



# TREATMENT OF UNCERTAINTIES IN MULTI-PHYSICS MODEL FOR WIND TURBINE ASSET MANAGEMENT

**Elias FEKHARI**

ÉLECTRICITÉ DE FRANCE R&D

*Chatou, France*

&

CÔTE D'AZUR UNIVERSITY

*Nice, France*

This dissertation is submitted for the degree of

*Doctor of Philosophy*

publicly defended on January xx, 2024 in front of the following jury:

Pr. Mireille BOSSY,	INRIA, Sophia-Antipolis	Examiner
Dr. Vincent CHABRIDON	EDF R&D, Chatou	Co-advisor
Dr. Sébastien DA VEIGA	ENSAI, Rennes	Examiner
Dr. Bertrand IOOSS	EDF R&D, Chatou	Thesis director
Dr. Anaïs LOVERA	EDF R&D, Saclay	Invite
Dr. Joseph MURÉ	EDF R&D, Chatou	Co-advisor
Pr. Franck SCHOEFS	Nantes Université, Nantes	Reviewer
Pr. Daniel STRAUB	TUM, Munich	Reviewer
Pr. Bruno SUDRET	ETH, Zürich	Examiner



# Table of contents

<b>List of figures</b>	<b>v</b>
<b>List of tables</b>	<b>vii</b>
<b>Introduction</b>	<b>1</b>
<b>I Introduction to uncertainty quantification and wind energy</b>	<b>7</b>
<b>1 Treatment of uncertainties in computer experiments</b>	<b>9</b>
1.1 Introduction . . . . .	10
1.2 Black-box model specification . . . . .	10
1.3 Enumerating and modeling the uncertain inputs . . . . .	11
1.3.1 Sources of the input uncertainties . . . . .	11
1.3.2 Modeling uncertain inputs with the probabilistic framework . . . . .	11
1.3.3 Joint input probability distribution . . . . .	13
1.4 Central tendency uncertainty propagation . . . . .	15
1.4.1 Numerical integration . . . . .	16
1.4.2 Numerical design of experiments . . . . .	19
1.4.3 Central tendency estimation . . . . .	20
1.5 Reliability-oriented uncertainty propagation . . . . .	20
1.5.1 Problem formalization . . . . .	20
1.5.2 Rare event estimation methods . . . . .	20
1.6 Sensitivity analysis . . . . .	21
1.6.1 Global sensitivity analysis . . . . .	21
1.6.2 Reliability-oriented sensitivity analysis . . . . .	21
1.7 Metamodeling . . . . .	21
1.7.1 Global metamodel . . . . .	21
1.7.2 Contour finding for rare-event estimation . . . . .	21
1.8 Conclusion . . . . .	21

<b>2</b>	<b>Introduction to wind turbine modeling and design</b>	<b>23</b>
2.1	Introduction . . . . .	24
2.2	Wind turbine modeling . . . . .	24
2.2.1	Synthetic wind generation [TurbSim, Kaimal spectrum] . . . . .	24
2.2.2	Synthetic wave generation . . . . .	24
2.2.3	Aerodynamic interactions . . . . .	24
2.2.4	Servo-Hydro-Aero-Elastic wind turbine simulation [DIEGO] . . . . .	24
2.2.5	Soil modeling . . . . .	24
2.2.6	Wake modeling [FarmShadow] . . . . .	24
2.3	Recommended design practices . . . . .	24
2.3.1	Design load cases . . . . .	24
2.3.2	Dynamic response design . . . . .	24
2.3.3	Fatigue response design . . . . .	24
2.4	Uncertain inputs . . . . .	24
2.4.1	Environmental inputs . . . . .	24
2.4.2	System inputs . . . . .	24
2.4.3	Probabilistic fatigue assessment . . . . .	24
2.5	Conclusion . . . . .	24
	<b>References</b>	<b>25</b>
	<b>Appendix A Univariate distribution fitting</b>	<b>27</b>
A.1	Main parametric methods . . . . .	27
A.2	Main nonparametric methods . . . . .	28
	<b>Appendix B Nonparametric copula estimation</b>	<b>31</b>
B.1	Empirical copula . . . . .	31
B.2	Empirical Bernstein & Beta copula . . . . .	31
B.3	Goodness-of-fit . . . . .	32
	<b>Appendix C Rare event estimation algorithms</b>	<b>35</b>
	<b>Appendix D Résumé étendu de la thèse</b>	<b>37</b>

# List of figures

1	General uncertainty quantification framework (adapted from <a href="#">Ajenjo (2023)</a> ) . . . . .	3
1.1	Samples of three joint distributions with identical marginals and different dependence structures . . . . .	13
1.2	Ranked samples represented in the Fig. <a href="#">1.1</a> . . . . .	14
1.3	Monte Carlo and quasi-Monte Carlo designs ( $n = 256$ ) . . . . .	19
1.4	Monte Carlo and quasi-Monte Carlo designs ( $n = 256$ ) . . . . .	19
A.1	Adequation of two different Weibull models using their likelihood with a sample of observations (black crosses). . . . .	28
A.2	Fit of a bimodal density by KDE using different tuning parameters. . . . .	29
A.3	QQ-plot between the data from Example <a href="#">2</a> and a KDE model. . . . .	30
B.1	Evolution of $m_{\text{IMSE}}$ for different dimensions and sample sizes. . . . .	33



## List of tables





# Introduction

## Industrial context and motivation

The shift in wind energy projects from limited onshore resources to the vast potential of offshore locations is a growing trend. Offshore wind energy offers several advantages, including more consistent winds and the ability to install larger turbines. Since the installation of the first offshore wind farm in Vindeby, Denmark, in 1991, the industry has experienced rapid growth, with a total capacity of 56GW exploited worldwide in 2021. Over time, offshore wind technology has matured, resulting in significant achievements such as securing projects in Europe through "zero-subsidy bids," where electricity generated by wind farms is sold at wholesale prices.

However, despite the progress of this sector, scaling limitations emerge and numerous scientific challenges. To meet ambitious national and regional development targets, the wind energy industry must address various scaling issues, including port logistics, the demand for critical natural resources, and sustainable end-of-life processes. Furthermore, the field presents various scientific challenges that often involve coupling data with numerical simulations of physical systems and their surrounding environment. The wind energy community is focused on several objectives, including enhancing the design of floating offshore wind turbines, refining wind resource estimation techniques, and optimizing maintenance operations. Additionally, the design, installation and exploitation of these industrial assets implicate several decision-making steps, considering limited access to information. Therefore, properly modeling and treating the various uncertainties along this process proved to be a key success factor in this highly competitive industry.

Overall, the industry needs methods and techniques for uncertainty management to optimize safety margins and asset management. As a wind farm project developer, the attention is first drawn to refining the wind potential of candidate sites by combining different sources of information and modeling the multivariate distribution of environmental conditions within a wind farm. In floating projects, the probabilistic design helps to define safer and more robust solutions. As a wind farm owner, another significant consideration revolves around end-of-life management. This involves evaluating three possible outcomes: extending the operating assets' lifetime, replacing current turbines with more advanced models, or dismantling and selling the wind farm. The first two solutions require assessing the current reliability of the structure

and its remaining useful life. These quantitative evaluations are studied by certification bodies and insurance providers to issue exploitation permits. To deliver rigorous risk assessments, the generic *uncertainty quantification methodology* may be adopted.

## Generic methodology for uncertainty quantification

Uncertainty Quantification (UQ) aims at modeling and managing uncertainties in complex systems. Over the year, generic UQ frameworks were proposed [add ref Deroc.] to quantify and analyze the relations between uncertain input factors and the systems' outcomes. UQ is particularly relevant in situations where experiments or direct observations are costly, time-consuming, or even impossible to conduct.

Computer experiments, also known as numerical experiments or simulations, play an important role in UQ. They involve the use of numerical models to simulate the behavior of a system under various conditions and parameter settings. These virtual experiments provide a cost-effective way to explore the behavior of complex systems and make robust and well-informed decisions. They enable researchers and decision-makers to gain a deeper understanding of the system dynamics, optimize designs, assess risk, and make robust predictions. As a result, uncertainty quantification has become an essential tool in wind energy, benefiting from the multiphysics numerical models simulating the behavior of wind farms interacting with their environment. Nevertheless, numerical models should be calibrated against measured data and pass validation, and verification processes [add ref] to minimize the residual modeling error. Figure 1 illustrates the UQ methodologies and the standardized usual steps encountered during a study, which are detailed hereafter:

- **Step A – Problem specification:** at this step, it is necessary to establish the system under study and construct a numerical model capable of precisely simulating its behavior. Specifying the problem also involves defining the complete set of parameters inherent to the computer model. This includes the input variables as well as determining the specific output quantity that will be generated by the numerical model;
- **Step B – Uncertainty modeling:** The objective of the second step is to identify all the sources of uncertainty impacting the input variables. Most of the time choosing a probabilistic framework, the modeling methods will depend on the available information (e.g., amount of data, input dimension);
- **Step C – Uncertainty propagation:** This step consists in propagating the uncertain inputs through the computer model, making the output uncertain. Then, the goal becomes the estimation of a quantity of interest (i.e., a statistic on the random output variable of interest). The uncertainty propagation method may differ depending on the quantity of interest targeted (e.g., central tendency, rare event);
- **Step C' – Inverse analysis:** In this additional step, a sensitivity analysis can be performed to study the role allocated to each uncertain input leading to the uncertain output;

- **Metamodeling:** Since this methodology is frequently used with computationally expensive numerical models, it becomes interesting to emulate these models using statistical models constructed from a limited number of simulations. The uncertainty quantification is then performed on the so-called “metamodel” (or surrogate model) at a reasonable computation cost.

