

# Uncertainty Quantification

Personal summary by Gilberto Lem

Summer 2018

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Uncertainty Quantification . . . . .	2
1.2	Basics of Probability Theory and Statistics . . . . .	2
<b>2</b>	<b>Sampling Methods</b>	<b>5</b>
2.1	Monte Carlo . . . . .	5
2.2	Variance-Reduction Techniques . . . . .	5
2.3	Quasi Monte-Carlo . . . . .	6
<b>3</b>	<b>Basic Interpolation and Quadrature</b>	<b>8</b>
3.1	Interpolation . . . . .	8
3.2	Quadrature . . . . .	8
<b>4</b>	<b>Polynomial Chaos</b>	<b>10</b>
4.1	Multivariate Polynomial Chaos Expansion . . . . .	11
<b>5</b>	<b>Sparse Grids</b>	<b>12</b>
<b>6</b>	<b>Sensitivity Analysis</b>	<b>12</b>
<b>7</b>	<b>Random Fields</b>	<b>12</b>
<b>8</b>	<b>Software</b>	<b>12</b>
<b>9</b>	<b>Bayesian Inverse Problems</b>	<b>12</b>

# 1 Introduction

## 1.1 Uncertainty Quantification

Uncertainty Quantification determines how likely certain outcomes are if some aspects of the system are not exactly known.

The number of uncertain parameters is called **stochastic dimensionality**.

UQ in general is slow, because one needs to run the model several times to get good results. To speed up the process there are three options:

- Make the deterministic model faster
- Use surrogate model instead of expensive high-fidelity model (tradeoff for accuracy)
- Use clever UQ method with low amount of evaluation points (maybe trade-off for accuracy)

## 1.2 Basics of Probability Theory and Statistics

A **probability space** is a triplet  $(\Omega, \mathcal{F}, \mathcal{P})$ , where:

- $\Omega$  is the sample space, the set of all possible single outcomes
- $\mathcal{F}$  ( $\sigma$ -algebra) is the combinatorial of all the elements of  $\Omega$
- $\mathcal{P}$  is the probability measure,  $\mathcal{F} \rightarrow [0, 1]$

### 1.2.1 Univariate Concepts

A **random variable** is a function  $X$  which maps  $\Omega \rightarrow \mathbb{R}$ . It represents an outcome of a random phenomenon as a number.

We can represent the probability of obtaining  $X = x$ , being  $x$  a real number, as  $f_X(x)$ . If  $X$  is continuous, this function is called the **probability density function (PDF)**, if  $X$  is discrete, is called the **probability mass function (PMF)**.

The value of  $f_X$  has to be always greater or equal than zero, and its accumulate value over all its domain has to be 1.

Every random variable has an associated **Cumulative Distribution Function (CDF)**, represented by  $F_X(x)$ . It represents the probability that the random variable  $X$  has a value less or equal than a specified  $x$ , i.e.  $F_X(x) = \mathcal{P}\{X \leq x\}$ . This can be represented as an integral:

$$F_X(x) = \int_{-\infty}^x f_X(s) ds$$

The **expectation** of a continuous random variable  $X$  is defined as:

$$\mu := \mathbb{E}[X] = \int x f_X(x) dx$$

The **variance** of  $X$  is the expected value of the square of the variation of  $X$  with respect to its mean:

$$\sigma^2 := \text{Var}(X) = \int (x - \mathbb{E}[X])^2 f_X(x) dx = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

If  $X$  is discrete, the **sample mean** and **sample variance** are defined:

$$\bar{X} = \frac{1}{n} \sum_i X_i, \quad S^2 = \frac{1}{n-1} \sum_i (X_i - \bar{X})^2$$

### 1.2.2 Multivariate Concepts

$X$  can be generalized to a  $n$ -dimensional random vector  $\mathbf{X}$ , which maps  $\Omega \rightarrow \mathbb{R}^n$ .  $\mathbf{X}$  consists of  $n$  random variables.

The **covariance** of two random variables  $X$  and  $Y$  is

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

Covariance is high when  $X$  and  $Y$  are “similar”.

