UnworthySegments))
224                print("")
225
226            return UnworthySegments
227
228        def TreatSegments(self,ReducedArray,verbose=1):
229            self.ReducedArray = ReducedArray
230            self.verbose = 0
231
232            if verbose == 1:
233                print("Run TreatSegments")
234
235            TreatedArray = self.ReducedArray
236            ListOfSegmentsToBeTreated = []
237            ListOfSegmentsToBeTreated = self.GetUnworthySegments(self.ReducedArray,verbose=0)
238
239            if verbose == 1:
240                print("ListOfSegmentsToBeTreated: " + str(ListOfSegmentsToBeTreated))
241
242            for i in ListOfSegmentsToBeTreated:
243                NewTime = int(((self.ReducedArray[i,0] - self.ReducedArray[i-1,0]) / 2) + self.ReducedArray[i-1,0])
244                NewValue = self.GetKaplanMeierValueSmooth(NewTime)
245
246                TreatedArray = np.vstack([TreatedArray, np.array([NewTime, NewValue])])
247
248                if verbose == 1:
249                    print("Segment to work on (i): " + str(i))
250                    print("NewTime: " + str(NewTime))
251                    print("NewValue: " + str(NewValue))
252                    print("Added to TreatedArray")
253                    print("")
254
255            TreatedArray = TreatedArray[TreatedArray[:,0].argsort()]
256
257            if verbose == 1:
258                print("This is the TreatedArray:")
259                print(TreatedArray)
260                print("")
261
262            return TreatedArray
263
264        def RemoveUnnecessaryPoints(self,ReducedArray,verbose=1):
265
266            self.ReducedArray = ReducedArray
267
268            if verbose == 1:
269                print("Run RemoveUnnecessaryPoints...")
270
271                print("Starting Array:")
272                print(self.ReducedArray)
273
274
275            AmountOfPoints = int(self.ReducedArray.shape[0])
276            if verbose == 1:
277                print("Amount of Points: " + str(AmountOfPoints))
278                print("\n")
279            ListOfIndicesToBeRemoved = []
```

```
280
281
282            for i in range(0, AmountOfPoints-1):
283                if verbose == 1:
284                    print("Iteration Number " + str(i) + ":")
285
286                if i > 0 & i != AmountOfPoints-1:
287                    # Removing that Index
288                    self.ReducedArray = ReducedArray
289                    ArrayWithIndexRemoved = np.delete(self.ReducedArray, i, axis=0)
290
291                    ListOfUnworthySegmentsAfterRemoval = self.GetUnworthySegments(ArrayWithIndexRemoved)
292
293                    if len(ListOfUnworthySegmentsAfterRemoval) != 0:
294                        # print("ListOfUnworthySegmentsAfterRemoval is not empty, so we keep that index")
295                        pass
296                    else:
297                        # print("ListOfUnworthySegmentsAfterRemoval is empty, so this index is not important. Will be added to the
       ListOfIndicesToBeRemoved")
298                        ListOfIndicesToBeRemoved.append(i)
299                        # print("")
300                else:
301                    if verbose == 1:
302                        print("Iteration Number is 0, do nothing")
303                        print("")
304                    pass
305
306            if verbose == 1:
307                print("ListOfIndicesToBeRemoved: " + str(ListOfIndicesToBeRemoved))
308
309            self.ReducedArray = np.delete(ReducedArray, (ListOfIndicesToBeRemoved), axis=0)
310
311            if verbose == 1:
312                print(self.ReducedArray)
313
314            return self.ReducedArray
315
316        def ReducedDataIntelligent(self, verbose=0):
317
318            base, ext = os.path.splitext(str(self.filename))
319
320            ThresholdName = str(self.ReductionThreshold).replace(".", "-")
321
322            FileNameForIntelligentlyReducedArray = str(base) + "_Intelligent_" + str(ThresholdName) + "_Array.txt"
323
324            try:
325                IntelligentlyReducedArray = np.loadtxt(FileNameForIntelligentlyReducedArray)
326                # print("Loaded IntelligentlyReducedArray from existing text file " + str(FileNameForIntelligentlyReducedArray))
327
328            except:
329                # print("Running ReducedDataIntelligent()...")
330
331                # Initial KM Array, smoothened version
332                KaplanMeierArray = self.KaplanMeierArrayWithoutCensoredSmooth()
333
334                # Create initial Reduced Array.
335                # First Time: 0                    First Value: 1
336                # Last Time: MaximumEventTime    Last Value: Last known Value
337                IntelligentlyReducedArray = np.array([0, 1])
338                LastTime = self.MaximumEventTime
339                LastValue = KaplanMeierArray[-1,-1]
340
341                IntelligentlyReducedArray = np.vstack([IntelligentlyReducedArray, np.array([LastTime, LastValue])])
342
343                if self.verbose == 1:
344                    print("Start creating reduced Array:")
345                    print(IntelligentlyReducedArray)
346                    print("")
347
348                AmountOfUnworhtySegments = len(self.GetUnworthySegments(IntelligentlyReducedArray, verbose=0))
349
350                if self.verbose == 1:
351                    print("AmountOfUnworhtySegments in the beginning (should be 1): " + str(AmountOfUnworhtySegments))
352                    print("")
353
354                while AmountOfUnworhtySegments > 0:
355                    IntelligentlyReducedArray = self.TreatSegments(IntelligentlyReducedArray, verbose=0)
356                    AmountOfUnworhtySegments = len(self.GetUnworthySegments(IntelligentlyReducedArray))
357
358                IntelligentlyReducedArray = self.RemoveUnnecessaryPoints(IntelligentlyReducedArray, verbose=0)
359
360                np.savetxt(FileNameForIntelligentlyReducedArray, IntelligentlyReducedArray)
361
362            self.AmountOfAutomaticallyFoundSegments = int(IntelligentlyReducedArray.shape[0]) - 1
363
364            return IntelligentlyReducedArray
365
366        def ReturnValueArray(self, TimeArray, Method):
367
368            self.TimeArray = TimeArray.reshape(-1,1)
369            self.method = Method
370
371            # Method: KM, EXP, REDUCED, DYNAMIC
372
373            ValueArray = np.zeros([TimeArray.shape[0],1])
374            ValueArray = np.hstack([self.TimeArray, ValueArray])
```

```
375
376              if self.method == "KM":
377                  for i in range (0,int(TimeArray.shape[0])):
378                      ValueArray[i,1] = self.GetKaplanMeierValueSmooth(int(ValueArray[i,0]))
379
380              if self.method == "EXP":
381                  for i in range (0,int(TimeArray.shape[0])):
382                      ValueArray[i,1] = self.GetExponentialValue(int(ValueArray[i,0]))
383
384              if self.method == "REDUCED":
385                  for i in range (0,int(TimeArray.shape[0])):
386                      ValueArray[i,1] = self.GetReducedValue(self.GlobalReducedArray,int(ValueArray[i,0]))
387
388              if self.method == "DYNAMIC":
389                  for i in range (0,int(TimeArray.shape[0])):
390                      ValueArray[i,1] = TrueDynamic.GetR(int(ValueArray[i,0]))
391
392              return ValueArray
393
394          def Kolmogorov(self,method1,method2):
395
396              # KM, REDUCED, EXP, SELFMADE
397              self.method1 = method1
398              self.method2 = method2
399
400              NumberOfKolmogorovCheckPoints = 30
401
402              print("\nRun Kolmogorov()...")