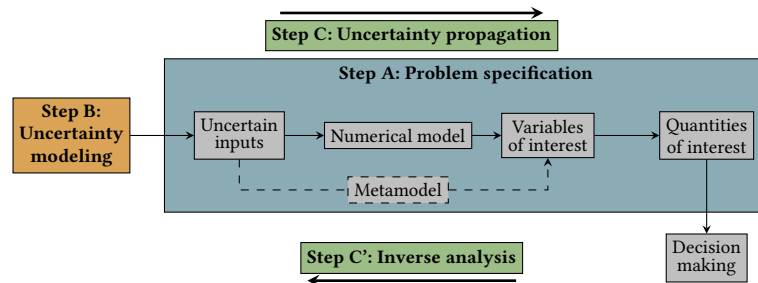


Fig. 1 General uncertainty quantification framework (adapted from Ajenjo (2023))

## Problem statement and outline of the thesis

[Rewrite para] A general topic of research for EDF R&D is to adapt the UQ methodology to offshore wind turbine industrial cases. However, this problem presents various specificities, raising scientific challenges. First, the numerical model studied is composed of a series of three codes, among which one is intrinsically stochastic (i.e., running twice the same numerical model with the same set of inputs results in different outputs). Second, the computational cost of these numerical models quickly requires the use of efficient techniques deployed on high-performance computers to perform UQ. Then, the probabilistic modeling tools available to model the uncertain inputs are challenged by a complex underlying dependence structure. In the presence of large amounts of data describing these complex inputs, different methods to quantify and propagate the uncertainties are needed. Finally, performing a risk assessment on this case study combines all the challenges previously stated. In order to adapt the UQ framework to this industrial case, this thesis aims at answering the following questions:

- Q1** *How to accurately model the complex dependence structure underlying the multivariate distribution of the environmental conditions?*
- Q2** *How to perform an efficient and accurate given-data uncertainty propagation on a costly and stochastic numerical model?*
- Q3** *How to couple rare event estimation with reliability-oriented sensitivity analysis?*

To intend at solving these problems, this thesis is divided into three parts. The first part gathers an introduction to UQ’s state-of-the-art and a specification of the offshore wind turbine problem. The second part presents the contributions to uncertainty quantification and propagation while the third part the contributions to rare event estimation. This manuscript is divided into seven chapters, which are summarized hereafter:

**Chapter 1** introduces the Treatment of uncertainties in computer experiments. Uncertainty quantification

**Chapter 2** Introduction to wind turbine modeling and design

**Chapter 3** Kernel-based uncertainty quantification

**Chapter 4** Kernel-based central tendency estimation

**Chapter 5** Kernel-based metamodel validation

**Chapter 6** Nonparametric rare event estimation

**Chapter 7** Sequential reliability oriented sensitivity analysis

## Numerical developments

In the vain of an open-data approach, this aims at sharing the implementations developed and allows the reader to reproduce numerical results. Along this thesis, the contributions to numerical developments are summarized below:

otkerneldesign <sup>1</sup>	<ul style="list-style-type: none"> <li>• This Python package generates designs of experiments based on kernel methods such as Kernel Herding and Support Points. A tensorized implementation of the algorithms was proposed, significantly increasing their performances. Additionally, optimal weights for Bayesian quadrature are provided.</li> <li>• This Python package, developed in collaboration with J.Muré, is available on the platform Pypi and fully documented.</li> </ul>
bancs <sup>2</sup>	<ul style="list-style-type: none"> <li>• This Python package proposes an implementation of the “Bernstein Adaptive Nonparametric Conditional Sampling” method for rare event estimation.</li> <li>• This Python package is available on the PyPI platform and is illustrated with examples and analytical benchmarks.</li> </ul>
ctbenchmark <sup>3</sup>	<ul style="list-style-type: none"> <li>• This Python package presents a standardized process to benchmark different sampling methods for central tendency estimation.</li> <li>• This Python package is available on a GitHub repository with analytical benchmarks.</li> </ul>
copulogram <sup>4</sup>	<ul style="list-style-type: none"> <li>• This Python package proposes an implementation of a synthetic visualization tool for multivariate distributions.</li> <li>• This Python package, developed in collaboration with V.Chabridon, is available on the Pypi platform and fully documented.</li> </ul>

<sup>1</sup>Documentation: <https://efekhari27.github.io/otkerneldesign/master/>

<sup>2</sup>Repository: <https://github.com/efekhari27/bancs>

<sup>3</sup>Repository: <https://github.com/efekhari27/ctbenchmark>

<sup>4</sup>Repository: <https://github.com/efekhari27/copulogram>

## Publications and communications

The research contributions in this manuscript are based on the following publications:

Book Chap.	<u>E. Fekhari</u> , B. Iooss, J. Muré, L. Pronzato and M.J. Rendas (2023). “Model predictivity assessment: incremental test-set selection and accuracy evaluation”. In: <i>Studies in Theoretical and Applied Statistics</i> , pages 315–347. Springer.
Jour Pap.	<u>E. Fekhari</u> , V. Chabridon, J. Muré and B. Iooss (2023). “Fast given-data uncertainty propagation in offshore wind turbine simulator using Bayesian quadrature”. In: <i>Data-Centric Engineering</i> .
Int. Conf	<p><u>E. Fekhari</u>, B. Iooss, V. Chabridon, J. Muré (2022). “Numerical Studies of Bayesian Quadrature Applied to Offshore Wind Turbine Load Estimation”. In: <i>SIAM Conference on Uncertainty Quantification (SIAM UQ22)</i>, Atlanta, USA. (Talk)</p> <p><u>E. Fekhari</u>, B. Iooss, V. Chabridon, J. Muré (2022). “Model predictivity assessment: incremental test-set selection and accuracy evaluation”. In: <i>22nd Annual Conference of the European Network for Business and Industrial Statistics (ENBIS 2022)</i>, Trondheim, Norway. (Talk)</p> <p><u>E. Fekhari</u>, B. Iooss, V. Chabridon, J. Muré (2022). “Efficient techniques for fast uncertainty propagation in an offshore wind turbine multi-physics simulation tool”. In: <i>Proceedings of the 5th International Conference on Renewable Energies Offshore (RENEW 2022)</i>, Lisbon, Portugal. (Paper &amp; Talk)</p> <p><u>E. Fekhari</u>, V. Chabridon, J. Muré and B. Iooss (2023). “Bernstein adaptive nonparametric conditional sampling: a new method for rare event probability estimation”. In: <i>Proceedings of the 13th International Conference on Applications of Statistics and Probability in Civil Engineering (ICASP 14)</i>, Dublin, Ireland. (Paper &amp; Talk)</p> <p>E. Vanem, <u>E. Fekhari</u>, N. Dimitrov, M. Kelly, A. Cousin and M. Guiton (2023). “A joint probability distribution model for multivariate wind and wave conditions”. In: <i>Proceedings of the ASME 2023 42th International Conference on Ocean, Offshore and Arctic Engineering (OMAE 2023)</i>, Melbourne, Australia. (Paper)</p> <p>A. Lovera, <u>E. Fekhari</u>, B. Jézéquel, M. Dupoirion, M. Guiton and E. Ardillon (2023). “Quantifying and clustering the wake-induced perturbations within a wind farm for load analysis”. In: <i>Journal of Physics: Conference Series (WAKE 2023)</i>, Visby, Sweden (Paper)</p>
Nat. Conf.	<p><u>E. Fekhari</u>, B. Iooss, V. Chabridon, J. Muré (2022). “Kernel-based quadrature applied to offshore wind turbine damage estimation”. In: <i>Proceedings of the Mascot-Num 2022 Annual Conference (MASCOT NUM 2022)</i>, Clermont-Ferrand, France (Poster)</p> <p><u>E. Fekhari</u>, B. Iooss, V. Chabridon, J. Muré (2023). “Rare event estimation using nonparametric Bernstein adaptive sampling”. In: <i>Proceedings of the Mascot-Num 2023 Annual Conference (MASCOT-NUM 2023)</i>, Le Croisic, France (Talk)</p>

