

SUSA

**Software for Uncertainty
and Sensitivity Analyses**

Classical Methods





Gesellschaft für Anlagen-
und Reaktorsicherheit
(GRS) gGmbH

SUSA

**Software for Uncertainty
and Sensitivity Analyses**

Classical Methods

Martina Kloos
Nadine Berner

March 2021

Remark:

This documentation is associated to the SUSA software.

The documentation refers to the research projects RS1528 and RS1559 which have been funded by the German Federal Ministry of Economic Affairs and Energy (BMWi).

The authors are responsible for the content of the report.

GRS - 631
ISBN 978-3-949088-20-9

Key-Words:

BEPU Analysis, Best Estimate Plus Uncertainty Analysis, Monte Carlo Simulation, Sensitivity Analysis, SUSA, Uncertainty Analysis

Contents

1	Probabilistic uncertainty and sensitivity analysis.....	1
2	Input Uncertainties	3
2.1	Distribution.....	4
2.1.1	Parametric Distribution.....	7
2.1.2	Nonparametric Distribution	30
2.2	Dependency	36
2.2.1	Population-related correlation	38
2.2.2	Sample-related correlation.....	44
2.2.3	Full Dependence.....	46
2.2.4	Conditional Distribution	48
2.2.5	Function of Parameters.....	48
2.2.6	Inequality	49
2.3	Proportions	51
3	Sample Generation	53
3.1	Pseudorandom number generators	53
3.2	Simple Random Sampling (SRS).....	55
3.2.1	SRS with consideration of population-related correlations	55
3.2.2	SRS with consideration of sample-related correlations	58
3.3	Latin Hypercube Sampling (LHS).....	60
3.3.1	LHS with consideration of population-related correlations.....	60
3.3.2	LHS with consideration of sample-related correlations.....	61
3.4	Sample generation with consideration of complete dependencies	62
3.5	Sample generation with consideration of conditional distributions.....	62
3.6	Sample generation with consideration of functional relationships.....	63
3.7	Sample generation with consideration of inequalities.....	63
3.8	Computer code runs	63

4	Uncertainty Analysis	65
4.1	Basic statistics	67
4.2	Tolerance limits	69
4.2.1	Wilks non-parametrical tolerance limits.....	69
4.2.2	Tolerance limits in case of a Normal or Lognormal distribution	75
4.2.3	Bootstrapped tolerance limits.....	76
4.3	Interval limits from Chebychev and Chebychev-Cantelli inequalities	78
4.4	Parametric distribution fitting.....	79
4.4.1	Kolmogorov-Smirnov goodness-of-fit test	80
4.4.2	Lillifors goodness-of-fit test	81
4.5	Construction and application of a surrogate model	82
4.6	Uncertainty quantifications for multiple variables.....	85
4.6.1	Simultaneous multiple tolerance limits	85
4.6.2	Probability of compliance of multiple limiting values.....	86
5	Sensitivity Analysis	89
5.1	Correlation based sensitivity indices	90
5.1.1	Pearson's correlation	90
5.1.2	Spearman's rank correlation	98
5.1.3	Blomqvist's medial correlation.....	99
5.1.4	Kendall's rank correlation.....	101
5.1.5	Partial correlation coefficient, standardized regression coefficient and coefficient of determination relating to Spearman's, Blomqvist's and Kendall's correlations.....	102
5.2	Multiple correlation coefficients	103
5.2.1	Pearson's multiple correlation	103
5.2.2	Spearman's, Blomqvist's and Kendall's multiple correlations	104
5.3	Correlation ratio	105
5.3.1	Correlation ratio on ranks.....	107
5.4	Sobol indices	109
5.5	Goodman/Kruskal coefficients from 2x2 contingency tables	112
5.6	Results from stepwise regression	113

References	115
List of Figures.....	121
List of Tables	125

1 Probabilistic uncertainty and sensitivity analysis

A very convenient method to get a quantification of the uncertainty of a computational result (figure of merit) is the Monte Carlo (MC) simulation method (/MCK 79/, /HEL 96/, /WIC 98/, /HEL 06/). It relies on the uncertainties quantified for relevant input parameters and the propagation of these uncertainties through the computer model. That means possible values of the uncertain input parameters are sampled based on the respective uncertainty quantifications and supplied as input to corresponding computer runs. The different values finally obtained for the computational result can then be analyzed by statistical methods in order to derive appropriate indicators for the uncertainty of the result.

To identify the most important uncertainty sources of a computational result, an additional (global) sensitivity analysis is useful (/HOF 99/, /SAL 00/). It can show where to improve the state of knowledge in order to reduce the uncertainty of the computational result most effectively.

To facilitate the performance of uncertainty and sensitivity analyses based on the MC simulation method, the tool SUSANA (Software for the Uncertainty and Sensitivity Analyses) was developed. SUSANA combines well established methods from probability calculus and statistics with a comfortable graphical user interface (GUI). The concept of SUSANA enables the user to fully concentrate on the analysis input including the identification of the input parameters which represent the main uncertainty sources of the computational result and the formulation of the corresponding uncertainties. After this is done, SUSANA provides support to quantify the uncertainties probabilistically and to perform the different steps of an uncertainty and sensitivity analysis.

The main steps of a probabilistic uncertainty and sensitivity analysis supported by SUSANA can be summarized as follows:

1. Identify the uncertain input parameters which may essentially contribute to the uncertainty of the computational result.
2. Document the parameters in SUSANA and use SUSANA to quantify the uncertainty of the parameters in terms of univariate probability distributions and dependences (e.g. association measures, conditional distributions, or functional relationships).
3. Prompt SUSANA to generate a sample of values for the parameters based on the quantifications made in step 2.
4. With the support of SUSANA, start the corresponding computer code runs for the sets of parameter values sampled in step 3 to get a sample of values for the computational result.
5. Prompt SUSANA to calculate statistics useful for quantifying the uncertainty of the computational result.
6. Prompt SUSANA to calculate sensitivity indices and to provide a ranking of the parameters with respect to their contribution to the uncertainty of the computational result.

In this manual, the methods integrated in SUSANA to perform the aforementioned steps 2 – 6 of an uncertainty and sensitivity analysis are described in detail. Section 2 provides a description of the methods available to quantify input uncertainties. The methods for generating the samples of parameter values and of computational results are outlined in Section 3. Subject of Section 4 are the methods for quantifying the uncertainty of the computational result. Section 5 deals with the sensitivity indices in SUSANA.

2 Input Uncertainties

In this Section, the basic principles and methods used within SUSAs to transfer experts' knowledge into appropriate (subjective) probability distributions are explained. The Section aims to assist the analyst to adequately design uncertainties by the optimal use of the resources provided by SUSAs. The process of transferring experts' knowledge into descriptive statistical measures such as central tendency (e.g. mean, median, or mode), dispersion (variability) and association (e.g. correlation) in order to derive a suitable subjective probability distribution is generally referred to as *elicitation of uncertainty* /GEL 13/. The concept of using a probability distribution as an expression of uncertainty essentially corresponds to a Bayesian or subjective interpretation of probability as a *degree of belief*. For the purpose of performing a probabilistic uncertainty and sensitivity analysis the specification of the parameter uncertainties is mainly referred to as *input specification* /MCK 95/. The main steps of the input specification may be considered as:

1st step: choice of an appropriate *univariate probability distribution for a parameter* in order to model its uncertainty (/KLO 91/).

2nd step: assignment of suitable *measures of association, conditional distributions, inequalities or other functional relationships between parameters* in order to model knowledge dependencies (see Section 2.2 and /KRZ 88/).

Based on the specified univariate distributions in tandem with the assigned dependencies, a sample of parameter values of size n in compliance with these properties is generated. The sampling procedures available in SUSAs are explained in detail in Section 3. Each element of the generated sample represents one realization of the set of uncertain parameters. It is used as part of the input to the computer code of interest in order to perform Monte Carlo simulation.

The first step of the input specification comprises the formulation of probability distributions for the uncertain parameters. This can be done either directly by the specification of an appropriate analytical formula or indirectly by the specification of distribution characteristics reflecting the experts' state of knowledge. In the latter case, SUSAs applies suitable approaches to deduce an appropriate distribution in compliance with the experts' knowledge. In the second step, (knowledge) dependencies between parameters are taken into account. They may be formulated in terms of measures of association, complete dependencies, conditional distributions, inequalities or other functional

relationships between parameters. To support the analyst during her/his task of uncertainty specification, the strategies available at each specification step are described in detail.

This Section is organized according to the structure of the graphical user interface (GUI) of SUSAs as follows:

- all distributions available in SUSAs are defined and their parameters are explained,
- the available options to quantify knowledge dependencies are defined and their intuitive interpretation is described, and
- the use of proportions to model further aspects of association are exemplified.

2.1 Distribution

The theoretical way to consider parameter uncertainty, indicating the likeliness that a parameter acquires alternative values within a certain range, is accomplished by assuming a suitable probability distribution. In general experts' beliefs are rarely provided in a convenient parameterized form. Therefore, an appropriate approximation to the uncertainty considered for an input parameter has to be derived based on scientific judgement using all of the relevant information available. This information may include:

- measurement data (i.e. the statistical analysis of a series of indication values),
- expert's knowledge about the behavior and properties of the relevant process or system,
- findings from previous uncertainty evaluations, and
- data provided in calibration studies and other technical reports.

Dependent on the range of alternative values and further probabilistic characteristics, appropriate distribution types can be modelled to best capture the uncertainty. Once a distribution type is selected, the distribution parameters can be calculated either analytically or by applying a random search iteration.

Each distribution can be truncated at a given minimum and/or maximum of the uncertainty range. Due to this combination of subjective information given by experts, the elicitation process of uncertainty quantification leads to descriptive and subjective probability distributions. In general, for each considered uncertain parameter at least one of the following characteristics should be known by the analyst:

- **maximum range** of possibly applicable alternative parameter values, i.e. the support of a hypothetical probability distribution
- **intermediate values to given degrees of belief** for a parameter, i.e. supporting points of a probability distribution in terms of at least two quantiles (percentiles).

One of these characteristics suffices to derive an appropriate approximation in terms of a probability distribution to the parameter uncertainty. By providing the distribution type (e.g. Normal distribution, Beta distribution, Gamma distribution, etc.) as well as the maximum range and/or some quantiles and the corresponding quantile probabilities the suitable parameterization can be derived *analytically* or by a *simple random search*. The full spectrum of analytical approaches to quantify the uncertainty of a parameter is carried out in detail in the context of the Normal distribution (Section 2.1.1.1). If the analytical determination of the distribution parameters is not feasible, a simple random search is applied as an iterative solution to a nonlinear optimization problem /KLO 91/. Thereby, the optimization problem to estimate the distributional parameters p_1 and p_2 is defined by the α_1 -quantile q_1 and the α_2 -quantile q_2 provided for the corresponding cumulative distribution function F as

$$\left(F_{p_1, p_2}(q_1) - \alpha_1\right)^2 + \left(F_{p_1, p_2}(q_2) - \alpha_2\right)^2 \stackrel{!}{=} \min \quad (2.1)$$

If further quantiles (up to $K=10$) are indicated, the optimization problem is defined as

$$\sum_{i=1}^K w(q_i) \cdot \left(F_{p_1, p_2}(q_i) - \alpha_i\right)^2 \stackrel{!}{=} \min. \quad (2.2)$$

with $w(q_i)$ representing the subjective weight assigned to quantile q_i .

At the starting point of the iteration process, the specified distribution is assumed to have default parameter values p_1 and p_2 . The optimization procedure, that is an iterative

variation of the parameter values such that the optimization problem becomes minimal, is performed as a random search method.

SUSA offers the direct visualization of the modelled probability distribution (density, cumulative distribution function, or complementary cumulative distribution function) and the comparison to alternative distributions in order to achieve the most suitable distributional shape representing the experts' belief.

As mentioned before, the range of alternative parameter values must not correspond to the original support of the selected distribution type (e.g. $(-\infty, +\infty)$ for the Normal distribution). It may be specified according to the experts' belief. That means distributions with infinite tails can be truncated at one or both tails.

A comprehensive summary of all available input combinations for each available distribution is given in /KLO 91/. The following Sections summarize and characterize the major probability distributions commonly employed to quantify the input uncertainty in the context of Monte Carlo simulation. For this purpose, each probability distribution is explained via its particular properties and its common field of application. For each distribution, the probability density function and the distribution function are provided and the available alternative input strategies besides the plain parametrization are briefly explained.

More information about uncertainty evaluation strategies to specify input uncertainties can be found, e.g., in the book /BED 01/ and online via /ITL 17/. A general but detailed overview of the probability distributions commonly used for uncertainty assignments is provided in /JOH 94/, /JOH 95/. A more practical overview in the context of uncertainty assignments in engineered systems can be found in /HAL 00/.

2.1.1 Parametric Distribution

2.1.1.1 Normal Distribution

The Normal distribution (*aka* Gaussian distribution) is used to describe natural as well as technical processes, such as measurement processes. Its importance basically originates from the Central Limit Theorem. This theorem says that the sum of a large number of independent random variables (uncertain parameters) asymptotically follows a Normal distribution regardless of the individual distributions of the random variables (uncertain parameters). For measurement processes, it might be simplified as follows: if a measurement result is linearly influenced by an infinitely large number of uncertainty sources, then the distribution of the measurement result approaches the Normal distribution regardless of the individual distributions describing the uncertainty sources. Even in more realistic situations, i.e. a limited set of uncertainty sources, the uncertainty of a measurement result may be adequately approximated via the Normal distribution /COL 09/.

The Normal distribution of variable X is defined via the density function as

$$f(x) = \frac{1}{\sqrt{2\pi} p_2} \cdot \exp\left(-\frac{1}{2} \left(\frac{x - p_1}{p_2}\right)^2\right), \quad x \in \mathbb{R} \quad (2.3)$$

and via the cumulative distribution function as

$$F(x) = \int_{-\infty}^x f(t) dt = \frac{1}{2} \cdot \left(1 + \operatorname{erf}\left(\frac{x - p_1}{p_2 \cdot \sqrt{2}}\right)\right), \quad x \in \mathbb{R} \quad (2.4)$$

It is parametrized by the mean p_1 and the standard deviation $p_2 > 0$ (or variance p_2^2) of X for the support range $S =] - \infty, \infty[$. The expression $\operatorname{erf}(\cdot)$ refers to the error function defined as

$$\operatorname{erf}(y) = \frac{2}{\sqrt{\pi}} \int_0^y e^{-t^2} dt \quad (2.5)$$

In general the Normal distribution has an unbounded support, i.e. the x values may range from minus to plus infinity. In this case, the most remarkable intuitive property of the Normal distribution depicts its symmetric bell-shape. If the distribution parameters p_1, p_2 are not explicitly known, the specific parameters can be analytically computed from two quantile-probability pairs (q_1, α_1) and (q_2, α_2) with $0 < \alpha_1 < \alpha_2 < 1$ and $q_1 < q_2$:

$$p_1 = \frac{q_1 \cdot F^{-1}(\alpha_2) - q_2 \cdot F^{-1}(\alpha_1)}{F^{-1}(\alpha_2) - F^{-1}(\alpha_1)} \quad (2.6)$$

$$p_2 = \frac{q_2 - q_1}{F^{-1}(\alpha_2) - F^{-1}(\alpha_1)} \quad (2.7)$$

The Normal distribution can be truncated, that is the support range may be restricted to a closed interval (i.e. two-sided truncation) or half closed interval (i.e. left- or right-sided truncation). The support of the untruncated Normal distribution $S =] - \infty, \infty [=] a, b [$ may be restricted to the closed interval $S = [a', b']$ with $-\infty < a' < b' < \infty$ representing a two-sided truncation, to the half-closed interval $S = [a', \infty [$ representing a left-sided truncation or to the half-closed interval $S =] - \infty, b']$ representing a right-sided truncation. The density function $f(x)$ and distribution function $F(x)$ of the truncated distribution can be generally derived as:

$$\text{two - sided: } \begin{cases} f_{a',b'}(x) = \frac{f(x)}{F(b') - F(a')} \\ F_{a',b'}(x) = \frac{F(x) - F(a')}{F(b') - F(a')} \end{cases}, \quad x \in [a', b'] \quad (2.8)$$

$$\text{left - sided: } \begin{cases} f_{a'}(x) = \frac{f(x)}{1 - F(a')} \\ F_{a'}(x) = \frac{F(x) - F(a')}{1 - F(a')} \end{cases}, \quad x \in [a', b[\quad (2.9)$$

$$\text{right - sided: } \begin{cases} f_{b'}(x) = \frac{f(x)}{F(b')} \\ F_{b'}(x) = \frac{F(x)}{F(b')} \end{cases}, \quad x \in]a, b'] \quad (2.10)$$

Besides the support range, the Normal distribution may be specified via the distribution parameters p_1 and p_2 or via a set of at least two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) . Note, the parameters of the Normal distribution are equivalent to the first and second statistical moment of the untruncated distribution, that is the mean or expectation value (i.e. $E(X) = p_1$) and the standard deviation (i.e. $SD(X) = p_2$). In case the Normal distribution is truncated the simple random search approach needs to be applied. The Normal distribution as obtained in the visual output of SUSAs is exemplified in the following Fig. 2.1.

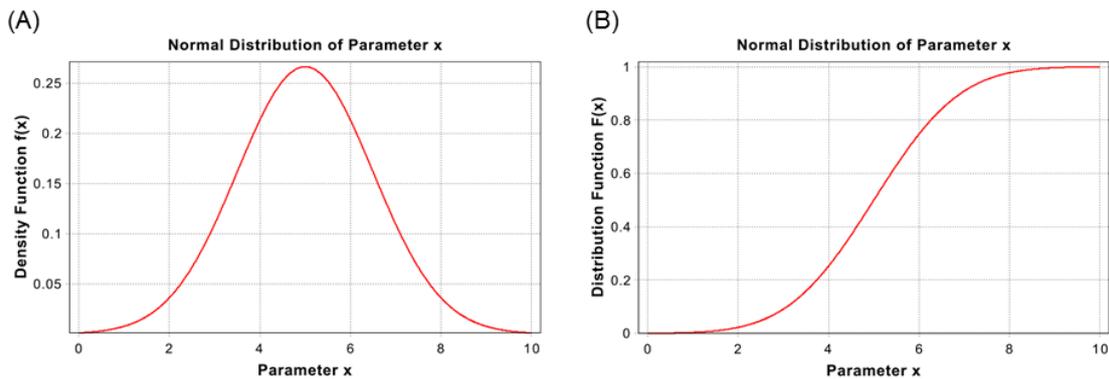


Fig. 2.1 The probability density function (A) and the cumulative distribution function (B) of a Normal distribution with mean $p_1 = 5$ and standard deviation $p_2 = 1$ over the support $[0, 10]$

2.1.1.2 Lognormal Distribution

The logarithmic Normal distribution (also Lognormal distribution) is the distribution of a variable X in case the transformation $\ln(X)$ (natural logarithm of X) is normally distributed. An intuitive comparison between Normal and Lognormal distributions and handy explanation about their deeper understanding is provided in /LIM 01/.

The Lognormal distribution of variable X is defined via the density function as

$$f(x) = \frac{1}{\sqrt{2\pi} p_2} \cdot \frac{1}{x} \cdot \exp\left(-\frac{1}{2} \left(\frac{\ln(x) - p_1}{p_2}\right)^2\right), \quad x \in \mathbb{R}^{>0} \quad (2.11)$$

and via the cumulative distribution function as

$$F(x) = \int_0^x f(t) dt = \frac{1}{2} + \frac{1}{2} \cdot \operatorname{erf}\left(\frac{\ln(x) - p_1}{\sqrt{2} \cdot p_2}\right), \quad x \in \mathbb{R}^{>0} \quad (2.12)$$

where $\mathbb{R}^{>0}$ means the set of positive real numbers (exceeding zero).

The Lognormal distribution is parametrized by the mean p_1 and the standard deviation $p_2 > 0$ (or variance p_2^2) of $\ln(X)$ (natural logarithm of X). The most remarkable intuitive property of the Lognormal distribution depicts its asymmetric, i.e. positively skewed shape, resulting from its bounded support, i.e. the distribution $f(x)$ is only non-zero for positive x values.

The Lognormal distribution can be truncated left-, right or two-sided and may be specified via the distribution parameters p_1 and p_2 or via a set of at least two quantiles each given by the corresponding quantile-probability pair (q_i, a_i) .

Another strategy to specify the Lognormal distribution is given by providing the first and second statistical moments as

$$E(x) = \exp\left(p_1 + \frac{p_2^2}{2}\right) \quad (2.13)$$

and

$$\operatorname{Var}(x) = \exp(2p_1 + p_2^2) \cdot (\exp(p_2^2) - 1) \quad (2.14)$$

Based on these formulas the parameters p_1 and p_2 can then be analytically computed as follows:

$$p_1 = \ln\left(\frac{E(x)}{\sqrt{\operatorname{Var}(x) + E(x)^2}}\right) \quad (2.15)$$

$$p_2^2 = \ln \left(1 + \frac{Var(x)}{E(x)^2} \right) \quad (2.16)$$

Another strategy arises from providing the median M and the factor k_{95} . The median M is the 50 %-quantile of the distribution, i.e. $F(M)=0.5$ with F representing the cumulative distribution function. The factor k_{95} is defined as the ratio of the 95 %-quantile to the 50 %-quantile.

$$p_1 = \ln(M) \quad (2.17)$$

$$p_2 = \frac{\ln(k_{95})}{1.645} \quad (2.18)$$

Any kind of truncation can be applied via the formulas provided in Eqs. (2.8) – (2.10).

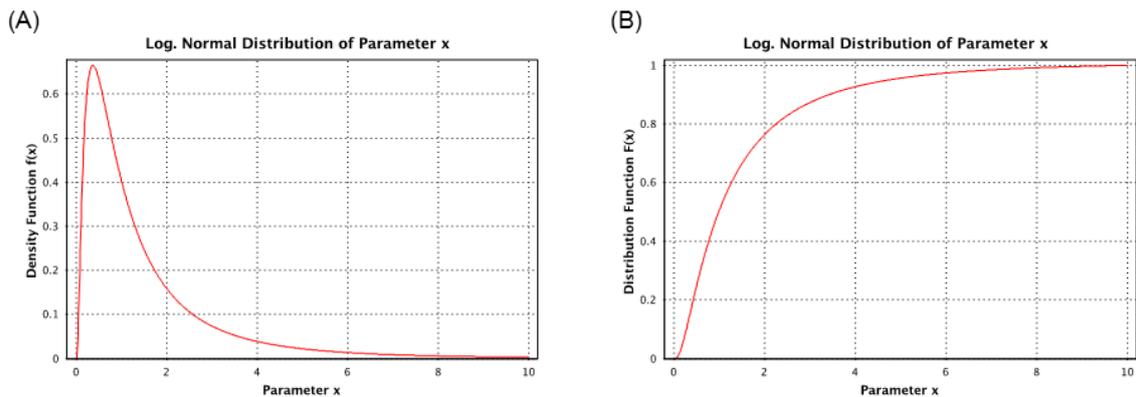


Fig. 2.2 The probability density function (A) and the cumulative distribution function (B) of a Lognormal distribution with parameters $p_1 = 0$ and $p_2 = 1$ over the support $[0, 10]$

2.1.1.3 Uniform Distribution

The Uniform distribution may be assumed in case no a priori information about the uncertainty of an input parameter is available. By providing the lower and upper bound of the support, i.e. $x \in [Min, Max]$, each value within this support is equally likely to occur as an alternative input value. Thus, a Uniform distribution is considered as the most conservative or least-informative uncertainty assumption. In other words, the Uniform

distribution is commonly used in case any value is as likely as any other within a support range.

The Uniform distribution of a variable X is defined via the density function as

$$f(x) = \frac{1}{p_2 - p_1}, \quad x \in [p_1, p_2] \in \mathbb{R} \quad (2.19)$$

and via the cumulative distribution function as

$$F(x) = \frac{x - p_1}{p_2 - p_1}, \quad x \in [p_1, p_2] \in \mathbb{R} \quad (2.20)$$

It is parametrized by the lower bound $p_1 = Min$ and the upper bound $p_2 = Max$ of the support of variable X , i.e. $Min < Max$.

The Uniform distribution may be specified via the distribution parameters p_1 (lower bound) and p_2 (upper bound) or via two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) . Since the support range can be flexibly adapted to any restriction of the support range no truncation option is provided in SUSAS.

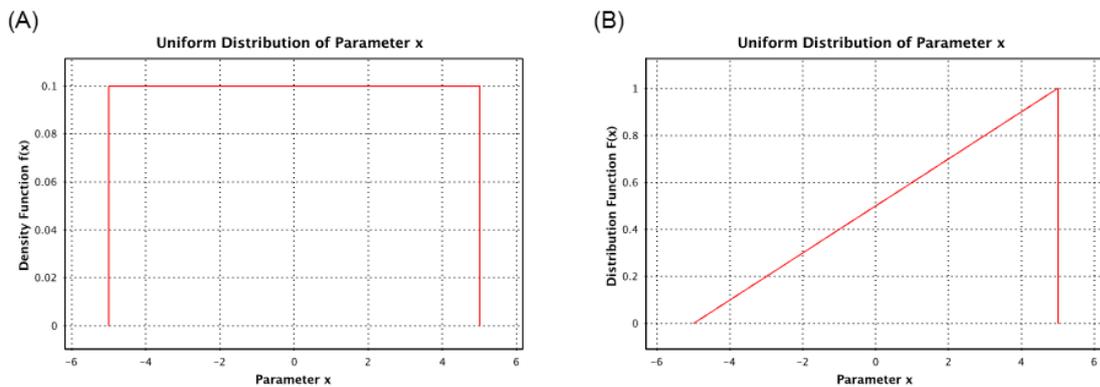


Fig. 2.3 The probability density function (A) and the cumulative distribution function (B) of a Uniform distribution over the support $[-5, 5]$

2.1.1.4 Loguniform Distribution

The logarithmic Uniform distribution (also Loguniform distribution) is the distribution of a variable X in case the transformation $\ln(X)$ (natural logarithm) is uniformly distributed.

The Loguniform distribution is an alternative conservative uncertainty assumption in case where inputs cover large ranges of values, but little is known about their underlying distribution. Thus, the log-transformed input parameter is assumed to be uniformly distributed over a positive support range, i.e. $x \in [Min, Max]$. Due to the log-transformation the support of the resulting Loguniform distribution encloses only positive values, i.e. $0 < Min < Max$, and is of an asymmetric, i.e. positively skewed, shape.

The Loguniform distribution of variable X is defined via the density function as

$$f(x) = \frac{1}{x} \cdot \frac{1}{\ln\left(\frac{p_2}{p_1}\right)}, \quad x \in [p_1, p_2] \in \mathbb{R}^{>0} \quad (2.21)$$

and via the cumulative distribution function as

$$F(x) = \frac{\ln\left(\frac{x}{p_1}\right)}{\ln\left(\frac{p_2}{p_1}\right)}, \quad x \in [p_1, p_2] \in \mathbb{R}^{>0} \quad (2.22)$$

It is parametrized by the positive lower bound $p_1 = Min$ and the positive upper bound $p_2 = Max$ of the support of variable X with $0 < Min < Max$.

The Loguniform distribution may be specified via the distribution parameters p_1 (Min) and p_2 (Max) or via two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) . Since the support range can be flexibly adapted to any restriction of the positive support range, no truncation option is provided in SUSAs.

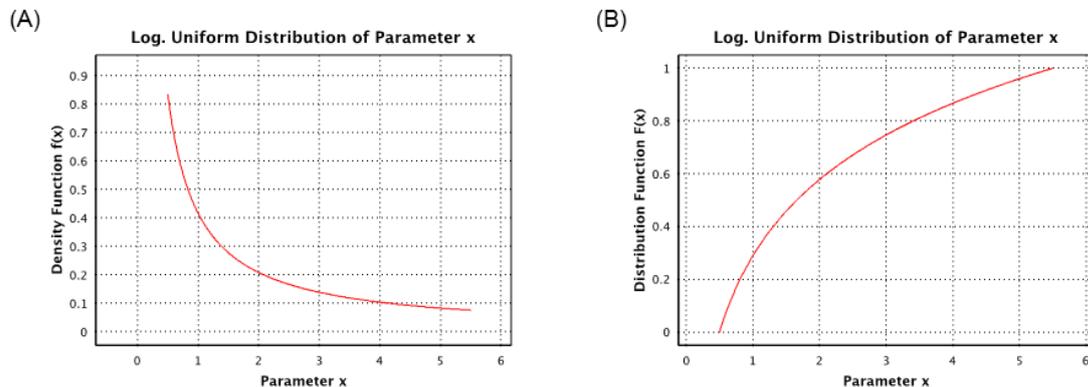


Fig. 2.4 The probability density function (A) and the cumulative distribution function (B) of a Loguniform distribution over the support $[0.5, 5.5]$

2.1.1.5 Triangular Distribution

The Triangular distribution is to a less conservative model of uncertainty, i.e. it encodes more information than the Uniform distribution. Besides the support range, the mode as the location of the highest probability density is used to describe a more informative distribution function. Since no further assumptions are made about the shape of the distribution, a triangular form is considered as the most conservative formulation in this setting.

The Triangular distribution of variable X is defined via the density function as

$$f(x) = \begin{cases} \frac{2 \cdot (x - Min)}{(Max - Min) \cdot (p_1 - Min)}, & \text{for } Min \leq x \leq p_1 \\ \frac{2 \cdot (Max - x)}{(Max - Min) \cdot (Max - p_1)}, & \text{for } p_1 < x \leq Max \end{cases}, x \in \mathbb{R} \quad (2.23)$$

and via the cumulative distribution function as

$$F(x) = \begin{cases} \frac{(x - Min)^2}{(Max - Min) \cdot (p_1 - Min)}, & \text{for } Min \leq x \leq p_1 \\ \frac{(Max - x)^2}{(Max - Min) \cdot (Max - p_1)}, & \text{for } p_1 < x \leq Max \end{cases}, x \in \mathbb{R} \quad (2.24)$$

It is parametrized by the lower bound Min and the upper bound Max of the support of variable X and the mode p_1 with $Min < p_1 < Max$. Depending on the mode value, the Triangular distribution can be symmetric or asymmetric, i.e. skewed with respect to its tail behavior.

The Triangular distribution may be specified via the support range Min , Max and the mode p_1 or via two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) and the mode p_1 . Since the support range can be flexibly adapted to any restriction of the support range, no truncation option is provided in SUSAs.

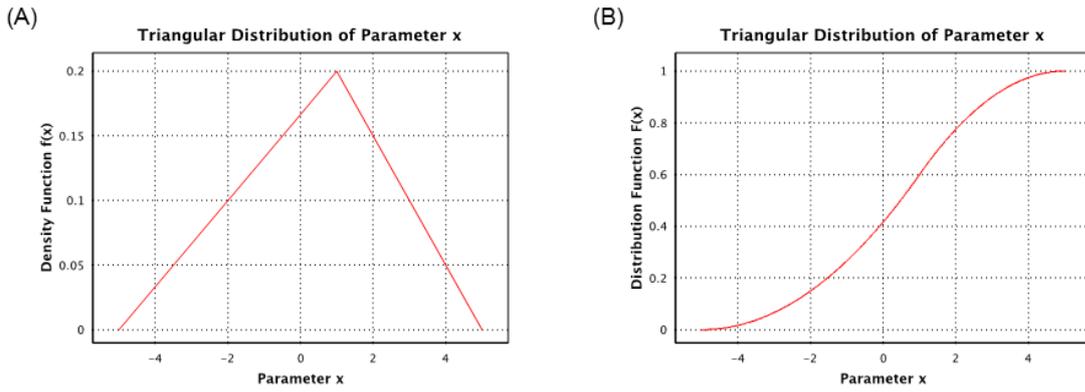


Fig. 2.5 The probability density function (A) and the cumulative distribution function (B) of a Triangular distribution over the support $[-5, 5]$ with mode 1

2.1.1.6 Logtriangular Distribution

The logarithmic Triangular distribution (also Logtriangular distribution) is the distribution of a variable X in case the transformation $\ln(x)$ is triangularly distributed. The Logtriangular distribution is – compared to the Loguniform distribution – a less conservative uncertainty model in case where inputs cover large ranges of values, but little is known about the underlying distribution shape except the mode. Due to the log-transformation, the support $[Min, Max]$ of the resulting Logtriangular distribution encloses only positive values, i.e. $0 < Min < Max$, and is of an asymmetric, i.e. positively skewed, shape.

The Logtriangular distribution of variable X is defined via the density function as

$$f(x) = \begin{cases} \frac{2 \cdot \ln\left(\frac{x}{Min}\right)}{\ln\left(\frac{Max}{Min}\right) \cdot \ln\left(\frac{p_1}{Min}\right)} \cdot \frac{1}{x}, & \text{for } Min \leq x \leq p_1 \\ \frac{2 \cdot \ln\left(\frac{Max}{x}\right)}{\ln\left(\frac{Max}{Min}\right) \cdot \ln\left(\frac{Max}{p_1}\right)} \cdot \frac{1}{x}, & \text{for } p_1 < x \leq Max \end{cases}, x \in \mathbb{R}^{>0} \quad (2.25)$$

and via the cumulative distribution function as

$$F(x) = \begin{cases} \frac{\ln^2\left(\frac{x}{Min}\right)}{\ln\left(\frac{Max}{Min}\right) \cdot \ln\left(\frac{p_1}{Min}\right)}, & \text{for } Min \leq x \leq p_1 \\ 1 - \frac{\ln^2\left(\frac{Max}{x}\right)}{\ln\left(\frac{Max}{Min}\right) \cdot \ln\left(\frac{Max}{p_1}\right)}, & \text{for } p_1 < x \leq Max \end{cases}, x \in \mathbb{R}^{>0} \quad (2.26)$$

It is parametrized by the lower bound Min and the upper bound Max of the support of variable X and, additionally, by the mode p_1 with $0 < Min < p_1 < Max$.