403              print("method1: " + str(self.method1))
404              print("method2: " + str(self.method2))
405
406              print("NumberOfKolmogorovCheckPoints: " +str(NumberOfKolmogorovCheckPoints))
407
408              if NumberOfKolmogorovCheckPoints > 12:
409                  CriticialValue = 1.36 * math.sqrt((NumberOfKolmogorovCheckPoints+NumberOfKolmogorovCheckPoints) / (
         NumberOfKolmogorovCheckPoints * NumberOfKolmogorovCheckPoints))
410              elif NumberOfKolmogorovCheckPoints == 10:
411                  CriticialValue = 0.7
412
413              ListOfKolmogorovValues = []
414
415              for c in range(0,100):
416                  print("Iteration Nr. " + str(c) + " / 100")
417
418                  CheckPointArray = np.sort(np.random.randint(0,self.MaximumEventTime,size=NumberOfKolmogorovCheckPoints)).reshape
         (-1,1)
419
420                  NumberOfRows = CheckPointArray.shape[0]
421                  ZeroArray = np.zeros([NumberOfRows,5])
422                  KolmogorovArray = np.hstack((CheckPointArray,ZeroArray))
423
424                  ## Fill KM and exponential values
425
426                  for i in range(0,NumberOfRows):
427
428                      TimeForPrediction = KolmogorovArray[i,0]
429
430                      if self.method1 == "KM":
431                          KolmogorovArray[i,1] = self.GetKaplanMeierValueSmooth(TimeForPrediction)
432                      elif self.method1 == "EXP":
433                          KolmogorovArray[i,1] = self.GetExponentialValue(TimeForPrediction)
434                      elif self.method1 == "REDUCED":
435                          KolmogorovArray[i,1] = self.GetReducedValue(self.GlobalReducedArray,TimeForPrediction)
436                      elif self.method1 == "SELFMADE":
437                          KolmogorovArray[i,1] = 1 - SelfmadeData.getF(int(TimeForPrediction))
438                      elif self.method1 == "POLY":
439                          KolmogorovArray[i,1] = self.GetPolyValue(TimeForPrediction)
440                      elif self.method1 == "DYNAMIC":
441                          KolmogorovArray[i,1] = Dynamic.GetR(TimeForPrediction)
442
443                      if self.method2 == "KM":
444                          KolmogorovArray[i,2] = self.GetKaplanMeierValueSmooth(TimeForPrediction)
445                      elif self.method2 == "EXP":
446                          KolmogorovArray[i,2] = self.GetExponentialValue(TimeForPrediction)
447                      elif self.method2 == "REDUCED":
448                          KolmogorovArray[i,2] = self.GetReducedValue(self.GlobalReducedArray,TimeForPrediction)
449                      elif self.method2 == "SELFMADE":
450                          KolmogorovArray[i,2] = 1 - SelfmadeData.getF(int(TimeForPrediction))
451                      elif self.method2 == "POLY":
452                          KolmogorovArray[i,2] = self.GetPolyValue(TimeForPrediction)
453                      elif self.method2 == "DYNAMIC":
454                          KolmogorovArray[i,2] = TrueDynamic.GetR(TimeForPrediction)
455
456                  for i in range(0,NumberOfRows):
457                      if i == 0:
458                          KolmogorovArray[i,3] = KolmogorovArray[i,1]
459                          KolmogorovArray[i,4] = KolmogorovArray[i,2]
460
461                      else:
462                          KolmogorovArray[i,3] = KolmogorovArray[i-1,3] + KolmogorovArray[i,1]
463                          KolmogorovArray[i,4] = KolmogorovArray[i-1,4] + KolmogorovArray[i,2]
464
465                  for i in range(0,NumberOfRows):
466                      KolmogorovArray[i,5] = abs(KolmogorovArray[i,3] - KolmogorovArray[i,4])
467
468                  KolmogorovValue = np.max(KolmogorovArray[:,5])
```

```
469                KolmogorovIndex = np.argmax(KolmogorovArray[:,5])
470                ListOfKolmogorovValues.append(KolmogorovValue)
471
472          MeanKolmogorovValue = mean(ListOfKolmogorovValues)
473
474          if MeanKolmogorovValue <= CriticialValue:
475              print(str(MeanKolmogorovValue) + " < " + str(CriticialValue) + ": Samples based on same distribution!")
476
477          if MeanKolmogorovValue > CriticialValue:
478              print(str(MeanKolmogorovValue) + " > " + str(CriticialValue) + ": Samples NOT based on same distribution!")
479
480          return KolmogorovArray, KolmogorovValue, KolmogorovIndex, MeanKolmogorovValue
481
482      def FindNearestIndex(self, array, value):
483          idx = (np.abs(array-value)).argmin()
484          return idx
485
486      def FindIndexLowerThanValue(self, array, value):
487          DifferenceArray = (array-value)
488
489          if value == 0:
490              idx = 0
491          else:
492              idx = (np.argwhere(DifferenceArray<0)).argmax()
493
494          return idx
495
496      def FitPolyToKaplanMeier(self):
497          KaplanMeierArray = self.KaplanMeier()
498          X = KaplanMeierArray[:,0]
499          Y = KaplanMeierArray[:,5]
500          Z = np.polyfit(X, Y, self.PolyFitDegrees)
501          p = np.poly1d(Z)
502          return p
503
504      def GetPolyValue(self, time):
505          self.time = time
506          Poly = self.FitPolyToKaplanMeier()
507          return Poly(self.time)
508
509      def KaplanMeier(self):
510
511          base, ext = os.path.splitext(str(self.filename))
512
513          FileNameForKaplanMeierArray = str(base) + "_KM_Array.txt"
514
515          # if self.IsNumpyArray == False:
516          try:
517              KaplanMeierArray = np.loadtxt(FileNameForKaplanMeierArray)
518              # print("Loaded KaplanMeierArray from existing text file " + str(FileNameForKaplanMeierArray))
519
520          except:
521              # Create an array with the Kaplan-Meier data inside
522
523              # print("Running KaplanMeier()...")