# PART I:

## INTRODUCTION TO UNCERTAINTY QUANTIFICATION AND WIND ENERGY

*Toute pensée émet un coup de dé.*

---

S. MALLARMÉ





# Treatment of uncertainties in computer experiments

---

1.1	Introduction . . . . .	10
1.2	Black-box model specification . . . . .	10
1.3	Enumerating and modeling the uncertain inputs . . . . .	11
1.3.1	Sources of the input uncertainties . . . . .	11
1.3.2	Modeling uncertain inputs with the probabilistic framework . . . . .	11
1.3.3	Joint input probability distribution . . . . .	13
1.4	Central tendency uncertainty propagation . . . . .	15
1.4.1	Numerical integration . . . . .	16
1.4.2	Numerical design of experiments . . . . .	19
1.4.3	Central tendency estimation . . . . .	20
1.5	Reliability-oriented uncertainty propagation . . . . .	20
1.5.1	Problem formalization . . . . .	20
1.5.2	Rare event estimation methods . . . . .	20
1.6	Sensitivity analysis . . . . .	21
1.6.1	Global sensitivity analysis . . . . .	21
1.6.2	Reliability-oriented sensitivity analysis . . . . .	21
1.7	Metamodeling . . . . .	21
1.7.1	Global metamodel . . . . .	21
1.7.2	Contour finding for rare-event estimation . . . . .	21
1.8	Conclusion . . . . .	21

---

## 1.1 Introduction

The progress of computer simulation gradually allows the virtual resolution of more complex problems in scientific fields such as physics, astrophysics, engineering, climatology, chemistry, or biology. This domain often provides a deterministic solution to complex problems depending on several inputs. Associating a UQ analysis with these possibly nonlinear numerical models is a key element to improving the understanding of the phenomena studied. A wide panel of UQ methods has been developed over the years to pursue these studies with a reasonable computational cost.

This chapter presents the standard tools and methods from the generic UQ framework [add ref to intro], exploited later in this thesis. It is structured as follows: Section [add ref] describes the context of the model specification step; Section [add ref] presents a classification of the inputs uncertainties and the probabilistic framework to model them; Section [add ref] and xx introduce various methods to propagate the input uncertainties through the numerical model for different purposes; finally, Section [add ref] presents the main inverse methods to perform sensitivity analysis in our framework.

## 1.2 Black-box model specification

The uncertainty quantification studies in our framework are performed around an input-output numerical simulation model. This numerical model, or code, is hereafter considered as *black-box* since the knowledge of the underlying physics doesn't inform the UQ methods. Alternatively, one could consider *intrusive* UQ methods, introducing uncertainties within the resolution of computer simulation (see e.g., [add ref]). In practice, the numerical model might be a sequence of codes executed in series to obtain a variable of interest.

Moreover, the simulation model is in most cases deterministic, otherwise, it is qualified as intrinsically stochastic (i.e., two runs of the same model taking the same inputs return different outputs). Then, most numerical simulation presents modeling errors. In the following, it will be assumed that the numerical models passed a *validation & verification* phase, to quantify their confidence and predictive accuracy.

Formally, part of the problem specification is the definition of the set of  $p$  input variables  $\mathbf{x} = (x_1, \dots, x_p)^\top$  considered uncertain (e.g., wind speed, wave period, etc.). In this thesis, the models considered will only present scalar outputs. UQ methods dedicated to other types of outputs exist (see e.g., for time series outputs Lataniotis (2019), [functional Alvaro?]). Let us then define the following numerical model:

$$\mathcal{M} : \left\{ \begin{array}{lcl} \mathcal{D}_x \subseteq \mathbb{R}^p & \longrightarrow & \mathcal{D}_y \subseteq \mathbb{R} \\ \mathbf{x} & \longmapsto & y. \end{array} \right. \quad (1.1)$$

Unlike the typical machine learning input-output dataset framework, the UQ analyst can simulate the output image of any inputs (in the input domain), using the numerical model.

However, numerical simulations often come with an important computational cost. Therefore UQ methods should be efficient and require as few simulations as possible. In this context, metamodels (or surrogate models) are statistical approximations of the costly numerical model, that can be used to perform tractable UQ. Metamodels are only built and validated on a limited number of simulations (in a *supervised learning* framework). In practice, the model specification step is often associated with the development of a *wrapper* of the code [explain wrapper], with can be deployed on a *high-performance computer*. Once the model is specified, a critical step of uncertainty quantification is enumerating the input uncertainties and building an associated mathematical model.

## 1.3 Enumerating and modeling the uncertain inputs

### 1.3.1 Sources of the input uncertainties

To ensure a complete risk assessment (e.g., associated with the exploitation of a wind turbine throughout its life span), the analyst should construct a list of uncertain inputs as exhaustive as possible. Even if these uncertainties might have different origins, they should all be considered jointly in the UQ study. The authors proposed to classify them for practical purposes into two groups:

- **aleatory uncertainty** regroups the uncertainties that arise from natural randomness (e.g., [add example]). From a risk management point of view, these uncertainties are qualified as *irreducible* since the industrials facing them will not be able to acquire additional information to reduce them (e.g., additional measures).
- **epistemic uncertainty** gathers the uncertainties resulting from a lack of knowledge. Contrarily to the aleatory ones, epistemic uncertainties might be reduced by investigating their origin.

Der Kiureghian and Ditlevsen (2009) offers a discussion on the relevance of this classification. They affirm that this split is practical for decision-makers to identify possible ways to reduce their uncertainties. However, this distinction should not affect the way of modeling or propagating uncertainties. [To illustrate the limits of this split, some uncertainties present both an aleatory and epistemic aspect.] In the following, the probabilistic framework is introduced to deal with uncertainties.

### 1.3.2 Modeling uncertain inputs with the probabilistic framework

Uncertainties are traditionally modeled with objects from the probability theory. In this thesis, the *probabilistic framework* is adopted. Alternative theories exist to mathematically model uncertainties. For example, imprecise probability theory allows more general modeling of the uncertainties. It becomes useful when dealing with very limited and possibly contradictory

information (e.g., expert elicitation). The core probabilistic tools and objects are introduced hereafter.

The *probability space* (i.e., a measure space with its total measure summing to one), also called probability triple and denoted  $(\Omega, \mathcal{A}, \mu)$ . This mathematical concept first includes a sample space  $\Omega$ , which contains a set of outcomes  $\omega \in \Omega$ . An *event* is defined as a set of outcomes in the sample space. Then, a  $\sigma$ -algebra  $\mathcal{A}$  (also called event space) is a set of events. Finally, a probability function  $\mu : \mathcal{A} \rightarrow [0, 1]$ , is a positive probability measure associated with an event. Most often, the choice of the probability space will not be specified. The main object will be functions defined over this probability space: random variables.

The *random vector*  $\mathbf{X}$  (i.e., multivariate random variable) is a measurable function defined as:

$$\mathbf{X} : \begin{cases} \Omega & \longrightarrow & \mathcal{D}_{\mathbf{x}} \subseteq \mathbb{R}^p \\ \omega & \longmapsto & \mathbf{X}(\omega) = \mathbf{x}. \end{cases} \quad (1.2)$$

In the following, the random vector  $\mathbf{X}$  will be considered to be a squared-integrable function against the measure  $\mu$  (i.e.,  $\int_{\Omega} |\mathbf{X}(\omega)|^2 d\mu(\omega) < \infty$ ). Moreover, this work will focus on continuous random variables.

The *probability distribution* of the random vector  $\mathbf{X}$  is the pushforward measure of  $\mu$  by  $\mathbf{X}$ . Which is a probability measure on  $(\mathcal{D}_{\mathbf{x}}, \mathcal{A})$ , denoted  $\mu_{\mathbf{X}}$  and defined by:

$$\mu_{\mathbf{X}}(B) = \mu(\mathbf{X} \in B) = \mu(\omega \in \Omega : \mathbf{X}(\omega) \in B), \quad \forall B \in \mathcal{A}. \quad (1.3)$$

The *cumulative distribution function* (CDF) is a common tool to manipulate random variables. It is a function  $F_{\mathbf{X}} : \mathcal{D}_{\mathbf{x}} \rightarrow [0, 1]$  defined for all  $\mathbf{x} \in \mathcal{D}_{\mathbf{x}}$  as:

$$F_{\mathbf{X}}(\mathbf{x}) = \mu(\mathbf{X} \leq \mathbf{x}) = \mu(X_1 \leq x_1, \dots, X_p \leq x_p) = \mu_{\mathbf{X}}([-\infty, x_1] \times \dots \times [-\infty, x_p]). \quad (1.4)$$

The CDF is a positive, increasing, right-continuous function, which tends to 0 as  $\mathbf{x}$  tends to  $-\infty$  and to 1 as  $\mathbf{x}$  tends to  $+\infty$ . In the continuous case, one can also define a corresponding *probability density function* (PDF)  $f_{\mathbf{X}} : \mathcal{D}_{\mathbf{x}} \rightarrow \mathbb{R}_+$  with  $f_{\mathbf{X}}(\mathbf{x}) = \frac{\partial^p F_{\mathbf{X}}(\mathbf{x})}{\partial x_1 \dots \partial x_p}$ .