If you normalize this covariance you get the **Pearson correlation coefficient**:

$$\rho_{XY} = \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

which have values  $[-1, 1]$ .

**Independence** Let  $A$  and  $B$  denote two events. If the probability for  $A$  to occur is the same if  $B$  happens or if  $B$  does not happen,  $A$  and  $B$  are **independent**. If additionally, both events have the same distribution of probability, they are **independent and identically distributed (iid)** random variables.

### 1.2.3 Biased vs. unbiased estimators

An **estimator** is a procedure to construct estimates for a quantity  $q$  based on random samples. An **estimate** is a realization of the estimator. Example estimators are the mean and the variance.

An estimator is called **biased** if its expected value is not the same as the parameter to be estimated.

Correspondingly, an **unbiased** estimator is when the expected value of the estimator of a sample is equal to the real value of the parameter we are estimating.

For an unbiased mean, if we take the mean of a sample, we would expect the same as the mean of the population.

In the following equations,  $\bar{X}$  and  $S^2$  represent unbiased estimators:

$$\mathbb{E}[\bar{X}] = \mu, \quad \mathbb{E}[S^2] = \sigma^2$$

If we take the following definitions of  $\bar{X}$  and  $S^2$  :

$$\bar{X} = \frac{1}{n} \sum_i X_i \quad S^2 = \frac{1}{n} \sum_i (X_i - \bar{X})^2$$

$\bar{X}$  is unbiased but  $S^2$  is biased. To correct for this, we multiply our definition of sample variance by  $n/(n-1)$  to end up with:

$$S^2 = \frac{1}{n-1} \sum_i (X_i - \bar{X})^2$$

Note that numpy's var function computes the biased version of the variance. When we use that function we have to correct it.

## 2 Sampling Methods

Sampling methods consist in select a subset of “individuals” from a statistical “population”, to estimate characteristics of the whole population.

### 2.1 Monte Carlo

**Monte Carlo Sampling** consists in generating samples with a pseudo-random number generator from a given probability distribution. After that compute deterministic computation with these samples and aggregate the results.

One application of Monte Carlo Sampling is Monte Carlo Integration. It computes the value of an integral of a function over a domain via sampling of the function. The error of the approximation is given by the **Root Mean Squared Error**:

$$RMSE := \sqrt{\mathbb{E}[(I - \hat{I}_f)^2]}$$

Or with the equivalent expression:

$$RMSE \approx \frac{\hat{\sigma}_f}{\sqrt{N}}$$

with

$$\hat{\sigma}_f = \frac{1}{N-1} \sum_i \left( f(U_i) - \hat{I}_f \right)^2$$

Monte Carlo Sampling is easy, robust and embarassingly parallel, but its convergence rate is  $\mathcal{O}(\frac{1}{\sqrt{N}})$ , although independent of the input dimension, is very slow.

So, how can we reduced the RMSE given by Monte Carlo? We could make our simulations faster, also we could increase N (which is not desirable because it would require more deterministic simulations), decrease  $\hat{\sigma}_f$  (variance-reduction techniques) or improve the random number generation (quasi-montecarlo).

### 2.2 Variance-Reduction Techniques

#### 2.2.1 Antithetic Sampling

We can use a different sampling (choosing of  $X$  values) to reduce the variance. This type of sampling assumes that a continuous PDF  $f_X(x)$  is symmetric around  $c$ . This implies that the symmetric counterpart of  $x$  is  $\tilde{x} = 2c - x$ , and that  $f_X(x) = f_X(\tilde{x})$ .

Thus, we can sample  $n/2$  times from  $f_X(x)$  and the remaining  $n/2$  samples are obtained via reflection.

For example for an uniform distribution in  $[0, 1]$ , we would get  $n/2$  values for  $X$ , and the rest of the  $n/2$  values would be obtained with the relation  $\tilde{U} = 1 - U$ .

Taking into account the symmetry of the function helps reducing the variance of the sampling.

### 2.2.2 Stratified Sampling

This type of sampling divides the interval of sampling to prevent clustering in a particular region of the interval. Suppose our interval is  $[x_i, x_f]$ . We would choose a  $\lambda \in (0, 1)$  and make two intervals. We take  $\lambda n$  samples in  $[x_i, \lambda(x_f - x_i)]$  and  $(1 - \lambda)n$  samples in  $[\lambda(x_f - x_i), x_f]$ .