524              # Number of events in RawDataArray. Used for counting
525
526              NumberOfEvents = int(self.GetNumberOfEvents())
527
528              # Sort the RawDataArray according to the event times - failure times and censoring times
529              # Save it in SortedArray
530
531              SortedArray = self.RawDataArray[self.RawDataArray[:,0].argsort()]
532
533              # Create RankArray from 1 to NumberOfEvents
534              # Make it vertically afterwards
535
536              RankArray = np.arange(1,NumberOfEvents+1)
537              RankArray = np.vstack(RankArray)
538
539              # Create a reversed version of the RankArray
540              ReverseRankArray = RankArray[::-1]
541
542              # Create array of ones
543              OnesArray = np.ones([NumberOfEvents,2])
544
545              # Stack all arrays horizontally
546
547              KaplanMeierArray = np.hstack([SortedArray, RankArray, ReverseRankArray, OnesArray])
548
549
550              # Calculate p in 4th column
551              KaplanMeierArray[:,4] = (KaplanMeierArray[:,3]-1)/(KaplanMeierArray[:,3])
552
553
554              # Ugly routine to set p to 1 for all censored sets
555              for i in range(0, NumberOfEvents):
556                  if KaplanMeierArray[i,1] == self.censor_indicator:
557                      KaplanMeierArray[i,4] = 1
558
559              # Build Product for Kaplan-Meier Estimator in last column
560
561              for i in range(0, NumberOfEvents):
562                  KaplanMeierArray[i,5] = np.product([KaplanMeierArray[0:i,4]])
563
564
```

```python
565                    if self.IsNumpyArray == False:
566                        np.savetxt(FileNameForKaplanMeierArray, KaplanMeierArray)
567
568            return KaplanMeierArray
569
570        def KaplanMeierArrayWithoutCensored(self):
571            base, ext = os.path.splitext(str(self.filename))
572
573            FileNameForKaplanMeierArrayWithoutCensored = str(base) + "_KM_Array_WO_Censored.txt"
574
575            try:
576                KaplanMeierArrayWithoutCensored = np.loadtxt(FileNameForKaplanMeierArrayWithoutCensored)
577                # print("Loaded KaplanMeierArrayWithoutCensored from existing text file " + str(
       FileNameForKaplanMeierArrayWithoutCensored))
578
579            except:
580                # print("Running KaplanMeierArrayWithoutCensored()...")
581                KaplanMeierArray = self.GlobalKaplanMeierArray
582
583                KaplanMeierArrayWithoutCensored = KaplanMeierArray[KaplanMeierArray[:,1] == 1]
584                if KaplanMeierArray[-1,1] == 0:
585                    KaplanMeierArrayWithoutCensored = np.vstack([KaplanMeierArrayWithoutCensored, KaplanMeierArray[-1:,]])
586
587                np.savetxt(FileNameForKaplanMeierArrayWithoutCensored, KaplanMeierArrayWithoutCensored)
588
589            return KaplanMeierArrayWithoutCensored
590
591        def GetKaplanMeierValue(self, TimeForPrediction):
592            self.TimeForPrediction = TimeForPrediction
593            KaplanMeierArray = self.KaplanMeier()
594
595            ResultArray = np.array([0])
596
597            for x in np.nditer(TimeForPrediction):
598                RowIndex = self.FindIndexLowerThanValue(KaplanMeierArray[:,0], x)
599                ResultArray = np.hstack((ResultArray, KaplanMeierArray[RowIndex+1,-1]))
600
601            ResultArray = np.delete(ResultArray, (0), axis=0)
602
603            return ResultArray
604
605
606        def GetKaplanMeierValueSmooth(self, TimeForPrediction):
607            self.TimeForPrediction = TimeForPrediction
608            KaplanMeierArray = self.GlobalKaplanMeierArrayWithoutCensored
609
610            ResultArray = np.array([0])
611            for x in np.nditer(TimeForPrediction):
612                Result = np.interp(self.TimeForPrediction, KaplanMeierArray[:,0], KaplanMeierArray[:,-1])
613                ResultArray = np.hstack((ResultArray, Result))
614
615            return Result
616
617        def GetReducedValue(self, ReducedArray, TimeForPrediction):
618
619            self.ReducedArray = ReducedArray
620            self.TimeForPrediction = TimeForPrediction
621
622            X = self.ReducedArray[:,0]
623            Y = self.ReducedArray[:,1]
624            Estimation = np.interp(self.TimeForPrediction, X, Y)
625
626            return Estimation
627
628        def GetExponentialValue(self, TimeForPrediction):
629            self.TimeForPrediction = TimeForPrediction
630            return math.e ** (-(self.EstimatedLambda*self.TimeForPrediction))
631
632        def CompareValues(self, TimeForPrediction, Method):
633
634            self.TimeForPrediction = TimeForPrediction
635            self.Method = Method
636
637            KaplanMeierValue = self.GetKaplanMeierValue(self.TimeForPrediction)
638            ExponentialValue = self.GetExponentialValue(self.TimeForPrediction)
639            KaplanMeierValueSmooth = self.GetKaplanMeierValueSmooth(self.TimeForPrediction)
640
641            if self.Method == "Exponential":
642                ResultValue = ((ExponentialValue - KaplanMeierValue) / KaplanMeierValue) * 100
643
644            if self.Method == "Reduced":
645                ResultValue = ReducedValue - KaplanMeierValue
646
647            if self.Method == "Smooth":
648                ResultValue = ((ExponentialValue - KaplanMeierValueSmooth) / KaplanMeierValueSmooth) * 100
649
650            return ResultValue
651
652        def GetUpperConfidenceInterval(self, EstimatedFailureProbability, NumberOfFailures):
653            self.EstimatedFailureProbability = EstimatedFailureProbability
654            self.NumberOfFailures = NumberOfFailures
655
656            UpperConfidenceInterval = self.EstimatedFailureProbability + 1.96 * ( self.EstimatedFailureProbability * ( 1 - self.
       EstimatedFailureProbability)/self.NumberOfFailures)**0.5
657
658            return UpperConfidenceInterval
```

```
659
660       def GetLowerConfidenceInterval(self, EstimatedFailureProbability, NumberOfFailures):
661           self.EstimatedFailureProbability = EstimatedFailureProbability
662           self.NumberOfFailures = NumberOfFailures
663
664           LowerConfidenceInterval = self.EstimatedFailureProbability - 1.96 * ( self.EstimatedFailureProbability * ( 1 - self.