The expected value of a random vector  $\mathbb{E}[\mathbf{X}]$ , also called the first moment, is a vector defined as:

$$\mathbb{E}[\mathbf{X}] = \int_{\Omega} \mathbf{X}(\omega) d\mu(\omega) = \int_{\mathcal{D}_{\mathbf{x}}} \mathbf{x} f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_p])^{\top}. \quad (1.5)$$

In addition, considering two random variables  $X_i$  and  $X_j$ , with  $i, j \in \{1, \dots, p\}$ , one can write their respective variance:

$$\text{Var}(X_i) = \mathbb{E}[X_i - \mathbb{E}[X_i]]^2, \quad (1.6)$$

and a covariance describing their joint variability:

$$\text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]. \quad (1.7)$$

The standard deviation  $\sigma_{X_j} = \sqrt{\text{Var}(X_j)}$  and coefficient of variation  $\delta_{X_j} = \frac{\text{Var}(X_j)}{|\mathbb{E}[X_j]|}$  are two quantities directly associated to the two first moments.

[Remark on high dimension: We call high dimension anything higher than  $p > 10$  and it creates issues.]

### 1.3.3 Joint input probability distribution

This section aims at presenting various techniques to model and infer a joint probability distribution (or multivariate distribution). It will first introduce the *copula*, a universal mathematical tool to model the dependence structure of a joint distribution. Then, a few methods to fit a joint distribution over a dataset will be mentioned. And finally, a panel of tools to evaluate the goodness of fit between a probabilistic model and a dataset will be recalled.

From a practical point of view, people tend to properly model the single effects of their input uncertainties. However, modeling the dependence structure unlying in a joint distribution is often overlooked. To illustrate the importance of this step, Fig. 1.1 represents three i.i.d samples from three bivariate distributions sharing the same single effects (e.g., here two exponential distributions) but different dependence structures. One can assume that the joint distribution is the composition of the single effects, also called marginals, and an application governing the dependence between them.

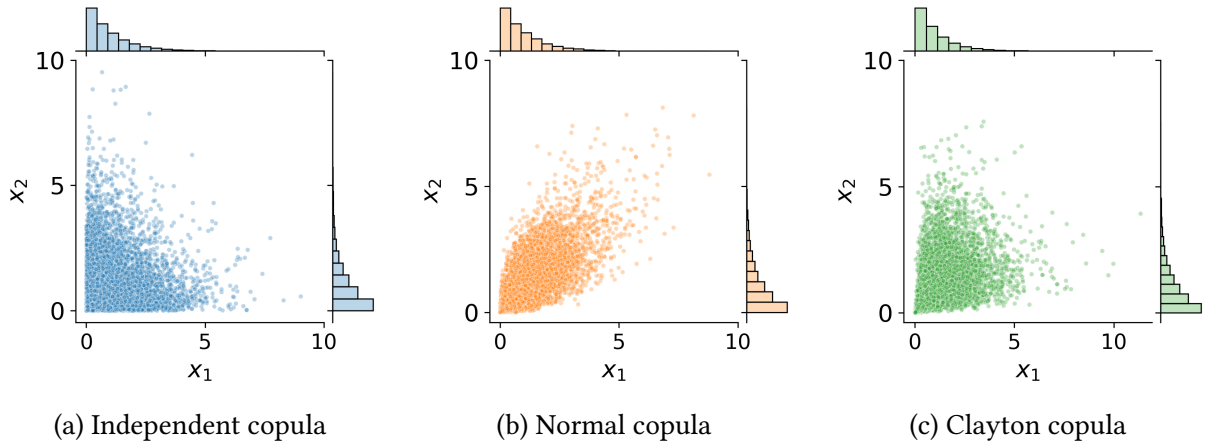


Fig. 1.1 Samples of three joint distributions with identical marginals and different dependence structures

An empirical way of isolating the three dependence structures from this example is to transform the samples in the ranked space. Let us consider a  $n$ -sized sample  $\mathbf{X}_n = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\} \in \mathcal{D}_{\mathbf{x}}^n$ . The corresponding ranked sample is defined as:  $\mathbf{R}_n = \{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(n)}\}$ , where  $r_j^{(l)} = \sum_{i=1}^n \mathbb{1}_{\{x_j^{(i)} \leq x_j^{(l)}\}}$ ,  $\forall j \in \{1, \dots, p\}$ . Ranking a multivariate dataset allows us to isolate the dependence structure witnessed empirically. Fig. 1.2 shows the same three samples from Fig. 1.1 in the ranked space. One can first notice that the marginals are uniform since each rank is uniformly distributed. Then, the scatter plot from the distribution with independent copula (left plot) is uniform while the two others present different patterns.

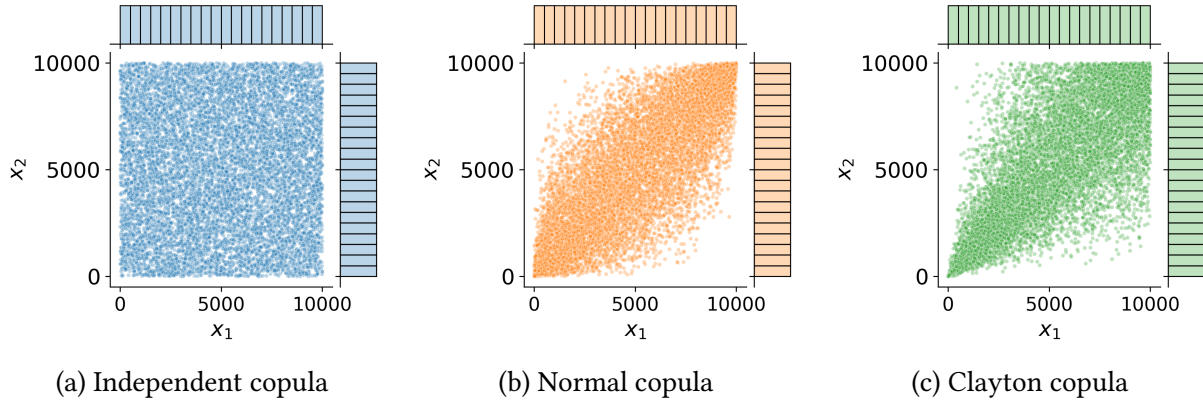


Fig. 1.2 Ranked samples represented in the Fig. 1.1

A theorem states that the multivariate distribution of any random vector can be broken down into two objects (Joe, 1997). First, a set of univariate marginal distributions describing the behavior of the individual variables; Second, a function describing the dependence structure between all variables: a copula.

**Theorem 1** (Sklar's theorem). *Let  $\mathbf{X} \in \mathbb{R}^p$  be a random vector and its joint CDF  $F_{\mathbf{X}}$  with marginals  $\{F_{X_j}\}_{j=1}^p$ , there exists a copula  $C : [0, 1]^p \rightarrow [0, 1]$ , such that:*

$$F_{\mathbf{X}}(x_1, \dots, x_p) = C\left(F_{X_1}(x_1), \dots, F_{X_p}(x_p)\right). \quad (1.8)$$

*If the marginals  $F_{X_i}$  are continuous, then this copula is unique.*

Theorem 1 expresses the joint CDF by combining marginal CDFs and a copula, which is practical for sampling joint distributions. Conversely, the copula can be defined by using the joint CDF and the marginal CDFs:

$$C(u_1, \dots, u_p) = F_{\mathbf{X}}(F_{X_1}^{-1}(u_1), \dots, F_{X_p}^{-1}(u_p)) \quad (1.9)$$

This equation allows us to extract a copula from a joint distribution by knowing its marginals. Additionally, copulas are invariant under creasing transformations. This property is important to understand the use of rank transformation to display the copula without the marginal effects.

[define the copula density and address the warning on the confusion.]

Identically to the univariate continuous distributions, a large catalog of families of copulas exists (e.g., independent, Normal, Clayton, Frank, Gumbel copula, etc.).

[define the independent copula, also called a product of marginals.]

To infer a joint distribution, this theorem divides the fitting problem into two independent problems: fitting the marginals and fitting the copula. Provided a dataset, this framework allows the combination of a parametric (or nonparametric) fit of marginals with a parametric (or nonparametric) fit of the copula.

Appendix A details the main techniques to estimate marginal distributions. Then, Appendix B introduces different nonparametric methods to infer a copula, including the empirical Bern-

stein copula and the Beta copula. The adequation between a fitted probabilistic model and a dataset should be validated, therefore, appendices A and B respectively present visual and quantitative tools for goodness-of-fit evaluation.

To infer a joint distribution over a dataset, the analyst should determine a fitting strategy. Smart data visualization helps to choose the fitting methods susceptible to be relevant to the problem. The following points can be checked at this early stage:

- Is the distribution unimodal? If not, mixtures methods or nonparametric models might be required;
- Is the validity domain restrictive? If so, specific families of parametric distributions can be chosen or truncations can be applied;
- Is the dimension high? If the dimension is too high: tensorize the distribution as much as possible.
- Is the dependence structure complex? Graphically, the dataset in the ranked space gives an empirical description but some independence tests exist as well.