The Logtriangular distribution may be specified via the support range Min , Max and the mode p_1 or via two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) and the mode p_1 . Since the support range can be flexibly adapted to any restriction to the positive support range, no truncation option is provided in SUSAS.

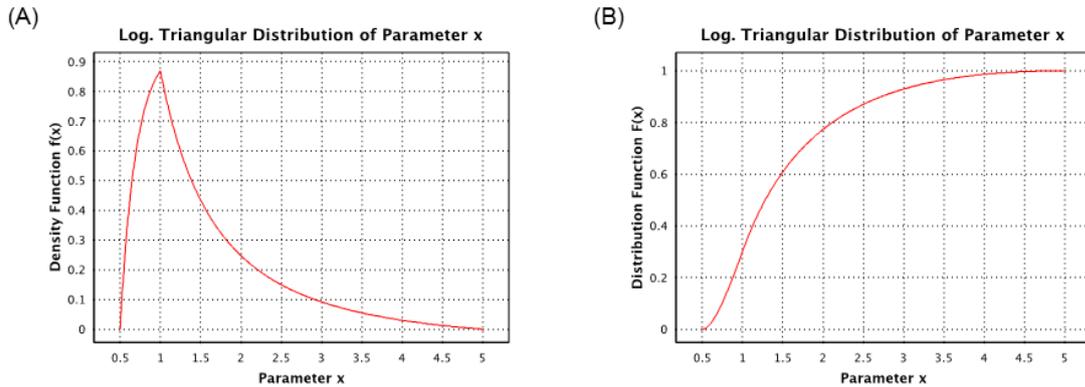


Fig. 2.6 The probability density function (A) and the cumulative distribution function (B) of a Logtriangular distribution over the support $[0.5, 5.0]$ and with mode 1

2.1.1.7 Weibull Distribution

The Weibull distribution is the limiting distribution of the smallest values among a sample of independent, identically distributed random variables (uncertain parameters). The Weibull distribution actually represents a family of distributions that can be flexibly modelled. Thereby, the location, scale and shape parameters allow to adapt the distributions shape in order to describe different uncertainty assumptions. A common application of the Weibull distribution is the modelling of time to failure or length of life of a component from a specified time to its failure as required e.g. in the field of reliability analysis. Moreover, due its versatility in can be easily used to capture the empirical behavior of many physical quantities. Depending on its specification, the Weibull distribution can be used to model a variety of life or occurrence behaviors.

The Weibull distribution of variable X is defined via the density function as

$$f(x) = \frac{p_1}{p_2} \cdot \left(\frac{x - Min}{p_2}\right)^{p_1-1} \cdot \exp\left(-\left(\frac{x - Min}{p_2}\right)^{p_1}\right), x \in \mathbb{R}^{\geq Min} \quad (2.27)$$

and via the cumulative distribution function as

$$F(x) = 1 - \exp\left(-\left(\frac{x - Min}{p_2}\right)^{p_1}\right), \quad x \in \mathbb{R}^{\geq Min} \quad (2.28)$$

The Weibull distribution is parametrized by the location parameter Min , the shape parameter $p_1 > 0$ and the scale parameter $p_2 > 0$. Thereby, the location parameter specifies the lower bound of the support of the distribution, i.e. $X \geq Min$, the shape parameter basically influences the mode of the distribution thus specifying the shape of the distribution and the scale parameter influences the width of the distribution.

The Weibull distribution may be specified via the parameters shape p_1 , scale p_2 , and the lower bound Min of the support range or by a set of at least two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) and the lower bound Min of the support range. The Weibull distribution can be only truncated right-sided.

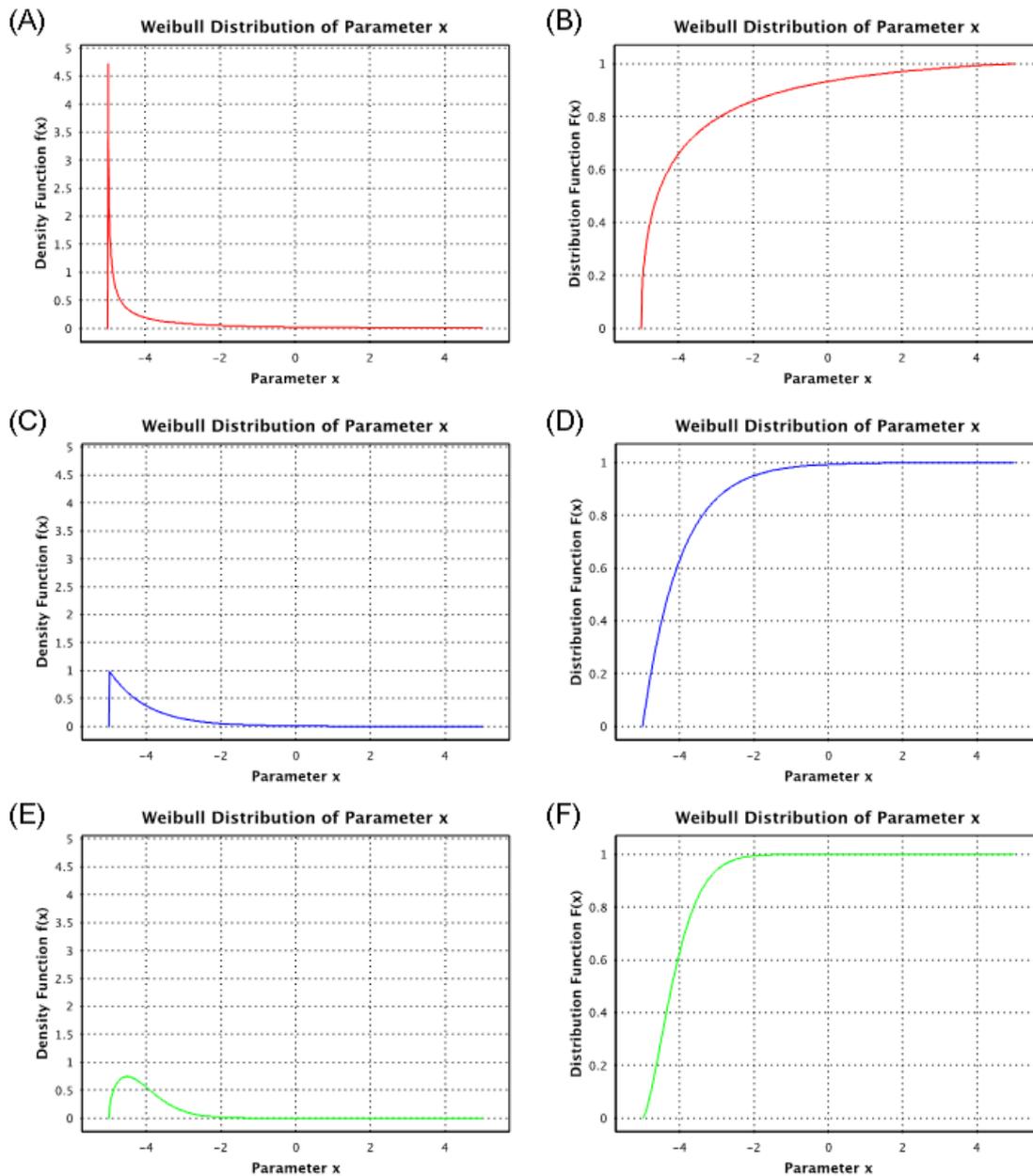


Fig. 2.7 The probability density function (A, C, E) and the cumulative distribution function (B, D, F) of Weibull distributions over the joint support $[-5, 5]$, with the scale parameter $p_2 = 1$ and varying shape parameter $p_1 = 0.5$ for (A, B), 1.0 for (C, D) and 1.5 for (E, F)

2.1.1.8 Beta Distribution

The Beta distribution actually represents a family of distributions that can flexibly model variability patterns over a fixed range. Due to its versatility the Beta distribution can be closely fitted to most classical distributions and, thus, is frequently used in many areas of application. Commonly, it is used to describe the uncertainty of proportions, fractions and percentages.

The Beta distribution of variable X is defined via the density function as

$$f(x) = \frac{(Max - Min)^{1-p_1-p_2}}{B_1(p_1, p_2)} \cdot (x - Min)^{p_1-1} \cdot (Max - x)^{p_2-1}, \quad (2.29)$$
$$x \in [Min, Max] \in \mathbb{R}$$

and via the cumulative distribution function as

$$F(x) = \frac{1}{B_1(p_1, p_2)} \cdot B_{\frac{x-Min}{Max-Min}}(p_1, p_2), \quad x \in [Min, Max] \in \mathbb{R} \quad (2.30)$$

$B_y(p_1, p_2)$ in Eq. (2.30) denotes the incomplete Beta function:

$$B_y(p_1, p_2) = \int_0^y t^{p_1-1} \cdot (1-t)^{p_2} dt \quad (2.31)$$

$B_1(p_1, p_2)$ in Eqs. (2.29) – (2.30) denotes the complete Beta function which is a special case of the incomplete Beta function (Eq. (2.31)) and is defined as

$$B_1(p_1, p_2) = B(p_1, p_2) = \frac{\Gamma(p_1) \cdot \Gamma(p_2)}{\Gamma(p_1 + p_2)} \quad (2.32)$$

where $\Gamma(x)$ denotes the Gamma function as defined in Eq. (2.39).

The Beta distribution is parametrized by the shape parameters $p_1 > 0$ and $p_2 > 0$. In general, the resulting distribution shape is symmetric in case $p_1 = p_2$ and skewed for $p_1 \neq p_2$.

The Beta distribution may be specified via the shape parameters p_1, p_2 , and the bounds Min and Max of the support range or by a set of at least two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) and the bounds Min and Max of the support range. Another strategy to specify the Beta distribution is given by providing the first and second statistical moments as

$$E(x) = Min + (Max - Min) \cdot \frac{p_1}{p_1 + p_2} \quad (2.33)$$

and

$$Var(x) = (Max - Min)^2 \cdot \frac{p_1 \cdot p_2}{(p_1 + p_2)^2 \cdot (p_1 + p_2 + 1)} \quad (2.34)$$

The shape parameters p_1 and p_2 can then be analytically computed as

$$p_1 = E(x) \cdot \left(\frac{E(x)}{Var(x)} \cdot (1 - E(x)) - 1 \right) \quad (2.35)$$

and

$$p_2 = (1 - E(x)) \cdot \left(\frac{E(x)}{Var(x)} \cdot (1 - E(x)) - 1 \right) \quad (2.36)$$

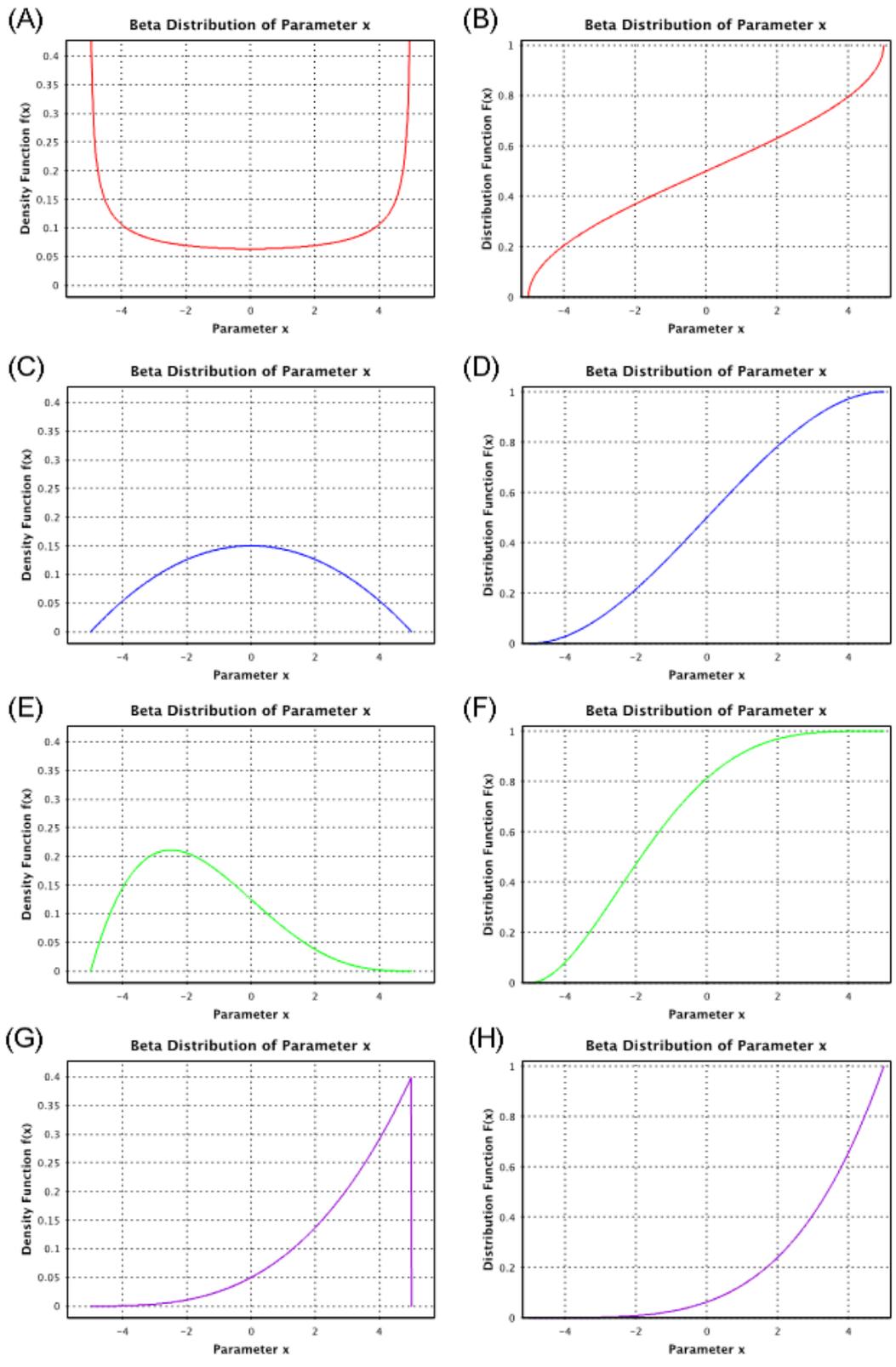


Fig. 2.8 The probability density function (A, C, E, G) and the cumulative distribution function (B, D, F, H) of Beta distributions over the joint support $[-5, 5]$, with the shape parameters $(p_1, p_2) = (0.5, 0.5)$ for (A, B), $(2, 2)$ for (C, D), $(2, 4)$ for (E, F) and $(4, 1)$ for (G, H)

2.1.1.9 Gamma Distribution

The shape of the Gamma distribution is similar to that of the Lognormal distribution but it is less positively skewed and less heavy-tailed. Thus, it may be used in situations similar to those where the Lognormal distribution would be appropriate. The Gamma distribution is particularly useful for describing the uncertainty on the time (or spaces) between events (items) that are not pure random processes. Quantities frequently exhibiting skewed distributions like the Gamma distribution are physical quantities as well as the time between malfunctions of components or the time required to complete maintenance of a component.

The Gamma distribution is defined via the density function as

$$f(x) = \frac{1}{\Gamma(p_1)} \cdot p_2^{p_1} \cdot x^{p_1-1} \cdot e^{-p_2 \cdot x}, \quad x \in \mathbb{R}^{>0} \quad (2.37)$$

and via the cumulative distribution function as

$$F(x) = \frac{\Gamma_{p_2 \cdot x}(p_1)}{\Gamma(p_1)}, \quad x \in \mathbb{R}^{>0} \quad (2.38)$$

It is parametrized by the shape parameter $p_1 > 0$, the rate (or inverse scale) parameter $p_2 > 0$ and the Gamma function $\Gamma(p_1)$. The incomplete Gamma function is denoted by

$$\Gamma_y(p_1) = \int_0^y t^{p_1-1} \cdot e^{-t} dt \quad (2.39)$$

with the complete Gamma function derived as $\Gamma_\infty(p_1) = \Gamma(p_1)$.

The Gamma distribution may be specified via the parameters p_1 (shape) and p_2 (rate) over the support range or by a set of at least two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) over the support range. Any kind of truncation can be applied via the formulas provided in Eqs. (2.8) – (2.10).

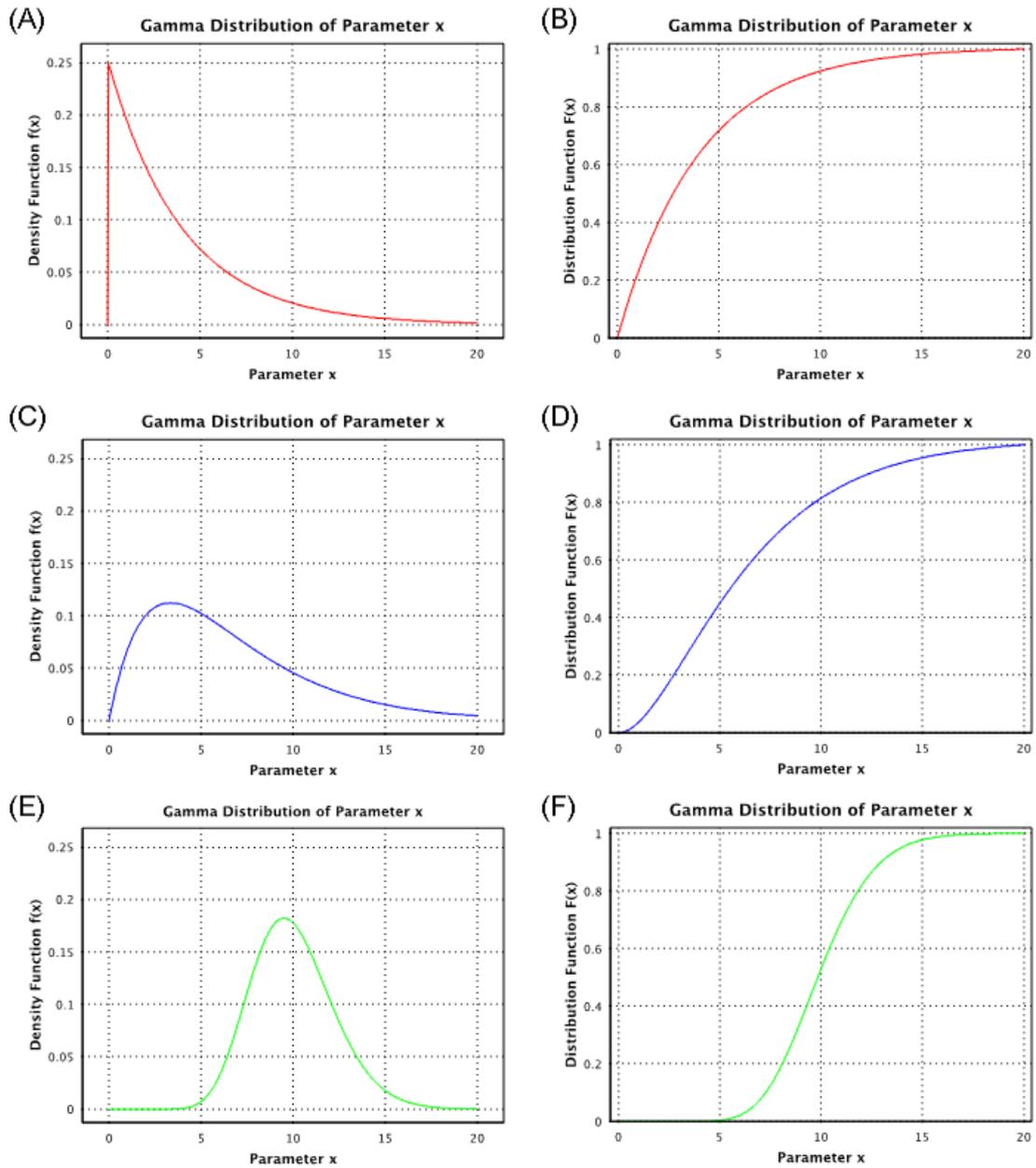


Fig. 2.9 The probability density function (A, C, E) and the cumulative distribution function (B, D, F) of Gamma distributions over the joint support $[0, 20]$ with the parameter settings $(p_1, p_2) = (1, 0.25)$ for (A, B), $(2, 0.3)$ for (C, D) and $(20, 2)$ for (E, F)

2.1.1.10 Extreme Value I Distribution

In many engineering applications, the extreme values of random variables (uncertain parameters) are of particular importance. The largest or smallest values may dictate a particular component or system design, e.g. wind speed, earthquake loads or flood levels. To construct an extreme value distribution, an underlying variable with a particular distribution is necessary. This underlying distribution governs the form of the corresponding extreme value distribution. Further reading about extreme values and their distributional patterns in the context of engineered systems can be found in /HAL 00/ and /JOR 05/.

The Extreme Value I distribution (*aka* Gumbel distribution) is a limiting distribution of the largest (or smallest) value among a sample of independent, identically distributed random variables (uncertain parameters). Thus, the Extreme Value I distribution is used to model the distribution of the maximum (or the minimum) of a number of samples. This distribution might be used to represent the uncertainty of the maximum measurement outcome or the uncertainty of an extreme event such as an earthquake, flood or other natural disaster. The potential applicability of the Gumbel distribution to represent the uncertainty of the extreme value of a sample originates from the extreme value theory, which indicates that the Gumbel distribution is likely to be useful if the distribution of the underlying sample data is of the Normal or Exponential type.

The Extreme Value I distribution is defined via the density function as

$$f(x) = \frac{1}{p_2} \cdot \exp\left(-\frac{x - p_1}{p_2}\right) \cdot \exp\left(-\exp\left(\frac{x - p_1}{p_2}\right)\right), \quad x \in \mathbb{R} \quad (2.40)$$

and via the cumulative distribution function as

$$F(x) = \exp\left(-\exp\left(\frac{x - p_1}{p_2}\right)\right), \quad x \in \mathbb{R} \quad (2.41)$$

It is parametrized by the location parameter p_1 and the scale parameter $p_2 > 0$.

The Extreme Value I distribution may be specified via the parameters p_1 (location) and p_2 (scale) or by a set of at least two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) . Any kind of truncation can be applied via the formulas provided in Eqs. (2.8) – (2.10).

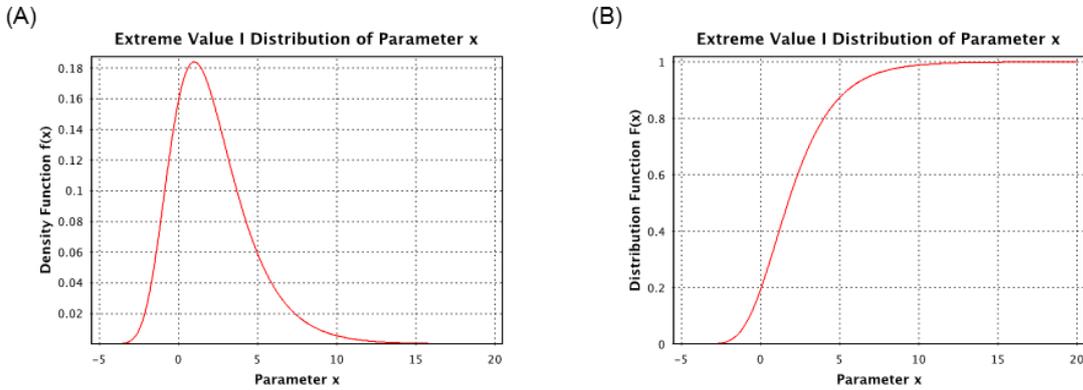


Fig. 2.10 The probability density function (A) and the cumulative distribution function (B) of an Extreme Value I distribution over the support $[-5, 20]$ for the location parameter $p_1 = 1$ and the scale parameter $p_2 = 2$

2.1.1.11 Extreme Value II Distribution

The Extreme Value II distribution (*aka* Fréchet distribution, *aka* inverse Weibull distribution) is a limiting distribution of the largest (or smallest) value among a sample of independent, identically distributed random variables (uncertain parameters). The Extreme Value II distribution is used to model the distribution of the maximum (or the minimum) of a number of samples. This distribution might be used to represent the uncertainty of the maximum measurement outcomes, since it captures the typical characteristics of asymmetric long tails. It is useful to represent the uncertainty of an extreme event such as an earthquake, flood or other natural disaster. The potential applicability of the Fréchet distribution to represent the distribution of the extreme value of a sample originates from the extreme value theory, which indicates that of the Fréchet distribution is likely to be useful, if the distribution of the underlying sample data is of the Cauchy or Lognormal type.

The Extreme Value II distribution of variable X is defined via the density function as

$$f(x) = \frac{p_1}{p_2} \cdot \left(\frac{x - Min}{p_2}\right)^{-(p_1+1)} \cdot \exp\left(-\left(\frac{x - Min}{p_2}\right)^{-p_1}\right), x \in \mathbb{R}^{>Min} \quad (2.42)$$

and via the cumulative distribution function as

$$F(x) = \exp\left(-\left(\frac{x - \text{Min}}{p_2}\right)^{-p_1}\right) \quad x \in \mathbb{R}^{>\text{Min}} \quad (2.43)$$

It is parametrized by the lower bound Min , i.e. $X > \text{Min}$, the shape parameter $p_1 > 0$ and the scale parameter $p_2 > 0$.

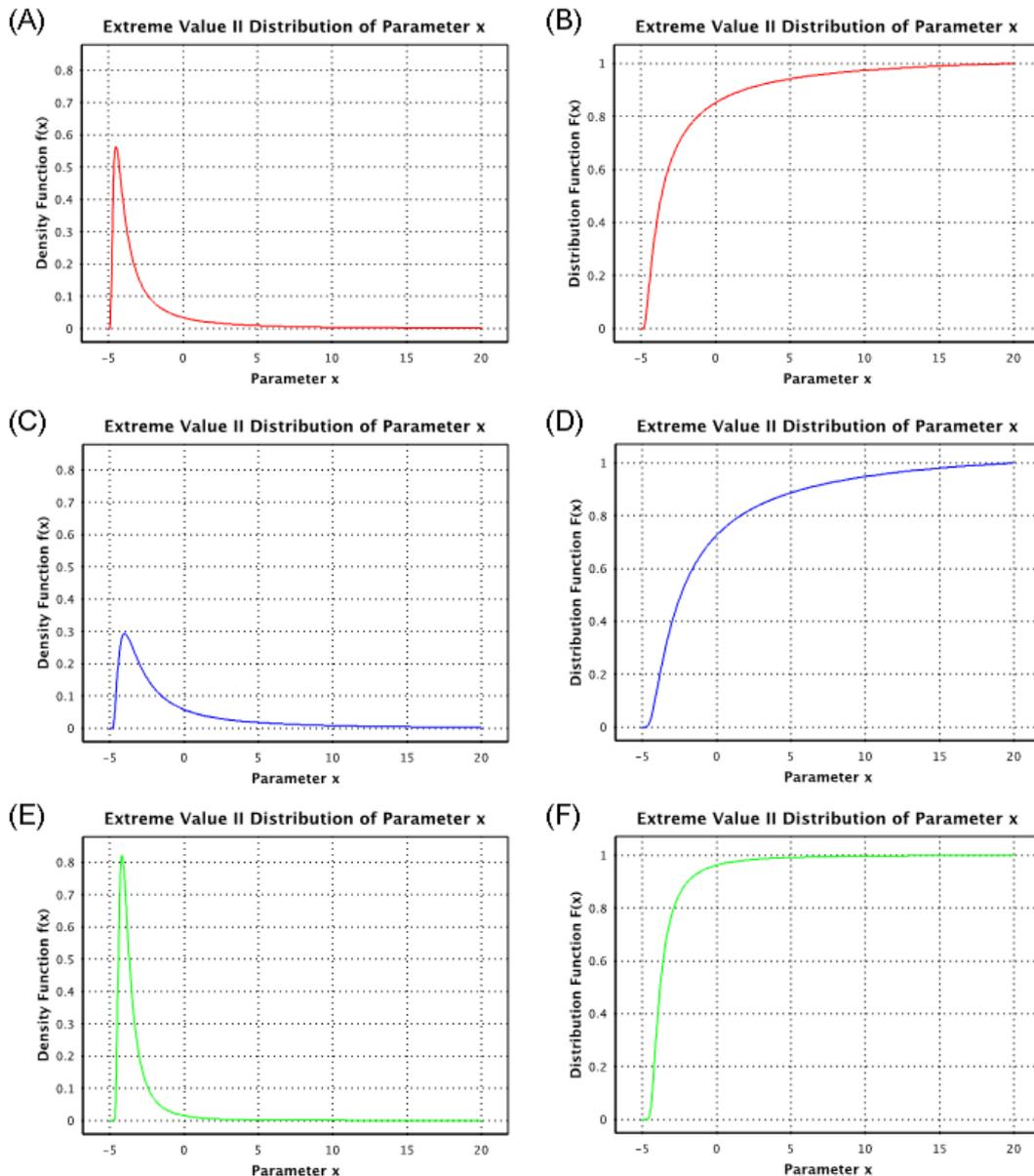


Fig. 2.11 The probability density function (A,C,E) and the cumulative distribution function (B,D,F) of Extreme Value II distributions over the joint support $[-5, 20]$, i.e. $\text{Min} = -5$, for the shape and scale parameters $(p_1, p_2) = (1, 1)$ for (A, B), $(1, 2)$ for (C, D) and $(2, 1)$ for (E, F)

The Extreme Value II distribution may be specified via the parameters p_1 (shape) and p_2 (scale) over the support range bounded by Min or by a set of at least two quantiles each given by the corresponding quantile-probability pair (q_i, α_i) over the support range bounded by Min . It can be only truncated right-sided.

2.1.1.12 Exponential Distribution

The Exponential distribution is often used in reliability theory and reliability engineering. The distribution describes the uncertainty of the time between events in a Poisson process, i.e. a process in which events occur continuously and independently at a constant average rate, e.g. failure rate of a system or component. Because of the memoryless property of this distribution, it is well-suited to model the constant hazard rate portion of the bathtub curve used in reliability theory. However, the exponential distribution is not appropriate to model the overall lifetime of technical devices, because realistic failure rates are not constant: more failures occur for very young and for very old systems. In physics, if you observe a gas at a fixed temperature and pressure in a uniform gravitational field, the heights of the various molecules also follow an approximate exponential distribution, known as the Barometric formula. This correspondence is a consequence of the entropy property and is often used to justify the use of the Exponential distribution for certain physical quantities.

The Exponential distribution is defined via the density function as

$$f(x) = p_1 \cdot e^{-p_1 \cdot x}, \quad x \in \mathbb{R}^{>0} \quad (2.44)$$

and via the cumulative distribution function as

$$F(x) = 1 - e^{-p_1 \cdot x}, \quad x \in \mathbb{R}^{>0} \quad (2.45)$$

It is parametrized by the rate (or inverse scale) parameter $p_1 > 0$ over a positive support range.

The Exponential distribution may be specified via the parameter p_1 (rate) or by a set of at least one quantile given by the corresponding quantile-probability pair (q_i, α_i) . Any kind of truncation can be applied via the formulas provided in Eqs. (2.8) – (2.10).