### 2.2.3 Control Variates

This does not reduce the variance through sampling, but through an auxiliary function. We take as an illustrative example the estimation of the integral of  $f(x)$  over a domain. Here we introduce a control function  $\phi$ , which can be easily integrated in that domain.

$$I = \int f(x)dx = \int (f - \phi)dx + \int \phi dx$$

We know the exact integral for  $\phi$ , so we would only need to sample the difference. If we look at the variance of the difference:

$$\text{Var}(f - \phi) = \text{Var}(f) + \text{Var}(\phi) - 2\text{cov}(f, \phi)$$

we note that if  $\text{cov}(f, \phi)$  is high, then  $\text{Var}(f - \phi) < \text{Var}(f)$ , and we would have reduced the variance. So we only need to find a “similar function” to  $f$ .

### 2.2.4 Importance Sampling

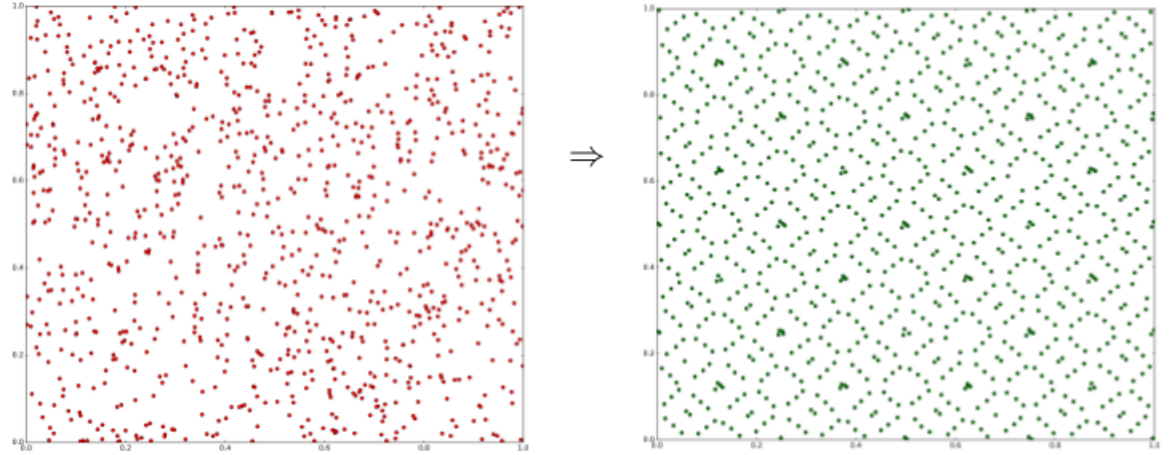
In a similar way to control variates, if we in Monte Carlo instead of sampling from the uniform distribution, sample from another function  $g_X$ , and  $f$  and  $g_X$  have similar shapes, then variance is reduced.

$$I = \int f(x)dx = \int \frac{f(x)}{g_X(x)} g_X(x) dx$$

In this case, the sampling would be done according to the  $g_X$ , and the function to evaluate would be  $f(x)/g_X(x)$ .

## 2.3 Quasi Monte-Carlo

Consists in alternative sampling techniques. Instead of using pseudo-random samples as in standard Monte-Carlo, which would result in a nonuniform sampling, we use deterministic samples. Examples are Fibonacci generators, latin hypercube sampling, Sobol’ sequences and Halton sequences.



In this course we will use techniques based on **low-discrepancy sequences**. For them, the error is  $\mathcal{O}\left(\frac{\log(N)^d}{N}\right)$ , where  $d$  is the dimension.

One example of low-discrepancy sequences are the Halton Sequences. They consist in choosing a prime number  $p$  and dividing the interval in  $p$  subintervals, and then each interval in other  $p$  subintervals, and so on until reaching the desired number of points.