## 1.4 Central tendency uncertainty propagation

The previous section aimed at building a probabilistic model of the uncertainties considering the knowledge available. This one will introduce diverse forward propagation of uncertainty through a numerical model. This step is hereafter qualified as “global” because the analysis of the resulting output random variable will particularly focus on its central tendency (i.e., expected value and variance). This approach contrasts with the uncertainty propagation dedicated to rare event estimation, which will be introduced in the next section (e.g., for a reliability or certification problem).

The difficulties related to any uncertainty propagation mostly arise from the practical properties of the numerical model. Its potential high dimension, low regularity and nonlinearities each represent a challenge. These studies rely on a finite number of observations which depends on the computational budget the analyst can afford. This forward propagation might be a finality of the uncertainty quantification, but keep in mind that it fully stands on an accurate uncertainty modeling. Uncertainty propagation should be perceived as a standardized process with modular bricks, on which the “garbage in, garbage out” concept fully applies.

This section introduces the main methods of global uncertainty propagation. Outlining the strong links between numerical integration (i.e., Lebesgue integration or central tendency estimation) and numerical design of experiments.



### 1.4.1 Numerical integration

Forward uncertainty propagation aims at integrating a measurable function  $g : \mathcal{D}_X \rightarrow \mathbb{R}$  with respect to a probability measure  $\mu$ . Numerical integration brings algorithmic tools to help the resolution of this probabilistic integration (i.e., Lebesgue integration).

In practice, this integral is approximated by summing a finite  $n$ -sized set of realizations  $\mathbf{y}_n = \{g(\mathbf{x}^{(1)}), \dots, g(\mathbf{x}^{(n)})\}$  from a set of input samples  $\mathbf{X}_n = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ . A *quadrature* establishes a rule to select the input samples  $\mathbf{X}_n$  (also called nodes), and an associated set of weights  $\mathbf{w}_n = \{w_1, \dots, w_n\} \in \mathbb{R}^n$ . The approximation given by a quadrature rule is defined as a weighted arithmetic mean of the realizations:

$$I_\mu(g) := \int_{\mathcal{D}_X} g(\mathbf{x}) d\mu(\mathbf{x}) \approx \sum_{i=1}^n w_i g(\mathbf{x}^{(i)}). \quad (1.10)$$

For a given sample size  $n$ , our goal is to find a set of tuples  $\{\mathbf{x}^{(i)}, w_i\}_{i=1}^n$  (i.e., quadrature rule), giving the best approximation of our quantity. Ideally, the approximation quality should be fulfilled for a wide class of integrands. Most quadrature rules only depend on the measure space  $(\Omega, \mathcal{A}, \mu)$ , regardless of the integrand values. In the context of a costly numerical model, this property allows the analyst to massively distribute the calls to the numerical model.

This section aims at presenting the main multivariate numerical integration techniques. These methods have very different properties: some are deterministic and some are aleatory; some are sequential (or nested) some are not; some are victims of the curse of dimensionality and some are not. A summary table of the different methods and their respective properties is proposed [\[add ref to table\]](#).

#### Classical multivariate deterministic quadrature

Historically, quadrature methods have been developed for univariate integrals. The Gaussian rule and the Fejér-Clenshaw-Curtis rule are two univariate deterministic quadratures that will be briefly introduced.

Gaussian quadrature is a powerful univariate quadrature building together a set of irregular nodes and a set of weights. The computed weights are positive, which ensures a numerically stable rule even for large sample sizes.

Different variants of rules exist, the most famous being the Gauss-Legendre quadrature. In this case, the function  $g$  to be integrated with respect to the uniform measure on  $[-1, 1]$  is approximated by Legendre polynomials. Considering the Legendre polynomial of order  $n$ , denoted  $l_n$ , the quadrature nodes  $x^{(i)}_{i=1}^n$  are given by the polynomial roots. The respective weights are given by the following formula:

$$w_i = \frac{2}{\left(1 - (x^{(i)})^2\right) \left(l'_n(x^{(i)})\right)^2}. \quad (1.11)$$



This rule guarantees a very precise approximation provided that the integrand is well-approximated by a polynomial of degree  $2n - 1$  or less on  $[-1, 1]$ . This rule is deterministic but not sequential, meaning that two rules with sizes  $n_1$  and  $n_2$ ,  $n_1 < n_2$  will not be nested. However, a sequential extension is proposed by the Gauss-Kronrod rule [add ref], offering lower accuracy.

To overcome this practical drawback, the Fejér-Clenshaw-Curtis is a fully nested rule with mostly equivalent accuracy as Gaussian quadrature. This method is usually presented to integrate a function with respect to the uniform measure on  $[-1, 1]$  and starts with a change of variables:

$$\int_{-1}^1 g(x) dx = \int_0^\pi g(\cos(\theta)) \sin(\theta) d\theta \quad (1.12)$$

This expression can be written as an expansion of the integrand using cosine series. Moreover, cosine series are closely related to the Chebyshev polynomials of the first kind. Féjer's "first rule" [LLOYD N. TREFETHEN 2008] relies on the Chebyshev polynomials roots as nodes  $x^{(i)} = \cos(\theta^{(i+1/2)})$ , and the following weights:

$$w_i = \frac{2}{n} \left( 1 - 2 \sum_{j=1}^{\lfloor n/2 \rfloor} \frac{1}{4j^2 - 1} \cos(j\theta^{(2i+1)}) \right) \quad (1.13)$$

These two univariate integration schemes are both very efficient on a wide panel of functions. Yet, Fejér-Clenshaw-Curtis is sequential and offers easy implementations, benefitting from powerful algorithms such as *fast Fourier transform*.

[add 1D Gaussian nodes and Féjer nodes and show that some are nested and some are not]

UQ problems are rarely unidimensional, but one can build a multivariate quadrature rule by defining the tensor product (also called full grids) of univariate rules. This exhaustive approach quickly shows its practical limits as the problem's dimension increases. Alternatively, sparse multivariate quadratures (e.g., Smolyak sparse grid) propose a more efficient exploration of the joint domain (see [add ref] for more details).

## Monte Carlo methods

Monte Carlo methods were initially developed in the 1940's to solve problems in neutronics. Ever since, this frequentist techniques have been applied to the resolution of the Lebesgue integral. To integrate a function  $g$  against a measure  $\mu$ , it randomly generates points following the input measure. The integral is estimated by taking the uniform arithmetic mean of the images of these nodes obtained by this random process.

This purely aleatory method requires to be able to generate points following a given distribution. To do so, the most common approach is to first generate a sequence of random points uniformly on  $[0, 1]$ . Then, inverse transform of the CDF is applied on marginals before adding the possible dependence effects using the Sklar theorem [add pointer]. These sequences mimic a actual uniform randomness but are in fact generated by deterministic algorithms (also called pseudorandom number generators). Pseudorandom algorithms generate a sequence of numbers with a very large, but finite length. This sequence can be exactly repeated by fixing the same

initial point, also called *pseudorandom seed*. Most programming languages use the Mersenne Twister pseudorandom generator [cite MATSUMOTO], offering a very long period (around  $4.3 \times 10^{6001}$  iterations).

Formally, the “Vanilla” Monte Carlo (sometimes called “crude” Monte Carlo) method uses a set of i.i.d samples  $\mathbf{X}_n = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$  following the joint distribution of  $\mu$ . The Monte Carlo estimator of the integral [add ref] is given by:

$$I_\mu(g) \approx \bar{y}_n^{MC} = \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}^{(i)}). \quad (1.14)$$

By construction, the law of large numbers makes this estimator unbiased, however it converges relatively slowly. Considering the images of the sample  $\mathbf{X}_n$ , one can also estimate the variance of the output random variable  $\hat{\sigma}_Y^2$ . The variance of the Monte Carlo estimator [add ref to above] results from a manipulation of the central tendency theorem:

$$\text{Var}(\bar{y}_n^{MC}) = \frac{1}{\sqrt{n}} \text{Var}(\hat{\sigma}_Y^2). \quad (1.15)$$

This estimator also comes with theoretical confidence intervals at  $\alpha\%$ , regardless of the output distribution:

$$I_\mu(g) \in \left[ \bar{y}_n^{MC} - q_\alpha \frac{\text{Var}(\hat{\sigma}_Y^2)}{\sqrt{n}}, \bar{y}_n^{MC} + q_\alpha \frac{\text{Var}(\hat{\sigma}_Y^2)}{\sqrt{n}} \right], \quad (1.16)$$

where  $q_\alpha$  is the  $\alpha$ -quantile of the standard normal distribution. Monte Carlo presents the advantage of being a universal method, with no bias and strong convergence guarantees. Moreover, it is worth noting that its convergence properties does not depend on the dimension of the input domain (i.e., no curse of dimension). The main limit of crude Monte Carlo is its convergence speed, more recent method aim at keeping the interesting properties if this technique while making it more efficient<sup>2</sup>. Among them, one can mention two variance-reduction techniques: importance sampling and multi-level Monte Carlo. [Quick definition and references].

[Remark: when it is hard to sample from the distribution, Markov Chain Monte Carlo is an option for UP.]