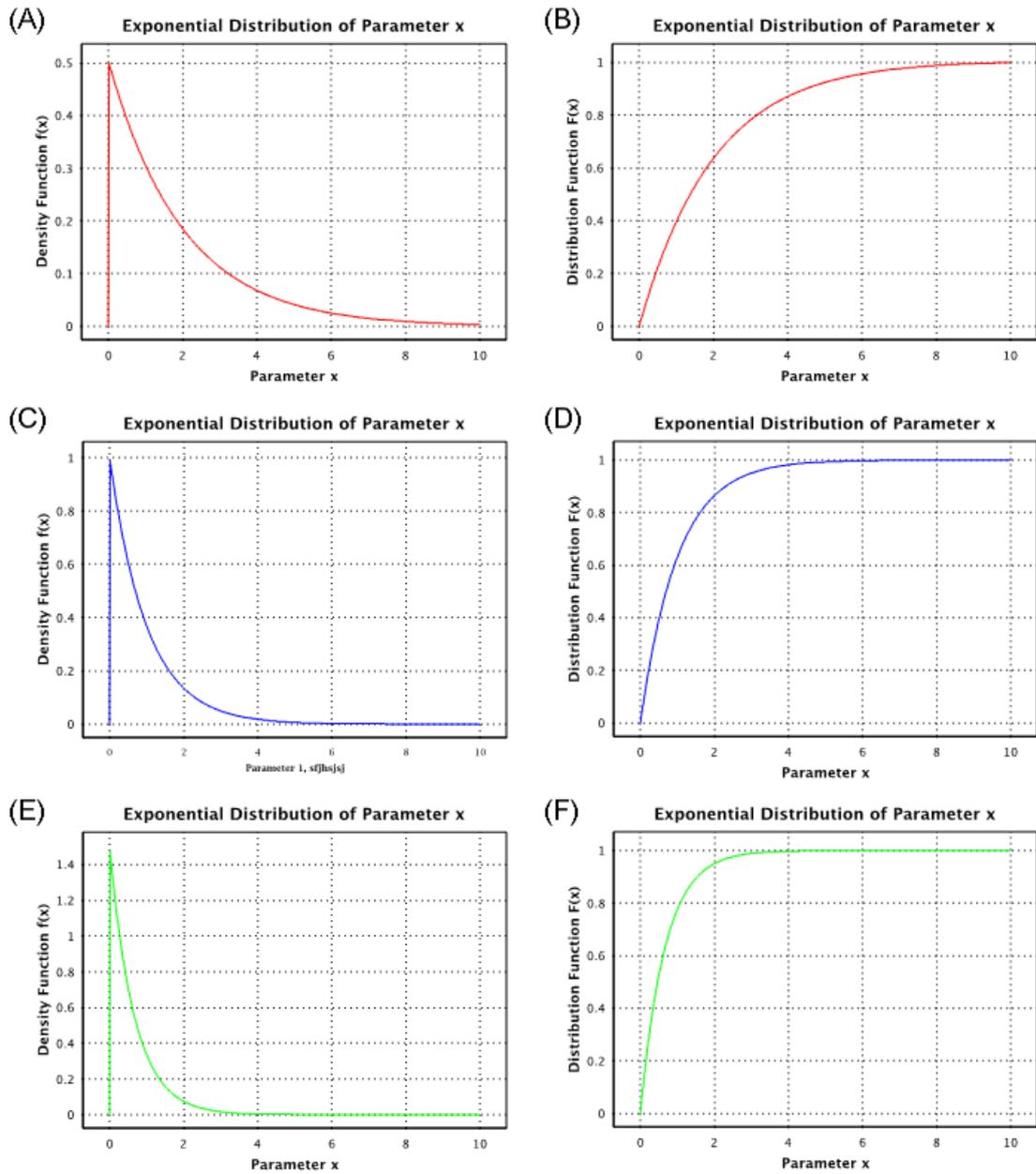


Fig. 2.12 The probability density function (A,C,E) and the cumulative distribution function (B,D,F) of Exponential distributions over the joint support $[0, 10]$ for the rate parameter $\rho_1 = 0.5$ for (A, B), 1.0 for (C, D) and 1.5 for (E, F)

2.1.1.13 ChiSquared Distribution

The ChiSquared (χ^2) distribution describes the uncertainty of a squared random variable which itself is normally distributed. The χ^2 -distribution is usually applied to the sum of square error between measured data and corresponding model predictions.

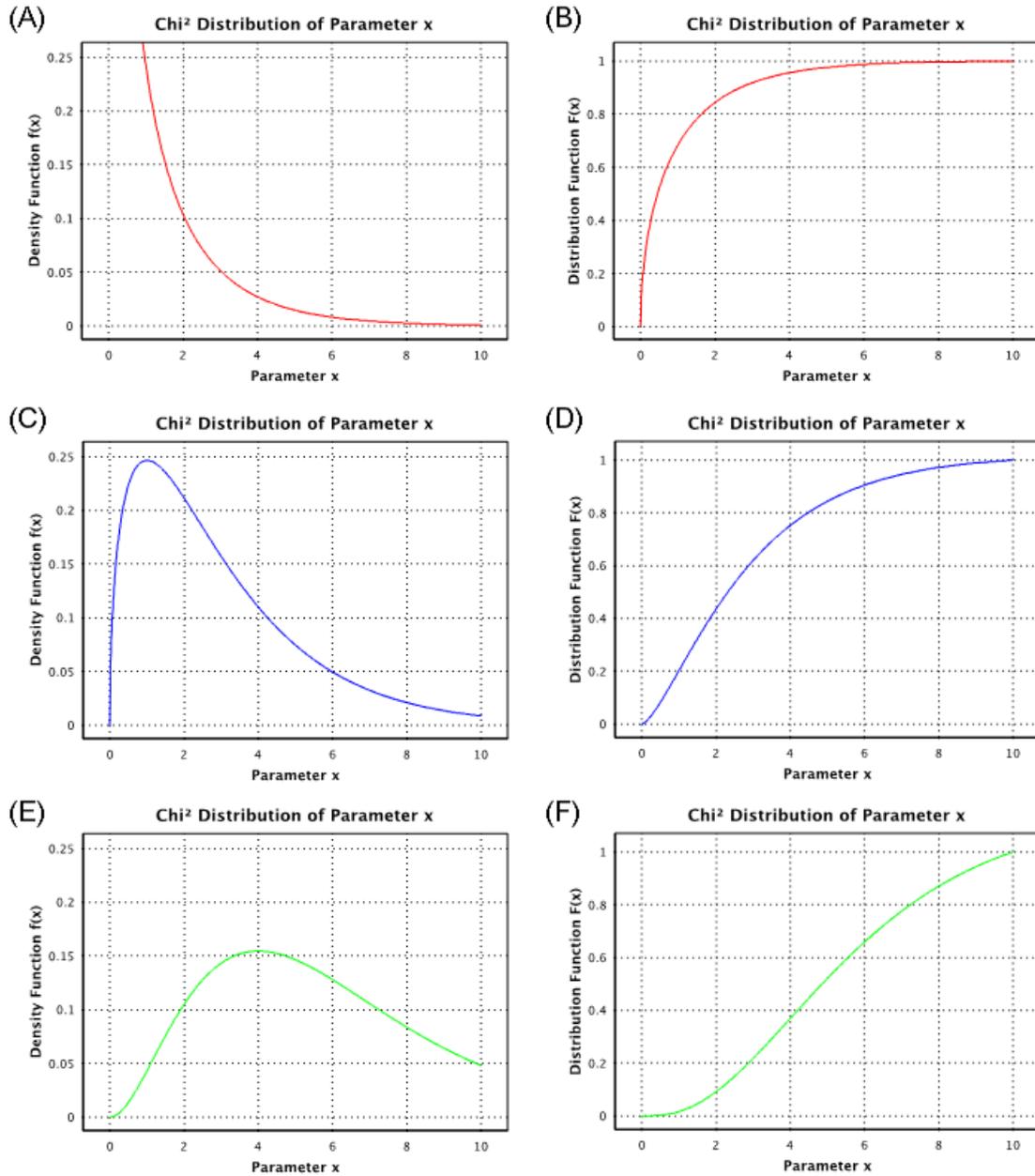


Fig. 2.13 The probability density function (A, C, E) and the cumulative distribution function (B, D, F) of ChiSquared distributions over the joint support [0, 10] for the degree of freedom $p_1 = 1$ for (A, B), 3 for (C, D) and 6 for (E, F)

The χ^2 -distribution is defined via the density function as

$$f(x) = \frac{1}{\Gamma\left(\frac{p_1}{2}\right)} \cdot 0.5^{\frac{p_1}{2}} \cdot x^{\frac{p_1}{2}-1} \cdot \exp\left(-\frac{x}{2}\right), \quad x \in \mathbb{R}^{>0} \quad (2.46)$$

and via the cumulative distribution function as

$$F(x) = \frac{\Gamma_x\left(\frac{p_1}{2}\right)}{\Gamma\left(\frac{p_1}{2}\right)}, \quad x \in \mathbb{R}^{>0} \quad (2.47)$$

It is parametrized by the degree of freedom $p_1 \in \mathbb{N}^{>0}$ over a positive support range.

The χ^2 -distribution may be specified via the degree of freedom parameter p_1 or by a set of at least one quantile given by the corresponding quantile-probability pair (q_i, α_i) . Any kind of truncation can be applied via the formulas provided in Eqs. (2.8) – (2.10).

2.1.2 Nonparametric Distribution

If, for any reason, it is not possible to represent the uncertainty of a variable by a parametric distribution, the specification of a nonparametric distribution may be appropriate.

2.1.2.1 Discrete Distribution

The selection of the Discrete distribution is appropriate, if the uncertainty can be modelled in a discrete manner by assigning a probability (degree of belief) to each possible value of a parameter X , i.e. by specifying value-probability pairs (x_i, p_i) for parameter X .

Based on the discrete support $S = \{z_1, z_2, \dots, z_n\} \in \mathbb{R}$ where each parameter value $z_i, i = 1, \dots, n$ and $n \in \mathbb{N}$, is associated with a probability $p_i \in [0,1]$, the Discrete distribution is defined via the density function as

$$f(x) = \begin{cases} p_i, & \text{for } x = z_i \ \forall \ i = 1, \dots, n \\ 0, & \text{else} \end{cases} \quad (2.48)$$

and via the cumulative distribution function as

$$F(x) = \begin{cases} 0, & \text{for } x < z_1 \\ \sum_{j=1}^i p_j, & \text{for } z_i \leq x < z_{i+1} \quad \forall i = 1, \dots, n-1 \\ 1, & \text{for } x \geq z_n \end{cases} \quad (2.49)$$

Note, in order that these functional expressions satisfy the basic properties of a probability function (namely a probability measure as formulated within probability theory), the discrete probability values p_i must lie in the interval $[0, 1]$ and their total sum needs to be equal to one, i.e.

$$\sum_{i=1}^n p_i = 1 \quad \forall p_i \in [0,1] \quad (2.50)$$

Moreover, for SUSAs it is important that the provided parameter values are ordered, that is $z_1 < z_2 < \dots < z_n$, and are assigned with probability values $p_i > 0$. Since the support range can be flexibly adapted to any restriction of the support range, no truncation option is provided in SUSAs.

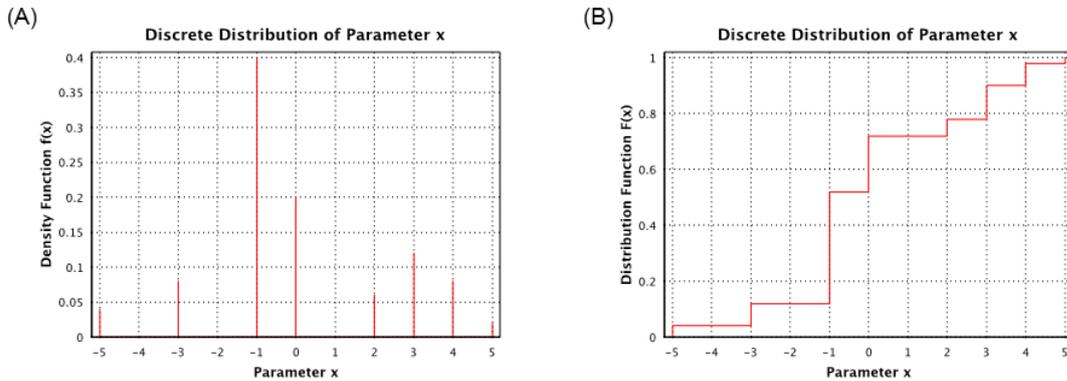


Fig. 2.14 The probability density function (A) and the cumulative distribution function (B) of a Discrete distribution for the value-probability pairs $\{(-5, 0.04), (-3, 0.08), (-1, 0.40), (0, 0.20), (2, 0.06), (3, 0.12), (4, 0.08), (5, 0.02)\}$ over the support $[-5, 5]$

2.1.2.2 Histogram Distribution

The Histogram distribution allows for modelling a piecewise Uniform distribution shape from a set of discrete value-probability pairs. This might be reasonable in case constant probability values between successive parameter values can be assumed.

Based on the discrete support $S = \{z_1, z_2, \dots, z_{n+1}\} \in \mathbb{R}$ where each parameter value z_i for $i = 1, \dots, n$ and $n \in \mathbb{N}$ is associated with a probability value $p_i \in [0,1]$ applicable to the interval $[z_i, z_{i+1}[$, the Histogram distribution can be defined via the density function as

$$f(x) = \begin{cases} \frac{p_i}{z_{i+1} - z_i}, & \text{for } x \in [z_i, z_{i+1}[\vee i = 1, \dots, n \\ \frac{p_n}{z_{n+1} - z_n}, & \text{for } x = z_{n+1} \end{cases} \quad (2.51)$$

and via the cumulative distribution function as

$$F(x) = \begin{cases} \frac{x - z_1}{z_2 - z_1} \cdot p_1, & \text{for } x \in [z_1, z_2[\\ \sum_{j=1}^{i-1} p_j + \frac{x - z_i}{z_{i+1} - z_i} \cdot p_i, & \text{for } x \in [z_i, z_{i+1}[\vee i = 2, \dots, n \\ 1, & \text{for } x = z_{n+1} \end{cases} \quad (2.52)$$

In order to satisfy the basic properties of a probability function (leading to the definition as a measure within probability theory), the discrete probability values p_i need to lie in the interval $[0, 1]$ and their total sum needs to be equal to one as formulated in Eq. (2.50). Since the support range can be flexibly adapted to any restriction of the support range, no truncation option is provided in SUSAS.

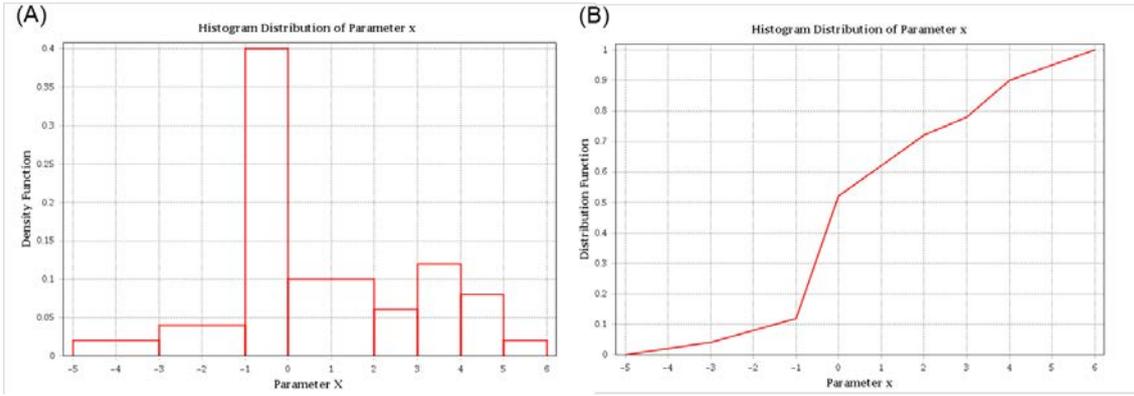


Fig. 2.15 The probability density function (A) and the cumulative distribution function (B) of a Histogram distribution for the value-probability pairs $\{(-5, 0.04), (-3, 0.08), (-1, 0.40), (0, 0.20), (2, 0.06), (3, 0.12), (4, 0.08), (5, 0.02)\}$ over the support $[-5, 6]$

2.1.2.3 Loghistogram Distribution

The logarithmic Histogram distribution (also Loghistogram distribution) allows for modeling a piecewise-logarithmic decreasing, continuous distribution shape from a set of discrete value-probability pairs.

Based on the discrete support $S = \{z_1, z_2, \dots, z_{n+1}\} \in \mathbb{R}^{>0}$ where each parameter value z_i for $i = 1, \dots, n$ and $n \in \mathbb{N}$ is associated with a probability value $p_i \in [0,1]$, the Loghistogram distribution can be defined via the density function as

$$f(x) = \begin{cases} \frac{p_i}{x \cdot \ln\left(\frac{z_{i+1}}{z_i}\right)}, & \text{for } x \in [z_i, z_{i+1}[\vee i = 1, \dots, n \\ \frac{p_n}{x \cdot \ln\left(\frac{z_n}{z_{n+1}}\right)}, & \text{for } x = z_{n+1} \end{cases} \quad (2.53)$$

and via the cumulative distribution function as

$$F(x) = \begin{cases} \frac{\ln\left(\frac{x}{z_1}\right)}{\ln\left(\frac{z_2}{z_1}\right)} \cdot p_1, & \text{for } x \in [z_1, z_2[\\ \sum_{j=1}^{i-1} p_j + \frac{\ln\left(\frac{x}{z_i}\right)}{\ln\left(\frac{z_{i+1}}{z_i}\right)} \cdot p_i, & \text{for } x \in [z_i, z_{i+1}[\quad \forall i = 2, \dots, n \\ 1, & \text{for } x = z_{n+1} \end{cases} \quad (2.54)$$

In order to satisfy the basic properties of a probability function (leading to the definition as a measure within probability theory), the discrete probability values p_i need to lie in the interval $[0, 1]$ and their total sum needs to be equal to one as formulated in Eq. (2.50). Since the support can be flexibly adapted to any restriction of the support range, no truncation option is provided in SUSAS.

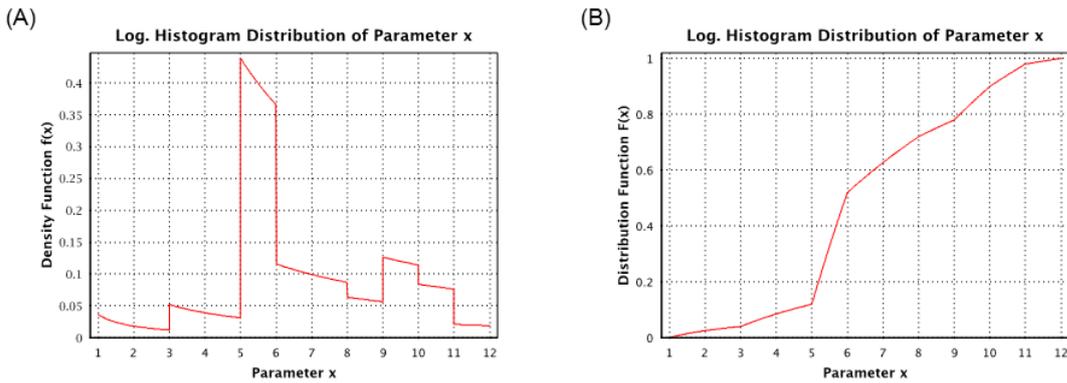


Fig. 2.16 The probability density function (A) and the cumulative distribution function (B) of a Loghistogram distribution for the value-probability pairs $\{(1, 0.04), (3, 0.08), (5, 0.40), (6, 0.20), (8, 0.06), (9, 0.12), (10, 0.08), (11, 0.02)\}$ over the support $[1, 12]$

2.1.2.4 Polygonal Line Distribution

In cases where it might be reasonable to assume a uniform or linear behavior of the probability density between successive discrete parameter values (i.e. piecewise linear density functions) a Polygonal Line distribution can be used. The Polygonal Line distribution particularly allows for flexibly modeling multimodal distribution patterns without indicating specific or complex shape patterns. In that sense, the value-probability pairs

to be specified for the Polygonal Line distribution represent (x, y) -coordinates as the base points of simple linear spline functions.

Based on the discrete support $S = \{z_1, z_2, \dots, z_{n+1}\} \in \mathbb{R}$ where each parameter value z_i for $i = 1, \dots, n + 1$ and $n \in \mathbb{N}$ is associated with a value $y_i \in [0,1]$ (representing the relative height of the density function at value z_i), the Polygonal Line distribution can be defined via the density function as

$$f(x) = p_{i-1} + (x - z_{i-1}) \frac{p_i - p_{i-1}}{z_i - z_{i-1}}, \quad \text{for } x \in [z_{i-1}, z_i[, i = 2, \dots, n \quad (2.55)$$

and via the cumulative distribution function as

$$F(x) = 0.5 \cdot \sum_{j=2}^{i-1} (z_j - z_{j-1}) \cdot (p_j + p_{j-1}) + (x - z_{i-1}) \cdot (p_i + p_{i-1}), \quad (2.56)$$

for $x \in [z_{i-1}, z_i[, i = 2, \dots, n$

where $p_i, i=1, \dots, n$, is calculated from the given z_i and y_i values as

$$p_i = \frac{y_i}{A} \quad (2.57)$$

with

$$A = \sum_{i=2}^n 0.5 \cdot (z_i - z_{i-1}) \cdot (y_i + y_{i-1}) \quad (2.58)$$

The discrete values y_i need to lie in the interval $[0, 1]$. In order that the area enclosed by the polygonal line may be considered as a probability distribution, SUSANA calculates the actual values $p_i, i=1, \dots, n$, of the density function according to Eqs. (2.57) – (2.58). Since the support range can be flexibly adapted to any restriction of the support range, a truncation option is not necessary.

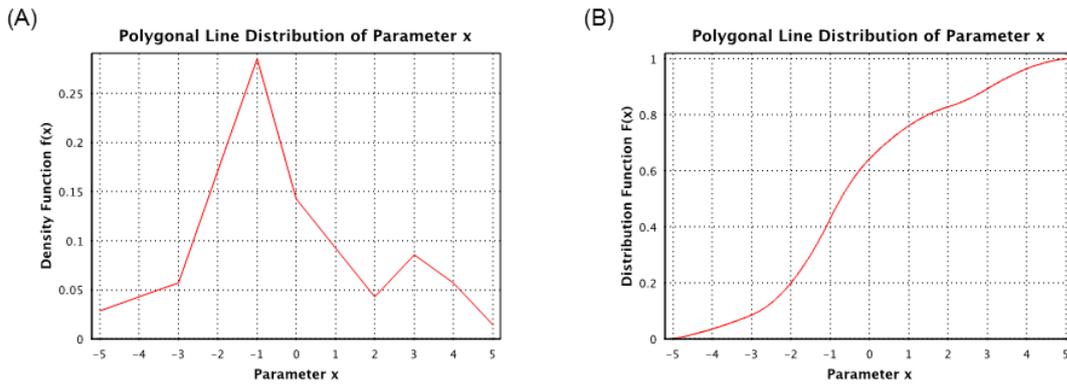


Fig. 2.17 The probability density function (A) and the cumulative distribution function (B) of a Polygonal Line distribution for the (x, y) -coordinates $\{(-5, 0.04), (-3, 0.08), (-1, 0.40), (0, 0.20), (2, 0.06), (3, 0.12), (4, 0.08), (5, 0.02)\}$ over the support $[-5, 5]$

2.2 Dependency

Besides the univariate probability distribution representing the experts' degree of belief on the possible values of a parameter, the specification of dependencies between different parameters may be important to appropriately model input uncertainties. In the context of epistemic uncertainties, dependency has to be interpreted as knowledge or epistemic dependency meaning the dependency between the states of knowledge on two or more parameters due to a common contribution. For instance, the failure rates of two components which are nominally identical but physically different are completely dependent, if the same data pool is used to estimate the failure rates /APO 81/. Two uncertain parameters are dependent, if they represent measurements from the same measuring instrument. These measurements are subject to an independent error and to a common uncertain bias term due to the measuring instrument.

For the representation of knowledge dependency, SUSA offers a variety of approaches to correlate uncertain parameters. The pairwise dependency between two uncertain parameters can be quantified by measures of association (Pearson's ordinary correlation, Spearman's rank correlation, Blomqvist's medial correlation, and Kendall's rank correlation), full dependency, conditional distributions or by inequalities. Additionally, a parameter can be specified as function of one or more other parameters.

In particular the measures of associations may be divided into two classes with respect to the sample of parameter values to be generated (experimental design /KRZ 88/):

- **population-related:** experimental design complies with specified properties of a desired multivariate population
- **sample-related:** experimental design complies with empirical properties

Most of the association measures have an intuitive interpretation which facilitates the transfer of experts' knowledge into knowledge dependency. In order to check whether the specified measure of association results in a bivariate sample that adequately meets the desired dependency pattern, the resulting bivariate sample in the plane of the support may be visualized as a scatter plot.

Besides the measures of association which are mainly non-parametric (distribution-free) approaches, further flexibility for the experimental design may be accomplished by considering association of uncertain parameters by the following approaches:

- **full dependency:** to model a completely positive or negative dependency between two parameters,
- **function of parameters:** to model an explicit deterministic functional relationship between multiple parameters,
- **conditional distributions:** to model a varying state of knowledge about one parameter for different parameter ranges of another parameter, or
- **inequality:** to model a deterministic boundary condition within the parameter space of two parameters.

All these approaches to model dependencies between the input uncertainties are motivated by different practical challenges the analyst faces when confronted with the task to transfer experts' knowledge into correlated mathematical/statistical expressions. To make optimal use of the provided dependency concepts, the basic mathematical/statistical principles are explained and the interpretation of assumptions are described within a hypothetical parameter plane determined by the support plane of the distributions as defined in Section 2.1.

2.2.1 Population-related correlation

Population-related correlation coefficients (measures of association) are used in case the sample of parameter values (experimental design) shall comply with the specified properties indicating the joint multivariate distribution of a population. These coefficients, here referred to as $corr(X, Y)$, between two parameters X and Y have at least the following basic properties:

- $-1 \leq corr(X, Y) \leq +1$, i.e. symmetric measure
- X, Y independent $\Rightarrow corr(X, Y) = 0$ (the reverse does not hold generally)
- $corr(X, Y) = 1 (-1) \Rightarrow Y$ is an increasing (decreasing) function of X

Apart from Pearson's ordinary correlation coefficient, all population-related measures are scale (or ordinally) invariant. This is explained in more detail in the context of Blomqvist's medial correlation coefficient (Section 2.2.1.2) but basically indicates, that the value of $corr(X, Y)$ is not changing in case a monotone transformation is applied to the parameters X and Y . Due to the intuitive interpretation of such ordinal measures, i.e. measure without any specific metric, they are often considered as the most conservative approach to encode subjective judgements of complex parameter dependency.

All population-related measures of association provided in SUSAs are outlined based on their intentional design in order to highlight their different concepts. Even though, there is hardly one 'correct' choice for a measure of association, but some measure which is most reasonable in the context of the experts' belief. SUSAs offers a variety of measures of association each suitable to model another aspect of knowledge dependency as presented in the following.

2.2.1.1 Pearson's (Ordinary) correlation coefficient

Pearson's correlation coefficient (or Ordinary correlation coefficient) characterizes the bivariate Normal distribution and, therefore, is often employed to model dependency between two parameters X and Y . It is defined as:

$$\rho(X, Y) = \frac{Cov(X, Y)}{\sqrt{Var(X) \cdot Var(Y)}} = \frac{E[(X - E[X]) \cdot (Y - E[Y])]}{\sqrt{Var(X) \cdot Var(Y)}} \quad (2.59)$$

with the variances $Var(X)$ and $Var(Y)$, the covariance $Cov(X, Y)$ and the expectations $E[X]$ and $E[Y]$ of the parameters X and Y .

Next to the aforementioned general properties of correlation coefficients, Pearson's correlation coefficient offers a fundamental relation between two parameters:

- $\rho(X, Y) = 1(-1) \Leftrightarrow Y$ is a *linearly* increasing (decreasing) function of X

However, the practical interpretation of Pearson's correlation coefficient apart from $\rho(X, Y) = -1, 0$ or 1 is only feasible for a very restricted setting:

- Let Z, X', Y' be independent normally distributed parameters of equal variance related by $X = Z + X'$ and $Y = Z + Y'$. Given this setting, Pearson's correlation coefficient is given by $\rho(X, Y) = \frac{1}{2}$.

Despite the fact that this interpretation as an additive structure does only hold for normally distributed parameters, the example finds an import correspondence in real settings:

- Let X and Y be measurements acquired by a common measuring instrument. That instrument may be biased, such that each unbiased measurement X', Y' is biased by an additive random term Z independent for each measurement. This dependency may be described by Pearson's correlation coefficient $\rho(X, Y) = \frac{1}{2}$.

But apart from this pragmatic use, the non-intuitive interpretation of arbitrary values of the measure of association makes Pearson's correlation coefficient often not feasible. Moreover, a further disadvantage depicts its invariance properties:

- Let $X' = g(X)$ and $Y' = h(Y)$ with g and h being monotone increasing or decreasing functions, then generally Pearson's correlation coefficient is *not scale (or ordinal) invariant*, i.e. $\rho(X', Y') \neq \rho(X, Y)$.

This is a big disadvantage of Pearson's correlation coefficient, since a simple monotone transformation of the parameters potentially falsifies the modelled correlation patterns and may not comply with the experts' belief anymore. Hence, SUSA offers further non-parametric (distribution-free) measures of association that allow a simple interpretation

of arbitrary correlation values and are scale invariant in order to reliably model correlation patterns based on subjective judgements.

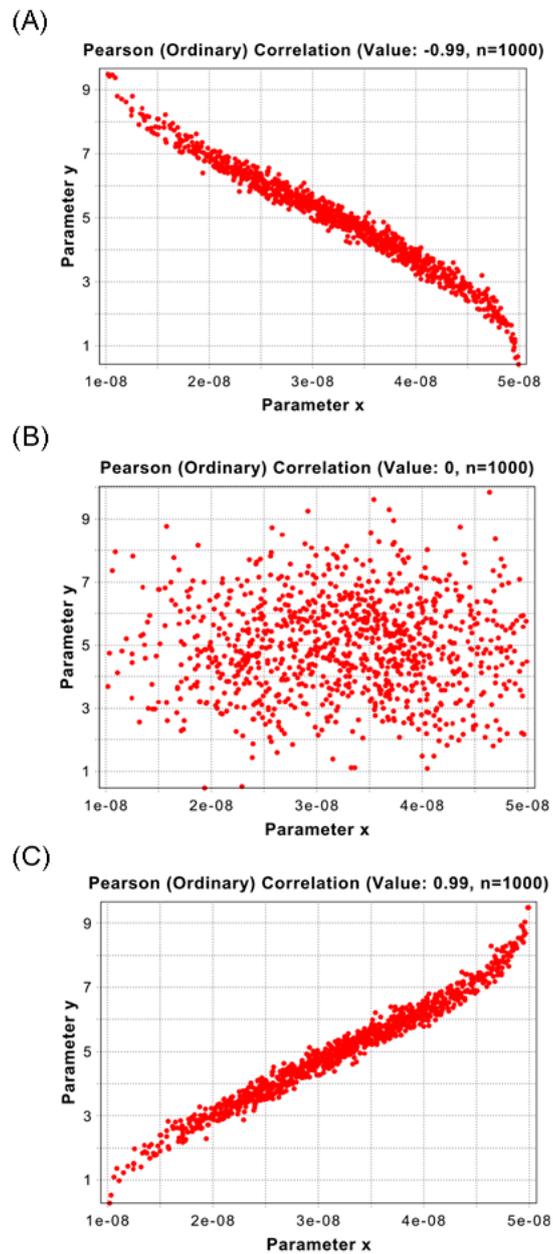


Fig. 2.18 The Pearson's (Ordinary) correlation coefficient for values of (A) strongly negative correlated $\rho=-0.99$, (B) uncorrelated $\rho=0$ and (C) strongly positive correlated $\rho=0.99$ parameters X and Y . These scatter plots exemplify the dependency assigned to the parameters X and Y

2.2.1.2 Blomqvist's medial correlation coefficient

Blomqvist's medial correlation coefficient (or Blomqvist's beta or population quadrant measure) is a practical approach to take into account a degree of association between parameters without structural information about the distribution of the corresponding parameters. This measure enables the analyst to practically take into account the experts' belief about the effect of an increasing parameter X on a parameter Y relative to the medians m_X and m_Y of the assigned distributions.

- A positive association in terms of decrease/increase in the parameter X goes with a decrease/increase in the parameter Y is also referred to as *concordance*.
- A negative association in terms of decrease/increase in the parameter X goes with an increase/decrease in the parameter Y is also referred to as *discordance*.

Given a fixed point in the support plane (X, Y) such as the median pair (m_X, m_Y) the concordance/discordance property of correlation provides an intuitive approach to make a subjective judgement about the association between two parameters as presented in the following Figure.

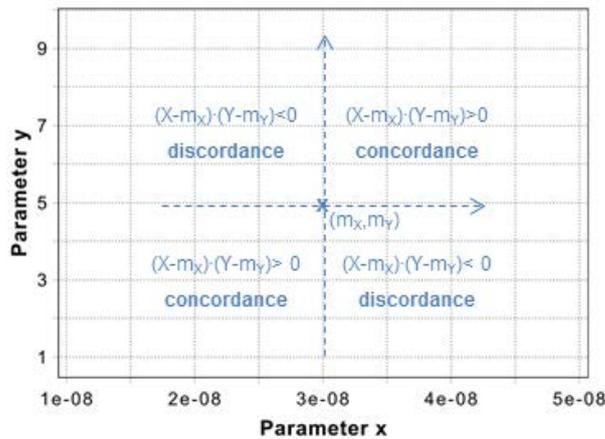


Fig. 2.19 The principle of concordance and discordance exemplified for the parameter pair (X, Y) . The reference point for the order is taken as the medians (m_X, m_Y) similar to the definition of Blomqvist's correlation coefficient

Based on the idea of concordance/discordance, Blomqvist's medial correlation coefficient $\beta(X, Y)$ offers an intuitive formulation of association between two parameters X and Y :

$$\beta(X, Y) = Prob((X - m_X) \cdot (Y - m_Y) > 0) - Prob((X - m_X) \cdot (Y - m_Y) < 0) \quad (2.60)$$

$\beta(X, Y)$ is the difference between the probabilities ($Prob(\cdot)$) of concordance and discordance of X and Y relative to the corresponding medians m_X and m_Y . Another formulation of this measure of association is given by the conditional probability of Y being smaller than the median m_Y given that X is smaller than the median m_X :

$$\beta(X, Y) = 2 \cdot Prob(Y < m_Y | X < m_X) - 1 \quad (2.61)$$

In contrast to Pearson's correlation coefficient Blomqvist's medial correlation coefficient offers an intuitive interpretation and, more importantly, a beneficial advantage due to its invariance properties:

- Let $X' = g(X)$ and $Y' = h(Y)$ with g and h being monotone increasing or decreasing functions, then Blomqvist's medial correlation coefficient is *scale (or ordinal) invariant* such that $\rho(X', Y') = \rho(X, Y)$ holds under monotone transformations.

Clearly, Blomqvist's medial correlation coefficient is a non-parametric (distribution-free) measure of association and therefore does not encode any structural information about the corresponding distributions of the parameters. Moreover, the property of concordance is considered relative to an *a priori* given reference point, i.e. the medians (m_X, m_Y), which may be regarded as a relatively arbitrary choice /KRU 58/.