In chaospy, it is very easy to generate QMC samples. If for example for generating 1000 Halton sequence-based samples on an uniform distribution:

```
distr = chaospy.Uniform()
samples = distr.sample(size=1000, rule='H')
```

## 3 Basic Interpolation and Quadrature

### 3.1 Interpolation

Interpolating consists in having some data points and estimating unknown values between those points. One way to do this is using polynomials, because they are analytically simple and computationally cheap.

Given a continuous function, and a grid where we know the value of the function, there exists a **unique** polynomial of degree  $N$  that satisfies the value of the function on those points. This polynomial can be expressed as a sum of **Lagrange cardinal polynomials**.

For a certain grid there is a unique polynomial, but we have freedom in choosing the grid. For defining the error of the interpolation, we can relate the error of our interpolation taking grid  $G$  with respect to the real function in terms of the error of an interpolation with the *best approximation polynomial* (with the same number of nodes).

$$Error_G \leq (1 + \Lambda_N(G)) (Error_{\text{Best grid}})$$

$\Lambda$  is the **Lebesgue constant** relative to the grid  $G$ , and contains all the information about the effect of the grid choice.

Lagrange interpolation on an uniform grid is not always a good idea, it can have strong variations near the limits of the interval (this is called the **Runge phenomenon**). Possible remedies for this is to use non-uniform grids (as Chebyscheff nodes) or different basis functions (maybe functions with local support, e.g. splines).

### 3.2 Quadrature

The idea of Quadrature is to approximate the definite integral of a function based on function evaluations. Quadrature uses rules that integrate polynomials exactly.

The general form of quadrature is a weighted sum of function values, however, as the function evaluation can be very computationally expensive, we need rules with small error of quadrature using few evaluations.

Gaussian Quadrature uses not only the weights of the quadrature as degrees of freedom, but also the nodes on the grid. If we make a quadrature considering  $N + 1$  points (quadrature of degree  $N$ ), this technique can integrate exactly polynomials of degree up to  $2N + 1$ . Gaussian quadrature has the following form:

$$\int f(x)w(x)dx = \sum_{i=0}^N w_i f(x_i) + \epsilon$$

Again,  $\epsilon = 0$  for polynomials up to degree  $2N + 1$ .

Note that we are integrating not only  $f$ , but a product with  $w$ . This is a weight function, and this makes it is useful for UQ when we take  $w$  as a



probability distribution. This would be equivalent to taking the expected value of  $f(x)$  with respect to the probability distribution  $w(x)$ . On the right side,  $w_i$  corresponds to the weights given to each point evaluation, and  $x_i$  correspond to the nodes.

**Getting the nodes** We can find the nodes if we find a family of polynomials  $p_i(x)$  that is orthogonal with respect to  $w(x)$  (i.e.  $\int p_i(x)p_j(x)w(x)dx = \delta_{ij}$ ). Then the nodes are simply the roots of the polynomial  $p_{N+1}$ . For  $w$  constant (as in the uniform probability distribution) our orthogonal family is the **Legendre polynomials**. For  $w$  gaussian (as in the normal distribution) our family is the **Hermite polynomials**.

**Getting the weights** The weights can be computed with the Lagrange cardinal polynomials evaluated at the nodes, integrated with respect to the weight function (at the end of the day we are integrating a  $(2N+1)$ -degree interpolation of  $f(x)$ ).

## 4 Polynomial Chaos

Forward UQ consists on having a stochastic input that goes into a model  $f$ , and a stochastic output which we want to characterize in the cheapest way possible. Montecarlo methods (and its improved versions) have slow convergence and need many samples. Another idea is to approximate our model  $f$  with a series of polynomials:

$$f(t, \omega) \approx \sum_{n=0}^{N-1} \hat{f}_n(t) \phi_n(\omega)$$

where  $t$  represent the independent variable(s) of our model,  $\omega$  is the stochastic input variable (with distribution  $\rho(\omega)$ ),  $\hat{f}_n(t)$  are the coefficients of the polynomials (note that are independent on the stochastic variable) and  $\phi_n(\omega)$  are orthogonal polynomials of degree  $n$  purely dependent on the stochastic variable.