### Quasi-Monte Carlo and Koksma-Hlawka inequality

- Numerical integration scheme over  $[0, 1]^p$  wrt uniform measure.
- Synthesis: deterministic quadratures are subject to the curse of dim while MC has a slow convergence rate.
- Quasi-MC combines the good properties from both deterministic and MC integration schemes by selecting nodes following a low discrepancy sequence.

## 1.4.2 Numerical design of experiments

### Space-filling metrics

[MinMax / PhiP / MaxMin / Discrepancies]

### “Good” properties

[Curse of dim / Projections in sub-spaces / Sequential / Deterministic]

### Monte Carlo and quasi-Monte Carlo designs

Monte Carlo sampling is the simplest way to propagate uncertainty (either for central tendency or probability estimation). Using a pseudo-random generator, one can generate  $n \in \mathbb{N}$  i.i.d realizations  $\{X^{(i)}\}_{i=1,\dots,n}$  of a random vector  $X$ .

[Remark on randomized quasi-Monte Carlo: paper on quantiles]

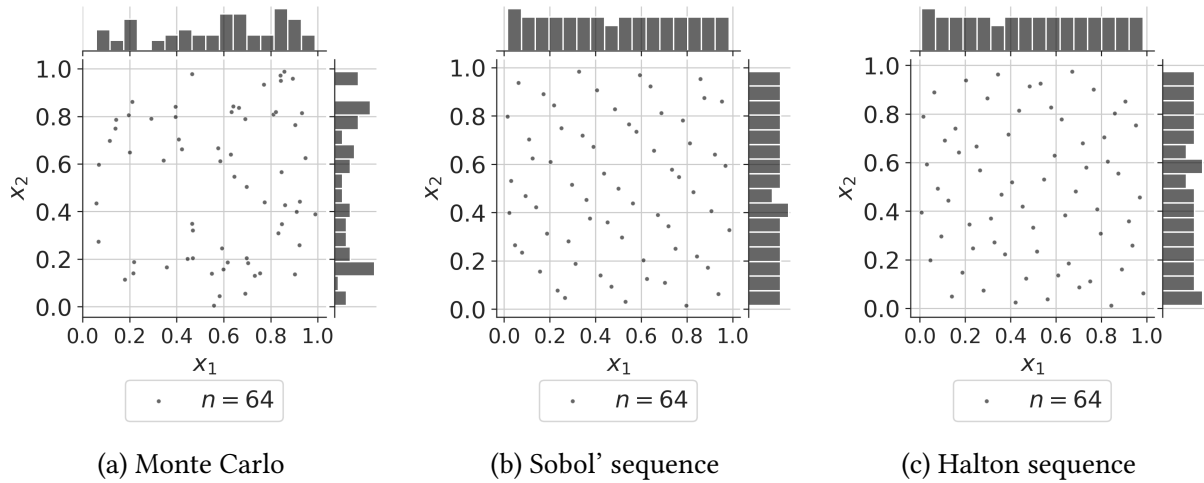


Fig. 1.3 Monte Carlo and quasi-Monte Carlo designs ( $n = 256$ )

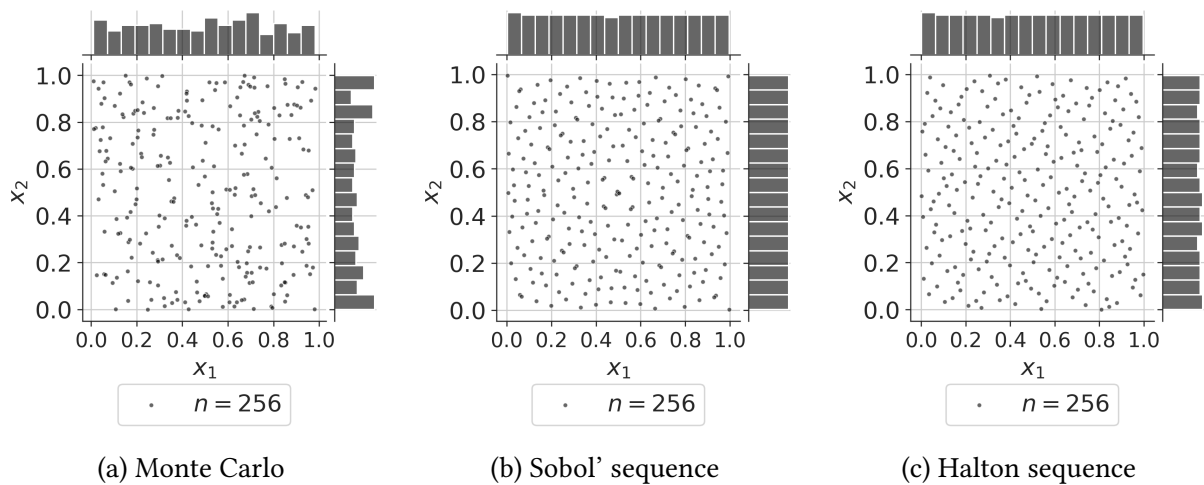


Fig. 1.4 Monte Carlo and quasi-Monte Carlo designs ( $n = 256$ )

## Latin hypercube sampling

The LHS is a method introduced in 1979 (Mckay et al., 1979), initially created to compute multivariate numerical integrals. The main idea is that to make sure that the information is not redundant, the distribution of each sub-projection of the domain should be as uniform as possible. To force this concept, each variable's domain is divided into  $n$  identical segments creating a square lattice over the domain. This lattice is composed of  $n^p$  squared elements. Among the  $n^p$  elements, the LHS samples  $n$  points (only one point by element) with the constraint of having one and only one sample for each segment of each variables. This ensures a good distribution over the sub-projections of the domain. Inside each selected element in the square lattice, the points can be placed in the center or randomly. [Illustration LHS vs. Monte Carlo]

[Illustration pathological LHS (diagonal design)]

## Optimized Latin hypercube sampling

Franco (2008) discusses how LHS can have very poor space-filling properties in some cases, which can be prevented by adding an optimization step. Knowing that an LHS can be permuted to another one, one can find the permutation on an LHS that would optimize a given space-filling criteria. The optimization according to the two following criteria is discussed in Damblin et al. (2013) and the sub-projections issue is also raised. Note that LHS designs are very efficient at exploring a domain but handle poorly some complex dependency models such as copulas.

[Present MaxPro, Uniform Projection designs. The way of optimizing is not the topic. It is more about what are the right criteria]

### 1.4.3 Central tendency estimation

#### Iso-probabilistic transformation

Central tendency estimation is a probabilistic integration

## 1.5 Reliability-oriented uncertainty propagation

### 1.5.1 Problem formalization

Limit-state function, failure event and domain

Risk measures [Failure probability, quantile, super-quantile]

### 1.5.2 Rare event estimation methods

[Why are the previous sampling methods not suited for rare events?]

**FORM/SORM**

**Monte Carlo**

**Importance sampling**

**Adaptive sampling (SS/NAIS/IS-CE/Moving particles)**

## **1.6 Sensitivity analysis**

### **1.6.1 Global sensitivity analysis**

### **1.6.2 Reliability-oriented sensitivity analysis**

## **1.7 Metamodeling**

Note that the calibration error is often larger than the metamodeling error.

### **1.7.1 Global metamodel**

### **1.7.2 Contour finding for rare-event estimation**

## **1.8 Conclusion**



# Introduction to wind turbine modeling and design

---

2.1	Introduction . . . . .	24
2.2	Wind turbine modeling . . . . .	24
2.2.1	Synthetic wind generation [TurbSim, Kaimal spectrum] . . . . .	24
2.2.2	Synthetic wave generation . . . . .	24
2.2.3	Aerodynamic interactions . . . . .	24
2.2.4	Servo-Hydro-Aero-Elastic wind turbine simulation [DIEGO] . . . . .	24
2.2.5	Soil modeling . . . . .	24
2.2.6	Wake modeling [FarmShadow] . . . . .	24
2.3	Recommended design practices . . . . .	24
2.3.1	Design load cases . . . . .	24
2.3.2	Dynamic response design . . . . .	24
2.3.3	Fatigue response design . . . . .	24
2.4	Uncertain inputs . . . . .	24
2.4.1	Environmental inputs . . . . .	24
2.4.2	System inputs . . . . .	24
2.4.3	Probabilistic fatigue assessment . . . . .	24
2.5	Conclusion . . . . .	24