2.2.1.3 Kendall's rank correlation coefficient

Kendall's rank correlation coefficient (or Kendall's tau) is a practical approach to take into account a degree of association between two parameters X and Y and may be considered as an extension of Blomqvist's medial correlation coefficient (Section 2.2.1.2). Both measures of association do have the same properties, only the choice of the reference point for the concordance deliberately differs. Instead of the fixed reference point in terms of the pair of medians (m_X, m_Y) for Blomqvist's measure, another bivariate parameter pair

(X_2, Y_2) is employed for Kendall's measure. Accordingly, Kendall's rank correlation coefficient is defined by using two independent pairs (X_1, Y_1) and (X_2, Y_2) derived from the *a priori* assumed bivariate distribution of (X, Y) :

$$\tau(X, Y) = Prob((X_1 - X_2) \cdot (Y_1 - Y_2) > 0) - Prob((X_1 - X_2) \cdot (Y_1 - Y_2) < 0) \quad (2.62)$$

that $\tau(X, Y)$ is the difference between the probabilities of concordance and discordance of X_1 and Y_1 relative to another independent pair X_2 and Y_2 . In this way, this measure of association is more flexible, i.e. less arbitrary and a more natural way of making use of a more detailed knowledge about the parameters (X, Y) . That is per design the Kendall rank correlation coefficient is influenced by the bivariate distribution structure F_{XY} /KRU 58/. Another formulation of this measure of association is given by the conditional probability of Y_2 being smaller than Y_1 given that X_2 is smaller X_1 , that is

$$\tau(X, Y) = 2 \cdot Prob(Y_2 < Y_1 | X_2 < X_1) - 1 \quad (2.63)$$

Note, the presented formulation of Kendall's rank correlation coefficient is the population analogue to the sample based Kendall's rank correlation coefficient taking into account the concordant and discordant bivariate observations $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ as given in the following equation.

$$\tau = \frac{n_c - n_d}{\frac{1}{2} \cdot n \cdot (n - 1)} \quad (2.64)$$

with the number of concordant pairs n_c , the number of discordant pairs n_d and the total number of pairs n . Since the denominator represents the total number of possible combinations of comparisons of the pairs (x_i, y_i) and (x_j, y_j) with $i = 1, \dots, n-1$ and $j = i+1, \dots, n$, the correspondence between the empirical to the population based measure of association is obvious.

2.2.1.4 Spearman rank correlation coefficient

Spearman's rank correlation coefficient (or Spearman's rho) is a practical approach to take into account a degree of association between two random parameters (X, Y) and may be considered as an extension of Kendall's rank correlation coefficient (Section 2.2.1.3). Both measures of association do have the same properties, only the choice of the reference point for the concordance deliberately differs. In contrast to Kendall's measure, instead of the two bivariate independent observations (X_1, Y_1) and (X_2, Y_2) obtained from an *a priori* assumed distribution of (X, Y) a third observation (X_3, Y_3) is employed for Spearman's measure. Accordingly, Spearman's rank correlation coefficient ρ_S is derived as indicated in the following equation.

$$\rho_S(X, Y) = Prob((X_2 - X_1) \cdot (Y_3 - Y_1) > 0) - Prob((X_2 - X_1) \cdot (Y_3 - Y_1) < 0) \quad (2.65)$$

ρ_S is the difference between the probabilities of concordance and discordance of X_2 and Y_3 , each from an independent pair (X_2, Y_2) and (X_3, Y_3) , relative to a third independent pair X_1 and Y_1 . In contrast to Kendall's rank correlation coefficient, Spearman's rank correlation coefficient ρ_S is per design influenced by the structures of the marginal distributions F_X and F_Y /KRU 58/. Another formulation of this measure of association is given by the conditional probability of Y_3 being smaller than Y_1 given that X_2 is smaller than X_1 .

$$\rho_S(X, Y) = 6 \cdot Prob(Y_3 < Y_1 | X_2 < X_1) - 3 \quad (2.66)$$

Note, the presented formulation of Spearman's rank correlation coefficient ρ_S is the population analogue to (the sample based) Spearman's sample rank correlation coefficient which is separately discussed in Section 2.2.2.1.

2.2.2 Sample-related correlation

As outlined in the previous Section, population-related measures of association are used in combination with a set of univariate distributions to design a population. The sample of parameter values which will be generated in this context can be considered as selected from an *a priori* defined multivariate distribution satisfying the univariate marginal distributions and the measures of association specified for the parameters. In contrast to that, sample-related measures of association are used to design a sample without

consideration of the corresponding population properties. They are attained by a specific transformation *a posteriori* applied to a sample selected from a multivariate distribution satisfying only the univariate marginal distributions of the parameters and not the measures of association. The practical consequence for the analyst is therefore:

- population-related measures of association: the input uncertainty specifications comply with *population properties*
- sample-related measures of association: the input uncertainty specifications comply with *empirical properties*

2.2.2.1 Spearman's sample rank correlation coefficient

The sample-related measure of association provided in SUSA is the empirical analogue to the Spearman's rank correlation coefficient. The so-called Spearman's sample rank correlation coefficient is formulated based on the ranking of two sets of parameter values (x_1, x_2, \dots, x_n) and (y_1, y_2, \dots, y_n) where each pair (x_i, y_i) is sampled from the according bivariate distribution of the parameter pair (X, Y) . The ranking assigns the ranks $r(x_i)$ and $r(y_i)$ to the values x_i and y_i , respectively. The ranks correspond to the rank order of the values of X and Y . A simple example shall serve to clarify the ranking process:

$$\begin{pmatrix} x_1 & y_1 \\ x_2 & y_2 \\ x_3 & y_3 \end{pmatrix} = \begin{pmatrix} 2.5 & 10 \\ 1.8 & 8 \\ 0.5 & 12 \end{pmatrix} \xrightarrow{\text{ranking}} \begin{pmatrix} r(x_1) & r(y_1) \\ r(x_2) & r(y_2) \\ r(x_3) & r(y_3) \end{pmatrix} = \begin{pmatrix} 3 & 2 \\ 2 & 1 \\ 1 & 3 \end{pmatrix} \quad (2.67)$$

The average of ranks is used as the rank of respective parameter values, if the values are tied (even on magnitude). Tied values may be obtained for discrete distributions.

Spearman's sample rank correlation coefficient corresponds to the empirical analogue of Pearson's correlation coefficient (Section 2.2.1.1) applied not to the values of (X, Y) themselves but to the assigned ranks $(r(X), r(Y))$.

$$\begin{aligned}
\hat{\rho}_S(X, Y) = \hat{\rho}(r(X), r(Y)) &= \frac{\hat{Cov}(r(X), r(Y))}{\sqrt{\hat{Var}(r(X))\hat{Var}(r(Y))}} \\
&= \frac{\sum_i (r(x_i) - \overline{r(X)})(r(y_i) - \overline{r(Y)})}{\sqrt{\sum_i (r(x_i) - \overline{r(X)})^2 \sum_i (r(y_i) - \overline{r(Y)})^2}} \\
&= \frac{\sum_i (r(x_i) - \frac{n+1}{2})(r(y_i) - \frac{n+1}{2})}{\frac{n(n^2 - 1)}{12}} \tag{2.68}
\end{aligned}$$

where $\hat{\cdot}$ = estimator derived from a sample or empirical analogue, $Var(\cdot)$ = variance, $Cov(\cdot)$ = covariance, n = sample size, and $\overline{r(\cdot)}$ = average of n ranks.

Similar to Pearson's correlation coefficient, there is no intuitive interpretation for Spearman's sample rank correlation coefficient as measure of association between two uncertain parameters X and Y . The analyst can only provide the following properties for the sample to be generated:

- The ranks and, consequently, the empirical distribution functions of the parameters X and Y shall be correlated according to Pearson's correlation coefficient $\hat{\rho}$ with $-1 \leq \hat{\rho}(r(X), r(Y)) \leq 1$.
- $\hat{\rho}(r(X), r(Y)) = 0$: The ranks (empirical distribution functions) of the parameters X and Y shall be independent.
- $\hat{\rho}(r(X), r(Y)) = 1$: The agreement between the ranks (empirical distribution functions) shall be perfect, i.e. the ranks of the parameters X and Y shall be identical.
- $\hat{\rho}(r(X), r(Y)) = -1$: The disagreement between the ranks (empirical distribution functions) shall be perfect, i.e. the ranks shall be reverse to each other.

The application of Spearman's sample rank correlation coefficient does not allow for deriving a sample in compliance with a joint multivariate distribution structure.

2.2.3 Full Dependence

The full (complete) positive (negative) dependency between the two parameters X and Y is interpreted in the sense that the uncertainty in Y completely derives from the

uncertainty in X , although an explicit functional relationship between X and Y is not known. Given the corresponding univariate distributions F_X and F_Y as representations of the uncertainties in X and Y , respectively, intuitive formulations of the strict monotone increasing (decreasing) functional relationship between X and Y are the following positive Eq. (2.69) and negative Eq. (2.70) dependencies.

$$Y = F_Y^{-1}(F_X(X)) \quad (2.69)$$

$$Y = F_Y^{-1}(1 - F_X(X)) \quad (2.70)$$

Completely dependent parameters X and Y fulfilling the relationships in Eq. (2.69) or Eq. (2.70) exhibit the following properties:

- X and Y are completely positively (negatively) dependent \Rightarrow all scale (or ordinal) invariant measures of association equal to 1 (-1) (the reverse does not hold generally)
- Pearson's correlation coefficient equals to 1 (-1) \Rightarrow X and Y are completely positively (negatively) dependent (the reverse does not hold generally)
- X and Y are completely positively (negatively) dependent \Rightarrow the sample rank correlation coefficient of any bivariate sample (X, Y) equals to 1 (-1)

The following comparative properties between the values of the completely positive (negative) dependent parameters X and Y may be of interest to the analyst:

- values of X and Y in (X, Y) are of the same (reverse) order
- values of X and Y have a strictly monotone increasing (decreasing) relationship
- X and Y do have the same (complement) quantiles

These rather descriptive properties are a direct consequence of the employed definition of complete dependence. The mathematical concept of full (complete) dependency as used in SUSA and consequences thereof are outlined in /KRZ 88/.

2.2.4 Conditional Distribution

In case the experts' knowledge about an uncertain parameter Y is different for ranges of the support $S_X = [\min(X), \max(X)]$ of uncertain parameter X , SUSAS offers the option to specify conditional probability distributions for parameter Y on condition on the different ranges of the support of parameter X .

As an example let X represent the (sub)model to be applied to simulate a specific physical phenomenon and let Y represent the uncertain correction factor which – multiplied with the corresponding model prediction – provides the true value. If two model alternatives M_a (indicated by model index 1) and M_b (indicated by model index 2) are available and it is uncertain, which of the models is the best to simulate the phenomenon, parameter X is an uncertain parameter for which two values (1 and 2) may be true. If the uncertainty on the correction factor Y is different for the two models, different distributions for Y may be specified on condition of the values of parameter X . For instance, the uncertainty on Y may be represented by a Triangular distribution with support $[0.8; 1.3]$ and mode = 1.0, if model M_a is the best model. A Uniform distribution with support $[0.85; 1.15]$ may be used to represent the uncertainty on Y , if model M_b is the best.

SUSAS requires that the marginal distribution of parameter X is completely specified as indicated in Section 2.1. Instead of the marginal distribution of Y which cannot be specified, the conditional distributions $F_{Y|X \in I_k}$ of Y on condition of $X \in I_k$, $k=1, \dots, K$ must be specified, where I_1, I_2, \dots, I_K represent the partition of the support S_X of X into K disjoint intervals.

2.2.5 Function of Parameters

An uncertain parameter Y may be associated to other uncertain parameters X_i, X_j, \dots by an explicit functional relationship. SUSAS offers the option to formulate such a relationship as a Fortran formula as exemplarily shown in Eq. (2.71).

$$Y = \frac{a \cdot X_i + b \cdot X_j}{\sqrt{X_k}} \quad (2.71)$$

The values of uncertain parameter Y are derived from the explicit functional dependency on the parameters X_i, X_j, \dots . They are affected by the *a priori* specified marginal distributions F_{X_i}, F_{X_j}, \dots

SUSA includes an internal compiler for formulas represented in Fortran language. This compiler can appropriately interpret all known arithmetic operators (+, -, *, /), the mathematical bracket '()' and the functions *int, abs, min, max, sqrt, log, log10, exp, sin, cos, and tan*.

2.2.6 Inequality

In some settings two uncertain parameters X and Y are associated to each other such that X or $a \cdot X$, with a being a real-valued factor, determines a lower threshold for Y . This relationship may be described by the following inequality

$$Y \geq a \cdot X \tag{2.72}$$

This inequality results in the restriction of the derived sample values (y_1, \dots, y_n) . The inequality modifies the joint support plane $S_Y \times S_X$ as illustrated in the following figure.

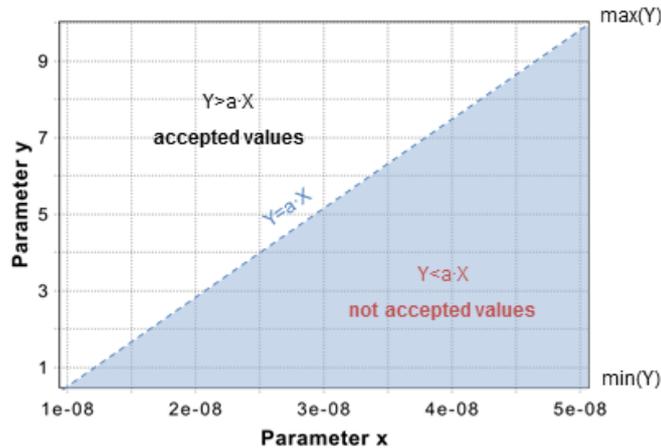


Fig. 2.20 The principle of dependency between uncertain parameters X and Y by the assumption of an inequality $Y \geq a \cdot X$. The blue shaded area marks the support plane in which the subjective association between X and Y implies that no sample pair (x_i, y_i) is accepted

In order to derive a sample element (x_i, y_i) in compliance with the above inequality, SUSAS offers the following two options:

- **independent repeated sampling** until a complete sample is generated that only contains values satisfying the inequality.

In general, the marginal distributions F_X and F_Y of the parameters X and Y are affected by this kind of modification. The computational time of this brute force approach may be high.

- **value modifications** according to the following equation

$$y_i' = a \cdot x_i + \frac{y_i - \min(Y)}{\max(Y) - \min(Y)} \cdot (\max(Y) - a \cdot x_i) \quad (2.73)$$

The modification of Y according to Eq. (2.73) is applied only to sample elements (x_i, y_i) which do not fulfil the inequality. To ensure that the modified values y_i' are larger than $a \cdot x_i$ within the limits of the support $S_Y = [\min(Y), \max(Y)]$, it is required that the relationship in Eq. (2.74) is fulfilled.

$$\max(Y) \geq a \cdot \max(X) \quad (2.74)$$

In general, the marginal distribution F_Y of Y is affected by this kind of modification.

Specification of the inequality Eq. (2.73) between the parameters X and Y is recommended only, if the following relationships hold:

- $\min(Y) < a \cdot \max(X)$
- $\max(Y) \geq a \cdot \max(X)$

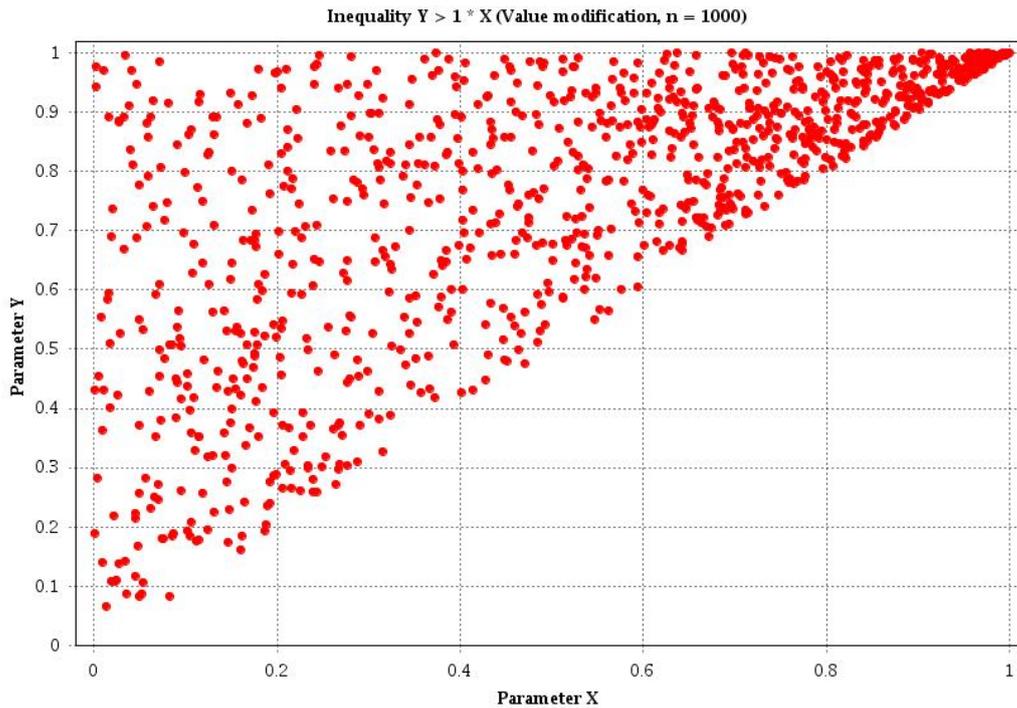


Fig. 2.21 Relationship between parameters X and Y defined by the inequality $Y \geq a \cdot X$ with $a = 1$ ($n = 1000$ data points (x, y))

2.3 Proportions

Another aspect of uncertainty specification relates to the association of multiple uncertain parameters which represent the proportions (percentages) of a whole and, therefore, must sum up to 1.0 (100 %). An example of a joint group of uncertain parameters representing proportions of a whole depicts the probabilities of the branches at a branching point of an event tree. Another example are the proportions of the different age groups in a population:

Age group (years)	Percentage
0 – 14	13 %
15 – 24	10 %
25 – 54	41 %
55 – 64	14 %
65 and over	22 %
Total	100 %

All uncertain proportions of a whole required to sum up to 1 have to be grouped in a so-called proportion group. The distribution to be quantified for an uncertain proportion in a proportion group must not be the distribution of the proportion itself but the distribution of a conditional proportion. This conditional proportion represents the corresponding proportion relative to the remaining total after the contribution of all proportions specified prior to the current one were excluded (e. g. the proportion of age group 25 – 54 relative to the total of the population without all lower age groups). During the sample generation (Section 3), the conditional proportions are finally transformed into the actual proportions.

3 Sample Generation

One major step of a probabilistic uncertainty and sensitivity analysis (Section 1) is the generation of a multivariate sample of values (experimental design) for the uncertain input parameters influencing the prediction of the applied computer code. The following two sampling procedures are implemented in SUSAs:

- Simple random sampling
- Latin Hypercube sampling

Both sampling procedures use a pseudorandom number generator providing values from a Uniform distribution. The pseudorandom number generators implemented in SUSAs are delineated in Section 3.1. Sections 3.2 and 3.3 give a description of the simple random and the Latin Hypercube sampling procedure, respectively. The sample generation algorithms accounting for specific dependences between uncertain parameters are described in Sections 3.4 - 3.7. Section 3.8 shortly describes how the sample of computational results is generated.

3.1 Pseudorandom number generators

The random number generators in SUSAs are pseudorandom number generators. The sequences of random numbers provided by these generators seem to be random although they are calculated by a deterministic algorithm. They are completely determined by the initial value (initial seed, e.g. 123457) specified as input. That means they always produce the same sequence of random numbers when initialized with that value. Pseudorandom number generators are often applied because of their speed in number generation and their reproducibility.

Two pseudorandom number generators in SUSAs are multiplicative congruential generators /HUL 62/ and recursively produce the sequence $\{x_i\}$ as indicated in the following equation.

$$x_{i+1} = (a \cdot x_i) \bmod m \quad i = 0, 1, 2, \dots \quad (3.1)$$

where *mod* means modulo, $m \in \{2, 3, 4, \dots\}$, $a \in \{1, 2, 3, \dots, m-1\}$ and $x_0 \in \{0, 1, 2, \dots, m-1\}$.

The period of these generators is at most m . The sequence $\{x_i/m\}$ is taken as the uniform random number sequence.

The first multiplicative congruential generator in SUSAN is characterized by

- multiplier $a = 16807$
- modulo $m = 2^{31}-1$

The second multiplicative congruential generator is characterized by

- multiplier $a = 48271$
- modulo $m = 2^{31}-1$

Both generators are frequently used in run-time libraries of various compilers, because they are fast and require minimal memory. Due to the serial correlation between successive values of the random number sequence, they should not be used for applications where high-quality randomness is required.

A high quality of randomness is provided by the Mersenne Twister - the other pseudorandom number generator in SUSAN. This generator was developed by Makoto Matsumoto and Takuji Nishimura /MAT 98/. The Mersenne Twister Fortran algorithm implemented in SUSAN is based on the algorithm MT19937 with improvements from 2002 considering Shawn Cokus' optimization, Matthew Bellew's simplification and Isaku Wada's real version. The following link gives more information:

<http://www.math.sci.hiroshima-u.ac.jp/~m-mat/MT/MT2002/emt19937ar.html>

The Mersenne Twister produces a sequence of 32-bit integers with the following properties:

- It has a very long period P of $2^{19937} - 1$.
- It passes numerous tests for statistical randomness.
- It is k -distributed to 32-bit accuracy for every $1 \leq k \leq 623$, i.e. each of the 2^{kv} possible combinations of bits occurs the same number of times in the period P , except for the all-zero combination that occurs once less often. v represents the number of most significant leading bits of the produced integers. From this property, it follows

that there is a very low correlation between successive values of the random number sequence.

3.2 Simple Random Sampling (SRS)

The SRS procedure selects each set of parameter values at random from a multivariate distribution defined by the marginal probability distributions, association measures and other dependency models specified as input to quantify the knowledge on the parameters and the corresponding (knowledge) dependencies (Section 2).

Section 3.2.1 describes the procedure for SRS, if population-related correlations between parameters have to be considered. Section 3.2.2 explains the procedure, if sample-related correlations between parameters have to be taken into account.

3.2.1 SRS with consideration of population-related correlations

A simple random sample of size n for $npar$ uncertain parameters $(X_1, X_2, \dots, X_{npar})$ is obtained by sampling the $npar$ values of the parameter vector $(X_1, X_2, \dots, X_{npar})$ independently n times according to the specified marginal distributions F_i . The sampling is performed by using a pseudorandom number generator providing values for uniformly distributed variables U_i , $i = 1, \dots, npar$, and by applying the inverse distribution function to these values Eq. (3.2).

$$X_i = F_i^{-1}(U_i) \quad (3.2)$$

If population-related correlation coefficients (Pearson, Blomqvist, Kendall or Spearman) between parameters have to be considered, the transformation of an appropriately chosen multivariate Normal distribution is carried out. Instead of the specified marginal distribution F_i of parameter X_i , the standard Normal distribution Φ of parameter Z_i is considered initially.

$$Z_i = \Phi^{-1}(U_i) \quad (3.3)$$

The originally specified correlation coefficient θ_{ij} between the parameters X_i and X_j is transformed to Pearson's ordinary correlation coefficient ρ_{ij} between the standard normally distributed parameters Z_i and Z_j . If Z'_i and Z'_j are two independent standard normally distributed parameters, Pearson's ordinary correlation coefficient ρ_{ij} between the parameters Z_i and Z_j is obtained, if the following relationships are fulfilled Eq. (3.4):

$$\begin{aligned} Z_i &= Z'_i \\ Z_j &= \rho_{ij} \cdot Z'_i + \sqrt{1 - \rho_{ij}^2} \cdot Z'_j \end{aligned} \tag{3.4}$$

For any correlation coefficient θ_{ij} between the parameters X_i and X_j , Pearson's ordinary correlation coefficient ρ_{ij} is mostly obtained by an appropriate iteration procedure applying the bisection method in combination with Monte Carlo simulation. In a few cases, an analytical formula is used to calculate ρ_{ij} /JOH 78/.

If the ordinary correlation coefficients ρ_{ij} between concerned standard normally distributed parameters Z_i and Z_j are determined, the multivariate Normal distribution is defined and the parameter X_i with $X_i \sim F_i$ can be obtained by Eq. (3.5).

$$X_i = F_i^{-1}(\Phi(Z_i)) \tag{3.5}$$

The steps of the iteration procedure to determine Pearson's ordinary correlation coefficient ρ_{ij} between the parameters Z_i and Z_j can be summarized as follows:

1. In stage k of the iteration process, $k = 0, 1, 2, \dots, K$, the value r_k of Pearson's ordinary correlation coefficient ρ_{ij} is specified to define the multivariate Normal distribution of the parameter vector (Z_i, Z_j) .

Stage 0:

If the correlation type of θ_{ij} is Pearson's ordinary correlation, r_0 is calculated as:

$$r_0 = \theta_{ij} \tag{3.6}$$

If the correlation type of θ_{ij} is Blomqvist's correlation or Kendall's correlation, r_0 is calculated as:

$$r_0 = \sin\left(\theta_{ij} \cdot \frac{\pi}{2}\right) \quad (3.7)$$

If the correlation type of θ_{ij} is Spearman's rank correlation, r_0 is calculated as:

$$r_0 = 2 \sin\left(\theta_{ij} \cdot \frac{\pi}{6}\right) \quad (3.8)$$

For stage $k > 0$, r_k is calculated according to the bisection method as:

$$r_k = \frac{a_k + b_k}{2} \quad (3.9)$$

a_k and b_k are determined as follows:

$k = 1$:

$a_k = r_{k-1}$, $b_k = 1$, if $r_{k-1} > 0$

$a_k = r_{k-1}$, $b_k = -1$, if $r_{k-1} < 0$

$k > 1$:

$a_k = r_{k-1}$, $b_k = b_{k-1}$, if $|\hat{\theta}_{ij}| < |\theta_{ij}|$

$a_k = a_{k-1}$, $b_k = r_{k-1}$, if $|\hat{\theta}_{ij}| \geq |\theta_{ij}|$

2. n realizations (z_i, z_j) are sampled from the multivariate Normal distribution of (Z_i, Z_j)
3. Each sample element (z_i, z_j) is transformed to (x_i, x_j) according to Eq. (3.5) and the sample correlation coefficient $\hat{\theta}_{ij}$ is calculated.
4. The following inequality is checked:

$$|\hat{\theta}_{ij} - \theta_{ij}| < \varepsilon \quad (3.10)$$

If the inequality is fulfilled, the iteration procedure is finished.

If the inequality is not fulfilled, the iteration continues with step 1.

The current settings for the aforementioned iteration procedure are as follows:

- Number K of iteration stages: 50
- Number n of sample elements: 2000
- Epsilon-value ε of the stopping criterion: 0.01

3.2.2 SRS with consideration of sample-related correlations

First, a simple random sample of size n for $npar$ uncertain parameters $(X_1, X_2, \dots, X_{npar})$ is generated by sampling the $npar$ values of the parameter vector $(X_1, X_2, \dots, X_{npar})$ independently n times according to the specified marginal distributions F_i and regardless of the specified rank correlations. Then, the n values obtained for each parameter X_i are permuted appropriately so that the (Spearman) rank correlation coefficients $\hat{\rho}_S$ between the parameters are very close to those specified as input.

The steps to obtain a simple random sample with sample-related correlations can be summarized as follows (see also /IMA 82/):

1. Generation of the matrix $X = (x_{ji})$ of parameter values with $i=1, \dots, npar, j=1, \dots, n$ and x_{ji} = element of the j^{th} row and i^{th} column of X . The rows of X are selected from a multivariate distribution which is only defined by the marginal distributions of the parameters without consideration of any correlation.
2. Calculation of the matrix of ranks $R^X = (r(x_{ji}))$, $i=1, \dots, npar, j=1, \dots, n$. The ranks are calculated separately for the values $x_{1i}, x_{2i}, \dots, x_{ni}$ of each parameter X_i in column i of matrix X .

Rank $r(x_{ji})$: If the j^{th} sampled value x_{ji} is the smallest value of X_i and there is no other sample element with this value, the corresponding rank is $r(x_{ji})= 1$. If the j^{th} sampled value x_{ji} is the highest value of X_i and there is no other sample element with this value, the corresponding rank is $r(x_{ji})= n$. The average of the respective ranks is used for tied values (even on magnitude).

3. Calculation of the empirical ordinary correlation matrix $C^S = (c_{ij}^S)$ of matrix R^X , i.e. calculation of Pearson's ordinary correlation coefficients between column i and column j of matrix R^X . This step is equivalent to the calculation of Spearman's sample rank correlation coefficients between the values sampled for parameters X_i and X_j , $i=1, \dots, npar, j=1, \dots, npar$

$$c_{ij}^S = \frac{\sum_k (r(x_{ki}) - \bar{r}_{x_i})(r(x_{kj}) - \bar{r}_{x_j})}{\sqrt{\sum_k (r(x_{ki}) - \bar{r}_{x_i})^2 \sum_k (r(x_{kj}) - \bar{r}_{x_j})^2}} \quad (3.11)$$

4. Calculation of the lower triangular matrix T^S with $C^S = T^S \cdot T^{S'}$ (= Cholesky decomposition of matrix C^S)
5. Calculation of the lower triangular matrix T of the matrix C of specified rank correlations with $C = T \cdot T'$ (= Cholesky decomposition of matrix C)
6. Calculation of the matrix $R' = R^X \cdot (T \cdot T^{S^{-1}})$. The ordinary correlation matrix of R' corresponds to the specified rank correlation matrix.
7. Calculation of the matrix $R^{R'}$ of ranks from matrix R' .
8. Permutation of the parameter values $x_{ji}, j=1, \dots, n$, in column i of matrix X according to column i of matrix $R^{R'}$, $i=1, \dots, npar$.

Remarks to the aforementioned procedure:

- The Cholesky decomposition of matrix C^S requires, that $n > npar$, i.e. the sample size n must be higher than the number $npar$ of uncertain parameters.
- The marginal distributions of the parameters are considered, however, a joint multivariate distribution of the uncertain parameters is not defined. Therefore, the sample elements finally provided cannot be considered as independently selected from such a multivariate distribution.
- To increase the sample size, a new matrix X with a higher number n of rows must be generated. Results of Monte-Carlo simulation runs based on the original matrix X with a smaller number n of rows cannot be used in general.

3.3 Latin Hypercube Sampling (LHS)

Latin Hypercube sampling may be considered as a special case of stratified sampling. First, the range of each parameter is divided into n distinct subintervals of equal probability $1/n$ which cover the parameter range. Then, one value is selected from each subinterval. Two options are available to select the value \tilde{x}_{ji} of parameter X_i :

- \tilde{x}_{ji} = median of the conditional distribution of X_i on condition of $X_i \in I_j$
- \tilde{x}_{ji} = value randomly selected from the conditional distribution of X_i on condition of $X_i \in I_j$

The n values finally obtained for each parameter are permuted randomly and combined to n parameter vectors of length $npar$. The value from each subinterval of each parameter is considered once and only once in the generated sample. If association measures have to be considered, the combinations of the n values of each parameter are modified appropriately. Section 3.3.1 describes the procedure for LHS, if population-related correlations between parameters have to be considered. Section 3.3.2 explains the procedure, if sample-related correlations between parameters have to be taken into account.

McKay et al. /MCK 79/ showed that LHS is better than SRS for estimating the mean and the population distribution function of the computational result Y , if Y is a monotonic function of the uncertain parameters.

3.3.1 LHS with consideration of population-related correlations

Let $\tilde{X} = (\tilde{x}_{ji})$, $i=1, \dots, npar$, $j=1, \dots, n$, be the matrix of parameter values selected from the subintervals of the range of each parameter. For each parameter X_i , the relationship between the selected parameter values is $\tilde{x}_{1i} < \tilde{x}_{2i} < \dots < \tilde{x}_{ni}$.

All parameter values \tilde{x}_{ji} , $j=1, \dots, n$, of parameter X_i are permuted according to the following procedure:

1. Generation of the matrix $X = (x_{ji})$ of parameter values according to the simple random sampling procedure described in Section 3.2.1.

2. Calculation of the matrix of ranks $R^X = (r(x_{ji}))$, $i=1, \dots, npar$, $j=1, \dots, n$. The ranks are calculated from the values $x_{1i}, x_{2i}, \dots, x_{ni}$ of parameter X_i in column i of matrix X , $i=1, \dots, npar$.
3. Permutation of the values $\tilde{x}_{1i} < \tilde{x}_{2i} < \dots < \tilde{x}_{ni}$ of parameter X_i in column i of matrix \tilde{X} , according to the ranks $r(x_{1i}), r(x_{2i}), \dots, r(x_{ni})$ in column i of matrix R^X , $i=1, \dots, npar$.

Since the LHS procedure uses the SRS procedure to generate corresponding values, even if instead of a random value, a conditional median is selected from each subinterval, a pseudorandom number generator is needed when the LHS procedure is applied.

3.3.2 LHS with consideration of sample-related correlations

Let $\tilde{X} = (\tilde{x}_{ji})$, $i=1, \dots, npar$, $j=1, \dots, n$, be the matrix of parameter values selected from the subintervals of the range of each parameter. For each parameter X_i , the relationship between the selected parameter values is $\tilde{x}_{1i} < \tilde{x}_{2i} < \dots < \tilde{x}_{ni}$.