              EstimatedFailureProbability)/ self.NumberOfFailures)**0.5
665
666           return LowerConfidenceInterval
667
668       def KaplanMeierWithConfidenceInterval(self):
669           KaplanMeierArray = self.KaplanMeier()
670           NumberOfEvents = self.GetNumberOfEvents()
671
672           ConfidenceArray = np.zeros([NumberOfEvents,2])
673
674           for i in range (0,NumberOfEvents):
675               ConfidenceArray[i,0] = self.GetLowerConfidenceInterval(KaplanMeierArray[i,5], NumberOfEvents)
676               ConfidenceArray[i,1] = self.GetUpperConfidenceInterval(KaplanMeierArray[i,5], NumberOfEvents)
677
678           KaplanMeierWithConfidenceIntervalArray = np.hstack([KaplanMeierArray, ConfidenceArray])
679
680           return KaplanMeierWithConfidenceIntervalArray
681
682       def KaplanMeierArrayWithoutCensoredSmooth(self, verbose=0):
683
684           KaplanMeierArrayWithoutCensored = self.KaplanMeierArrayWithoutCensored()
685
686           KaplanMeierArrayWithoutCensored_X = KaplanMeierArrayWithoutCensored[:,0]
687           KaplanMeierArrayWithoutCensored_Y = KaplanMeierArrayWithoutCensored[:,-1]
688
689           NumberOfRows = KaplanMeierArrayWithoutCensored_X.shape[0]
690
691           KaplanMeierArrayWithoutCensoredSmooth_Y = np.copy(KaplanMeierArrayWithoutCensored_Y)
692
693           for i in range (0,NumberOfRows-1):
694
695               First_Value = KaplanMeierArrayWithoutCensored_Y[i]
696               Second_Value = KaplanMeierArrayWithoutCensored_Y[i+1]
697               Difference = First_Value - Second_Value
698               New_Value = Second_Value + (Difference / 2.0)
699               KaplanMeierArrayWithoutCensoredSmooth_Y[i] = New_Value
700
701               if verbose == 1:
702                   print("Working on Line number " + str(i))
703                   print("\tTime: " + str(KaplanMeierArrayWithoutCensored_X[i]))
704                   print("\tR-Value of Line number " + str(i) + ": " + str(First_Value))
705                   print("\tR-Value of Line number " + str(i+1) + ": " + str(Second_Value))
706                   print("\tDifference: " + str(Difference))
707                   print("\tNew Value: " + str(New_Value))
708                   print("\n")
709
710           KaplanMeierArrayWithoutCensoredSmooth = np.column_stack([KaplanMeierArrayWithoutCensored_X,
              KaplanMeierArrayWithoutCensoredSmooth_Y])
711           FirstLine = np.array([0,1])
712           KaplanMeierArrayWithoutCensoredSmooth = np.vstack([FirstLine, KaplanMeierArrayWithoutCensoredSmooth])
713           return KaplanMeierArrayWithoutCensoredSmooth
714
715       def Plot(self, PlotKaplanMeier, PlotKaplanMeierSmooth, PlotKaplanMeierConfidenceInterval, PlotEstimatedExponential, PlotPolyFit,
              PlotReduced, PlotDifferences, ShowPrediction, PlotKolmogorov, PlotSelfmade, PlotDynamic, ShowPlot, SavePlot, PlotFile):
716
717           self.PlotKaplanMeier = bool(PlotKaplanMeier)
718           self.PlotKaplanMeierConfidenceInterval = bool(PlotKaplanMeierConfidenceInterval)
719           self.PlotKaplanMeierSmooth = bool(PlotKaplanMeierSmooth)
720           self.PlotEstimatedExponential = bool(PlotEstimatedExponential)
721           self.PlotPolyFit = bool(PlotPolyFit)
722           self.PlotReduced = bool(PlotReduced)
723           self.PlotDifferences = bool(PlotDifferences)
724           self.ShowPrediction = bool(ShowPrediction)
725           self.PlotKolmogorov = bool(PlotKolmogorov)
726           self.PlotSelfmade = bool(PlotSelfmade)
727           self.PlotDynamic = bool(PlotDynamic)
728           self.ShowPlot = bool(ShowPlot)
729           self.SavePlot = bool(SavePlot)
730           self.PlotFile = str(PlotFile)
731
732           # Plots various graphs, according to parameters.
733           # KaplanMeier - Plot KM if set to 1
734           # Show - Show the plot
735           # Save - Save the plot
736           # PlotFile - Filename for the saved plot
737
738           # Set generic plot options
739           # Set LaTeX font (Computer Modern)
740           rc('font', **{'family': 'serif', 'serif': ['Computer Modern']})#, 'size'   : 16})
741           rc('text', usetex=True)
742
743           # Set-up figure and axes
744           fig, ax1 = plt.subplots()#, sharex=True)
745
746           # Define Labels for x and y axis
747           plt.xlabel("Time t [h]")
748           plt.ylabel("Reliability R [-]")
749
750           # Axis range
751           ymin = 0
```

```python
752          xmax_value = int(self.MaximumEventTime)
753          plt.xlim(xmin=0, xmax=xmax_value)
754          plt.ylim(ymin=ymin, ymax=1.0)
755
756          # Set title. All "_" have to be replaced by "\_" for LaTeX escaping
757          TitleName = str(self.filename).replace("_", "\_")
758          Subtitle = "Number of Events: " + str(self.GetNumberOfEvents()) + " (" + str(self.GetNumberOfFailures()) + " Failures,
         " + str(self.GetNumberOfCensored()) + " Censored)"
759          ax1.set_title("Data from " + str(TitleName) + "\n" + Subtitle)
760
761          if PlotKaplanMeier == True:
762              # Run the KaplanMeier() function
763              KaplanMeierArray = self.GlobalKaplanMeierArray
764
765              # Define x and y for the plot
766              x_KM = KaplanMeierArray[:,[0]]
767              y_KM = KaplanMeierArray[:,[-1]]
768              ax1.step(x_KM,y_KM, color="red", linewidth=1, alpha=1, label="Kaplan-Meier Estimator")
769              ax1.legend(loc=3)
770
771
772          if PlotKaplanMeierSmooth == True:
773
774              KaplanMeierArrayWithoutCensoredSmooth = self.KaplanMeierArrayWithoutCensoredSmooth()
775
776              x_KM_Smooth = KaplanMeierArrayWithoutCensoredSmooth[:,0]
777              y_KM_Smooth = KaplanMeierArrayWithoutCensoredSmooth[:,-1]
778              ax1.plot(x_KM_Smooth,y_KM_Smooth, color = '#8c00ff', linewidth=1, alpha=1, label="Kaplan-Meier Estimator (Smooth)")
779
780          if PlotKaplanMeierConfidenceInterval == True:
781              KaplanMeierWithConfidenceIntervalArray = self.KaplanMeierWithConfidenceInterval()
782              X = KaplanMeierWithConfidenceIntervalArray[:,0]
783              Y_Lower = KaplanMeierWithConfidenceIntervalArray[:,6]
784              Y_Upper = KaplanMeierWithConfidenceIntervalArray[:,7]
785
786              ax1.step(X,Y_Upper, color="red", alpha=0.8, linestyle="-.", label="Upper KM Confidence Interval")
787              ax1.step(X,Y_Lower, color="red", alpha=0.8, linestyle="-.", label="Lower KM Confidence Interval")
788
789          if PlotEstimatedExponential == True:
790
791              # Set X-Axis to range
792              upper_limit = xmax_value
793              x_Lambda = np.arange(0,upper_limit)
794
795              EstimatedLambda = self.EstimateLambda()
796              EstimatedLambdaSci = '%.2E' % Decimal(EstimatedLambda)
797              ax1.plot(x_Lambda, math.e ** (-(EstimatedLambda*x_Lambda)), color="blue", linestyle="--", linewidth=2, alpha=1,
         label="Exponential Failure Probability ($\lambda$ = " + str(EstimatedLambdaSci) + ")")
798              ax1.legend(loc=3)
799
800          if PlotPolyFit == True:
801
802              X = self.KaplanMeier()[:,[0]]
803              Poly = self.FitPolyToKaplanMeier()
804              print(Poly)
805              ax1.plot(X, Poly(X), color="orange", linewidth=2, alpha=1, label="Polynomial fit (" + str(self.PolyFitDegrees) + "
         degrees)")
806
807          if PlotReduced == True:
808
809              ReducedArray = self.ReducedDataIntelligent()
810              X = ReducedArray[:,0]
811              Y = ReducedArray[:,1]
812
813              ax1.scatter(X, Y, color="#5492f7", linewidth=1.5, alpha=1)
814              ax1.plot(X, Y, color="#5492f7", linewidth=2, alpha=0.8, label="Reduced Data Set (" + str(self.