---

## **2.1 Introduction**

## **2.2 Wind turbine modeling**

### **2.2.1 Synthetic wind generation [TurbSim, Kaimal spectrum]**

### **2.2.2 Synthetic wave generation**

### **2.2.3 Aerodynamic interactions**

### **2.2.4 Servo-Hydro-Aero-Elastic wind turbine simulation [DIEGO]**

### **2.2.5 Soil modeling**

### **2.2.6 Wake modeling [FarmShadow]**

## **2.3 Recommended design practices**

### **2.3.1 Design load cases**

### **2.3.2 Dynamic response design**

### **2.3.3 Fatigue response design**

## **2.4 Uncertain inputs**

### **2.4.1 Environmental inputs**

### **2.4.2 System inputs**

### **2.4.3 Probabilistic fatigue assessment**

## **2.5 Conclusion**



# References

- Ajenjo, A. (2023). *Info-gap robustness assessment of reliability evaluations for the safety of critical industrialsystems*. PhD thesis, Université Bourgogne Franche-Comté.
- Chabridon, V., Balesdent, M., Perrin, G., Morio, J., Bourinet, J.-M., and Gayton, N. (2021). Global reliability-oriented sensitivity analysis under distribution parameter uncertainty. *Mechanical Engineering under Uncertainties: From Classical Approaches to Some Recent Developments*, pages 237–277.
- Damblin, G., Couplet, M., and Iooss, B. (2013). Numerical studies of space-filling designs: Optimization of Latin Hypercube Samples and subprojection properties. *Journal of Simulation*, 7.
- Der Kiureghian, A. and Ditlevsen, O. (2009). Aleatory or epistemic? Does it matter? *Structural Safety*, 31(2):105–112.
- Franco, J. (2008). *Planification d’expériences numériques en phase exploratoire pour la simulation des phénomènes complexes*. PhD thesis, Ecole Nationale Supérieure des Mines de Saint-Etienne.
- Joe, H. (1997). *Multivariate Models and Multivariate Dependence Concepts*. Chapman and Hall.
- Joe, H. and Kurowicka, D. (2011). *Dependence modeling: vine copula handbook*. World Scientific.
- Lasserre, M. (2022). *Apprentissages dans les réseaux bayésiens à base de copules non-paramétriques*. PhD thesis, Sorbonne Université.
- Lataniotis, C. (2019). *Data-driven uncertainty quantification for high-dimensional engineering problems*. PhD thesis, ETH Zürich.
- Mckay, M., Beckman, R., and Conover, W. (1979). A Comparison of Three Methods for Selecting Vales of Input Variables in the Analysis of Output From a Computer Code. *Technometrics*, 21:239 – 245.
- Sancetta, A. and Satchell, S. (2004). The Bernstein copula and its applications to modeling and approximations of multivariate distributions. *Econometric Theory*, 20(3):535–562.
- Segers, J., Sibuya, M., and Tsukahara, H. (2017). The empirical beta copula. *Journal of Multivariate Analysis*, 155:35–51.



## Univariate distribution fitting

This appendix recalls the main methods to infer a univariate distribution considering a  $n$ -sized i.i.d sample  $X_n = \{x^{(1)}, \dots, x^{(n)}\} \in \mathbb{R}^n$ . The goal is to use this finite set of observations of the random variable  $X$  to approach its underlying distribution by an estimated distribution. The inference techniques are split into two main groups, the methods assuming that the underlying distribution belongs to a family of parametric distributions are called parametric. Otherwise, the fitting method falls into the nonparametric group. Nonparametric methods often require a larger amount of data but allow more flexibility. In fact, nontrivial distributions (e.g., multimodal) might be easier to model using nonparametric approaches. To assess the quality of this estimation regarding the sample, a panel of goodness-of-fit methods are proposed [\[add ref\]](#), this appendix recalls a few of them. Note that the following tools can be used to estimate the marginals of a multivariate distribution.

### A.1 Main parametric methods

#### Moments method

The moments method aims at looking for a parametric distribution with density  $f_X(\theta)$ , whose first moments (e.g.,  $m(\theta)$  and  $\sigma^2(\theta)$ ) match the empirical moments of the sample  $X_n$  (e.g.,  $\widehat{m}_{X_n}$  and  $\widehat{\sigma}^2$ ). After computing the empirical moments:

$$\widehat{m}_{X_n} = \frac{1}{n} \sum_{i=1}^n x^{(i)}, \quad \widehat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n \left( x^{(i)} - \widehat{m}_{X_n} \right)^2, \quad (\text{A.1})$$

one can solve the system of equations  $(m(\theta) = \widehat{m}_{X_n}; \sigma^2(\theta) = \widehat{\sigma}^2)$  to determine the optimal set of parameters  $\theta$  in this situation. Some families of distributions are more suited to this method (i.e., ) because of the analytical expression of their moments. Moreover, this technique is sensitive to the possible biases in the estimation of the sample moments.

## Maximum likelihood estimation

Maximum likelihood estimation (MLE) is a popular alternative the moments method. Similarly, it aims at maximizing a given correspondence metric between the dataset  $X_n$  and a parametric distribution with density  $f_X(\theta)$ . This metric is the *likelihood* function, defined as:

$$\mathcal{L}(\theta|X_n) = \prod_{i=1}^n f_X(x^{(i)}; \theta), \quad (\text{A.2})$$

with the PDF taking the set of parameters  $\theta$  written:  $f_X(x^{(i)}; \theta)$ . For numerical reasons, the optimization is often performed on the natural logarithm of the likelihood function, called *log-likelihood*. The goal is then finding the optimal vector  $\hat{\theta}^*$  of parameters minimizing the following expression:

$$\hat{\theta}^* = \arg \min_{\theta \in \mathcal{D}_\theta} \left( - \sum_{i=1}^n \ln(f_X(x^{(i)}; \theta)) \right). \quad (\text{A.3})$$

Remark that the quick analytical results from the moment method can be used as a starting point of the MLE optimization. [Asymptotic behaviors of this method are decribed in: add ref]  
[This method can be applied to censored data in the field of survival analysis. Add ref]

**Example 1.** Considering a small set of observations  $X_n = \{1, 2, 3, 4, 6\}$ , the following figure xx represents

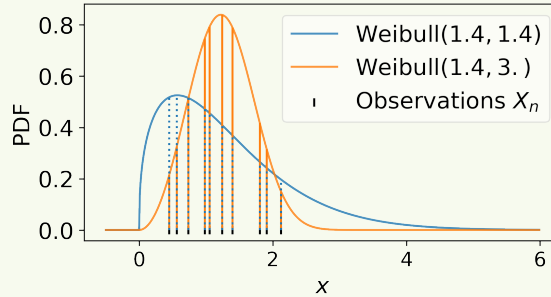


Fig. A.1 Adequation of two different Weibull models using their likelihood with a sample of observations (black crosses).

## A.2 Main nonparametric methods

### Empirical CDF and histogram

The empirical CDF is a cumulative stair-shaped representation of the sorted sample  $X_n$ :

$$\hat{F}_X(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x \geq x^{(i)}\}}. \quad (\text{A.4})$$

A histogram consists in sorting and gathering the observations in a sample  $X_n$  into a finite number of categories. These categories are called bins and each regroups the same number of observations (identical binwidth). The number of bins is the only tuning parameter of this method. Its definition has a great impact on the visual consistency of the plot, therefore, many rules exist to define it. Note that the empirical CDF can be seen as a cumulative histogram with the number of bins equal to the number of observations.

## Kernel density estimation

Kernel density estimation (KDE) is a nonparametric method, it estimates a PDF by weighing a sample of observations  $X_n$  with kernels. After setting a kernel  $k : \mathbb{R} \rightarrow \mathbb{R}_+$  and a scaling parameter  $h > 0$ , also called bandwidth, the kernel density estimator is defined as:

$$\hat{f}_X(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x - x^{(i)}}{h}\right) \quad (\text{A.5})$$

Different types of kernels are used for KDE, such as the uniform, triangular, squared exponential or Epanechnikov. The choice of bandwidth results in a bias-variance trade-off, that has been extensively discussed in the literature [\[add ref\]](#).

**Example 2.** Considering a small set of observations  $X_n = \{1, 2, 3, 4, 6\}$ , the following figure xx represents three fits obtained by.

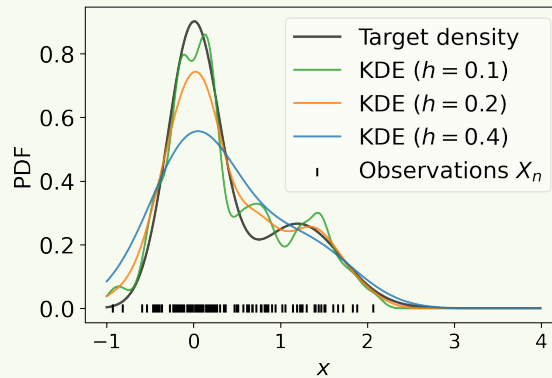


Fig. A.2 Fit of a bimodal density by KDE using different tuning parameters.

## Main goodness-of-fit methods

### Penalized likelihood criteria

Two quantitative goodness-of-fit criteria are commonly used to assess parametric inference: the *Akaike information criterion* (AIC) and the *Bayesian information criterion* (BIC). The likelihood as a goodness-of-fit criterion should only be applied to the same family of distributions. Otherwise, the comparison would unfairly advantage distributions with a large number of degrees of

freedom. The two following criteria are metrics based on the likelihood with a correction related to the number of degrees of freedom of the distribution. Moreover, let us remind that more flexible models will require more data to provide a robust estimation.