All parameter values \tilde{x}_{ji} , $j=1, \dots, n$, of parameter X_i are permuted according to the following procedure (see also /IMA 82/):

1. Generation of a random rank matrix $R = (r_{ji})$, $i=1, \dots, npar$, $j=1, \dots, n$. Each column i of R is a random permutation of $\{1, \dots, n\}$. The permutations in the columns are independent and have the same probability $prob = \frac{1}{n!}$.
2. Calculation of the correlation matrix $C^R = (c_{ij}^R)$ of R , $i=1, \dots, npar$, $j=1, \dots, npar$, comprising the ordinary correlations between the columns of R .
3. Calculation of the lower triangular matrix T^R of C^R with $C^R = T^R \cdot T^{R'}$ (= Cholesky decomposition of matrix C^R).
4. Calculation of the lower triangular matrix T of matrix C of specified rank correlations with $C = T \cdot T'$ (= Cholesky decomposition of matrix C).
5. Calculation of the matrix $R' = R \cdot (T \cdot T^{R'})^{-1}$. The correlation matrix of matrix R' comprising the ordinary correlations between the columns of R' corresponds to specified rank correlation matrix.
6. Calculation of the matrix $R^{R'}$ of ranks from matrix R' .

7. Permutation of the values $\tilde{x}_{1i} < \tilde{x}_{2i} < \dots < \tilde{x}_{ni}$ of parameter X_i in column i of matrix \tilde{X} , according to the ranks $r_{1i}, r_{2i}, \dots, r_{ni}$ in column i of the rank matrix $R^{R'}$, $i=1, \dots, npar$.

Remarks to the aforementioned procedure:

- The Cholesky decomposition of C^S and C requires, that each matrix is positive definite.
- The sample size n must be higher than the number $npar$ of uncertain parameters
- A joint multivariate distribution of the uncertain parameters is not defined. Therefore, the sample elements finally provided cannot be considered as independently selected from such a multivariate distribution.
- To increase the sample size, a new matrix X with a higher number n of rows must be generated. Results of Monte-Carlo simulation runs based on the original matrix X with a smaller number n of rows cannot be used in general.

3.4 Sample generation with consideration of complete dependencies

Complete dependency between two parameters X_i and X_j is realized by simply taking the same uniformly distributed random number U_i (or $(1.0 - U_i)$ for negative dependency) from the pseudorandom number generator when generating the values for X_i and X_j according to the SRS procedure (Eq. (3.2)).

3.5 Sample generation with consideration of conditional distributions

The conditional distributions of parameter X_2 on condition of the values of parameter X_1 (Section 2.2.4) are considered as follows:

1. The sample $X=(x_{ji})$ with $i=1, \dots, npar$ and $j=1, \dots, n$ is generated based on the marginal distributions and the correlations specified for the uncertain parameters $X_1, X_2, \dots, X_{npar}$. Since the marginal distribution of parameter X_2 is not specified, its values are initially set to zero.
2. The sample $V=(v_{jk})$ with $k=1, \dots, K$ and $j=1, \dots, n$ is generated for K variables (V_1, V_2, \dots, V_k) distributed according to the conditional distributions $F_{X_j|X_i \in I_k}$ of parameter X_j on condition of $X_i \in I_k$, $k=1, \dots, K$.

3. The values of X_j are changed appropriately according to the values sampled for X_i :
If $x_{li} \in I_k$, then $x_{lj} = v_{lk}$, $l=1, \dots, n$, $k= 1, \dots, K$.

If the LHS procedure is applied, the Latin Hypercube structure cannot generally be reached for the values of parameter X_j .

3.6 Sample generation with consideration of functional relationships

The functional relationship between uncertain parameter X_i and the uncertain parameters X_j, X_k, \dots (Section 2.2.5) is considered by applying the function indicated for parameter X_i to the values generated for parameters X_j, X_k, \dots . If the LHS procedure is applied, the Latin Hypercube structure cannot be reached for the values of parameter X_i .

3.7 Sample generation with consideration of inequalities

In order to consider a specific inequality relationship between two parameters (Section 2.2.6), SUSA offers the following two options:

- **independent repeated sampling** until a complete sample satisfying the specified inequalities is generated
- **value modifications** according to equation Eq. (2.73)

With the first option, the sampling step according to the SRS or LHS procedure may be repeated very often until the inequalities are fulfilled. The specified marginal distributions of the involved parameters may be changed. The latter drawback is also associated with the second option. If the LHS procedure is applied, another drawback of the second option is the generally unfulfillable Latin Hypercube structure of the parameter values derived from the modification function in Eq. (2.73).

3.8 Computer code runs

Each set of values sampled for the total of uncertain input parameters is supplied as input to a computer code run. When all runs are finished, a sample of values from the unknown probability distribution of the computational result is available. This sample can be analyzed by statistical methods in order to obtain indicators (indices) of the uncertainty and sensitivity of the computational result.

4 Uncertainty Analysis

The uncertainty of the computational result derives from the propagation of the uncertainties of the input parameters through the computer model. Computer code runs (Monte Carlo simulation runs) each performed with a set of values sampled for the uncertain input parameters provide a sample from the unknown probability distribution of the computational result. The statistical analysis of this sample provides estimators of the distribution and of its properties.

Estimators which may be used to quantify the uncertainty of the computational result are the empirical cumulative distribution function, the empirical mean, the unbiased empirical standard deviation or variance and empirical quantiles (Section 4.1). Very useful estimators especially for complex applications which don't allow performing many Monte Carlo simulation runs are tolerance limits (Section 4.2). Interval limits estimated from the inequalities of Chebychev and Cantelli may be adequate as well (Section 4.3). Another alternative of uncertainty quantification is the indication of a parametrical probability distribution well fitted to the empirical distribution of the computational result (Section 4.4). A further approach to derive uncertainty quantifications is the construction of a surrogate model and the application of the surrogate instead of the original complex model in order to be able to perform a huge number of Monte Carlo simulation runs (Section 4.5). From the large sample of values finally provided for the computational result, estimators of high accuracy can be derived to quantify the resulting uncertainty.

Estimators available in SUSA to quantify the uncertainty for multiple computational results are described in Section 4.6.

SUSA can perform the uncertainty analysis for a scalar as well as for a time/index-dependent result. A result is denoted as scalar, if it has one single value per run. A result is time/index-dependent, if it has values at different points in time, space, etc. for each run (Fig. 4.1). A time/index-dependent result at a specific point in time, space, etc. may be considered as a scalar result. So, a time/index-dependent uncertainty analysis is equivalent to a scalar uncertainty analysis at each time/index step.

Since a time/index-dependent analysis is computationally intensive, the number of options available for a time/index-dependent uncertainty analysis is smaller than that available for a scalar uncertainty analysis. The fitting of parametrical distributions (Section 4.4) and the construction and application of a surrogate model (Section 4.5) can only

be performed in the scalar analysis. Also the calculation of tolerance limits based on the assumption of a Normal or Lognormal distribution for the computational result (Section 4.2.2) is restricted to the scalar analysis.

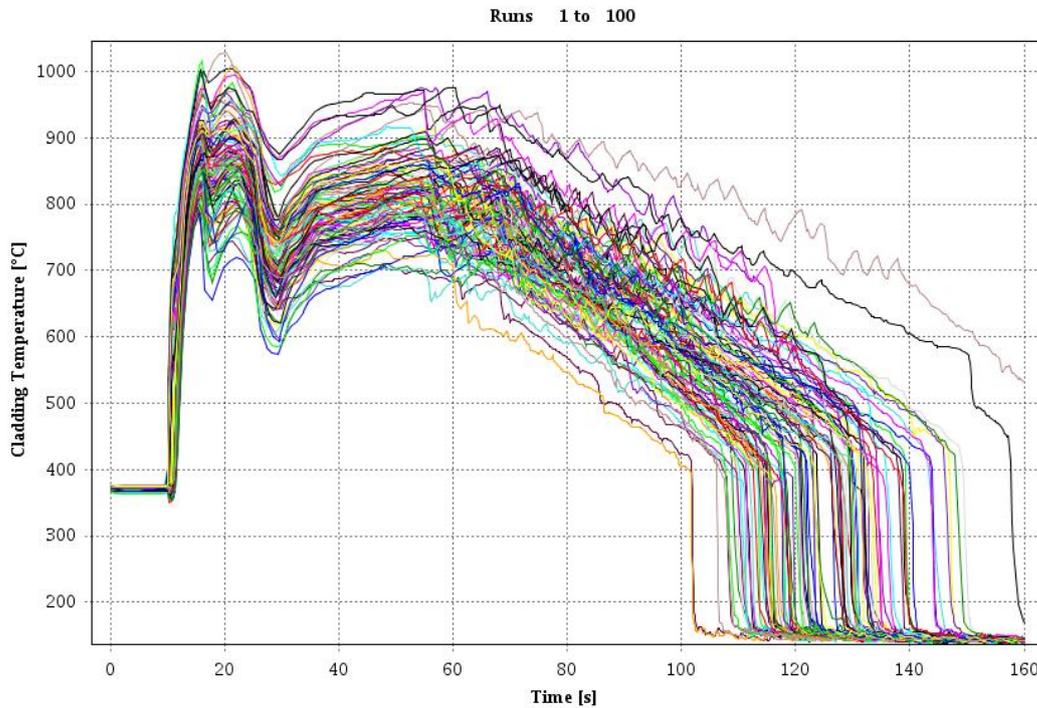


Fig. 4.1 Uncertainty of a time-dependent result represented by the time histories obtained from 100 Monte Carlo simulation runs

To perform the time/index-dependent analysis, the computational result must be available at the same series of time/index steps for all computer code runs. Since this requirement is usually not fulfilled, the following analysis steps have to be performed:

- SUSA either generates equidistant time/index steps of the minimal time range common to all computer code runs or the user defines the time/index steps of interest.
- SUSA performs stepwise linear interpolations with respect to two successive time/index steps t_{ik} and $t_{i_{k+1}}$ and the corresponding results $y(t_{ik})$ and $y(t_{i_{k+1}})$ provided by each run $i, i=1, \dots, n$ to derive $y(t_i)$ at each equidistant or user-defined time/index step t_i with $t_k < t_i \leq t_{k+1}$.

In the following, variable Y represents an uncertain scalar computational result and (y_1, y_2, \dots, y_n) denotes a sample of values (i.e. realizations) provided for Y via n (e.g. $n = 100$) computer code runs.

4.1 Basic statistics

From the sample (y_1, y_2, \dots, y_n) obtained for variable Y , SUSAS calculates the cumulative empirical distribution function $\hat{F}(y)$ as follows:

$$\hat{F}(y) = \frac{\#\{y_i | y_i \leq y\}}{n} = \frac{1}{n} \sum_{i=1}^n 1_{\{y_i \leq y\}} \quad (4.1)$$

where $\#$ is the cardinality symbol meaning the number of elements of a set $\{\}$.

Additionally, SUSAS provides the following basic statistics as indicators of the uncertainty of Y .

- Minimum(y_1, y_2, \dots, y_n)
- Maximum(y_1, y_2, \dots, y_n)
- Empirical percentiles (i.e. simple estimators of 1 %-quantile, 2 %-quantile, ..., 99 %-quantile): $y_{[n \cdot 0.01]:n} \leq y_{[n \cdot 0.02]:n} \leq \dots \leq y_{[n \cdot 0.99]:n}$ with $[\]$ being the floor function and $y_{j:n}$ indicating value No. j of the ordered sample $y_{1:n} \leq y_{2:n} \leq \dots \leq y_{n:n}$.
- Empirical (sample) mean:

$$\bar{y} = \frac{1}{n} \sum y_i \quad (4.2)$$

- Empirical median:

$$m = y_{[n \cdot 0.50]:n}$$

- Unbiased empirical (sample) variance s^2 and standard deviation s :

$$s^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2 \quad (4.3)$$

$$s = \sqrt{s^2}$$

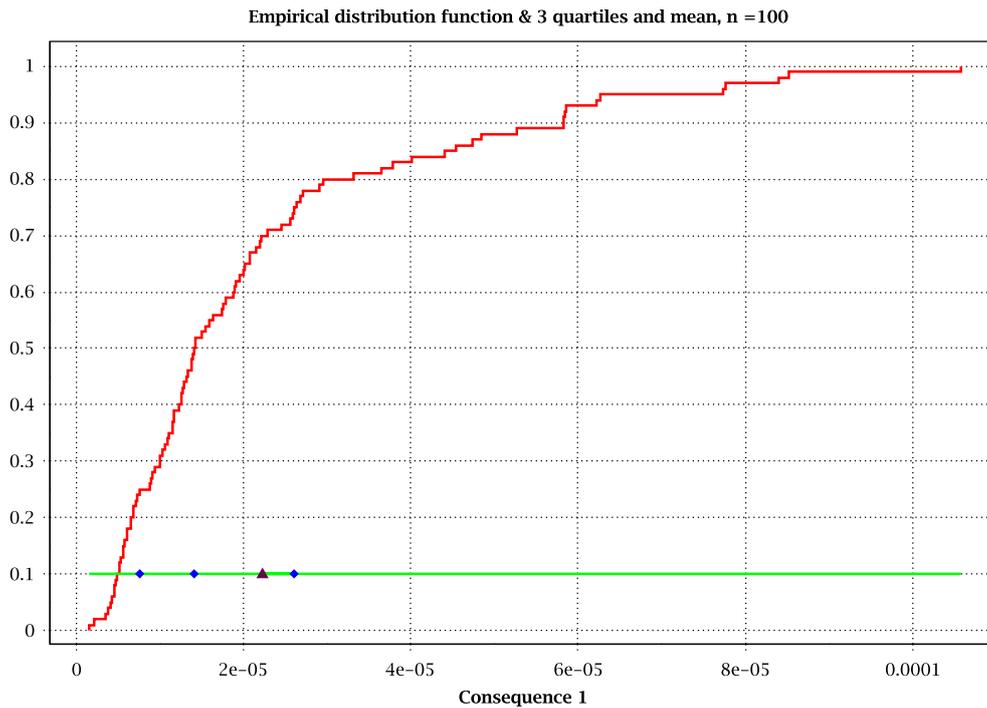


Fig. 4.2 Cumulative empirical distribution function of a scalar computational result (Consequence 1) together with the distribution support (cyan-colored horizontal line), the 3 empirical quartiles $y_{[n \cdot 0.25]:n} \leq y_{[n \cdot 0.50]:n} \leq y_{[n \cdot 0.75]:n}$ (blue diamonds) and the empirical mean (maroon triangle)

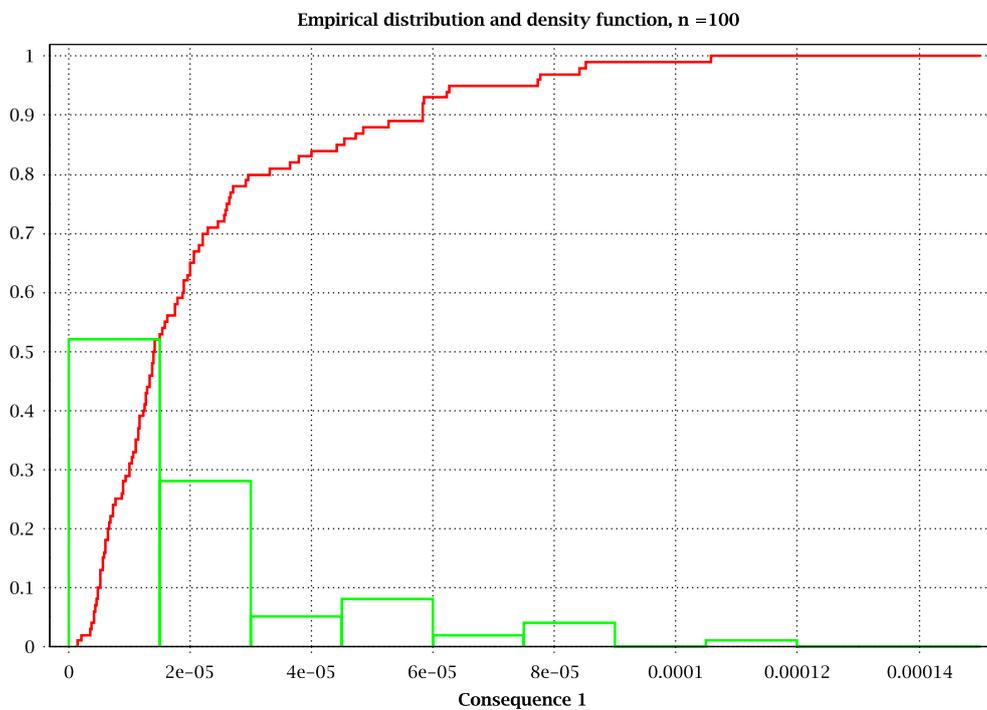


Fig. 4.3 Cumulative empirical distribution (red) and density (green) function of a scalar computational result (Consequence 1)

4.2 Tolerance limits

$\beta \cdot 100 \%$ / $\gamma \cdot 100 \%$ tolerance limits are estimators of the left and/or right endpoint of a two-sided, left-sided or right-sided closed interval (tolerance interval) covering a proportion of at least $\beta \cdot 100 \%$ (e.g., $\beta = 0.95$) of the values of a variable Y at a confidence level of at least $\gamma \cdot 100 \%$ (e.g., $\gamma = 0.95$). β is called the minimal coverage (probability) or the tolerance proportion and γ is the statistical confidence level. Usually, high values ≥ 0.90 are chosen for β and γ .

The confidence level γ accounts for the variability of tolerance limits from sample to sample and gives the probability (percentage) of all samples of the same size n for which the respective tolerance interval covers a proportion of at least $\beta \cdot 100 \%$.

The upper (lower) $\beta \cdot 100 \%$ / $\gamma \cdot 100 \%$ tolerance limit TL_u (TL_l) is a one-sided upper (lower) statistical confidence limit for the $\beta \cdot 100 \%$ ($(1-\beta) \cdot 100 \%$) quantile at a confidence level of at least $\gamma \cdot 100 \%$. The probability for the upper (lower) $\beta \cdot 100 \%$ / $\gamma \cdot 100 \%$ tolerance limit to be at least the $\beta \cdot 100 \%$ ($(1-\beta) \cdot 100 \%$) quantile is $\gamma \cdot 100 \%$ or higher.

The general formula for the $\beta \cdot 100 \%$ / $\gamma \cdot 100 \%$ tolerance limits is given in Eq. (4.4).

$$Prob(Prob(Y \in [TL_l, TL_u] \geq \beta)) \geq \gamma \quad (4.4)$$

where $Prob(\cdot)$ means probability, i.e. coverage probability or probability in the sense of confidence level. If $TL_l = -\infty$, then TL_u is the one-sided upper $\beta \cdot 100 \%$ / $\gamma \cdot 100 \%$ tolerance limit. Analogously, if $TL_u = +\infty$, then TL_l is the one-sided lower $\beta \cdot 100 \%$ / $\gamma \cdot 100 \%$ tolerance limit.

4.2.1 Wilks non-parametrical tolerance limits

The approach of Wilks [WIL 41/, [WIL 42/ to calculate tolerance limits does not require any assumption on the distribution of the considered variable Y and, therefore, is a non-parametrical approach. Wilks' lower and upper tolerance limits TL_l and TL_u are determined by appropriately chosen order statistics.

Let $Y_{1:n} < Y_{2:n} < \dots < Y_{n:n}$ be the order statistics (ordered by increasing size) associated with the sample (Y_1, Y_2, \dots, Y_n) of size n of a random variable Y and let $f(y)$ be the density function of Y . Then, the coverage (probability) of the region between the order statistics $Y_{r:n}$ and $Y_{s:n}$ with $0 \leq r < s \leq n+1$ has a Beta distribution with parameters $(s-r)$ and $n-(s-r)+1$ and, therefore, the following relationship in Eq. (4.5) is true for any $\beta \in (0, 1)$.

$$Prob\left(\int_{Y_{r:n}}^{Y_{s:n}} f(y)dy \geq \beta\right) = \frac{\int_{\beta}^1 u^{s-r-1} (1-u)^{n-(s-r)} du}{B(s-r, n-(s-r)+1)} \quad (4.5)$$

where $B(s-r, n-(s-r)+1)$ represents the complete Beta function.

The Beta distribution of the coverage probability is independent of the distribution of Y and only depends on the orders r and s .

Between the Beta distribution and the Binomial distribution, the following well-known relationship holds: 1. - $Beta_{k,n-k+1}(\beta) = Binomial_{n,\beta}(k-1)$. That means

$$\frac{\int_{\beta}^1 u^{k-1} (1-u)^{n-k} du}{B(k, n-k+1)} = \sum_{i=0}^{k-1} \binom{n}{i} \beta^i (1-\beta)^{n-i} \quad (4.6)$$

Due to Eqs. (4.5) – (4.6), the relationship in Eq. (4.7) can be applied to determine the order statistics $Y_{r:n}$ and $Y_{s:n}$ representing the lower and upper tolerance limit TL_l and TL_u , respectively.

$$Prob\left(\int_{Y_{r:n}}^{Y_{s:n}} f(y)dy \geq \beta\right) = \sum_{i=0}^{s-r-1} \binom{n}{i} \beta^i (1-\beta)^{n-i} \geq \gamma \quad (4.7)$$

To obtain a one-sided tolerance limit, the following definitions are made:

$$Y_{0:n} := -\infty \quad Y_{n+1:n} := +\infty$$

From these definitions, it follows that

- s must be set to $(n+1)$ in order to determine the one-sided lower tolerance limit TL_l
- r must be set to zero in order to determine the one-sided upper tolerance limit TL_u

Fig. 4.4 compares Wilks' one-sided upper 95 %/95 % tolerance limit with the empirical 95 %-quantile and the true 95 %-quantile. Both estimators, i.e. the tolerance limit and the empirical quantile are calculated from each of 1000 different samples of size $n = 100$. Each of the 1000 samples is selected from the standard Normal distribution with the 95 %-quantile being equal to 1.6448. While in most cases ($\geq \gamma \cdot 100 \% = 95 \%$) the one-sided upper tolerance limit exceeds the true 95 %-quantile, the empirical quantile is below or above the true 95 %-quantile with a probability of 0.5. In most cases, the tolerance limit is very conservative compared to the true 95 %-quantile as well as to the empirical quantile. But there is a chance of at most 5 % ($(1-\gamma) \cdot 100 \%$), that the tolerance limit remains below the true 95 %-quantile.

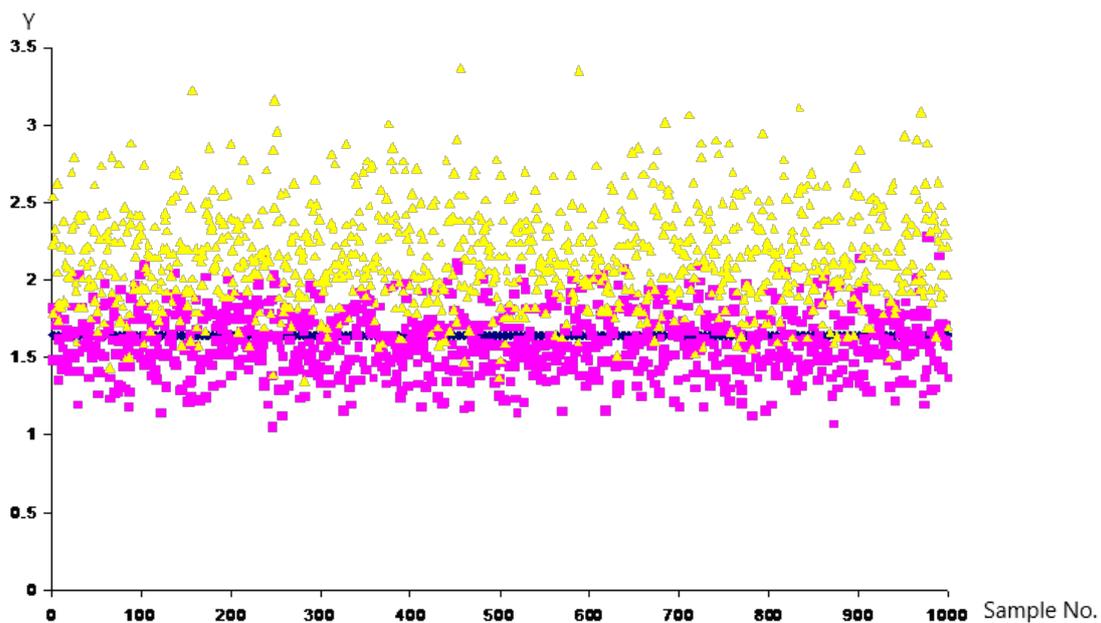


Fig. 4.4 Wilks' one-sided upper 95 %/ 95 % tolerance limit (yellow triangles) compared with the empirical 95 %-quantile (magenta) calculated in each of 1000 different samples of size $n=100$ from a standard Normal distribution with the 95 %-quantile = 1.6448 (blue horizontal line)

The minimum sample size to determine the $\beta \cdot 100 \% / \gamma \cdot 100 \%$ tolerance interval can be derived from Eq. (4.8) for the one-sided lower or upper tolerance interval or from Eq. (4.9) for the two-sided tolerance interval. Both equations follow from Eq. ((4.7).

$$1 - \beta^n \geq \gamma \quad (4.8)$$

$$1 - \beta^n - n \cdot (1 - \beta) \cdot \beta^{n-1} \geq \gamma \quad (4.9)$$

Tab. 4.1 Minimum sample size to determine the $\beta \cdot 100\%$ / $\gamma \cdot 100\%$ tolerance interval for selected coverage probabilities β and confidence levels γ

One-sided statistical tolerance limits				
$\gamma \backslash \beta$	0.90	0.95	0.99	
0.90	22	45	230	
0.95	29	59	299	
0.99	44	90	459	
Two-sided statistical tolerance limits				
$\gamma \backslash \beta$	0.90	0.95	0.99	
0.90	38	77	388	
0.95	46	93	473	
0.99	64	130	662	

Fig. 4.5 shows different cumulative Beta distribution functions of the coverage probability P of the interval $(-\infty, TL_u]$ right-sided closed by Wilks' upper 95 %/95 % tolerance limit TL_u for different sample sizes n . TL_u is identical to an appropriately chosen order statistic $Y_{s:n}$ where s depends on $\beta (=0.95)$, $\gamma (=0.95)$, and on the sample size n . Each distribution is an indicator of the conservativeness of the 95 %/95 % tolerance limit TL_u . As can be seen, the probability is 0.95 (i.e. the confidence level is 95 %), that the coverage of the interval $(-\infty, Y_{s:n}]$ with $s=n=59$ exceeds a proportion of 95 % (i.e. a probability P of 0.95). Simultaneously, the coverage of this interval exceeds

- a proportion of 96 % with a confidence level of 91 %
- a proportion of 97 % with a confidence level of 83 %
- a proportion of 98 % with a confidence level of 76 %
- a proportion of 99 % with a confidence level of 45 %

On the other side, the confidence level is 5 %, that the coverage of the interval $(-\infty, Y_{s:n}]$ with $s=n=59$ does not exceed a proportion of 95 % (i.e. a probability P of 0.95). The coverage of this interval does not exceed even

- a proportion of 94 % with a confidence level of 2.6 %
- a proportion of 93 % with a confidence level of 1.4 %

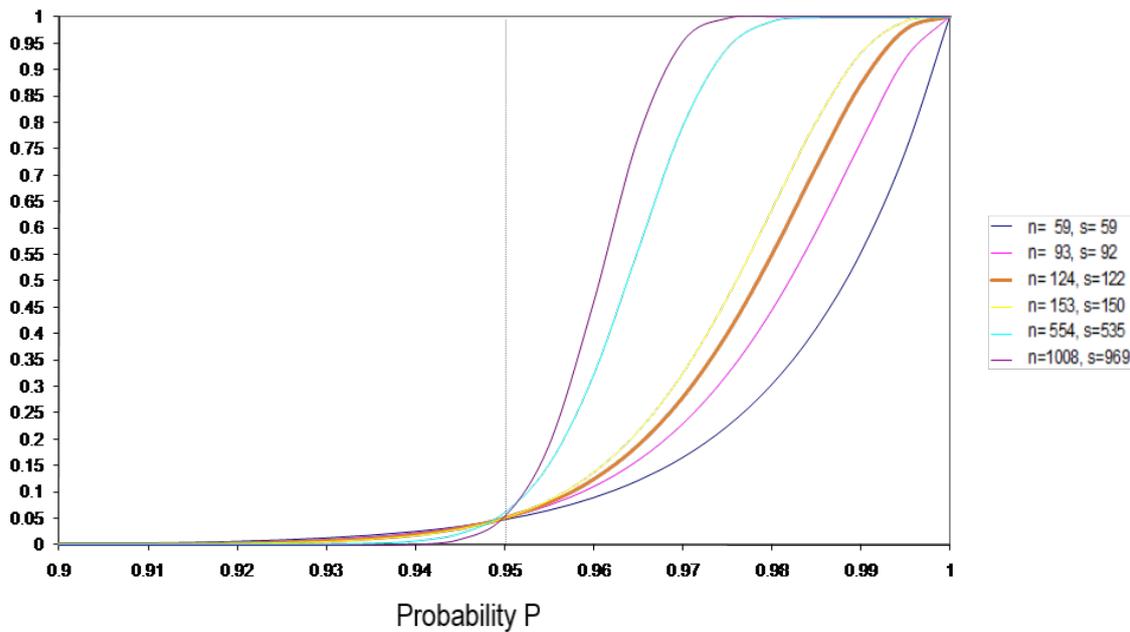


Fig. 4.5 Cumulative Beta distribution functions of the coverage probability P of the interval $(-\infty, Y_{s:n}]$ right-sided closed by Wilks' upper 95 %/ 95 % tolerance limit $Y_{s:n}$ for different orders s and different sample sizes n

Fig. 4.5 shows, that the quality of the 95 %/95 % tolerance interval is getting better with increasing sample size n and an appropriately adapted order s . That means the degree of conservativeness decreases with increasing sample size. While the coverage of the interval $(-\infty, Y_{s:n}]$ with $s=n=59$ exceeds a proportion of 99 % with a confidence level of 45 %, this confidence level is

- 24 % with $n= 93$ and $s= 92$
- 12 % with $n=124$ and $s=122$
- 7 % with $n=153$ and $s=150$

With $n \geq 234$, the confidence level is getting vanishingly low.

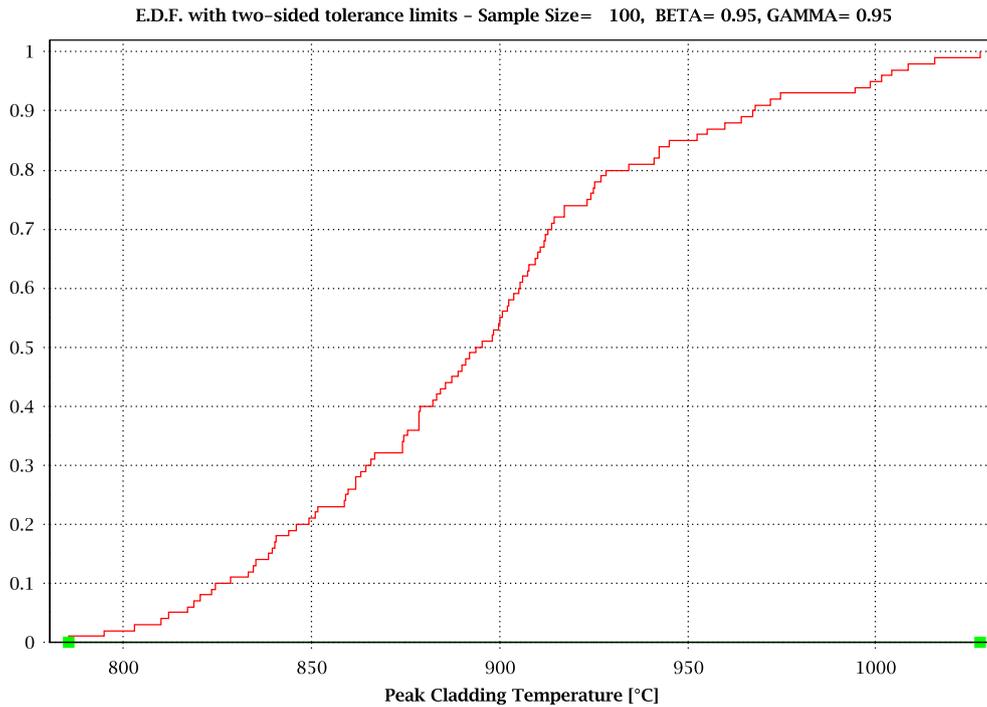


Fig. 4.6 Cumulative empirical distribution function of a scalar computational result (Peak Cladding Temperature) and Wilks' two-sided tolerance limits (green squares on the x-axis)

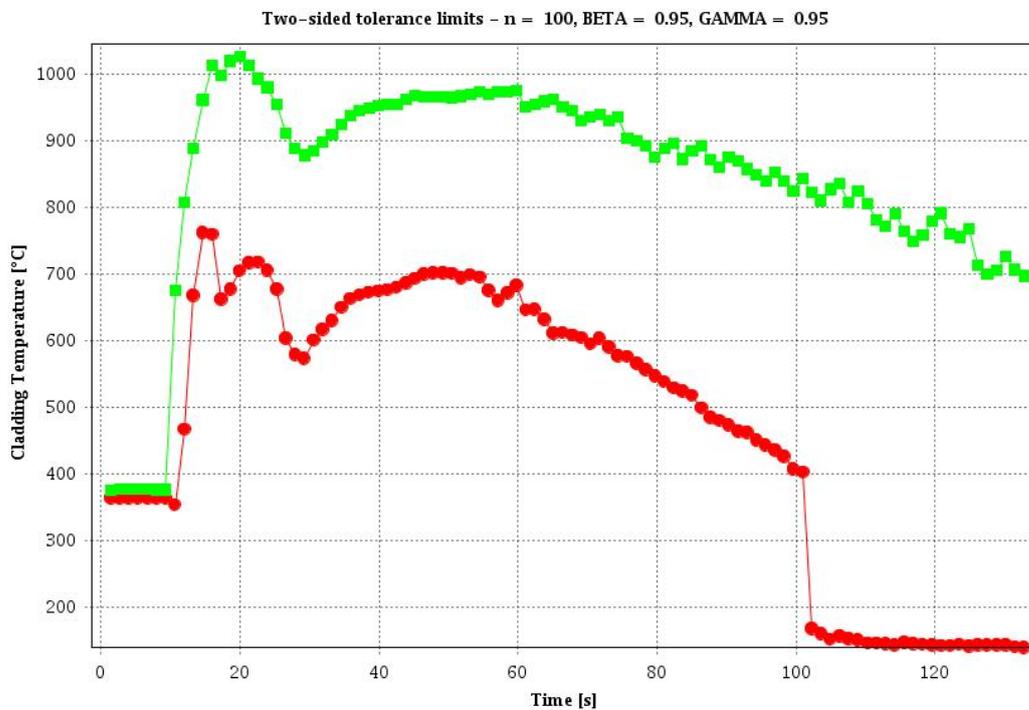


Fig. 4.7 Wilks' two-sided tolerance limits of a time-dependent computational result (Cladding Temperature)

4.2.2 Tolerance limits in case of a Normal or Lognormal distribution

If it can be assumed that the considered variable Y has a Normal distribution, then the lower and upper tolerance limits TL_l and TL_u can be derived from Eq. (4.10) and Eq. (4.11), respectively.

$$TL_l = \bar{y} - k \cdot s \quad (4.10)$$

$$TL_u = \bar{y} + k \cdot s \quad (4.11)$$

\bar{y} is the empirical mean (Eq. (4.2)) and s is the unbiased empirical standard deviation (Eq. (4.3)). k is appropriately determined for one and two-sided tolerance limits.