In order to make this approximation we have to define the family of polynomials, compute the coefficients, choose the maximum order of the polynomials and finally characterize the output using this approximation.

### Choosing the type of polynomials

The orthogonal polynomials  $\phi_n(\omega)$  would be chosen according to the input distribution of the stochastic variable ( $\rho(\omega)$ ). For example for a uniform distribution we would choose Legendre polynomials, for a normal distribution we would choose Hermite polynomials, and so on. Computation of other families of polynomials can be made with the Stieltjes' **three-term recursion relation**, which is satisfied by all orthogonal polynomials, and can be gotten with a numerically stable method, or we could use other schemes as a Gram-Schmidt algorithm.

### Computing the coefficients

As our polynomials represent an orthonormal basis, we could compute our coefficients with a simple inner product with respect to  $\rho(\omega)$ :

$$\hat{f}_n(t) = \langle f(t, \omega), \phi_n(\omega) \rangle_\rho = \int f(t, \omega) \phi_n(\omega) \rho(\omega) d\omega$$

This way of getting the coefficients is called the **pseudo-spectral approach**.

That inner product is computed through an integral, and as  $f$  can be expensive or even a black box, we rather use quadrature to solve it. Gaussian quadrature is specially adequate for this. We just need to get the nodes and weights (that depend solely on the probability distribution  $\rho$ ) and then we would only have to evaluate  $f$  once at every node. This would yield:

$$\hat{f}_n = \sum_{k=0}^{K-1} w_k f(t, \omega_k) \phi_n(\omega_k)$$

Note that the integral is with respect to  $\omega$ , and thus the grid  $\omega_k$  that we get as the roots of the polynomials, represent values of our stochastic variable  $\omega$ .

## Choosing maximum order of polynomials and quadrature

Consider that the number of nodes to evaluate the integral (order of the quadrature) is completely independent of the number of terms used to represent our model (order of the interpolation):

- Using  $K$  nodes to evaluate the quadrature would be as integrating exactly a polynomial of degree up to  $2K-1$  interpolating  $f(t, \omega)\phi_n(\omega)$  with respect to  $\omega$ .
- Independent of that, we have the polynomial of degree  $N-1$  that interpolates  $f(t, \omega)$  with respect to  $\omega$ .

It may be very expensive to evaluate  $f(t, \omega_k)$ , but once these values are known, with them we can compute  $\hat{f}_n$ , and with them we can evaluate our interpolation of  $f(t, \omega)$  with low computational effort. Thus, our  $K$  determines the computational effort. On the other side, generally  $\hat{f}_n(t)$  decay exponentially with respect to  $n$ , so few coefficients (low  $N$ ) are sufficient. As a rule of thumb we can use  $N \approx K/2$ .

## Statistical properties of the approximation

Once we have  $\hat{f}_n(t)$ , we can evaluate statistical moments of our resulting function  $f(t, \omega) \approx \sum_{n=0}^{N-1} \hat{f}_n(t)\phi_n(\omega)$  pretty easily.

- **Expectation value:**  $\mathbb{E}[f(t, \omega)] = \hat{f}_0(t)$  (just like in fourier series)
- **Variance:**  $\text{Var}[f(t, \omega)] = \sum_{n=1}^{N-1} \hat{f}_n^2(t)$  (all the mixed terms vanish because of the orthogonality of our polynomials)

### 4.1 Multivariate Polynomial Chaos Expansion

This approach can be generalized for  $d$  dimensions. Now instead of  $n$  we have a vector  $\mathbf{n} = (n_1, \dots, n_d) \in \mathbb{N}^d$ , and our polynomials are product of one-dimensional polynomials (orthogonal with respect to the combination of all their indices). Generally for this approach one chooses a given  $N$  as a global order, and then we have a finite number of combinations of  $(n_1, \dots, n_d)$ , such that  $n_1 + \dots + n_d \leq N$ . This number of combinations is notated by  $P := \binom{d+N}{d}$ . The computational cost, which is  $K$  is only indirectly related to  $P$ .

- 5 Sparse Grids
- 6 Sensitivity Analysis
- 7 Random Fields
- 8 Software
- 9 Bayesian Inverse Problems