         AmountOfAutomaticallyFoundSegments) + " Sections, " + str(self.ReductionThreshold) + " Threshold)")
815
816          if PlotDifferences == True:
817
818              ax2 = ax1.twinx()
819
820              Timerange = np.arange(0,self.MaximumEventTime,100)
821
822              if PlotReduced == True:
823                  DifferenceArrayReduced = self.CompareValues(Timerange,"Reduced")
824                  ax2.plot(Timerange, DifferenceArrayReduced, color="cyan", alpha=0.7, linewidth=0.8, label="Value difference (
         Reduced to KM)")
825
826              if PlotKaplanMeierSmooth == True:
827                  DifferenceArraySmooth = self.CompareValues(Timerange,"Smooth")
828                  ax2.plot(Timerange, DifferenceArraySmooth, color="green", alpha=0.7, linewidth=0.8, label="Value difference (
         Exponential to smooth Kaplan-Meier)")
829                  ax2.set_ylabel("Difference [\%]")
830                  ax2.tick_params('y', color="green")
831                  ax2.axhline(y=0, color="green", alpha=0.7, linewidth=0.5, linestyle="-")
832
833                  ax2.legend(loc=0)
834
835          if ShowPrediction == True:
836
837              # Expand the X axis. Use self.ExpansionFactor as a factor
838              plt.xlim(xmin=0, xmax=int(self.MaximumEventTime)*self.ExpansionFactor)
839
840              # Plot exponential for that area
841              x_Lambda = np.arange(int(self.MaximumEventTime),int(self.MaximumEventTime)*self.ExpansionFactor)
```

```
842
843              EstimatedLambda = self.EstimateLambda()
844              EstimatedLambdaSci = '%.2E' % Decimal(EstimatedLambda)
845
846              KaplanMeierArray = self.KaplanMeier()
847              Offset = KaplanMeierArray[-1,-1] - math.e ** (-(EstimatedLambda*self.MaximumEventTime))
848              OffsetSci = round(Offset, 2)
849
850              SmoothKaplanMeierValueAtMaximumEventTime = self.GetKaplanMeierValueSmooth(self.MaximumEventTime)
851              PredictLambda = -1 * (math.log(SmoothKaplanMeierValueAtMaximumEventTime) / self.MaximumEventTime)
852              PredictLambdaSci = '%.2E' % Decimal(PredictLambda)
853
854              if PlotEstimatedExponential == True:
855                  if Offset >= 0:
856                      ax1.plot(x_Lambda, math.e ** (-(EstimatedLambda*x_Lambda)) , color="blue", linestyle="--", linewidth=2,
         alpha=1)
857                  else:
858                      ax1.plot(x_Lambda, math.e ** (-(EstimatedLambda*x_Lambda)) , color="blue", linestyle="--", linewidth=2,
         alpha=0.5)
859              else:
860                  if Offset >= 0:
861                      ax1.plot(x_Lambda, math.e ** (-(EstimatedLambda*x_Lambda)) , color="blue", linestyle="--", linewidth=2,
         alpha=1, label="Original Exponential Failure Probability ($\lambda$ = " + str(EstimatedLambdaSci) + ")")
862                  else:
863                      ax1.plot(x_Lambda, math.e ** (-(PredictLambda*x_Lambda)) , color="green", linestyle="--", linewidth=2,
         alpha=1, label="Adapted Exponential Failure Probability ($\lambda_{Prediction}$ = " + str(PredictLambdaSci)+ ")")
864                      ax1.plot(x_Lambda, math.e ** (-(EstimatedLambda*x_Lambda)) , color="blue", linestyle="--", linewidth=2,
         alpha=1, label="Original Exponential Failure Probability ($\lambda$ = " + str(EstimatedLambdaSci) + ")")
865
866              if PlotPolyFit == True:
867                  X = np.arange(int(self.MaximumEventTime),int(self.MaximumEventTime)*self.ExpansionFactor)
868                  Poly = self.FitPolyToKaplanMeier()
869
870
871                  ax1.plot(X, Poly(X), color="green", linestyle="--", linewidth=2, alpha=1)
872
873              # Add hatching for predicted time frame
874
875              patterns = ['-', '+', 'x', 'o', 'O', '.', '*', '|', '/']
876              ax1.add_patch(
877                  patches.Rectangle(
878                      (self.MaximumEventTime, -0.2),   # (x,y)
879                      self.MaximumEventTime*self.ExpansionFactor,          # width
880                      1.3,               # height
881                      # hatch='/',
882                      hatch = patterns[8],
883                      edgecolor="red",
884                      # alpha=0.5,
885                      fill=False
886                      )
887                  )
888
889
890      if PlotKolmogorov == True:
891          KolmogorovArray, KolmogorovValue, KolmogorovIndex, MeanKolmogorovValue, SameDistribution = self.Kolmogorov(1)
892          KolmogorovTimes = KolmogorovArray[:,0]
893          KolmogorovKMValues = KolmogorovArray[:,1]
894          KolmogorovExpValues = KolmogorovArray[:,2]
895          KolmogorovKM_Cum_Values = KolmogorovArray[:,3]
896          KolmogorovExp_Cum_Values = KolmogorovArray[:,4]
897
898          plt.ylim(ymin=ymin, ymax=5)
899
900          ax1.plot(KolmogorovTimes, KolmogorovKM_Cum_Values, color='#8c00ff', alpha=1, linewidth=1, label="Cumulative Values
     for smoothened Kaplan-Meier")
901          ax1.plot(KolmogorovTimes, KolmogorovExp_Cum_Values, color='blue', alpha=1, linewidth=1, label="Cumulative Values
     for exponential distribution")
902          ax1.set_ylabel("Cumulative Values [-]")
903
904
905          for i in range(0,len(KolmogorovTimes)):
906              x = [KolmogorovTimes[i], KolmogorovTimes[i]]
907              y = [KolmogorovKM_Cum_Values[i], KolmogorovExp_Cum_Values[i]]
908
909              ax1.plot(x, y, marker = '', color = "green")
910
911          # Maximum Kolmogorov line
912          x_max = [KolmogorovTimes[KolmogorovIndex], KolmogorovTimes[KolmogorovIndex]]
913          y_max = [KolmogorovKM_Cum_Values[KolmogorovIndex], KolmogorovExp_Cum_Values[KolmogorovIndex]]
914          ax1.plot(x_max, y_max, marker = '', color = "red", label="Kolmogorov-Smirnov test value (" + str(round(
     KolmogorovValue,3)) + ")")
915
916
917      if PlotSelfmade == True:
918
919          upper_limit = xmax_value
920          x_Selfmade = np.arange(0,10100,100)
921          y_Selfmade = np.arange(0,10100,100)
922
923          y_list = []
924          for i in range(0,len(y_Selfmade)):
925              Time = int(y_Selfmade[i])
926              R_Value = 1 - SelfmadeData.getF(Time)
927
928              y_list.append(R_Value)
929
```

```
930
931                  ax1.plot(x_Selfmade, y_list, color="green", linestyle="-.", linewidth=2, alpha=1, label="True Selfmade distribution
             reliability")
932
933             if PlotDynamic == True:
934
935                  x_Dynamic = np.arange(0,26280,100)
936                  ax1.plot(x_Dynamic, TrueDynamic.GetR(x_Dynamic), color="green", linestyle="-.", linewidth=2, alpha=1, label="True
             dynamic reliability model")
937
938
939             if ShowPlot == True:
940
941                  plt.legend(loc="upper right")
942                  plt.show()
943
944
945             if SavePlot == True:
946                  plt.legend(loc="upper right")
947                  plt.savefig(PlotFile)
948
949  def KolmogorovArrayCheck(Array1, Array2):
950
951
952       Array1 = Array1.reshape(-1,1)
953       Array2 = Array2.reshape(-1,1)
954       if Array1.shape[0] != Array2.shape[0]:
955            print("The two input arrays do not have the same size! ABORT!")