For one-sided tolerance intervals and $n \leq 200$, $k = k_1$ is determined according to Eq. (4.12) /GUT 70/.

$$k_1 = \frac{t_{\gamma, n-1, \delta}}{\sqrt{n}} \quad (4.12)$$

where $t_{\gamma, n-1, \delta}$ is the $\gamma \cdot 100$ %-quantile of the non-central Student t distribution with parameter $(n-1)$ and non-centrality parameter δ defined as

$$\delta = z_\beta \cdot \sqrt{n} \quad (4.13)$$

For one-sided tolerance intervals and $n > 200$, $k = k_1$ is determined according to Eq. (4.14) /NAT 63/.

$$k_1 = \frac{z_\beta + \sqrt{z_\beta^2 - a \cdot b}}{a} \quad (4.14)$$

where z_β denotes the $\beta \cdot 100$ %-quantile of the standard Normal distribution and

$$a = 1 - \frac{z_{\gamma}^2}{2 \cdot (n - 1)} \quad (4.15)$$

$$b = z_{\beta}^2 - \frac{z_{\gamma}^2}{n} \quad (4.16)$$

For two-sided tolerance intervals, $k = k_2$ is determined according to Eq. (4.17) /HOW 69/.

$$k_2 = \frac{\frac{(n - 1)(n + 1)}{n} + z_{(1-\beta)/2}^2}{x_{1-\gamma, n-1}^2} \quad (4.17)$$

where $z_{(1-\beta)/2}^2$ is the $(1 - \beta)/2 \cdot 100$ %-quantile of the standard Normal distribution and $x_{1-\gamma, n-1}$ is the $(1 - \gamma) \cdot 100$ %-quantile of the X^2 distribution with parameter $(n - 1)$.

If it can be assumed that the considered variable Y has a Lognormal distribution, then the lower and upper tolerance limits TL_l and TL_u can be derived by the following steps:

- In-transformation (natural logarithm) of Y .
- Calculation of the tolerance limits TL'_l and TL'_u for the normally distributed variable $\ln(Y)$
- Retransformation of the tolerance limits TL'_l and TL'_u , i.e.

$$TL_l = \exp(TL'_l) \text{ and } TL_u = \exp(TL'_u)$$

4.2.3 Bootstrapped tolerance limits

Another approach to compute a tolerance interval of the computational result Y offers the nested (or double) bootstrapping of confidence intervals, in the following referred to as BTI. This approach may be selected, if

- no information about the distribution of the computational result is available
- not as many computations as required for Wilks' distribution-free/non-parametric approach can be afforded

- Wilks approach is considered to be too conservative /FER 01/

An intuitive motivation of the BTI idea provides the underlying approach of k-factors. Let $\vec{y} = (y_1, \dots, y_n) \in \mathbb{R}^n$ be a random sample following a distribution function F such that $y_i \stackrel{iid}{\sim} F$. An interval $[TL_l, TL_u]$ is called a β -content, γ -confidence tolerance interval for a cumulative distribution function F , if the statistics TL_l and TL_u satisfies $Prob(F(TL_u) - F(TL_l) \geq \beta) \geq \gamma$ while the statistics TL_l, TL_u are called tolerance limits.

Let's assume that a (non-symmetric) tolerance interval TI can be generally formulated as $TI = [\bar{y} + k_1 s; \bar{y} + k_2 s]$ with the sample mean \bar{y} (Eq. (4.2)), the sample standard deviation s (Eq. (4.3)) and the factors $k_1 < k_2$. \bar{y} and s are estimated from the random sample $\vec{y} = (y_1, \dots, y_n)$. Based on this vague formulation, the definition of a tolerance interval can be reformulated with respect to the predictor variable Z (any future observation from the underlying population) as

$$\begin{aligned} & Prob(F(\bar{y} + k_2 s) - F(\bar{y} + k_1 s) \geq \beta) \\ & = Prob(Prob(k_1 < \frac{Z - \bar{y}}{s} < k_2 | \{y_i\}_{i=1}^n) \geq \beta) \\ & \geq \gamma \end{aligned} \tag{4.18}$$

Via the formulation in Eq. (4.18), it becomes clear that two probabilistic expressions need to be iterated to derive a tolerance interval via a common bootstrapping approach.

The idea to derive a tolerance interval via nested bootstrapping is based on the work of /EFR 86/, /SHO 05/ and /REB 07/.

The basic steps of the algorithm can be summarized as described in the following:

- Bootstrap (i.e. resample with replacement) the complete sample of the computational result Y and generate B independent bootstrap samples $\vec{y}^{(1)*}, \dots, \vec{y}^{(B)*}$ with $y_i^{(b)*} \stackrel{iid}{\sim} \hat{F}$ and \hat{F} being the empirical distribution function derived from $\vec{y} = (y_1, \dots, y_n)$. For the sake of clarity, $(\cdot)^*$ indicates a bootstrap sample and any estimator derived from it;
- For each $b, b = 1, \dots, B$, repeat the following steps:
 - Sample P future observations $z^{(1)*}, \dots, z^{(P)*}$ with $z^{(i)*} \stackrel{iid}{\sim} \hat{F}$ with replacement from the sample $\vec{y} = (y_1, \dots, y_n)$

- For each $p, p = 1, \dots, P$, calculate predictor statistic $T^{(b,p)*} = \frac{z^{(p)*} - \bar{y}^{(b)*}}{s^{(b)*}}$ with the sample mean $\bar{y}^{(b)*}$ and the unbiased standard deviation $s^{(b)*}$ derived from the corresponding bootstrap sample $\vec{y}^{(b)*}$
- Localize the quantiles $\hat{t}^b(\beta_1)$ and $\hat{t}^b(\beta_2)$ of $\{T^{(b,p)*}\}_{p=1}^P$ with $\beta_2 - \beta_1 = \beta$; here, a simple symmetric algorithm is employed that iteratively removes the interval with the smallest lowest limit, then with largest upper limit, etc.;
- Find k_1 and k_2 such that a γ -portion of the mean coverage intervals of $\{[\hat{t}^b(\beta_1), \hat{t}^b(\beta_2)]\}_{b=1}^B$ are completely included in $[k_1; k_2]$.
- Calculate the bootstrapped $\beta \cdot 100\% / \gamma \cdot 100\%$ tolerance interval as $[\bar{y} + k_1 s; \bar{y} + k_2 s]$.

Another strategy to calculate BTI depicts the well-known content-corrected tolerance intervals which try to relax the assumption of normality (proof that asymptotic normality is sufficient). However, this strategy only ensures that the specified confidence level γ holds, but not the specified content or coverage probability β . Nevertheless, it serves as the only setting in which the performance of a BTI can be compared to its analytical/theoretical counterpart for multiple (non-normal as well as asymmetric) distributional shapes.

4.3 Interval limits from Chebychev and Chebychev-Cantelli inequalities

The Chebychev and/or Chebychev-Cantelli inequalities may be applied to estimate a two-sided interval and/or one-sided closed intervals covering at least a proportion of $\beta \cdot 100\%$. Since these inequalities are applicable to any distribution, the interval estimators are conservative in general.

The application of the Chebychev-Cantelli inequality provides one-sided left- or right-closed intervals covering at least a proportion of $\beta \cdot 100\%$ (Eq. (4.19)).

$$Prob(Y \geq E(Y) + t) \leq \frac{Var(Y)}{Var(Y) + t^2}, \quad t \geq 0 \quad (4.19)$$

$$Prob(Y \leq E(Y) - t) \leq \frac{Var(Y)}{Var(Y) + t^2}, \quad t \geq 0$$

where $E(Y)$ denotes the expected value and $Var(Y)$ denotes the variance of Y .

With $Prob(Y \geq E(Y) + t) \leq 1 - \beta$, the upper limit of a right-closed interval covering a proportion of at least $\beta \cdot 100\%$ is given by $E(Y) + \sqrt{\frac{\beta}{1-\beta}} \cdot \sqrt{VarY}$.

With $Prob(Y \leq E(Y) - t) \leq 1 - \beta$, the lower limit of a left-closed interval covering a proportion of at least $\beta \cdot 100\%$ is given by $E(Y) - \sqrt{\frac{\beta}{1-\beta}} \cdot \sqrt{VarY}$.

The application of the Chebychev inequality provides a two-sided closed interval covering a proportion of at least $\beta \cdot 100\%$ (Eq. (4.20)).

$$Prob(|Y - E(Y)| \geq t) \leq \frac{Var(Y)}{t^2}, \quad t \geq 0 \quad (4.20)$$

With $Prob(|Y - E(Y)| \geq t) \leq 1 - \beta$, the limits of a two-sided closed interval covering a proportion of at least $\beta \cdot 100\%$ are given by $E(Y) \pm \sqrt{\frac{1}{1-\beta}} \cdot \sqrt{VarY}$.

If $E(Y)$ and \sqrt{VarY} are estimated by the empirical mean \bar{y} (Eq. (4.2)) and the unbiased empirical standard deviation s (Eq. (4.3)), respectively, appropriate estimators of the intervals can be derived. However, these interval estimators are not associated with a statistical confidence level.

4.4 Parametric distribution fitting

The uncertainty of a variable Y may be quantified by an appropriate parametric probability distribution. Uncertainty quantifications using quantiles or intervals covering a proportion of $\beta \cdot 100\%$ can easily be derived from such a distribution.

If the sample values y_1, \dots, y_n are available for variable Y , an appropriate parametric distribution may be found by the following steps:

- Selection of a distribution type (Normal, Uniform, Beta, etc.)
- Estimation of the parameters of the selected distribution from the sample values y_1, \dots, y_n either by the maximum likelihood /KOT 88/ or by the moment-matching method (method of moments) /KOT 88/
- Goodness-of-fit test (Kolmogorov-Smirnov, Lilliefors) for the selected distribution defined by the estimated parameters

4.4.1 Kolmogorov-Smirnov goodness-of-fit test

The Kolmogorov-Smirnov goodness-of-fit test /KOT 88/ is a non-parametrical statistical test of the equality of two continuous probability distributions. The corresponding one-sample goodness-of-fit test can be used to compare the empirical distribution derived from the sample values y_1, \dots, y_n with a reference parametrical distribution (e.g. Normal distribution or any other continuous distribution). In this case, the hypothesis is tested, that the sample values are selected from the indicated reference distribution.

The one-sample Kolmogorov-Smirnov test statistic D_n quantifies the distance between the empirical cumulative distribution function $F_n(y)$ and the cumulative distribution function $F(y)$ of the reference parametrical distribution.

$$D_n = \sup_y |F_n(y) - F(y)| \quad (4.21)$$

where \sup_y is the supremum of the set of distances obtained for variable Y .

If $F_n(y)$ is the empirical distribution function of a sample selected from the distribution $F(y)$, then D_n converges to 0 almost surely with n approaching infinity. The distribution of D_n is called Kolmogorov-Smirnov distribution.

Let d denote a realization of the maximum distance D_n between the empirical cumulative distribution function $F_n(y)$ and the cumulative distribution function $F(y)$. If the sample is selected from the reference distribution, then the probability to exceed d is very small, if

d is high. This probability is higher, if d is smaller. Therefore, the following conclusions may be drawn:

- If $Prob(D_n > d) \leq 0.05$ (or 0.01), the null hypothesis is rejected at the significance level of 0.05 (or 0.01), that the considered sample of variable Y is selected from the reference parametrical distribution.
- If $Prob(D_n > d) > 0.05$ (or 0.01), the null hypothesis cannot be rejected at the significance level of 0.05 (or 0.01), that the considered sample of variable Y is selected from the reference parametrical distribution. But a large probability $Prob(D_n > d)$ may justify assuming, that the sample is selected from the reference distribution.

In practice, the Kolmogorov-Smirnov goodness-of-fit test requires a relatively large sample size in order to properly reject the null hypothesis. If either the distribution type or the distribution parameters are determined from the sample, the Kolmogorov-Smirnov test result may not be reliable, especially for small sample sizes.

4.4.2 Lilliefors goodness-of-fit test

The Lilliefors goodness-of-fit test /LIL 67/, /LIL 69/ uses the same test statistic D_n (Eq. (4.21)) as the Kolmogorov-Smirnov goodness-of-fit test. Whereas the Kolmogorov-Smirnov test is applicable to any parametrical continuous probability distribution, the Lilliefors test is applicable only to the Normal, Lognormal and Exponential distribution. Therefore, the corresponding Lilliefors distribution of D_n is stochastically smaller than the Kolmogorov-Smirnov distribution of D_n . That means the Lilliefors goodness-of-fit test would reject the null hypothesis, that the considered variable has a Normal distribution (or Lognormal or Exponential distribution), more likely than the Kolmogorov-Smirnov goodness-of-fit test.

Tables of the critical values for selected levels of significance may be found in /LIL 67/, /LIL 69/. For SUSAs, these tables were extended to cover sample sizes n with $30 < n \leq 60$ and more levels of significance. For sample sizes $n > 60$, the approximation of the critical value c_n in Eq. (4.22) is used:

$$c_n = c_{60} \cdot \sqrt{\frac{60}{n}} \quad (4.22)$$

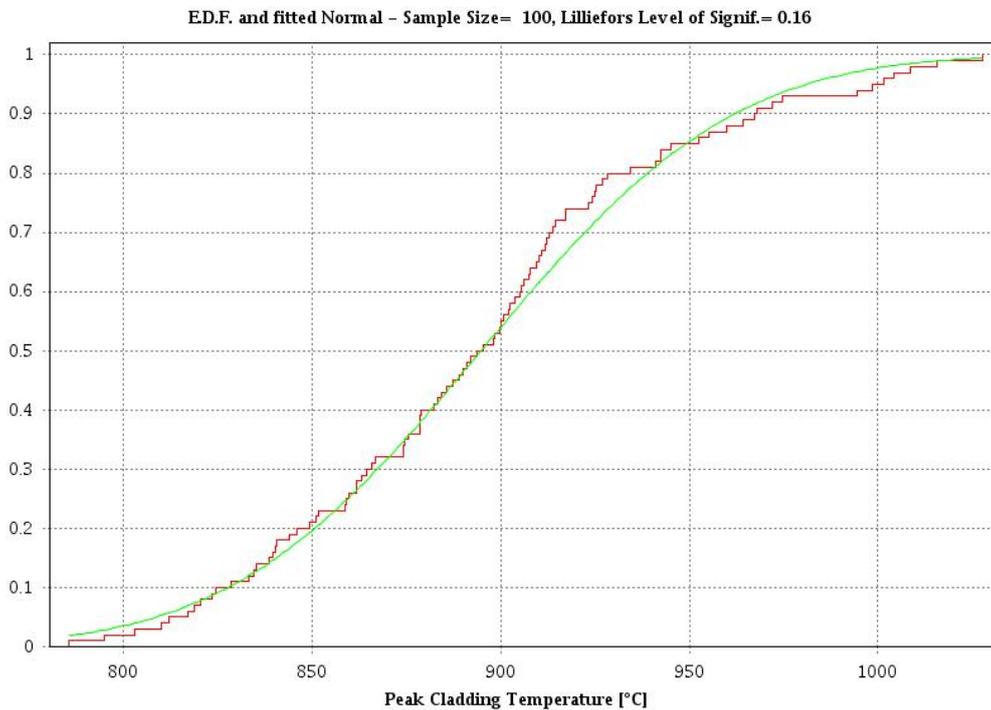


Fig. 4.8 Cumulative empirical distribution function of a scalar computational result (Peak Cladding Temperature) and fitted Normal distribution. The level of significance of the Lilliefors test is 0.16, i.e. the hypothesis cannot be rejected that Y has a Normal distribution.

4.5 Construction and application of a surrogate model

If the computer code is complex and slowly-running, the number of Monte-Carlo simulation runs, which can be performed within an acceptable time period, may be relatively small. Consequently, the sample of values y_1, y_2, \dots, y_n which can be obtained for the computational result is relatively small and does not allow for estimating distribution parameters and other characteristics useful for uncertainty quantifications with high

accuracy. Even the calculation of Wilks' tolerance limits may not be possible, because the minimum sample size to calculate these limits cannot be reached.

For complex computer codes, one way to obtain uncertainty quantifications for the computational result is the replacement of the code by a fast-running simpler surrogate code and the use of this fast-running code to perform many Monte-Carlo simulation runs. From the large sample of values finally available for the computational result, more accurate estimators including Wilks' tolerance limits can be calculated.

SUSA allows constructing a surrogate model and assessing its goodness-of-fit based on a comparison of its results with the results of the original model. The surrogate model is constructed from the vectors $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$ (each of length $npar$) sampled for the uncertain input parameters $X_1, X_2, \dots, X_{npar}$ and the corresponding results y_1, y_2, \dots, y_n provided by the n runs with the original model. For the construction process, SUSA applies a forward stepwise regression algorithm which successively inserts uncertain parameters into or removes them from the regression function according to the results of partial F-tests with the level of significance $\alpha = 0.05$. The user may control the parameter selection process by forcing parameters into the regression function or by excluding parameters from the analysis when the input of the uncertainty analysis is prepared. Parameters forced into the regression cannot be removed during the selection process.

The additional forward stepwise rank regression algorithm uses, instead of the original values, the rank transformed values $r(x_{j1}), r(x_{j2}), \dots, r(x_{jn})$ of each parameter $X_j, j = 1, \dots, npar$, and the rank transformed values $r(y_1), r(y_2), \dots, r(y_n)$ of the result Y (Eq. (2.67)).

After the model construction process, an ordinary regression model as indicated in Eq. (4.23) or a rank regression model as indicated in Eq. (4.24) is available.

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k \quad (4.23)$$

$$r_{\hat{Y}} = \beta'_0 + \beta'_1 r(X_1) + \dots + \beta'_l r(X_l) \quad (4.24)$$

With the constructed ordinary or rank regression model, new Monte-Carlo simulation runs may be performed. To this purpose, a new sample of parameter values must be

generated first. Since the corresponding runs are very fast, the size n_l of the new parameter sample may be larger than the original parameter sample size n .

The results $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{n_1}$ of an ordinary regression model are obtained by simply using the new parameter values as input parameters of the regression model. In case of a rank regression model, the results $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{n_1}$ are obtained by the following steps:

- The vectors $\underline{rx}_1, \underline{rx}_2, \dots, \underline{rx}_{n_1}$ (each of length n_{par}) of rank transformed parameter values are successively provided as input to the rank regression model to calculate n_l $r_{\hat{y}}$ -values (Eq. (4.24)).
- To obtain the actual result \hat{y} for the calculated $r_{\hat{y}}$ -value, piecewise linear interpolations (extrapolations) are performed with respect to two successive ranks $r(y_k)$ and $r(y_l)$ with $r(y_k) < r_{\hat{y}} \leq r(y_l)$ and the corresponding values y_k and y_l derived from the runs with the original code.

The Kolmogorov-Smirnov two-sample goodness-of-fit test /KOT 88/ is applied to test the hypothesis that the two samples y_1, y_2, \dots, y_n and $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{n_1}$ come from the same probability distribution (see Section 4.4.1). The corresponding test statistic D_{n_q} quantifies the distance between the two empirical cumulative distribution functions F_n and F_{n_1} at $n_q = 13$ selected empirical quantiles q_1, \dots, q_{n_q} (1 %, 5 %, 10 %, ..., 90 %, 95 %, 99 %).

$$D_{n_q} = \sup_{q_1 \dots q_{n_q}} |F_n(q_i) - F_{n_1}(q_i)| \quad (4.25)$$

Let d denote a realization of the maximum distance D_{n_q} between the two empirical cumulative distribution functions $F_n(y)$ and $F_{n_1}(\hat{y})$. If the two samples are from the same distribution, then the probability to exceed d is very small, if d is high. This probability is higher, if d is smaller. Therefore, the following conclusions may be drawn:

- If $Prob(D_{n_q} > d) \leq 0.05$ (or 0.01), the null hypothesis is rejected at the significance level of 0.05 (or 0.01), that the considered samples are selected from the same distribution.
- If $Prob(D_{n_q} > d) > 0.05$ (or 0.01), the null hypothesis cannot be rejected at the significance level of 0.05 (or 0.01), that the considered samples are selected from

the same distribution. But a large probability $Prob(D_{n_q} > d)$ may justify assuming, that the samples are selected from the same distribution.

4.6 Uncertainty quantifications for multiple variables

SUSA offers two options to derive uncertainty quantifications simultaneously relating to two and more computational results (figures of merit). These uncertainty quantifications are useful, for instance, to prove the simultaneous compliance of several safety limits in the design and licensing process of nuclear power plants.

The first option provides simultaneous multiple tolerance limits with a common confidence level. The second option provides an estimator of the probability of compliance of multiple limiting values.

4.6.1 Simultaneous multiple tolerance limits

Simultaneous multiple tolerance limits can be derived based on the inequality of Bonferroni (Eq. (4.26)). This inequality gives the lower limit of the probability that the events E_1, E_2, \dots, E_m from any finite or countable set of events occur simultaneously.

$$Prob\left(\bigcap_{k=1}^m E_k\right) \geq 1 - \sum_{k=1}^m (1 - Prob(E_k)) \quad (4.26)$$

If E_k denotes the event that the coverage (probability) of the tolerance interval for variable Y_k is at least β , simultaneous occurrence of the events E_1, E_2, \dots, E_m ($\bigcap_{k=1}^m E_k$) means that the coverages of the multiple tolerance intervals for the variables Y_1, Y_2, \dots, Y_m are simultaneously at least β . The probability $Prob(E_k)$ is identical to the confidence level γ_k of the individual tolerance interval of variable $Y_k, k = 1, \dots, m$, and the probability $Prob(\bigcap_{k=1}^m E_k)$ is the so-called common confidence level γ' of the simultaneous multiple tolerance intervals.

Based on the inequality of Bonferroni, simultaneous multiple tolerance intervals characterized by an identical coverage probability β (e.g. $\beta = 0.95$) and a common confidence level γ' (e.g. $\gamma' = 0.95$) are given by the individual $\beta \cdot 100\% / \gamma$ 100% tolerance intervals of the variables Y_1, Y_2, \dots, Y_m with $\gamma = \gamma_1 = \dots = \gamma_m$ derived from Eq. (4.27).

$$\gamma = 1 - \frac{1 - \gamma'}{m} \quad (4.27)$$

The individual $\beta \cdot 100\% / \gamma$ 100% tolerance limits may be derived according to Wilks approach or by assuming a Normal or Lognormal distribution (see Sections 4.2.1 and 4.2.2).

4.6.2 Probability of compliance of multiple limiting values

Let Y_1, Y_2, \dots, Y_m be variables which are required to fulfil the respective conditions C_1, C_2, \dots, C_m . $C_k, k = 1, \dots, m$, may be one of the following conditions:

- $Y_k \geq y_k^l$
- $Y_k \leq y_k^u$
- $y_k^l < Y_k \leq y_k^u$

Furthermore, let $\underline{y}_1, \underline{y}_2, \dots, \underline{y}_n$ be n vectors of length m each simultaneously sampled for the variables Y_1, Y_2, \dots, Y_m . Vector \underline{y}_i may represent the values of the different computational results (figures of merit) Y_1, Y_2, \dots, Y_m simultaneously obtained via the Monte Carlo simulation run No. $i, i = 1, \dots, n$.

The probability $Prob(\bigcap_{k=1}^m C_k)$ that the variables Y_1, Y_2, \dots, Y_m simultaneously fulfil the required conditions can be estimated from the sample $\underline{y}_1, \underline{y}_2, \dots, \underline{y}_n$. An appropriate estimator is the lower $\gamma \cdot 100\%$ confidence limit p_l of the probability. According to the approach of Clopper and Pearson /CLO 34/, p_l is derived by determining the number s of runs (of the total of n runs) where the corresponding vector \underline{y}_i fulfils the required conditions and by solving Eq. (4.28) for an appropriately selected γ -value (e.g. $\gamma=0.95$).

$$\sum_{i=s}^n \binom{n}{i} p_l^i (1-p_l)^{n-i} = 1 - \gamma \quad (4.28)$$

$$\sum_{i=0}^{s-1} \binom{n}{i} p_l^i (1-p_l)^{n-i} = 1 - \text{Prob}\left(F \leq \frac{n-s+1}{s} \frac{p_l}{1-p_l}\right) \quad (4.29)$$

Based on Eq. (4.29) with F being a random variable distributed according to the F -distribution with $2 \cdot s$ and $2 \cdot (n - s + 1)$ degrees of freedom, p_l can be determined as follows:

$$p_l = \frac{s \cdot F_{2s, 2(n-s+1), 1-\gamma}}{(n-s+1) + s \cdot F_{2s, 2(n-s+1), 1-\gamma}} \quad (4.30)$$

with $F_{2s, 2(n-s+1), 1-\gamma}$ being the $(1-\gamma) \cdot 100$ %-quantile of the F -distribution with $2 \cdot s$ and $2 \cdot (n - s + 1)$ degrees of freedom.

Since $F_{2s, 2(n-s+1), 1-\gamma} = \frac{1}{F_{2(n-s+1), 2s, \gamma}}$, p_l can also be determined via Eq. (4.31):

$$p_l = \frac{s}{(n-s+1) \cdot F_{2(n-s+1), 2s, \gamma} + s} \quad (4.31)$$

5 Sensitivity Analysis

A sensitivity analysis or, more precisely, an uncertainty importance analysis helps to identify those uncertain input parameters which mainly contribute to the uncertainty of the computational result /HOF 99/, /SAL 00/. Improvements of the state of knowledge on these parameters may help to reduce the uncertainty of the result most effectively.

Like the uncertainty analysis, the sensitivity analysis can be performed for a scalar as well as for a time/index-dependent computational result (see introduction of Section 4).

Sensitivity indices appropriate for computationally intensive computer codes which don't allow performing many Monte Carlo simulation runs are those related to statistical correlations. These sensitivity indices can be calculated from the same sample data as already generated for the uncertainty analysis. The correlation related indices applicable to individual parameters are described in Section 5.1 and those applicable to parameter groups are explained in Section 5.2.

Besides correlation related sensitivity measures, the classical correlation ratio from original and rank transformed data may serve as sensitivity index (Section 5.3). The square of the correlation ratio is equivalent to the variance based first order sensitivity index also known as Sobol's first order index. The procedure to calculate this index is described in Section 5.4.

Other types of sensitivity indices implemented in SUSA are the association measures from 2x2 contingency tables (Section 5.5) and the regression coefficients derived from a stepwise (rank) regression (Section 5.6).

In the following, variable Y represents an uncertain scalar computational result and (y_1, y_2, \dots, y_n) denotes a sample of values (i.e. realizations) provided for variable Y via n (e.g. $n = 100$) computer code runs. Input to each computer code run i is a vector $\underline{x}_i = (x_{i1}, x_{i2}, \dots, x_{i npar})$, $i = 1, \dots, n$, sampled for the $npar$ uncertain input parameters $X_1, X_2, \dots, X_{npar}$.

5.1 Correlation based sensitivity indices

The following types of correlation based sensitivity indices are implemented in SUSANA:

- Pearson's correlation (Section 5.1.1)
- Spearman's rank correlation (Section 5.1.2)
- Blomqvist's medial correlation (Section 5.1.3)
- Kendall's rank correlation (Section 5.1.4)

For each correlation type, statistical estimators of the ordinary and partial correlation coefficient as well as of the standardized regression coefficient can be calculated. An estimator of the coefficient of determination is additionally provided to inform on the usefulness of these correlation and regression coefficients as sensitivity indices /HOF 99/, /KLO 12/.

Since the ordinary and partial correlation coefficient and the standardized regression coefficient apply only to individual parameters, the multiple correlation coefficient is calculated, when the sensitivity indices shall be applied to parameter groups. This coefficient can be provided for each correlation type.

5.1.1 Pearson's correlation

5.1.1.1 Ordinary correlation coefficient

Pearson's ordinary correlation coefficient $CC_P(X_j, Y)$ between uncertain parameter X_j , $j = 1, \dots, npar$, and the computational result Y is defined as

$$CC_P(X_j, Y) = \rho(X_j, Y) = \frac{Cov(X_j, Y)}{\sqrt{Var(X_j) \cdot Var(Y)}} \quad (5.1)$$

with the variances $Var(X_j)$ and $Var(Y)$ and the covariance $Cov(X_j, Y)$ of X_j and Y .

Pearson's ordinary correlation coefficient $CC_P(X_j, Y)$ has the following well-known properties:

- $-1 \leq CC_P(X_j, Y) \leq 1$
- X_j and Y are independent $\rightarrow CC_P(X_j, Y) = 0$
- $CC_P(X_j, Y)$ measures the degree of linear dependency between X_j and Y
- $|CC_P(X_j, Y)| = 1 \leftrightarrow$ complete linear dependency between X_j and Y
- $CC_P(X_j, Y)$ is not invariant under monotone transformations, i.e. $CC_P(X_j, Y)$ is not ordinally invariant.

$CC_P(X_j, Y)$ may not necessarily reveal the true degree of sensitivity of Y with respect to the individual parameter X_j . This may happen, if Y is affected by uncertain parameters correlated with X_j . In this case, $CC_P(X_j, Y)$ quantifies the degree of sensitivity of Y with respect to the individual parameter X_j plus the contributions of the parameters correlated with X_j .

$CC_P^2(X_j, Y)$ represents the fraction of the variability of Y explained by X_j , if

- the functional relationship between X_j and Y can be approximated by a linear function
- X_j is not correlated with another uncertain parameter $X_k, k \neq j$.

Based on the sample (y_1, y_2, \dots, y_n) of the computational result Y and the corresponding sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ of the uncertain parameter $X_j; j = 1, \dots, npar$, $CC_P(X_j, Y)$ is estimated as

$$\widehat{CC}_P(X_j, Y) = \frac{\sum_{i=1}^n (x_{ij} - \bar{x}_{.j}) (y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_{.j})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (5.2)$$

with $\bar{x}_{.j}$ and \bar{y} representing the empirical means (Section 4.1) derived from $(x_{1j}, x_{2j}, \dots, x_{nj})$ and (y_1, y_2, \dots, y_n) , respectively.