The AIC and BIC are expressed as follows:

$$\text{AIC} = \frac{-2\ln(\mathcal{L}(\theta|X_n))}{n} + \frac{2q}{n}, \quad \text{BIC} = \frac{-2\ln(\mathcal{L}(\theta|X_n))}{n} + \frac{q\ln(n)}{n}, \quad (\text{A.6})$$

with the likelihood  $\mathcal{L}(\theta|X_n)$  and the number of distribution's number degrees of freedom denoted  $q$ . The second term adds a penalty depending on the number of parameters. The best inference will be given by the model with the smallest AIC or BIC. Note that an additional correction can be applied in a small data context.

## Kolmogorov-Smirnov adequacy test

### Quantile-quantile plot

The quantile-quantile plot (also called QQ-plot) is a graphical tool providing a qualitative check of the goodness of fit. It compares the CDF of the fitted model with the empirical CDF of the sample  $X_n$ . To do so, it represents a scatterplot of the empirical quantiles (i.e., the ranked observations), against the quantiles of the fitted model at the levels  $\{\alpha^{(i)}\}_{i=1}^n = \{\widehat{F}_X(x^{(i)})\}_{i=1}^n$ . The following [figure xx] is a QQ-plot of the model fitted in [Example xx]. The closer the scatter plot gets to the first bisector line the best the fit is.

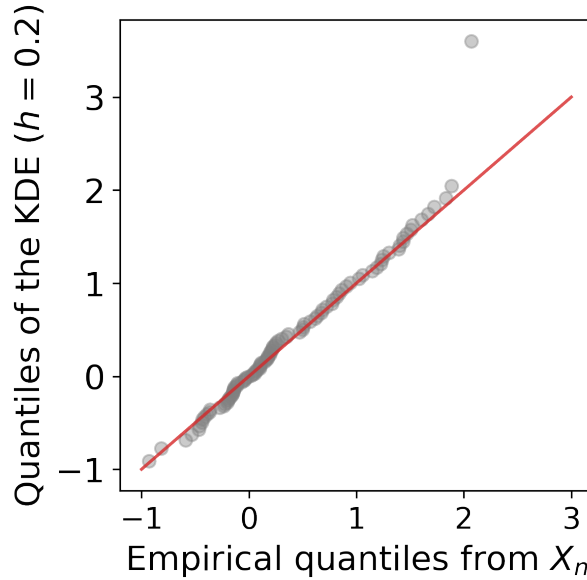


Fig. A.3 QQ-plot between the data from Example 2 and a KDE model.

## Nonparametric copula estimation

[update CDF notations]

[Change EBC notations using  $h$  for the polynomial orders]

When the distribution's dimension is higher than two, one can perform a parametric fit using vine copulas (Joe and Kurowicka, 2011), implying the choice of multiple types of parametric copulas. Otherwise, nonparametric fit by multivariate kernel density estimation (KDE) presents a computational burden as soon as the dimension increases (Chabridon et al., 2021). Since univariate marginals are usually well-fitted with nonparametric tools (e.g., KDE), let us introduce an effective nonparametric method for copula fitting.

### B.1 Empirical copula

In practice, considering a sample  $\mathbf{X}_n = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\} \in \mathbb{R}^{np}$  and the associated ranked sample  $\mathbf{R}_n = \{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(n)}\}$ , the corresponding empirical copula writes:

$$C_n(\mathbf{u}) := \frac{1}{n} \sum_{i=0}^n \prod_{j=1}^p \mathbb{1} \left\{ \frac{r_j^{(i)}}{n} \leq u_j \right\}, \quad (\text{B.1})$$

with  $\mathbf{u} = (u_1, \dots, u_d) \in [0, 1]^d$ . In the following, the polynomial order is set as equal in each dimension:  $\{m_i = m\}_{i=1}^d$ .

### B.2 Empirical Bernstein & Beta copula

Copulas are continuous and bounded functions defined on a compact set (the unit hypercube). Bernstein polynomials allow us to uniformly approximate as closely as desired any continuous and real-valued function defined on a compact set (Weierstrass approximation theorem). Therefore, they are good candidates to approximate unknown copulas. This concept was introduced as *empirical Bernstein copula* (EBC) by Sancetta and Satchell (2004) for applications in economics and risk management. Later on, Segers et al. (2017) offered further asymptotic studies. Formally,

the multivariate Bernstein polynomial for a function  $C : [0, 1]^d \rightarrow \mathbb{R}$  on a grid over the unit hypercube  $G := \left\{ \frac{0}{m_1}, \dots, \frac{m_1}{m_1} \right\} \times \dots \times \left\{ \frac{0}{m_d}, \dots, \frac{m_d}{m_d} \right\}$ ,  $\mathbf{m} = (m_1, \dots, m_d) \in \mathbb{N}^d$ , writes:

$$B_{\mathbf{m}}(C)(\mathbf{u}) := \sum_{t_1=0}^{m_1} \dots \sum_{t_d=0}^{m_d} C\left(\frac{t_1}{m_1}, \dots, \frac{t_d}{m_d}\right) \prod_{j=1}^d P_{m_j, t_j}(u_j), \quad (\text{B.2})$$

with  $\mathbf{u} = (u_1, \dots, u_d) \in [0, 1]^d$ , and the Bernstein polynomial  $P_{m,t}(u) := \frac{t!}{m!(t-m)!} u^m (1-u)^{t-m}$ . Notice how the grid definition implies the polynomial's order. When  $C$  is a copula, then  $B_{\mathbf{m}}(C)$  is called “Bernstein copula”. Therefore, the empirical Bernstein copula is an application of the Bernstein polynomial in Eq. (B.2) to the so-called “empirical copula”. [add a sentence to mean to refer to the previous subsection]

Theoretically, the tuning parameter can be optimized to minimize an “Mean Integrated Squared Error” (MISE), leading to a bias-variance tradeoff. Formally, the MISE of the empirical Bernstein copula  $B_{\mathbf{m}}(C_n)$  is defined as follows:

$$\mathbb{E}[\|B_{\mathbf{m}}(C_n) - C\|_2^2] = \mathbb{E}\left[\int_{\mathbb{R}^d} (B_{\mathbf{m}}(C_n)(\mathbf{u}) - C(\mathbf{u}))^2 d\mathbf{u}\right]. \quad (\text{B.3})$$

Then, [Sancetta and Satchell \(2004\)](#) prove in their Theorem 3 that:

- $B_{\mathbf{m}}(C_n)(\mathbf{u}) \rightarrow C(\mathbf{u})$  for any  $u_j \in ]0, 1[$  if  $\frac{m^{d/2}}{n} \rightarrow 0$ , when  $m, n \rightarrow \infty$ .
- The optimal order of the polynomial in terms of MISE is:  $m \lesssim m_{\text{MISE}} = n^{2/(d+4)}, \forall u_j \in ]0, 1[$ . The sign  $\lesssim$  means “less than or approximately”.

Let us remark that in the special case  $m = n$ , also called the “Beta copula” in [Segers et al. \(2017\)](#), the bias is very small while the variance gets large. To illustrate the previous theorem, [Lasserre \(2022\)](#) represents the evolution of the  $m_{\text{MISE}}$  for different dimensions and sample sizes (see Fig. B.1). In high dimension, the values of  $m_{\text{MISE}}$  tend towards one, which is equivalent to the independent copula. Therefore, high-dimensional problems should be divided into a product of smaller problems on which the EBC is tractable. Provided a large enough learning set  $\mathbf{X}_n$ , KDE fitting of marginals combined with EBC fitting of the copula delivers good results even on complex dependence structures. Moreover, EBC provides an explicit expression, making a Monte Carlo generation of i.i.d. samples simple.

### B.3 Goodness-of-fit

[Mention the vine copulas and how we want to only use nonparametric methods here.]

[Tails correlation / Kendall plot]



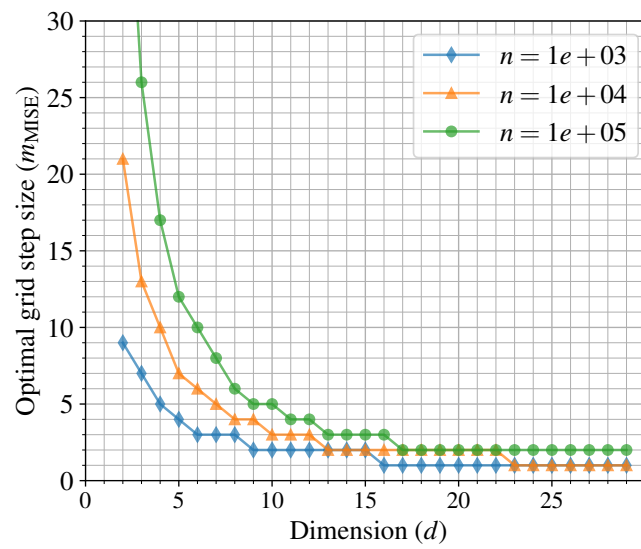


Fig. B.1 Evolution of  $m_{\text{IMSE}}$  for different dimensions and sample sizes.



# Appendix C

## Rare event estimation algorithms



Appendix	<b>D</b>
----------	----------

## Résumé étendu de la thèse