956
957
958       NumberOfKolmogorovCheckPoints = Array1.shape[0]
959
960       if NumberOfKolmogorovCheckPoints <= 12:
961            Critical_Values = [0,0,0,0,1,1,0.83,0.857,0.75,0.667,0.7,0.636,0.667]
962            CriticialValue = Critical_Values[NumberOfKolmogorovCheckPoints]
963
964       else:
965            CriticialValue = 1.36 * math.sqrt((NumberOfKolmogorovCheckPoints+NumberOfKolmogorovCheckPoints) / (
             NumberOfKolmogorovCheckPoints * NumberOfKolmogorovCheckPoints))
966
967       ListOfKolmogorovValues = []
968
969       for c in range(0,100):
970            NumberOfRows = NumberOfKolmogorovCheckPoints
971            ZeroArray = np.zeros([NumberOfRows,3])
972            ZeroArray2 = np.zeros([NumberOfRows,1])
973            KolmogorovArray = np.hstack((ZeroArray2, Array1, Array2, ZeroArray))
974
975            for i in range(0,NumberOfRows):
976                 if i == 0:
977                      KolmogorovArray[i,3] = KolmogorovArray[i,1]
978                      KolmogorovArray[i,4] = KolmogorovArray[i,2]
979
980                 else:
981                      KolmogorovArray[i,3] = KolmogorovArray[i-1,3] + KolmogorovArray[i,1]
982                      KolmogorovArray[i,4] = KolmogorovArray[i-1,4] + KolmogorovArray[i,2]
983
984            for i in range(0,NumberOfRows):
985                 KolmogorovArray[i,5] = abs(KolmogorovArray[i,3] - KolmogorovArray[i,4])
986
987
988            KolmogorovValue = np.max(KolmogorovArray[:,5])
989            KolmogorovIndex = np.argmax(KolmogorovArray[:,5])
990            ListOfKolmogorovValues.append(KolmogorovValue)
991
992       MeanKolmogorovValue = mean(ListOfKolmogorovValues)
993
994       if MeanKolmogorovValue <= CriticialValue:
995            print(str(MeanKolmogorovValue) + "\t<\t " + str(CriticialValue) + ":\tSamples based on same distribution!")
996            SameDistribution = 1
997
998       if MeanKolmogorovValue > CriticialValue:
999            print(str(MeanKolmogorovValue) + "\t>\t" + str(CriticialValue) + ":\tSamples NOT based on same distribution!")
1000            SameDistribution = 0
1001
1002       return KolmogorovArray, KolmogorovValue, KolmogorovIndex, MeanKolmogorovValue, SameDistribution
1003
1004  def ReturnRandomSample(InputArray,Number):
1005
1006       NumberOfEvents = int(InputArray.shape[0])
1007       ListOfIndices = random.sample(range(0, NumberOfEvents), Number)
1008       RandomSampleArray = np.sort(InputArray.take(ListOfIndices, axis=0), axis=0)
1009
1010       return RandomSampleArray
1011
1012
1013
1014
1015  def CensorArray(Array):
1016
1017       Length = Array.shape[0]
1018       Ones = np.ones([Length,1])
1019       CensoredArray = np.hstack((Array,Ones))
1020
1021       return CensoredArray
1022
```

```python
1023   def HowManyPoints(Montecarlo_Counter):
1024
1025       DataSetToBeChecked = Selfmade
1026
1027       # This is the Raw Data Array to use as source. Used for Exprosoft and Backblaze
1028       FullDataArray = DataSetToBeChecked.RawDataArray
1029       filename_prefix = "how_many_points_temp_data_selfmade/Selfmade"
1030       NumberOfKolmogorovCheckPoints = 30
1031
1032       # This is the number of entries in that Raw Data Array
1033       NumberOfMaximumEntries = FullDataArray.shape[0]
1034
1035       DictionaryOfSuccess = {}
1036       ListOfSuccess = []
1037       ListOfSamplesize = []
1038
1039       ResultArray = np.array([0,999])
1040
1041       for c in range(5,15,2):
1042           print("Monte Carlo Number " + str(Montecarlo_Counter) + ", Iteration with " + str(c) + " samples...")
1043
1044           NumberOfSamples = c
1045
1046
1047           # What indices of the Raw Data Array are to be considered? For Wellmaster and Backblaze
1048           # ListOfIndicesToBeConsidered = random.sample(range(0, NumberOfMaximumEntries), NumberOfSamples)
1049           # TestingArray = np.sort(FullDataArray.take(ListOfIndicesToBeConsidered, axis=0), axis=0)
1050
1051           TestingArray = np.array([])
1052           for i in range(0,c):
1053               TestingArray = np.append(TestingArray, SelfmadeData.GetSingleFailureTime())
1054               TestingArray = TestingArray.reshape(-1,1)
1055
1056           TestingArray = CensorArray(TestingArray)
1057
1058           # print("Our TestingArray:")
1059
1060           FileNameForTempArray = filename_prefix + "_Raw_Data_Iteration_" + str(c).zfill(4) + "_Samples.txt"
1061           np.savetxt(FileNameForTempArray, TestingArray, delimiter=";")
1062
1063           # print("\tArray " + str(FileNameForTempArray) + " saved...")