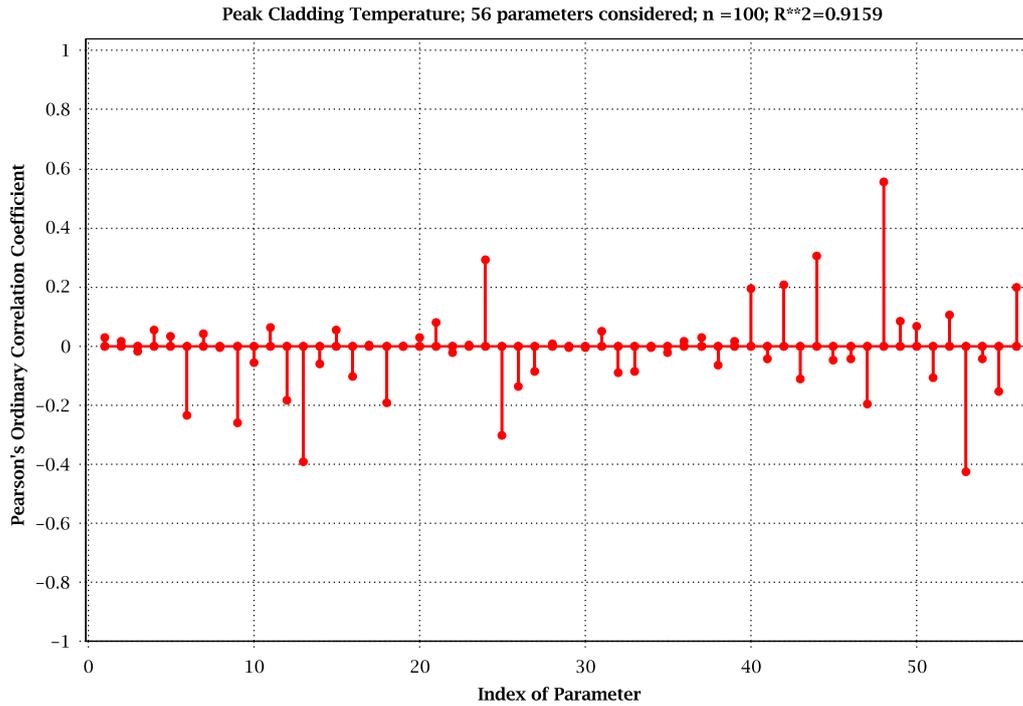


Fig. 5.1 Pearson's ordinary correlation coefficients as sensitivity indices of a scalar computational result (Peak Cladding Temperature) with respect to 56-uncertain input parameters

5.1.1.2 Partial correlation coefficient

Let $\hat{Y}_{|k \neq j}$ denote the linear regression of Y on all uncertain parameters but parameter X_j and let \hat{X}_j denote the linear regression of X_j on all uncertain parameters different from $X_j, j = 1, \dots, npar, \text{ i.e.}:$

$$\hat{Y}_{|k \neq j} = \sum_{k \neq j} a_k^{(j)} X_k \quad (5.3)$$

$$\hat{X}_j = \sum_{k \neq j} b_k^{(j)} X_k \quad (5.4)$$

The $a_k^{(j)}$ -values in Eq. (5.3) and the $b_k^{(j)}$ values in Eq. (5.4) with $k = 1, \dots, npar, k \neq j$, are the regression coefficients which minimize the corresponding mean square errors.

If $R_{\hat{Y}_{|k \neq j}}$ and $R_{\hat{X}_j}$ denote the residuals $(Y - \hat{Y}_{|k \neq j})$ and $(X_j - \hat{X}_j)$, respectively, then the partial correlation coefficient $PCC_P(X_j, Y)$ relating to Pearson's correlation between parameter X_j , $j = 1, \dots, npar$, and computational result Y is defined as the ordinary correlation coefficient (Eq. (5.1)) between these residuals, i.e.:

$$PCC_P(X_j, Y) = \rho(R_{\hat{X}_j}, R_{\hat{Y}_{|k \neq j}}) \quad (5.5)$$

$PCC_P(X_j, Y)$ measures the degree of linear association between X_j and Y after having removed, from both variables, all linear effects of the uncertain parameters different from X_j .

From Eq. (5.5) it can be concluded that $-1 \leq PCC_P(X_j, Y) \leq 1$.

Since the partial correlation coefficient $PCC_P(X_j, Y)$ measures the degree of linear dependency between the residuals $R_{\hat{Y}_{|k \neq j}}$ and $R_{\hat{X}_j}$, it may give misleading information. A large $PCC_P(X_j, Y)$ does not mean a high sensitivity of Y with respect to X_j but a high sensitivity of $R_{\hat{Y}_{|k \neq j}}$ with respect to $R_{\hat{X}_j}$. That means $PCC_P(X_j, Y)$ may be large, even if only a small residual $R_{\hat{Y}_{|k \neq j}}$ of Y is well explained by the residual $R_{\hat{X}_j}$ of X_j .

The partial correlation coefficient is calculated from the inverse CC_P^{-1} of Pearson's (population) correlation matrix CC_P of the entire parameter vector $(X_1, X_2, \dots, X_{npar})$ plus the corresponding computational result Y . CC_P and CC_P^{-1} are defined as follows:

$$CC_P = \begin{pmatrix} CC_P(X_1, X_1) \dots & CC_P(X_1, X_{npar}) & CC_P(X_1, Y) \\ \vdots & \vdots & \vdots \\ CC_P(X_{npar}, X_1) \dots & CC_P(X_{npar}, X_{npar}) & CC_P(X_{npar}, Y) \\ CC_P(Y, X_1) \dots & CC_P(Y, X_{npar}) & CC_P(Y, Y) \end{pmatrix} \quad (5.6)$$

$$CC_P^{-1} = \begin{pmatrix} IC_P(X_1, X_1) \dots & IC_P(X_1, X_{npar}) & IC_P(X_1, Y) \\ \vdots & \vdots & \vdots \\ IC_P(X_{npar}, X_1) \dots & IC_P(X_{npar}, X_{npar}) & IC_P(X_{npar}, Y) \\ IC_P(Y, X_1) \dots & IC_P(Y, X_{npar}) & IC_P(Y, Y) \end{pmatrix} \quad (5.7)$$

Based on the definition in Eq. (5.7), the partial correlation coefficient $PCC_P(X_j, Y)$ between X_j , $j = 1, \dots, npar$, and Y is calculated as:

$$PCC_p(X_j, Y) = -\frac{IC_p(X_j, Y)}{(IC_p(X_j, X_j) \cdot IC_p(Y, Y))^{1/2}} \quad (5.8)$$

The estimator $\widehat{PCC}_p(X_j, Y)$ of the partial correlation coefficient is derived from the inverse of the sample correlation matrix which includes the corresponding estimators $\widehat{CC}_p(\cdot, \cdot)$ calculated according to Eq. (5.2) instead of Pearson's population correlation coefficients. To be able to calculate $\widehat{PCC}_p(X_j, Y)$, the sample size n must exceed $npar + 1$. Otherwise, the sample correlation matrix is not positive definite and, therefore, its inverse cannot be derived. For sample sizes n close to $(npar + 1)$, the accuracy of $\widehat{PCC}_p(X_j, Y)$ is not satisfactory.

5.1.1.3 Standardized regression coefficient

Let X'_j and Y' denote the standardized version of the uncertain parameter X_j , $j = 1, \dots, npar$, and the computational result Y . For instance, X'_j is defined as:

$$X'_j = \frac{X_j - E(X_j)}{\sqrt{Var(X_j)}} \quad (5.9)$$

with the expectation $E(X_j)$ and the variance $Var(X_j)$ of X_j .

If \widehat{Y}' denotes the linear regression of the standardized variable Y' on all standardized parameters $X'_1, X'_2, \dots, X'_{npar}$ (Eq. (5.10)), then the standardized regression coefficient $SRC_p(X_j, Y)$ between uncertain parameter X_j , $j = 1, \dots, npar$, and the computational result Y is given by the regression coefficient a_j associated with X'_j :

$$\widehat{Y}' = \sum_{j=1}^{npar} a_j X'_j \quad (5.10)$$

$SRC_p(X_j, Y)$ indicates how many standard deviation changes of Y correspond to one standard deviation change of X_j , all other $X_k, k \neq j$, kept constant.

Based on the specific elements $IC_P(X_j, Y)$ and $IC_P(Y, Y)$ of the inverse of the population correlation matrix (Eq. (5.7)), the standardized regression coefficient $SRC_P(X_j, Y)$ between $X_j, j = 1, \dots, npar$, and Y can be calculated as

$$SRC_P(X_j, Y) = -\frac{IC_P(X_j, Y)}{IC_P(Y, Y)} \quad (5.11)$$

The estimator $\widehat{SRC}_P(X_j, Y)$ of the standardized regression coefficient is derived from the inverse of the sample correlation matrix instead of the population correlation matrix. To be able to calculate $\widehat{SRC}_P(X_j, Y)$, the sample size n must exceed $npar + 1$ (see Section 5.1.1.2 for more information).

It can be shown that

$$SRC_P(X_j, Y) = CC_P(X_j, Y) - \sum_{k \neq j} SRC_P(X_k, Y) \cdot CC_P(X_k, Y) \quad (5.12)$$

Comparing Eq. (5.8) and (5.11), it can be concluded that

$$SRC_P(X_j, Y) = PCC_P(X_j, Y) \cdot \sqrt{\frac{IC_P(Y, Y)}{IC_P(X_j, X_j)}} \quad (5.13)$$

Eq. (5.13) indicates that $SRC_P(X_j, Y)$ is not restricted to values between -1 and 1.

$$|SRC_P(X_j, Y)| > 1, \text{ if } \frac{IC_P(Y, Y)}{IC_P(X_j, X_j)} > 1 \text{ and } PCC_P(X_j, Y) > \sqrt{\frac{IC_P(X_j, X_j)}{IC_P(Y, Y)}}.$$

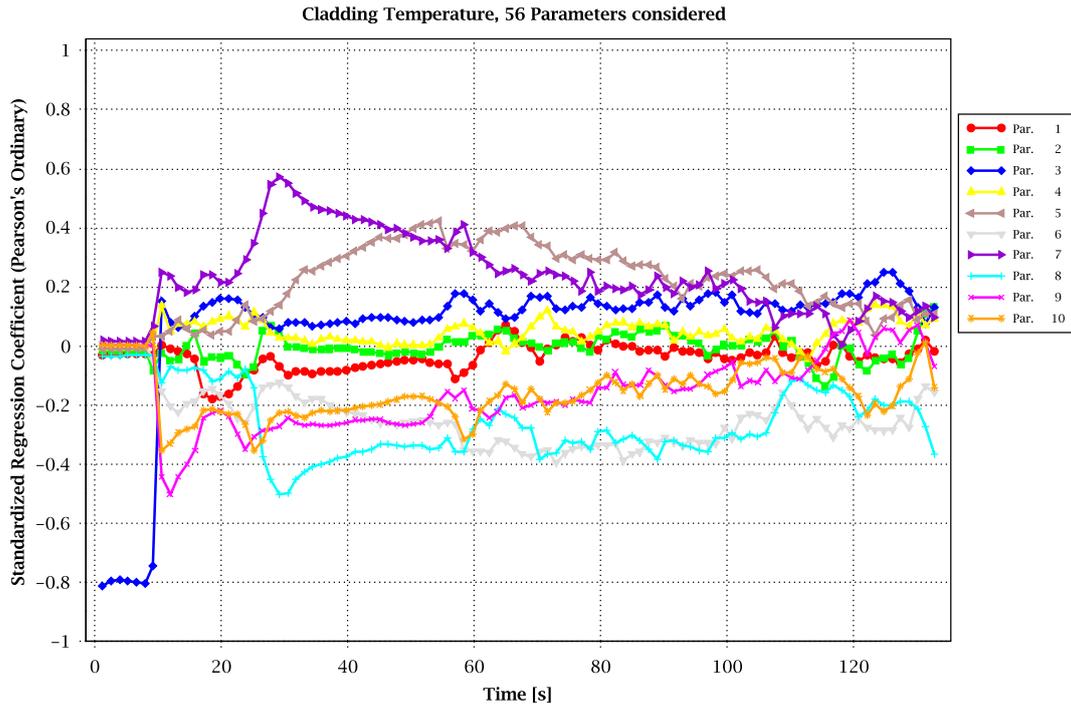


Fig. 5.2 Standardized regression coefficients (with respect to Pearson's ordinary correlation) as sensitivity indices of a time-dependent computational result (Cladding Temperature) with respect to 56 uncertain parameters

5.1.1.4 Coefficient of determination

The coefficient of determination R_p^2 is defined as the square of the population correlation coefficient between the code result Y and the variable \hat{Y} obtained from linear regression of Y on the parameters $X_1, X_2, \dots, X_{npar}$:

$$R_p^2 = \rho^2(Y, \hat{Y}) \quad (5.14)$$

It can be shown that R_p^2 is the proportion of the total variation of code result Y explained by the overall influence of the uncertain parameters $X_1, X_2, \dots, X_{npar}$ as modelled by a linear regression of Y on the parameters, i.e.:

$$R_p^2 = \frac{Var(\hat{Y})}{Var(Y)} \quad (5.15)$$

R_p^2 is useful to assess the quality of the ordinary and partial correlation coefficient as well as the standardized regression coefficient as sensitivity indices.

R_p^2 ranges between 0 and 1.

If $IC_p(Y, Y)$ denotes the last element of the inverse matrix CC_p^{-1} of the correlation matrix (Eq. (5.7)), R_p^2 can be calculated as

$$R_p^2 = 1 - \frac{1}{IC_p(Y, Y)} \quad (5.16)$$

The estimator \hat{R}_p^2 of the coefficient of determination is derived from the inverse of the sample correlation matrix instead of the population correlation matrix. To be able to calculate \hat{R}_p^2 , the sample size n must exceed $n_{par} + 1$ (see Section 5.1.1.2 for more information).

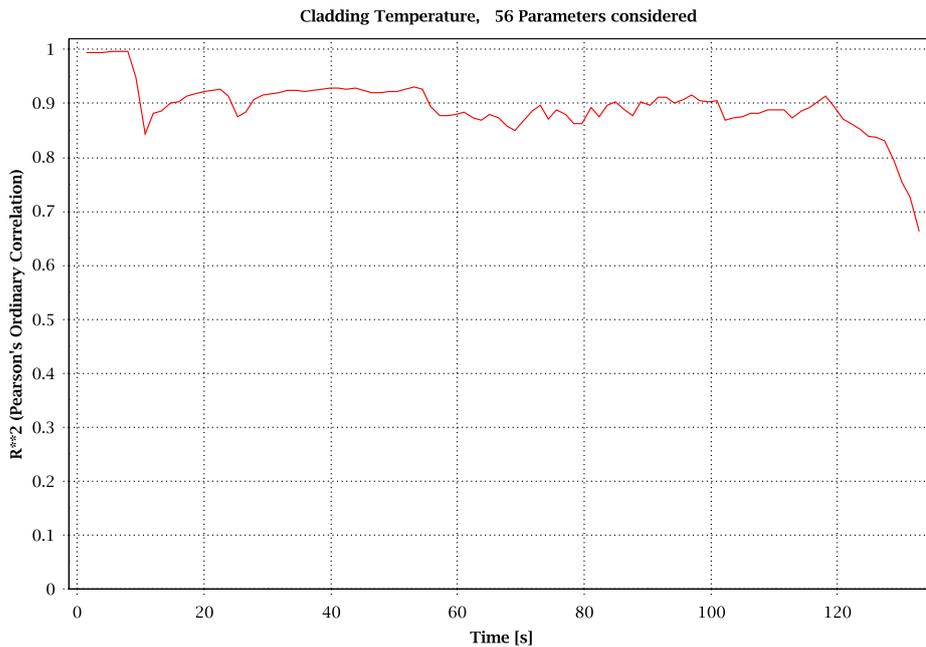


Fig. 5.3 Coefficients of determination with respect to Pearson's ordinary correlation for a time-dependent computational result (Cladding Temperature) influenced by 56 uncertain parameters

5.1.2 Spearman's rank correlation

Spearman's rank correlation coefficient $CC_S(X_j, Y)$ between uncertain parameter X_j with $j = 1, \dots, npar$ and the computational result Y is defined as

$$CC_S(X_j, Y) = \rho(F_{X_j}, F_Y) \quad (5.17)$$

where F_{X_j} and F_Y denote the cumulative distribution functions of X_j and Y , respectively, and ρ denotes Pearson's ordinary correlation coefficient defined in Eq. (5.1).

Eq. (5.17) indicates that Spearman's rank correlation coefficient between the two variables X_j and Y is equivalent to Pearson's ordinary correlation coefficient applied on the distribution functions of the two variables.

Spearman's rank correlation coefficient $CC_S(X_j, Y)$ has the following well-known properties:

- $-1 \leq CC_S(X_j, Y) \leq 1$
- X_j and Y are independent $\rightarrow CC_S(X_j, Y) = 0$
- $CC_S(X_j, Y)$ only detects monotonic dependence structures and measures the degree of monotonic dependency between X_j and Y
- $CC_S(X_j, Y)$ indicates whether upper (lower) quantiles of X_j lead in tendency to upper (lower) quantiles of Y
- $|CC_S(X_j, Y)| = 1 \leftrightarrow$ complete dependency between X_j and Y
- $CC_S(X_j, Y)$ is ordinally invariant
- $CC_S(X_j, Y)$ is not affected by outliers

Based on the sample (y_1, y_2, \dots, y_n) of the computational result Y and the corresponding sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ of the uncertain parameter X_j with $j = 1, \dots, npar$, $CC_S(X_j, Y)$ is estimated as follows

$$\begin{aligned} \widehat{CC}_S(X_j, Y) &= \frac{\sum_{i=1}^n (r(x_{ij}) - \bar{r}_{x_j}) (r(y_i) - \bar{r}_y)}{\sqrt{\sum_{i=1}^n (r(x_{ij}) - \bar{r}_{x_j})^2 \cdot \sum_{i=1}^n (r(y_i) - \bar{r}_y)^2}} \\ &= \frac{\sum_{i=1}^n (r(x_{ij}) - \frac{n+1}{2}) (r(y_i) - \frac{n+1}{2})}{\frac{n(n^2-1)}{12}} \end{aligned} \quad (5.18)$$

with $r(\cdot)$ representing the rank of the corresponding value in a sample of size n (see Eq. (2.67), and \bar{r}_{x_j} and \bar{r}_y representing the empirical means (see Section 4.1) derived from $(r(x_{1j}), r(x_{2j}), \dots, r(x_{nj}))$ and $(r(y_1), r(y_2), \dots, r(y_n))$, respectively.

5.1.3 Blomqvist's medial correlation

Blomqvist's medial correlation coefficient $CC_B(X_j, Y)$ between uncertain parameter X_j , $j = 1, \dots, npar$, and the computational result Y is defined as

$$CC_B(X_j, Y) = \rho(\text{sgn}(X_j - M_{X_j}), \text{sgn}(Y - M_Y)) \quad (5.19)$$

where M_{X_j} and M_Y are the medians of the distributions of X_j and Y , respectively. ρ denotes Pearson's ordinary correlation coefficient defined in Eq. (5.2) and sgn denotes the signum function defined as:

$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases} \quad (5.20)$$

Eq. (5.19) indicates that Blomqvist's medial correlation coefficient between the two variables X_j and Y is equivalent to Pearson's ordinary correlation coefficient applied on the transformations $\text{sgn}(X_j - M_{X_j})$ and $\text{sgn}(Y - M_Y)$.

Blomqvist's medial correlation coefficient $CC_B(X_j, Y)$ has the following well-known properties:

- $-1 \leq CC_B(X_j, Y) \leq 1$
- X_j and Y are independent $\rightarrow CC_B(X_j, Y) = 0$
- $CC_B(X_j, Y)$ provides the difference between the probabilities of concordance and discordance of X_j and Y relative to the corresponding medians M_{X_j} and M_Y (Fig. 2.19)
- $|CC_B(X_j, Y)| = 1 \leftrightarrow$ complete concordance or complete discordance between X_j and Y relative to the corresponding medians M_{X_j} and M_Y
- $CC_B(X_j, Y)$ is ordinally invariant
- $CC_B(X_j, Y)$ is not affected by outliers

Based on the sample (y_1, y_2, \dots, y_n) of the computational result Y and the corresponding sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ of the uncertain parameter X_j with $j = 1, \dots, npar$, $CC_B(X_j, Y)$ is estimated as follows

$$\widehat{CC}_B(X_j, Y) = \frac{1}{n} \sum_{i=1}^n \text{sgn}(x_{ij} - m_{x_j}) \text{sgn}(y_i - m_y) \quad (5.21)$$

with m_{x_j} and m_y representing the empirical medians (see Section 4.1) derived from $(x_{1j}, x_{2j}, \dots, x_{nj})$ and (y_1, y_2, \dots, y_n) , respectively.

5.1.4 Kendall's rank correlation

Kendall's rank correlation coefficient $CC_K(X_j, Y)$ between the uncertain parameter X_j , $j = 1, \dots, npar$, and the computational result Y is defined as

$$CC_K(X_j, Y) = \rho(\text{sgn}(X_j - X'_j), \text{sgn}(Y - Y')) \quad (5.22)$$

where (X'_j, Y') is another pair of variables distributed like the pair (X_j, Y) . ρ is defined in Eq. (5.2). sgn denotes the signum function defined in Eq. (5.6).

Eq. (5.22) indicates that Kendall's medial correlation coefficient between the two variables X_j and Y is equivalent to Pearson's ordinary correlation coefficient applied on the the transformations $\text{sgn}(X_j - X'_j)$ and $\text{sgn}(Y - Y')$.

Kendall's rank correlation coefficient $CC_K(X_j, Y)$ has the following well-known properties:

- $-1 \leq CC_K(X_j, Y) \leq 1$
- X_j and Y are independent $\rightarrow CC_K(X_j, Y) = 0$
- $CC_K(X_j, Y)$ provides the difference between the probabilities of concordance and discordance of X_j and Y relative to the variables X'_j and Y' , respectively, with the pair (X'_j, Y') distributed like the pair (X_j, Y)
- $|CC_K(X_j, Y)| = 1 \leftrightarrow$ complete concordance or complete discordance between X_j and Y relative to the variables X'_j and Y'
- $CC_K(X_j, Y)$ is ordinally invariant
- $CC_K(X_j, Y)$ is not affected by outliers

Based on the sample (y_1, y_2, \dots, y_n) of the computational result Y and the corresponding sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ of the uncertain parameter X_j with $j = 1, \dots, npar$, $CC_K(X_j, Y)$ is estimated as follows

$$\widehat{CC}_K(X_j, Y) = \frac{1}{\binom{n}{2}} \sum_{i=1}^{n-1} \sum_{k=i+1}^n \text{sgn}(x_{ij} - x_{kj}) \text{sgn}(y_i - y_k) \quad (5.23)$$

5.1.5 Partial correlation coefficient, standardized regression coefficient and coefficient of determination relating to Spearman's, Blomqvist's and Kendall's correlations

The correlation coefficients of Spearman (Eq. (5.17)), Blomqvist (Eq. (5.19)) and Kendall (Eq. (5.22)) can be formulated as Pearson's ordinary correlation coefficient applied on appropriate transformations of the considered uncertain parameter X_j , $j = 1, \dots, npar$, and the computational result Y . That means the matrix CC consisting of the population correlation coefficients of any of the three correlation types is equivalent to the ordinary population correlation matrix CC_P of appropriately transformed variables X_j and Y (Eq. (5.6)). Therefore, the partial correlation coefficient $PCC(X_j, Y)$, the standardized regression coefficient $SRC(X_j, Y)$ and the coefficient of determination R^2 relating to any of the three correlation types can be derived from elements of the respective inverse matrix CC^{-1} . Eqs. (5.8), (5.11), and (5.16) give the basic formulas used to calculate $PCC(X_j, Y)$, $SRC(X_j, Y)$ and R^2 . Of course, the expressions in the equations which refer to Pearson's population correlation matrix must be appropriately replaced by the corresponding expressions referring either to Spearman's, Blomqvist's or to Kendall's population correlation matrix.

$PCC(X_j, Y)$ and $SRC(X_j, Y)$ with respect to Spearman's, Blomqvist's or Kendall's correlation are ordinally invariant.

The estimators $\widehat{PCC}(X_j, Y)$, $\widehat{SRC}(X_j, Y)$ and \widehat{R}^2 relating to any of the correlation types are derived from the inverse of the respective sample correlation matrix. Dependent on the correlation type, this sample correlation matrix includes the corresponding estimators of the (population) correlation coefficients (Eqs. (5.18), (5.21), and (5.23)).

To be able to calculate $\widehat{PCC}(X_j, Y)$, $\widehat{SRC}(X_j, Y)$ and \widehat{R}^2 for any of the correlation types, the sample size n must exceed $n_{par} + 1$. Otherwise, the sample correlation matrix is not positive definite and, therefore, its inverse cannot be derived.

5.2 Multiple correlation coefficients

5.2.1 Pearson's multiple correlation

The multiple correlation coefficient $R_p^2(X^{(G)}, Y)$ between the group $X^{(G)}$ of n_G uncertain parameters - without loss of generality $X^{(G)} = X_1, X_2, \dots, X_{n_G}$ - and the computational result Y is defined as

$$R_p^2(X^{(G)}, Y) = (\rho(X_1, Y), \dots, \rho(X_{n_G}, Y)) \cdot CC_{P_{X_{n_G}}}^{-1} \cdot \begin{pmatrix} \rho(X_1, Y) \\ \vdots \\ \rho(X_{n_G}, Y) \end{pmatrix} \quad (5.24)$$

where ρ denotes Pearson's ordinary correlation coefficient (Eq. (5.1)) and $CC_{P_{X_{n_G}}}^{-1}$ denotes the inverse of the ordinary $n_G \times n_G$ correlation matrix of the parameters X_1, X_2, \dots, X_{n_G} of the group $X^{(G)}$.

$R_p^2(X^{(G)}, Y)$ has the following well-known properties:

- $0 \leq R_p^2(X^{(G)}, Y) \leq 1$
- Each parameter X_j , $j = 1, \dots, n_G$, of the group $X^{(G)}$ and Y are independent $\rightarrow R_p^2(X^{(G)}, Y) = 0$
- $R_p^2(X^{(G)}, Y) = 1 \leftrightarrow$ complete linear dependency of Y from $X^{(G)} = X_1, X_2, \dots, X_{n_G}$
- $CC_p(X_j, Y)$ measures the degree of multiple linear dependency between Y and the parameters of the group $X^{(G)}$.
- If the parameter group $X^{(G)}$ includes all parameters $X_1, X_2, \dots, X_{n_{par}}$, then the multiple correlation coefficient $R_p^2(X^{(G)}, Y)$ corresponds to the coefficient of determination R_p^2 (Section 5.1.1.4), i.e. $X^{(G)} = X_1, X_2, \dots, X_{n_{par}} \rightarrow R_p^2(X^{(G)}, Y) = R_p^2$.

The estimator $\hat{R}_P^2(X^{(G)}, Y)$ of the multiple correlation coefficient between the group $X^{(G)}$ of uncertain parameters and the computational result Y is obtained by considering in Eq. (5.24)

- the sample correlation coefficients $\widehat{CC}_P(X_1, Y), \dots, \widehat{CC}_P(X_{n_G}, Y)$ according to Pearson (Eq. (5.2)) instead of Pearson's ordinary correlation coefficients $\rho(X_1, Y), \dots, \rho(X_{n_G}, Y)$
- the inverse of Pearson's $n_G \times n_G$ sample correlation matrix of the parameters X_1, X_2, \dots, X_{n_G} of the group $X^{(G)}$ instead of the inverse of Pearson's correlation matrix $CC_{P_{X_{n_G}}}^{-1}$

To be able to calculate $\hat{R}_P^2(X^{(G)}, Y)$, the sample size n must exceed the number n_G of parameters of the group $X^{(G)}$. Otherwise, the $n_G \times n_G$ sample correlation matrix is not positive definite and, therefore, its inverse cannot be derived.

5.2.2 Spearman's, Blomqvist's and Kendall's multiple correlations

The correlation coefficients of Spearman (Eq. (5.17)), Blomqvist (Eq. (5.19)) and Kendall (Eq. (5.22)) can be formulated as Pearson's ordinary correlation coefficient applied on appropriate transformations of the considered uncertain parameters X_j and the computational result Y . So, the formula of the multiple correlation coefficient $R^2(X^{(G)}, Y)$ according to any of the three correlation types is obtained, if in Eq. (5.24) Pearson's ordinary correlation coefficient is applied on the appropriately transformed variables. Also the $n_G \times n_G$ correlation matrix in Eq. (5.24) must be applied on the respective transformations of the parameters of the group $X^{(G)}$.

The estimator $\hat{R}^2(X^{(G)}, Y)$ relating to any of the three correlation types is derived by considering in Eq. (5.24)

- the sample correlation coefficients $\widehat{CC}(X_1, Y), \dots, \widehat{CC}(X_{n_G}, Y)$ according to Spearman, Blomqvist or Kendall (Eqs. (5.18), (5.21), or (5.23)) instead of Pearson's ordinary correlation coefficients $\rho(X_1, Y), \dots, \rho(X_{n_G}, Y)$
- the inverse of the $n_G \times n_G$ sample correlation matrix (according to Spearman, Blomqvist or Kendall) of the parameters of the group $X^{(G)}$ instead of the inverse of Pearson's correlation matrix $CC_{P_{X_{n_G}}}^{-1}$.

To be able to calculate $\hat{R}^2(X^{(G)}, Y)$ for any of the correlation types, the sample size n must exceed $npar + 1$.

5.3 Correlation ratio

The correlation ratio $CR(X_j, Y)$ between uncertain parameter $X_j, j = 1, \dots, npar$, and the computational result Y is defined as

$$CR(X_j, Y) = \sqrt{\frac{Var(E(Y|X_j))}{Var(Y)}} \quad (5.25)$$

with $Var(\)$ denoting the variance of a variable and $E(Y|X_j)$ denoting the conditional expectation of Y conditioned on X_j .

$CR(X_j, Y)$ is based on the following well-known variance decomposition /MCK 96/:

$$Var(Y) = E(Var(Y|X_j)) + Var(E(Y|X_j)) \quad (5.26)$$

where $E(\)$ and $Var(\)$ denote the expectation and variance of a variable; $E(Y|X_j)$ and $Var(Y|X_j)$ are the conditional expectation and variance of Y conditioned on X_j .

Since

$$CR^2(X_j, Y) = \frac{Var(Y) - E(Var(Y|X_j))}{Var(Y)} \quad (5.27)$$

$CR(X_j, Y)$ is an indicator of the expected reduction in the variance of Y , if X_j could be fixed.

It can be shown that

$$CR(X_j, Y) = \rho(Y, E(Y|X_j)) \quad (5.28)$$

That means $CR(X_j, Y)$ is equivalent to Pearson's ordinary correlation coefficient between Y and its conditional expectation $E(Y|X_j)$ conditioned on X_j .

$CR(X_j, Y)$ has the following properties:

- $0 \leq CR(X_j, Y) \leq 1$
- X_j and Y are independent $\rightarrow CR(X_j, Y) = 0$
- $CR(X_j, Y) = 1 \leftrightarrow Y$ is a function of X_j

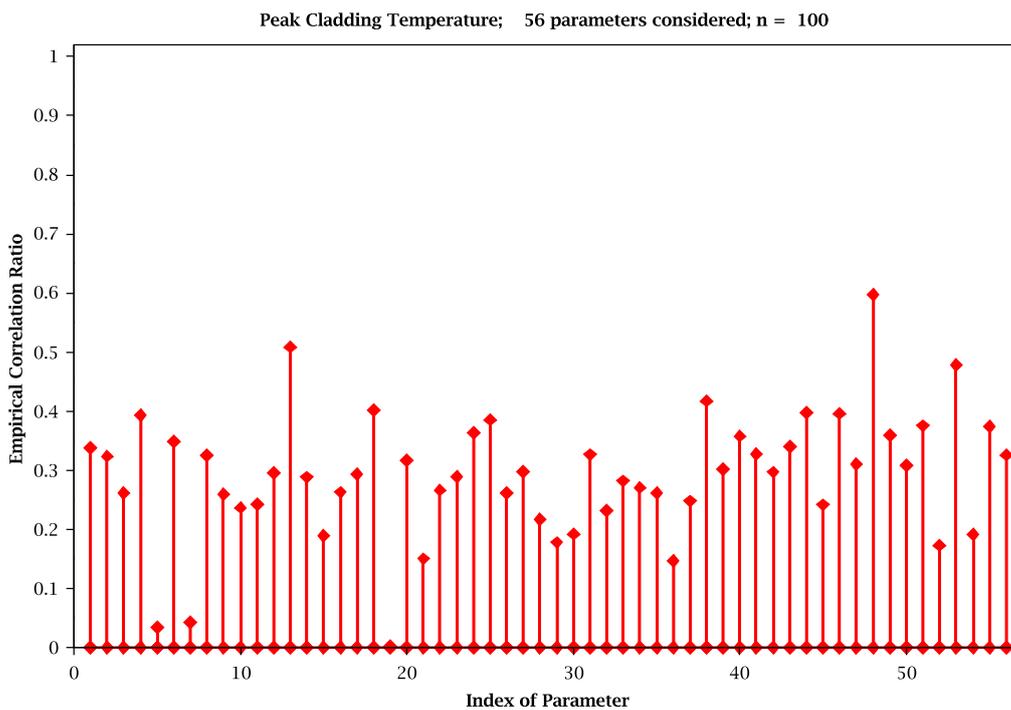


Fig. 5.4 Correlation ratios as sensitivity indices of a scalar computational result (Peak Cladding Temperature) with respect to 56 uncertain input parameters

Following procedure is applied to estimate $CR(X_j, Y)$ from the sample (y_1, y_2, \dots, y_n) of the computational result Y and the sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ of the uncertain parameter X_j ; $j = 1, \dots, npar$ /KEN 73/:

- The sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ is ordered by increasing size.
- The range of the ordered sample $(x_{1:nj}, x_{2:nj}, \dots, x_{n:nj})$ is divided into $n_X = \lfloor \sqrt{n} \rfloor$ subsets where each subset $I_k, k = 1, \dots, n_X$, consists of at least $n_k = n_X$ successive values of X_j ($\lfloor \sqrt{n} \rfloor$ means the greatest integer smaller than or equal to \sqrt{n}).
- The mean \bar{y} of all y -values and the means $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_{n_X}$ of the y -values corresponding to the x_j -values of the subsets I_1, I_2, \dots, I_{n_X} , respectively, are calculated.
- The estimator of the correlation ratio $CR(X_j, Y)$ is determined as

$$\widehat{CR}(X_j, Y) = \sqrt{\frac{\sum_{k=1}^{n_X} n_k \cdot (\bar{y}_k - \bar{y})^2}{\sum_{k=1}^{n_X} \sum_{l=1}^{n_k} (y_{kl} - \bar{y})^2}} \quad (5.29)$$

$\widehat{CR}(X_j, Y)$ can be determined, even if the number of parameters exceeds the sample size.