1064
1065           TempDataset = Dataset( FileNameForTempArray,
1066                                  timebase="hours",
1067                                  failure_indicator = 1,
1068                                  censor_indicator = 0,
1069                                  NumberOfReducedSections = 5,
1070                                  PolyFitDegrees = 4,
1071                                  ExpansionFactor = 3,
1072                                  ReductionThreshold = 0.0005,
1073                                  IsNumpyArray = False,
1074                                  ArrayName = "blank")
1075
1076           # Create a random time array with NumberOfKolmogorovCheckPoints number of times between 0 and the maximum time in the
1077           current iteration
1077           CheckPointArray = np.sort(np.random.randint(0, TempDataset.MaximumEventTime, size=NumberOfKolmogorovCheckPoints)).reshape
1077           (-1,1)
1078
1079           # Create array with KM-smooth values of this sample array
1080           KM_Sample_Array = np.zeros([NumberOfKolmogorovCheckPoints,1])
1081           for i in range(0,NumberOfKolmogorovCheckPoints):
1082               KM_Sample_Array[i] = TempDataset.GetKaplanMeierValueSmooth(CheckPointArray[i])
1083
1084           # Create array with Reliability values of true array
1085           KM_Full_Array = np.zeros([NumberOfKolmogorovCheckPoints,1])
1086           for i in range(0,NumberOfKolmogorovCheckPoints):
1087               # KM_Full_Array[i] = DataSetToBeChecked.GetReducedValue(DataSetToBeChecked.GlobalReducedArray, CheckPointArray[i])
1088               KM_Full_Array[i] = 1 - SelfmadeData.getF(int(CheckPointArray[i]))
1089
1090           KolmogorovArray, KolmogorovValue, KolmogorovIndex, MeanKolmogorovValue, SameDistribution = KolmogorovArrayCheck(
1090           KM_Sample_Array, KM_Full_Array)
1091
1092           # Setting first column in ResultArray
1093           ResultArrayLine = np.array([c, KolmogorovValue])
1094           ResultArray = np.vstack([ResultArray, ResultArrayLine])
1095           print("")
1096
1097       np.savetxt("how_many_points_temp_data_selfmade_results/Selfmade_HowManyPoints_Result_" + str(Montecarlo_Counter).zfill(4) +
1097       "_Samples.txt", ResultArray)
1098
1099       return ResultArray
1100
1101   def MonteCarlo(runs):
1102
1103       for i in range(0,runs):
1104           print("Iteration " + str(i))
1105           Resultarray = HowManyPoints(i)
1106
1107           if i == 0:
1108               # print("First run!")
1109               MonteCarloArray = Resultarray
1110
1111           if i > 0:
1112               # print("Not first run")
1113               MonteCarloArray = np.hstack([MonteCarloArray, Resultarray[:,-1].reshape(-1,1)])
1114
```

```
1115
1116            # print("Array so far: ")
1117            # print(MonteCarloArray)
1118
1119          np.savetxt("how_many_points_temp_data_selfmade_results/Selfmade_MonteCarlo_Result_" + str(i).zfill(4) + "th_Run.txt",
           MonteCarloArray)
1120
1121          files_to_be_deleted  = glob.glob('how_many_points_temp_data_selfmade\*')
1122          for f in files_to_be_deleted:
1123              os.remove(f)
1124
1125      print("Finished")
1126      np.savetxt("how_many_points_temp_data_selfmade_results/Selfmade_MonteCarlo_FinalArray.txt", MonteCarloArray)
1127
1128
1129
1130  A = Dataset(   "failure_dates/A-Data.txt",
1131                 timebase="hours",
1132                 failure_indicator = 1,
1133                 censor_indicator = 0,
1134                 NumberOfReducedSections = 5,
1135                 PolyFitDegrees = 4,
1136                 ExpansionFactor = 3,
1137                 ReductionThreshold = 0.0004,
1138                 IsNumpyArray = False,
1139                 ArrayName = "blank")
1140
1141
1142
1143  A.Plot(      PlotKaplanMeier=0,
1144               PlotKaplanMeierSmooth=0,
1145               PlotKaplanMeierConfidenceInterval=0,
1146               PlotEstimatedExponential=1,
1147               PlotPolyFit=1,
1148               PlotReduced=1,
1149               PlotDifferences=0,
1150               ShowPrediction=0,
1151               PlotKolmogorov=0,
1152               PlotSelfmade=0,
1153               PlotDynamic=1,
1154               ShowPlot=1,
1155               SavePlot=0,
1156               PlotFile="C_KM_Exponential_Comparison.pdf")
```

Listing A.2: statistic_tools.py Python source code

```
1   import numpy as np
2   import math
3   import sys
4   import os
5
6   class Dynamic:
7       def __init__(self, ParameterArray):
8           self.ParameterArray = ParameterArray
9
10      def GetF(self, TimeForPrediction):
11
12          self.TimeForPrediction = TimeForPrediction
13          Result = np.interp(self.TimeForPrediction, self.ParameterArray[:,0], self.ParameterArray[:,-1])
14          return Result
15
16      def GetR(self, TimeForPrediction):
17
18          self.TimeForPrediction = TimeForPrediction
19          Result = 1 - self.GetF(self.TimeForPrediction)
20          return Result
21
22      def GetTime(self, FailureProbability):
23          self.FailureProbability = FailureProbability
24          ResultTime = np.interp(self.FailureProbability, self.ParameterArray[:,-1], self.ParameterArray[:,0])
25          return int(ResultTime)
26
27  def CreateFailureTimes(Filename, Modelname, Number):
28
29      try:
30          output = open(Filename, "w")
31
32      except:
33          sys.stderr.write('Unable to open file "%s"\n' % Filename)
34          sys.stderr.flush()
35          exit
36
37      for i in range(0, Number):
38          output.write(str(Modelname.GetTime(np.random.uniform(0,1.0))) + ";1\n")
39
40      output.close()
```

Listing A.3: dynamic_reliability.py Python source code

```python
1   # Import numerical python module with short handle np
2   import numpy as np
3
4   # Import math module
5   import math
6   import os
7
8   # Generate a Treatment class
9
10  class Selfmade:
11      def __init__(self, t1, t2, t3, t4, t5, t6, t7, t8, t9, t10, t11, z1, z2, z3, z4, z5, z6, z7):
12          self.t1 = t1
13          self.t2 = t2
14          self.t3 = t3
15          self.t4 = t4
16          self.t5 = t5
17          self.t6 = t6
18          self.t7 = t7
19          self.t8 = t8
20          self.t9 = t9
21          self.t10 = t10
22          self.t11= t11
23          self.z1 = z1
24          self.z2 = z2
25          self.z3 = z3
26          self.z4 = z4
27          self.z5 = z5
28          self.z6 = z6
29          self.z7 = z7
30
31      def getz(self, t):
32          self.t = t
33          if self.t < 0:
34              t = 0
35              print("Negative time entered - t set to zero")
36
37          elif self.t > 0 and self.t <= self.t1:
38              z = ((self.z2 - self.z1)/self.t1)*self.t+self.z1
39              return z
40
41          elif self.t > self.t1 and self.t <= self.t2:
42              z = self.z2
43              return z
44
45          elif self.t > self.t2 and self.t <= self.t3:
46              z = ((self.z3 - self.z2)/(self.t3 - self.t2))*(self.t-self.t2)+self.z2
47              return z
48
49          elif self.t > self.t3 and self.t <= self.t4:
50              z = self.z3
51              return z
52
53          elif self.t > self.t4 and self.t <= self.t5:
54              z = ((self.z4 - self.z3)/(self.t5 - self.t4))*(self.t-self.t4)+self.z3
55              return z
56
57          elif self.t > self.t5 and self.t <= self.t6:
58              z = self.z4
59              return z
60
61          elif self.t > self.t6 and self.t <= self.t7:
62              z = ((self.z5 - self.z4)/(self.t7 - self.t6))*(self.t-self.t6)+self.z4
63              return z
64
65          elif self.t > self.t7 and self.t <= self.t8:
66              z = self.z5
67              return z
68
69          elif self.t > self.t8 and self.t <= self.t9:
70              z = ((self.z6 - self.z5)/(self.t9 - self.t8))*(self.t-self.t8)+self.z5
71              return z
72
73          elif self.t > self.t9 and self.t <= self.t10:
74              z = self.z6
75              return z
76
77          elif self.t > self.t10 and self.t <= self.t11:
78              z = ((self.z7 - self.z6)/(self.t11 - self.t10))*(self.t-self.t10)+self.z6
79              return z
80
81      def getInt(self, t):
82          self.t = t
83          integral = 0
84          if self.t < 0:
85              t = 0
86              print("Negative time entered - t set to zero")
87
88          elif self.t >= 0 and self.t <= self.t1:
89              integral_temp = np.trapz([self.z1, self.z2], x=[0, self.t])
90              integral = integral_temp
91              return integral
92
93          elif self.t > self.t1 and self.t <= self.t2:
94              integral_temp = np.trapz([self.z2, self.z2], x=[self.t1, self.t])
95              integral = integral_temp + self.getInt(self.t1)
96              return integral
```

```
97
98          elif self.t > self.t2 and self.t <= self.t3:
99              integral_temp = np.trapz([self.z2,self.z3], x=[self.t2,self.t])
100             integral = integral_temp +  self.getInt(self.t2)
101             return integral
102
103         elif self.t > self.t3 and self.t <= self.t4:
104             integral_temp = np.trapz([self.z3,self.z3], x=[self.t3,self.t])
105             integral = integral_temp + self.getInt(self.t3)
106             return integral
107
108         elif self.t > self.t4 and self.t <= self.t5:
109             integral_temp = np.trapz([self.z3,self.z4], x=[self.t4,self.t])
110             integral = integral_temp + self.getInt(self.t4)
111             return integral
112
113         elif self.t > self.t5 and self.t <= self.t6:
114             integral_temp = np.trapz([self.z4,self.z4], x=[self.t5,self.t])
115             integral = integral_temp + self.getInt(self.t5)
116             return integral
117
118         elif self.t > self.t6 and self.t <= self.t7:
119             integral_temp = np.trapz([self.z4,self.z5], x=[self.t6,self.t])
120             integral = integral_temp + self.getInt(self.t6)
121             return integral
122
123         elif self.t > self.t7 and self.t <= self.t8:
124             integral_temp = np.trapz([self.z5,self.z5], x=[self.t7,self.t])
125             integral = integral_temp + self.getInt(self.t7)
126             return integral
127
128         elif self.t > self.t8 and self.t <= self.t9:
129             integral_temp = np.trapz([self.z5,self.z6], x=[self.t8,self.t])
130             integral = integral_temp + self.getInt(self.t8)
131             return integral
132
133         elif self.t > self.t9 and self.t <= self.t10:
134             integral_temp = np.trapz([self.z6,self.z6], x=[self.t9,self.t])
135             integral = integral_temp + self.getInt(self.t9)
136             return integral
137
138         elif self.t > self.t10 and self.t <= self.t11:
139             integral_temp = np.trapz([self.z6,self.z7], x=[self.t10,self.t])
140             integral = integral_temp + self.getInt(self.t10)
141             return integral
142
143
144     def getF(self,t):
145         self.t = t
146         return 1 - math.e ** (-1 * self.getInt(t))
147
148     def getpdf(self,t):
149         self.t = t
150         z = self.getz(t)
151         pdf = math.e ** (-1 * self.getInt(t)) * z
152         return pdf
153
154     def array(self,t):
155         self.t = t
156         data = []
157
158         for i in range (1,int(self.t),1):
159             temp = np.array([i, self.getz(i), self.getF(i), self.getpdf(i)])
160             data.append(temp)
161
162         return np.array(data)
163
164     def save_parameters(self):
165         np.savetxt("S_10000_Parameters.txt", try2.array(t), delimiter=";")
166
167
168
169     def GetSingleFailureTime(self):
170         dataset = np.loadtxt("S_10000_Parameters.txt", delimiter=";")
171         times = dataset[:, [0]]
172
173         # Third column are the F values
174         failure = np.around(dataset[:, [2]], decimals=3)
175         max_F = np.amax(failure)
176
177         value = np.random.uniform(low=0.0001, high=max_F)
178         random_index = np.searchsorted(np.ravel(failure), value, side="left")
179
180         if random_index > 0 and (random_index == len(failure) or math.fabs(value - failure[random_index-1]) < math.fabs(value -
            failure[random_index])):
181             T = times[random_index-1]
182         else:
183             T = times[random_index]
184
185         return T
186
187
188
189
190
191
```

```python
192        def failure_times(self, count, repeat, random):
193            self.count = count
194            self.repeat = repeat
195            self.random = random
196
197            mu = 1
198            sigma = 0.2
199
200            t = 10000
201            # Load the generated parameters from the S function
202            dataset = np.loadtxt("S_10000_Parameters.txt", delimiter=";")
203
204            # First column are the times
205            times = dataset[:, [0]]
206
207            # Third column are the F values
208            failure = np.around(dataset[:, [2]], decimals=3)
209            max_F = np.amax(failure)
210
211            for r in range(1, self.repeat+1):
212                failure_times = np.array([])
213
214                for i in range(0, self.count):
215                    value = np.random.uniform(low=0.0001, high=max_F)
216                    random_index = np.searchsorted(np.ravel(failure), value, side="left")
217
218                    if random == True:
219                        if random_index > 0 and (random_index == len(failure) or math.fabs(value - failure[random_index-1]) < math.fabs(value - failure[random_index])):
220                            T = int(times[random_index-1] * np.absolute(np.random.normal(mu, sigma)))
221                            failure_times = np.append(failure_times, T)
222                        else:
223                            T = int(times[random_index] * np.absolute(np.random.normal(mu, sigma)))
224                            failure_times = np.append(failure_times, T)
225                    else:
226                        if random_index > 0 and (random_index == len(failure) or math.fabs(value - failure[random_index-1]) < math.fabs(value - failure[random_index])):
227                            T = times[random_index-1]
228                            failure_times = np.append(failure_times, T)
229                        else:
230                            T = times[random_index]
231                            failure_times = np.append(failure_times, T)
232
233                if random == True:
234                    np.savetxt("failure_times/SN_" + str(self.count) + "_" + str(r) + ".txt", np.sort(failure_times), delimiter=";")
235                else:
236                    np.savetxt("failure_times/S_" + str(self.count) + "_" + str(r) + ".txt", np.sort(failure_times), delimiter=";")
```

Listing A.4: generator.py Python source code