5.3.1 Correlation ratio on ranks

Instead on the original variables $X_j, j = 1, \dots, npar$, and Y , the correlation ratio on ranks $CR_R(X_j, Y)$ is applied on the corresponding distribution functions F_{X_j} and F_Y , i.e.:

$$CR_R(X_j, Y) = \sqrt{\frac{Var\left(E\left(F_Y \mid F_{X_j}\right)\right)}{Var\left(F_{X_j}\right)}} \quad (5.30)$$

Following procedure is applied to estimate $CR_R(X_j, Y)$ from the sample (y_1, y_2, \dots, y_n) of the computational result Y and the sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ of the uncertain parameter X_j ; $j = 1, \dots, npar$:

- The sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ is ordered by increasing size.

- The range of the ordered sample $(x_{1:n_j}, x_{2:n_j}, \dots, x_{n:n_j})$ is divided into $n_X = \lfloor \sqrt{n} \rfloor$ subsets where each subset $I_k, k = 1, \dots, n_X$, consists of n_k successive values of X_j
- The ranks of the y -values corresponding to the ordered sample $(x_{1:n_j}, x_{2:n_j}, \dots, x_{n:n_j})$ are determined and appropriately assigned to the subsets I_1, I_2, \dots, I_{n_X} . For instance, $r_{kl}(y)$ is the rank of the y -value corresponding to the l^{th} largest x_j -value in subset I_k .
- The mean rank $\bar{r} = \frac{n(n+1)}{2}$ and the mean ranks $\bar{r}_1(y), \bar{r}_2(y), \dots, \bar{r}_{n_X}(y)$ of the y -values corresponding to the x_j -values in the subsets I_1, I_2, \dots, I_{n_X} , respectively, are calculated.
- The estimator of the correlation ratio $CR_R(X_j, Y)$ is determined as

$$\widehat{CR}_R(X_j, Y) = \sqrt{\frac{\sum_{k=1}^{n_X} n_k \cdot (\bar{r}_k(y) - \bar{r})^2}{\sum_{k=1}^{n_X} \sum_{l=1}^{n_k} (r_{kl}(y) - \bar{r})^2}} \quad (5.31)$$

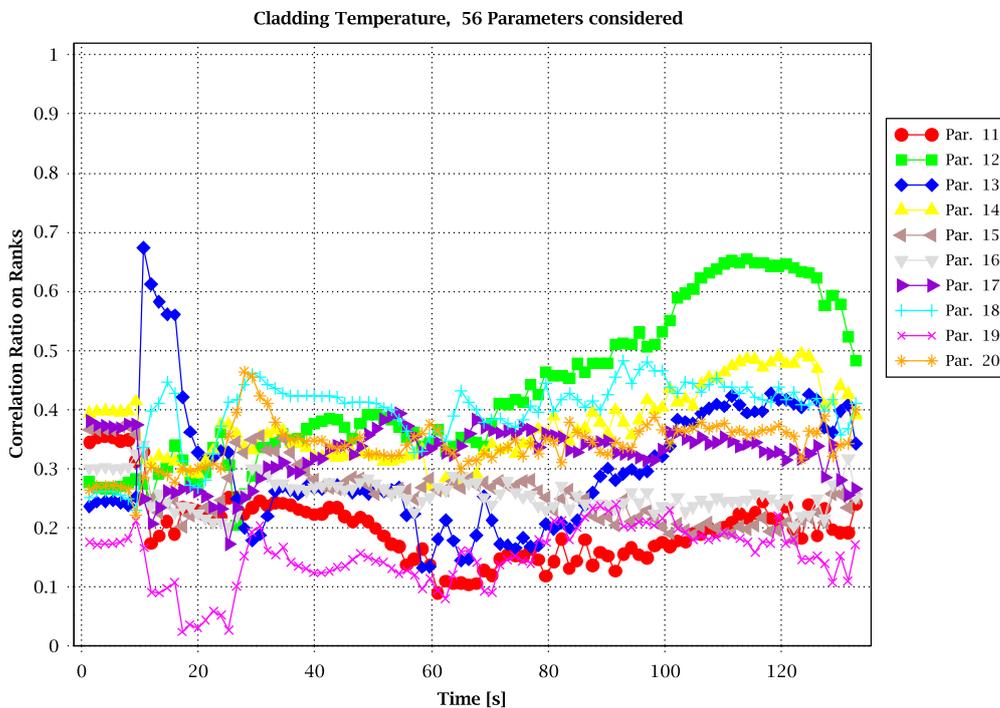


Fig. 5.5 Correlation ratios on ranks as sensitivity indices of a time-dependent computational result (Cladding Temperature) with respect to 56 uncertain parameters

5.4 Sobol indices

A well-known approach to calculate sensitivity indices offers the framework of the Sobol indices (SI) as variance-based sensitivity measures /SOB 99/. The SIs describe the sensitivity patterns of a model via the full decomposition of the variance of the model response into terms depending on the model input parameters and their interactions.

Given the model $Y = f(X_1, X_2, \dots, X_{npar})$ with the uncertain input parameters X_1, \dots, X_{npar} , it can be shown that the variance $Var(Y)$ of the response Y can be decomposed as indicated in Eq. (5.26). The functional decomposition is based on the ANOVA framework and given in condensed form in e.g. /SAL 10/.

Sobol's first order sensitivity index or uncertainty importance measure with respect to parameter X_j is defined as

$$S_j = \frac{Var(E(Y|X_j))}{Var(Y)} = \frac{D_j}{D} \quad (5.32)$$

with the conditional expectation $E(Y|X_j)$ of Y conditioned on X_j , the partial variance D_j and the total variance $Var(Y) = D$. $Var(E(Y|X_j))$ can be interpreted as the expected reduction in the variance $Var(Y)$, if X_j could be fixed.

The square root of Sobol's first order sensitivity index is equivalent to the correlation ratio defined in Section 5.3.

S_j has the following properties:

- $0 \leq S_j \leq 1$
- A high S_j value indicates that parameter X_j strongly influences the variance $Var(Y)$.

Sobol's total sensitivity index (or total effect index) with respect to parameter X_j is defined as

$$\begin{aligned}
S_{Tj} &= \frac{E\left(\text{Var}(Y|X_{\sim j})\right)}{\text{Var}(Y)} \\
&= 1 - \frac{\text{Var}\left(E(Y|X_{\sim j})\right)}{\text{Var}(Y)} \\
&= S_j + \sum_{1 \leq j < k \leq p} S_{jk} = \frac{D_j}{D} + \sum_{1 \leq j < k \leq p} \frac{D_{jk}}{D}
\end{aligned} \tag{5.33}$$

with $X_{\sim j}$ denoting all parameters but parameter X_j . $E\left(\text{Var}(Y|X_{\sim j})\right)$ can be interpreted as the expected remaining variance, if all parameters but parameter X_j could be fixed.

S_{Tj} quantifies the total effect of parameter X_j on the variance $\text{Var}(Y)$.

The computation of Sobol's first order sensitivity index S_j requires the analytical or numerical approximation of $\text{Var}\left(E(Y|X_j)\right)$. The implemented algorithm introduced in /SAL 02/ and based on the original approach of /SOB 99/, may basically be described via the following steps:

- Create a $2n \times n_{par}$ sample matrix and define two sample (reference) matrices A, B , and a (composite) matrix C_j , $j = 1, \dots, n_{par}$, based on reference matrix B :

$$A = \begin{pmatrix} x_1^{(1)} & \dots & x_{n_{par}}^{(1)} \\ \vdots & & \vdots \\ x_1^{(n)} & \dots & x_{n_{par}}^{(n)} \end{pmatrix} \quad B = \begin{pmatrix} x_1^{(n+1)} & \dots & x_{n_{par}}^{(n+1)} \\ \vdots & & \vdots \\ x_1^{(2n)} & \dots & x_{n_{par}}^{(2n)} \end{pmatrix} \tag{5.34}$$

$$C_j = B_A^{(j)} = \begin{pmatrix} x_1^{(n+1)} \dots x_j^{(1)} \dots x_k^{(n+1)} \\ \vdots & \vdots & \vdots \\ x_1^{(2n)} \dots x_j^{(n)} \dots x_k^{(2n)} \end{pmatrix}$$

The matrix C_j is identical to the matrix B with the exception that the j^{th} column including the values sampled for parameter X_j is taken from matrix A . Note, the formulation of the matrix triplet should be consistent, but can be arbitrarily chosen (i.e.

$C_j := A_B^{(j)}$) in case pure Monte Carlo samples are employed.

- Compute the n model responses $\vec{y}_A = f(A)$, $\vec{y}_B = f(B)$ and $\vec{y}_{C_j} = f(C_j)$ by evaluating the model at the input values from A, B , and $C_j, j = 1, \dots, npar$. These computations require a total of $n(npar + 2)$ simulation runs.
- Calculate the estimates for all terms of S_j and S_{Tj} , $j = 1, \dots, npar$ (with $\Sigma := \sum_{i=1}^n$):
 - Sobol's first order sensitivity index:

$$\hat{S}_j = \frac{\frac{1}{n} \sum y_{A,i} \cdot y_{C_j,i} - \frac{1}{n} \sum y_{A,i} \cdot \frac{1}{n} \sum y_{B,i}}{\frac{1}{n} \sum y_{A,i} \cdot y_{A,i} - \frac{1}{n} \sum y_{A,i} \cdot \frac{1}{n} \sum y_{B,i}} \quad (5.35)$$

- Sobol's total sensitivity index (or total effect index):

$$\hat{S}_{Tj} = 1 - \frac{\frac{1}{n} \sum y_{B,i} \cdot y_{C_j,i} - \frac{1}{n} \sum y_{A,i} \cdot \frac{1}{n} \sum y_{B,i}}{\frac{1}{n} \sum y_{A,i} \cdot y_{A,i} - \frac{1}{n} \sum y_{A,i} \cdot \frac{1}{n} \sum y_{B,i}} \quad (5.36)$$

Alternative estimates are based on Jansen's proposal /JAN 99/:

$$\hat{S}_j = 1 - \frac{\frac{1}{2n} \sum (y_{A,i} - y_{C_j,i})^2}{\frac{1}{n} \sum y_{A,i} \cdot y_{A,i} - \frac{1}{n} \sum y_{A,i} \cdot \frac{1}{n} \sum y_{B,i}} \quad (5.37)$$

- Sobol's total sensitivity index (or total effect index):

$$\hat{S}_{Tj} = \frac{\frac{1}{2n} \sum (y_{B,i} - y_{C_j,i})^2}{\frac{1}{n} \sum y_{A,i} \cdot y_{A,i} - \frac{1}{n} \sum y_{A,i} \cdot \frac{1}{n} \sum y_{B,i}} \quad (5.38)$$

5.5 Goodman/Kruskal coefficients from 2x2 contingency tables

The Goodman/Kruskal coefficient $\gamma(X_j, Y)$ /GOO 63/ is obtained as follows:

- The ranges of the parameter $X_j, j = 1, \dots, npar$, and of the computational result Y are divided into two disjoint intervals each:
 - $X_j: I_{j1} = (-\infty, a_j]$ and $I_{j2} = (a_j, +\infty)$.
 - $Y: I_{Y1} = (-\infty, a_Y]$ and $I_{Y2} = (a_Y, +\infty)$.
- The following 2x2 contingency table is built, and the probabilities p_{11}, \dots, p_{22} of the four table cells are calculated:

X_j/Y	I_{Y1}	I_{Y2}
I_{j1}	p_{11}	p_{12}
I_{j2}	p_{21}	p_{22}

- $\gamma(X_j, Y)$ is calculated as

$$\gamma(X_j, Y) = \frac{p_{11}p_{22} - p_{12}p_{21}}{p_{11}p_{22} + p_{12}p_{21}} \quad (5.39)$$

The Goodman/Kruskal coefficient $\gamma(X_j, Y)$ has the following properties:

- $-1 \leq \gamma(X_j, Y) \leq 1$
- X_j and Y are independent $\rightarrow \gamma(X_j, Y) = 0$
- Monotone increasing relationship $\rightarrow \gamma(X_j, Y) = +1$
- Monotone decreasing relationship $\rightarrow \gamma(X_j, Y) = -1$

Based on the sample (y_1, y_2, \dots, y_n) of the computational result Y and the corresponding sample $(x_{1j}, x_{2j}, \dots, x_{nj})$ of the uncertain parameter X_j ; $j = 1, \dots, npar$, $\gamma(X_j, Y)$ is estimated as

$$\hat{\gamma}(X_j, Y) = \frac{n_{11}n_{22} - n_{12}n_{21}}{n_{11}n_{22} + n_{12}n_{21}} \quad (5.40)$$

- n_{11} = number of data pairs (x_{ij}, y_i) in the cell $I_{j1} \times I_{Y1}$
- n_{12} = number of data pairs (x_{ij}, y_i) in the cell $I_{j1} \times I_{Y2}$
- n_{21} = number of data pairs (x_{ij}, y_i) in the cell $I_{j2} \times I_{Y1}$
- n_{22} = number of data pairs (x_{ij}, y_i) in the cell $I_{j2} \times I_{Y2}$

The estimator $\hat{\gamma}(X_j, Y)$, $j = 1, \dots, npar$, is calculated for three different 2x2 contingency tables. These tables are defined by the different values considered for a_j and a_Y which divide the ranges of X_j and Y into two disjoint intervals. In the first 2x2 contingency table, a_j and a_Y correspond to the lower sample quartile; in the second table, a_j and a_Y correspond to the sample median, and in the last table, a_j and a_Y correspond to the upper sample quartile of the respective variable.

5.6 Results from stepwise regression

To identify only those uncertain parameters which contribute most significantly to the uncertainty of the computational result Y , a forward stepwise regression may be performed. The corresponding regression algorithm implemented in SUSA successively inserts uncertain parameters into or removes them from the regression function according to the results of partial F-tests with the level of significance $\alpha = 0.05$. The user may control the parameter selection process by forcing parameters into the regression function or by excluding parameters from the analysis when the input of the sensitivity analysis is prepared. Parameters forced into the regression cannot be removed during the selection process.

The additional forward stepwise rank regression algorithm uses, instead of the original values, the rank transformed values $r(x_{j1}), r(x_{j2}), \dots, r(x_{jn})$ of each parameter

$X_j, j = 1, \dots, npar$, and the rank transformed values $r(y_1), r(y_2), \dots, r(y_n)$ of the result Y (see Eq. (2.67)).

For each uncertain parameter X_j identified as important via the ordinary forward stepwise regression, the estimators of the ordinary (Eq. (5.10)) and standardized (Eq. (5.11)) regression coefficient are provided. Additionally, the estimator of the coefficient of determination R_p^2 is calculated (Eq. (5.16)). All estimators are derived from the reduced ordinary correlation matrix considering only those parameters identified as most important (Eqs. (5.6) – (5.7)).

If the forward stepwise rank regression algorithm is applied, the estimators of the ordinary and standardized rank regression coefficient are calculated. Additionally, the estimator of the coefficient of determination is calculated. All estimators are derived from the reduced correlation matrix with respect to Spearman (Section 5.1.5).

For the regression model constructed by the ordinary forward stepwise regression (not rank!), *PRESS* (**p**redicted **r**esidual **e**rror **s**um of **s**quares) statistics can be calculated to get information on the stability of the model. They can help to identify data points influencing the model.

The $PRESS_k$ statistics, $k = 1, \dots, n$, is calculated as follows:

- Based on the regression model constructed by the ordinary forward stepwise regression, new regression coefficients are calculated from a reduced sample where the k^{th} sample element $\underline{x}_k = (x_{k1}, x_{k2}, \dots, x_{k npar})$ of the $npar$ uncertain input parameters $X_1, X_2, \dots, X_{npar}$ and the corresponding k^{th} sample element y_k of the computational result Y are not considered.
- The regression function with the new regression coefficients is applied on each parameter vector \underline{x}_i to calculate $\hat{y}_{ki}, i = 1, \dots, n$.
- $PRESS_k = \sum_i (y_i - \hat{y}_{ki})^2$ (*PRESS* value of the k^{th} sample element)

The total *PRESS* statistics is calculated as

$$PRESS = \sum_{k=1}^n PRESS_k \tag{5.41}$$

References

- /APO 81/ Apostolakis G., Kaplan S., Pitfalls in risk calculation. Reliability Engineering, 2, 135-145, 1981.
- /BED 01/ Bedford, T. and Cooke, R., Probabilistic Risk Analysis: Foundations and Methods, Cambridge University Press, 2001.
- /CLO 34/ Clopper, C.J., Pearson E.S., The use of confidence or fiducial limits illustrated in the case of the binomial, Biometrika, Vol. 26, 404-413, 1934.
- /COL 09/ Coleman, H. W., and Steele, W. G., Experimentation, validation, and uncertainty analysis for engineers, pp 31, 3rd edition, John Wiley & Sons, 2009.
- /EFR 86/ Efron, B. and R. Tibshirani, Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. Statistical science, 54-75, 1986.
- /FER 01/ Fernholz, L. T., and J. A. Gillespie, Content-corrected tolerance limits based on the bootstrap. Technometrics 43.2, 147-155, 2001.
- /GEL 13/ Gelman A., et al., Bayesian Data Analysis, 3rd Edition, Chapman & Hall / CRC Texts in Statistical Science, 2013.
- /GOO 63/ Goodman L. A., Kruskal W. H., Measures of Association for Cross Classification, Journal American Statistical Association, 58, 310-364, 1963.
- /GUT 70/ Guttman, I., Statistical tolerance regions, classical and Bayesian. Griffin, London, 1970.
- /HAL 00/ Haldar, A. and Mahadevan, S., Probability, Reliability and Statistical Methods in Engineering Design, John Wiley & Sons, 2000.
- /HEL 96/ Helton J.C. et al., Uncertainty and sensitivity analysis results obtained in the 1992 performance assessment for the waste isolation pilot plant. Reliab Eng Syst Safety, 51, 53-100, 1996.

- /HEL 06/ Helton J.C., et al., Survey of sampling-based methods for uncertainty and sensitivity analysis. *Reliab Eng Syst Safety*; 91, 1175–1209, 2006.
- /HOF 99/ Hofer E., Sensitivity analysis in the context of uncertainty analysis for computationally intensive models. *Computer Physics Communications*, 117, 21-34, 1999.
- /HOW 69/ Howe, W. G., Two-sided Tolerance Limits for Normal Populations - Some Improvements, *Journal of the American Statistical Association*, 64, 610-620, 1969.
- /HUL 62/ Hull, T. E., Dobell, A. R., Random number generators. *SIAM Review*, Vol. 4, No. 3, 1962.
- /ITL 17/ NIST/SEMATECH e-Handbook of Statistical Methods, <http://www.itl.nist.gov/div898/handbook/>, 19.07.2017.
- /IMA 82/ Iman, R., Conover W. J., A distribution free approach to inducing rank correlations among input variables, *Communication in Statistics - Simulation and Computation*, 11(3), 311-334, 1982.
- /JAN 99/ Jansen, M.J.W., Analysis of variance designs for model output, *Computer Physics Communications*, 117, 35–43, 1999.
- /JOH 78/ Johnson, M.E. ; Ramberg, J.S., Transformations of the multivariate normal distribution with applications to simulation. Johnson Transformation System. 11th International Conference on Systems Sciences, Honolulu, HI, USA, LA-UR-77-2595; CONF-780103-3, 1978.
- /JOH 94/ Johnson, N. L., Kotz, S., Balakrishnan, N., Continuous univariate distributions, Vol. 1 of Wiley series in Probability and Mathematical Statistics: Applied Probability and Statistics, 1994.
- /JOH 95/ Johnson, N. L., Kotz, S., Balakrishnan, N., Continuous univariate distributions, Vol. 2 of Wiley series in Probability and Mathematical Statistics: Applied Probability and Statistics, 1995.

- /JOR 05/ Jordaan, I., Decisions under uncertainty: probabilistic analysis for engineering decisions, pp.452, Cambridge University Press, 2005.
- /KEN 73/ Kendall M. G., Stuart A., The advanced theory of statistics, Vol. 2: Inference and relationship, 3rd ed., MacMillan Publishing Co., New York, 1973.
- /KLO 91/ Kloos, M., et al., DIVIS: An Interactive Software Package to Support the Probabilistic Modelling of Parameter Uncertainties, GRS-A-1760, Gesellschaft für Anlagen- und Reaktorsicherheit, Garching, Germany, 1991.
- /KLO 12/ Kloos M., Sensitivity analyses supplemented to epistemic uncertainty analyses for PSA results. Proceedings of PSAM 11, Helsinki, Finland, 2012.
- /KOT 88/ Kotz S., Johnson N. L. (eds.), Encyclopedia of statistical sciences, Volumes 1-9, John Wiley & Sons, 1982-1988.
- /KRZ 88/ Krzykacz, B. and Hofer, E., The Generation of Experimental Designs for Uncertainty and Sensitivity Analysis of Model Predictions with Emphasis on Dependences between Uncertain Parameters, Reliability Of Radioactive Transfer Models (Desmet, G. Ed), Elsevier Applied Science Publishers, London New York, 1988.
- /KRU 58/ Kruskal, W. H., Ordinal Measures of Association, Journal of the American Statistical Association, 53(284), 814-861, 1958.
- /LIL 67/ Lilliefors, H. W., On the Kolmogorov-Smirnov Test for Normality with Mean and Variance Unknown. Journal of American Statistician Assoc. 62, 399-402, 1967.
- /LIL 69/ Lilliefors, H. W., On the Kolmogorov-Smirnov Test for the Exponential Distribution with Mean Unknown. Journal of American Statistician Assoc. 64, 387-389, 1969.
- /LIM 01/ Limpert, E., Stahel, W.A. and Abbt, M., Log-normal Distributions across the Sciences: Keys and Clues, AIBS Bulletin, 51(5), 341-352, 2001.

- /MAT 98/ Matsumoto, M., Nishimura, T., Mersenne Twister: A 623-dimensionally equidistributed uniform pseudorandom number generator. *ACM Transactions on Modeling and Computer Simulation (TOMACS)* - Special issue on uniform random number generation. Volume 8, Issue 1, 3-30, 1998.
- /MCK 79/ McKay MD, Beckman RJ, Conover WJ., A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21, 239-245, 1979.
- /MCK 95/ McKay, M.D., et al., Evaluating prediction uncertainty, US Nuclear Regulatory Commission, 1995.
- /MCK 96/ McKay M. D., Variance-based methods for assessing uncertainty importance, NUREG-1150 analyses. LA-UR-96-2695; 1–27, 1996.
- /NAT 63/ Natrella, M. G., *Experimental Statistics*, NBS Handbook 91, US Department of Commerce, 1963.
- /REB 07/ Rebafka, T., S. Cléménçon and M. Feinberg, Bootstrap-based tolerance intervals for application to method validation. *Chemometrics and Intelligent Laboratory Systems* 89.2, 69-81, 2007.
- /SAL 00/ Saltelli A., K. Chan, E. Scott (Eds.), *Sensitivity analysis*, Wiley Series in Probability and Statistics, Wiley, 2000.
- /SAL 02/ Saltelli A., Making best use of model valuations to compute sensitivity indices, *Computer Physics Communications* 145, 280–297, 2002.
- /SAL 08/ Saltelli, A., M. Ratto et al., *Global sensitivity analysis: the primer*. John Wiley & Sons, 2008.
- /SAL 10/ Saltelli, A., et al., Variance based sensitivity analysis of model output. Design and estimator for the total sensitivity index. *Computer Physics Communications* 181.2, 259-270, 2010.
- /SHO 05/ Shoung, J-M., S. Altan, J. Cabrera, Double bootstrapping a tolerance limit. *Journal of biopharmaceutical statistics* 15.2, 367-373, 2005.

- /SOB 99/ Sobol, I.M., Levitan, Yu L., On the use of variance reducing multipliers in Monte Carlo computations of a global sensitivity index. Computer Physics Communications 117. 1-2, 52-61, 1999.
- /WIC 98/ Wickett, T., et. al., Report of the uncertainty methods study for advanced best estimate thermal hydraulic code applications, Vol. 1 (Comparison) and Vol. 2 (Re-port by the participating institutions), NEA/CSNI/R(97)35, 1998.
- /WIL 41/ Wilks, S.S., Determination of sample sizes for setting tolerance limits, Annals of Mathematical Statistics 1 (1), 91-96, 1941.
- /WIL 42/ Wilks, S.S., Statistical prediction with special reference to the problem of tolerance limits, Annals of Mathematical Statistics 13 (4), 400-409, 1942.

List of Figures

Fig. 2.1	The probability density function (A) and the cumulative distribution function (B) of a Normal distribution with mean $\mu = 5$ and standard deviation $\sigma = 1$ over the support $[0, 10]$	9
Fig. 2.2	The probability density function (A) and the cumulative distribution function (B) of a Lognormal distribution with parameters $\mu = 0$ and $\sigma = 1$ over the support $[0, 10]$	11
Fig. 2.3	The probability density function (A) and the cumulative distribution function (B) of a Uniform distribution over the support $[-5, 5]$	12
Fig. 2.4	The probability density function (A) and the cumulative distribution function (B) of a Loguniform distribution over the support $[0.5, 5.5]$	13
Fig. 2.5	The probability density function (A) and the cumulative distribution function (B) of a Triangular distribution over the support $[-5, 5]$ with mode 1	15
Fig. 2.6	The probability density function (A) and the cumulative distribution function (B) of a Logtriangular distribution over the support $[0.5, 5.0]$ and with mode 1	16
Fig. 2.7	The probability density function (A, C, E) and the cumulative distribution function (B, D, F) of Weibull distributions over the joint support $[-5, 5]$, with the scale parameter $\sigma = 1$ and varying shape parameter $\mu = 0.5$ for (A, B), 1.0 for (C, D) and 1.5 for (E, F).....	18
Fig. 2.8	The probability density function (A, C, E, G) and the cumulative distribution function (B, D, F, H) of Beta distributions over the joint support $[-5, 5]$, with the shape parameters $(\mu_1, \mu_2) = (0.5, 0.5)$ for (A, B), (2, 2) for (C, D), (2, 4) for (E, F) and (4, 1) for (G, H).....	21
Fig. 2.9	The probability density function (A, C, E) and the cumulative distribution function (B, D, F) of Gamma distributions over the joint support $[0, 20]$ with the parameter settings $(\mu_1, \mu_2) = (1, 0.25)$ for (A, B), (2, 0.3) for (C, D) and (20, 2) for (E, F)	23
Fig. 2.10	The probability density function (A) and the cumulative distribution function (B) of an Extreme Value I distribution over the support $[-5, 20]$ for the location parameter $\mu = 1$ and the scale parameter $\sigma = 2$	25
Fig. 2.11	The probability density function (A,C,E) and the cumulative distribution function (B,D,F) of Extreme Value II distributions over the joint support $[-5, 20]$, i.e. $\text{Min} = -5$, for the shape and scale parameters $(\mu_1, \mu_2) = (1, 1)$ for (A, B), (1, 2) for (C, D) and (2, 1) for (E, F)	26

Fig. 2.12	The probability density function (A,C,E) and the cumulative distribution function (B,D,F) of Exponential distributions over the joint support $[0, 10]$ for the rate parameter $p_1 = 0.5$ for (A, B), 1.0 for (C, D) and 1.5 for (E, F).....	28
Fig. 2.13	The probability density function (A, C, E) and the cumulative distribution function (B, D, F) of ChiSquared distributions over the joint support $[0, 10]$ for the degree of freedom $p_1 = 1$ for (A, B), 3 for (C, D) and 6 for (E, F).....	29
Fig. 2.14	The probability density function (A) and the cumulative distribution function (B) of a Discrete distribution for the value-probability pairs $\{(-5, 0.04), (-3, 0.08), (-1, 0.40), (0, 0.20), (2, 0.06), (3, 0.12), (4, 0.08), (5, 0.02)\}$ over the support $[-5, 5]$	31
Fig. 2.15	The probability density function (A) and the cumulative distribution function (B) of a Histogram distribution for the value-probability pairs $\{(-5, 0.04), (-3, 0.08), (-1, 0.40), (0, 0.20), (2, 0.06), (3, 0.12), (4, 0.08), (5, 0.02)\}$ over the support $[-5, 6]$	33
Fig. 2.16	The probability density function (A) and the cumulative distribution function (B) of a Loghistogram distribution for the value-probability pairs $\{(1, 0.04), (3, 0.08), (5, 0.40), (6, 0.20), (8, 0.06), (9, 0.12), (10, 0.08), (11, 0.02)\}$ over the support $[1, 12]$	34
Fig. 2.17	The probability density function (A) and the cumulative distribution function (B) of a Polygonal Line distribution for the (x, y)-coordinates $\{(-5, 0.04), (-3, 0.08), (-1, 0.40), (0, 0.20), (2, 0.06), (3, 0.12), (4, 0.08), (5, 0.02)\}$ over the support $[-5, 5]$	36
Fig. 2.18	The Pearson's (Ordinary) correlation coefficient for values of (A) strongly negative correlated $\rho=-0.99$, (B) uncorrelated $\rho=0$ and (C) strongly positive correlated $\rho=0.99$ parameters X and Y . These scatter plots exemplify the dependency assigned to the parameters X and Y	40
Fig. 2.19	The principle of concordance and discordance exemplified for the parameter pair (X, Y) . The reference point for the order is taken as the medians (m_x, m_y) similar to the definition of Blomqvist's correlation coefficient.....	41
Fig. 2.20	The principle of dependency between uncertain parameters X and Y by the assumption of an inequality $Y \geq a \cdot X$. The blue shaded area marks the support plane in which the subjective association between X and Y implies that no sample pair (x_i, y_i) is accepted.....	49

Fig. 2.21	Relationship between parameters X and Y defined by the inequality $Y \geq a \cdot X$ with $a = 1$ ($n = 1000$ data points (x, y))	51
Fig. 4.1	Uncertainty of a time-dependent result represented by the time histories obtained from 100 Monte Carlo simulation runs.....	66
Fig. 4.2	Cumulative empirical distribution function of a scalar computational result (Consequence 1) together with the distribution support (cyan-colored horizontal line), the 3 empirical quartiles $y_n \cdot 0.25: n \leq y_n \cdot 0.50: n \leq y_n \cdot 0.75: n$ (blue diamonds) and the empirical mean (maroon triangle)	68
Fig. 4.3	Cumulative empirical distribution (red) and density (green) function of a scalar computational result (Consequence 1)	68
Fig. 4.4	Wilks' one-sided upper 95 %/ 95 % tolerance limit (yellow triangles) compared with the empirical 95 %-quantile (magenta) calculated in each of 1000 different samples of size $n=100$ from a standard Normal distribution with the 95 %-quantile = 1.6448 (blue horizontal line).....	71
Fig. 4.5	Cumulative Beta distribution functions of the coverage probability P of the interval $(-\infty, Y_{s:n}]$ right-sided closed by Wilks' upper 95 %/ 95 % tolerance limit $Y_{s:n}$ for different orders s and different sample sizes n	73
Fig. 4.6	Cumulative empirical distribution function of a scalar computational result (Peak Cladding Temperature) and Wilks' two-sided tolerance limits (green squares on the x-axis)	74
Fig. 4.7	Wilks' two-sided tolerance limits of a time-dependent computational result (Cladding Temperature)	74
Fig. 4.8	Cumulative empirical distribution function of a scalar computational result (Peak Cladding Temperature) and fitted Normal distribution. The level of significance of the Lilliefors test is 0.16, i.e. the hypothesis cannot be rejected that Y has a Normal distribution.....	82
Fig. 5.1	Pearson's ordinary correlation coefficients as sensitivity indices of a scalar computational result (Peak Cladding Temperature) with respect to 56 uncertain input parameters.....	92
Fig. 5.2	Standardized regression coefficients (with respect to Pearson's ordinary correlation) as sensitivity indices of a time-dependent computational result (Cladding Temperature) with respect to 56 uncertain parameters	96

Fig. 5.3	Coefficients of determination with respect to Pearson's ordinary correlation for a time-dependent computational result (Cladding Temperature) influenced by 56 uncertain parameters.....	97
Fig. 5.4	Correlation ratios as sensitivity indices of a scalar computational result (Peak Cladding Temperature) with respect to 56 uncertain input parameters.....	106
Fig. 5.5	Correlation ratios on ranks as sensitivity indices of a time-dependent computational result (Cladding Temperature) with respect to 56 uncertain parameters.....	108

List of Tables

Tab. 4.1 Minimum sample size to determine the $\beta \cdot 100\%$ / $\gamma \cdot 100\%$ tolerance interval for selected coverage probabilities β and confidence levels γ 72

**Gesellschaft für Anlagen-
und Reaktorsicherheit
(GRS) gGmbH**

Schwertnergasse 1
50667 Köln
Telefon +49 221 2068-0
Telefax +49 221 2068-888

Forschungszentrum
Boltzmannstraße 14
85748 Garching b. München
Telefon +49 89 32004-0
Telefax +49 89 32004-300

Kurfürstendamm 200
10719 Berlin
Telefon +49 30 88589-0
Telefax +49 30 88589-111

Theodor-Heuss-Straße 4
38122 Braunschweig
Telefon +49 531 8012-0
Telefax +49 531 8012-200

www.grs.de