



Sequential experimental design based generalised ANOVA



Souvik Chakraborty*, Rajib Chowdhury

Department of Civil Engineering, Indian Institute of Technology Roorkee, Roorkee, India

ARTICLE INFO

Article history:

Received 14 November 2015

Received in revised form 21 February 2016

Accepted 21 April 2016

Available online 25 April 2016

Keywords:

Sequential experimental design

Generalised ANOVA

Polynomial chaos expansion

Distribution adaptive

ABSTRACT

Over the last decade, surrogate modelling technique has gained wide popularity in the field of uncertainty quantification, optimization, model exploration and sensitivity analysis. This approach relies on experimental design to generate training points and regression/interpolation for generating the surrogate. In this work, it is argued that conventional experimental design may render a surrogate model inefficient. In order to address this issue, this paper presents a novel distribution adaptive sequential experimental design (DA-SED). The proposed DA-SED has been coupled with a variant of generalised analysis of variance (G-ANOVA), developed by representing the component function using the generalised polynomial chaos expansion. Moreover, generalised analytical expressions for calculating the first two statistical moments of the response, which are utilized in predicting the probability of failure, have also been developed. The proposed approach has been utilized in predicting probability of failure of three structural mechanics problems. It is observed that the proposed approach yields accurate and computationally efficient estimate of the failure probability.

© 2016 Elsevier Inc. All rights reserved.

1. Introduction

Analysis and design of structural systems by considering the presence of uncertainties involve evaluation of structural reliability of the structure. Structural reliability is a theoretical framework for computing the probability of failure, and hence the reliability index, of structural systems. Determination of probability of failure depends on the probability density functions (PDF) of the input variables that have to be identified as preponderant.

The easiest way to determine probability of failure is the Monte Carlo simulation (MCS) [1–3]. Within this framework, large number of deterministic analyses are performed at some randomly generated sample points. Probability of failure is defined as the ratio of number of exceedance (i.e., number of times a given threshold is exceeded) to the total number of analyses. However, this method is computationally demanding, specifically for structures that require expensive finite element simulations. A number of improvements to the conventional MCS, such as the stratified sampling [4,5], importance sampling [6–8], directional simulation [9–11] and subset simulation [12,13] have also been proposed.

An alternate approach for structural reliability analysis is based on the Taylor's series expansion. The basic goal is to determine the point in standard normal space which is nearest to the origin. This point is often known as the design point or most probable point [14]. If first order Taylor's series expansion is used, the method is known as first order reliability method (FORM) [15]. On contrary if second order Taylor's series is used, the method is known as second order reliability method (SORM) [14,16]. However, the classical FORM/SORM often yields erroneous results for high dimensional nonlinear

* Corresponding author.

E-mail addresses: csouvik41@gmail.com (S. Chakraborty), rajibfce@iitr.ac.in (R. Chowdhury).

problems. It is worth mentioning that researchers have attempted to improve the classical FORM/SORM. For instance, variants of SORM based on parabolic failure surface [17] and Asymptotic distributions [18] have been proposed. Both these methods, to some extent, eliminate the limitation of SORM regarding high dimensional problems. However, these methods approximate the failure surface near the design point/most probable point by using a quadratic hypersurface (Type I error). As a consequence, results obtained for highly nonlinear problem are often erroneous. Additionally, these methods are applicable for systems involving only Gaussian variables (Type II error). In case, the system involves non-Gaussian variables, it is necessary to transform the variables to equivalent Gaussian variables. This in turn introduces some error into the system.

Over the last two decades, surrogate based structural reliability analysis has become quite popular among researchers. This method is mostly suitable for the case where the limit-state function has no known closed form and needs to be evaluated point-wise by numerical methods such as finite element (FE) method. In a sense, surrogate model is a system identification procedure, in which a transfer function relating the input parameters (loading and system conditions) to the output parameters (response in terms of displacements, stress, etc.) needs to be found in a suitable way. The observations required for the identification are usually taken from systematic numerical experiments with the full mechanical model, and the approximate transfer function obtained is termed as the surrogate. Surrogate models that are popular in literature include least square based response surface method (RSM) [19–21], high-order stochastic response surface method (HO-SRSM) [22], moving least square based RSM [23–25], Kriging [26,27], radial basis function [28,29], polynomial chaos expansion (PCE) [30–32], spectral functions [33,34], artificial neural network [35,36] and support vector machine [37]. However, most of the available surrogate models suffer from the following issues:

- Firstly, for most of the available surrogate models, the number of unknown coefficients increases factorially with the increase in number of variables. Naturally to avoid an underdetermined system, the number of training points needs to be increased. This is known as the *curse of dimensionality*. Due to this effect, most of the available surrogate models are not applicable for high dimensional problems.
- Secondly, almost all the available surrogate models yield accurate results only for problems that are either linear or weakly nonlinear in nature. If higher order polynomial bases are incorporated to account for the higher order of non-linearity, the common issue of over fitting comes into picture. As a consequence, the surrogate models yield erroneous results specifically at points that are located away from the training points.
- Thirdly, almost all the available surrogate models are applicable only for systems where the random variables are independent. If correlated random variables are present, some *ad hoc* transformations are employed to transform the correlated random variables into uncorrelated random variables. These transformations, which are based on some prior assumptions, introduces some unwanted error into the system. For example, Nataf transformation [38] assumes the dependent structure to be Gaussian and therefore, yields erroneous results for non-Gaussian variables. Similarly, orthogonal transformation [39] yields accurate result for lower value of correlation coefficient. However, if the correlation coefficient is higher, result obtained is inaccurate.
- Finally, in structural reliability problems, the training points for the construction of surrogate model are mostly selected empirically. Not much attention has been paid on development of a systematic framework for the selection of training points.

Of late, analysis of variance decomposition (ANOVA) [40–42] has been proposed for structural reliability analysis. This method, to some extent, addresses the first two limitations. As pointed out by [43], the number of training points required in ANOVA, with the increase in dimension, increases in a tabular manner. Therefore, this method is suitable for problems involving high number of random variables [44]. However, the third and fourth issues are not resolved by using the conventional ANOVA.

Hooker [45] and Li & Rabitz [46] extended the concepts of conventional ANOVA for treating systems having correlated random variables. This method, referred to as generalised ANOVA (G-ANOVA), represents the component functions of ANOVA by using the extended basis. As already demonstrated in [47–50], G-ANOVA yields excellent result for problems involving both uncorrelated and correlated random variables. However, the fourth issue is yet to be addressed.

The primary goal of this paper is to present a generalised framework for structural reliability analysis that addresses all the four limitations mentioned above. For that purpose, a novel sequential experimental design (SED) approach has been proposed and coupled with a newly developed variant of generalised ANOVA. The novelty of the proposed approach is mainly twofold:

- Firstly, a novel distribution adaption SED has been proposed. By utilizing the proposed SED, the optimum number of training points required for a given problem is determined in a systematic manner.
- Secondly, the proposed variant of generalised ANOVA presented in this paper expresses the component functions by using PCE. It is worth mentioning that PCE can be viewed as a special case of extended bases. However, utilizing PCE reduces the computational cost.

The rest of the paper is organised as follows. After defining the use of surrogate modelling technique for structural reliability analysis in Section 2, Section 3 present the concepts of G-ANOVA. In Section 4, the proposed sequential experimental design has been presented. Section 5 presents the proposed framework for sequential experimental design based

G-ANOVA. Generalised analytical formulas for the first two statistical moments, that are utilized for computing the probability of failure, have also been developed. Implementation of the proposed approach for structural reliability analysis has been presented in Section 6. Section 7 presents the concluding remarks.

2. Surrogate modelling technique for structural reliability analysis

Suppose, $\mathbf{x} = (x_1, x_2, \dots, x_N) : \Omega \rightarrow \mathbb{R}^N$ be a N -dimensional random vector with cumulative distribution function (CDF) $F_{\mathbf{x}}(\mathbf{X}) = \mathbb{P}(\mathbf{x} \leq \mathbf{X})$, where \mathbb{P} denotes probability, Ω is the probability space and $\mathbf{x} \in \mathbb{R}^N$. Traditionally, a reliability problem is defined by a performance function, often known as limit-state function $\mathcal{J}(\mathbf{x})$, where $\mathcal{J}(\mathbf{x}) < 0$ denotes the failure domain Ω_F and $\mathcal{J}(\mathbf{x}) \geq 0$ denotes the safe region, i.e.,

$$\Omega_F \triangleq \{\mathbf{x} : \mathcal{J}(\mathbf{x}) < 0\} \quad (1)$$

The failure probability P_f is defined as:

$$P_f = P(\mathbf{x} \in \Omega_F) = \int_{\Omega_F} dF_{\mathbf{x}}(\mathbf{X}) = \int_{\Omega} \xi_{\Omega_F}(\mathbf{x}) dF_{\mathbf{x}}(\mathbf{X}) \quad (2)$$

where $P(\bullet)$ denotes probability of an event and $\xi_C(\mathbf{x})$ is a characteristics function which satisfies:

$$\xi_C(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in C \\ 0 & \text{if } \mathbf{x} \notin C \end{cases} \quad (3)$$

where C in Eq. (3) denotes the domain. It is apparent that the limit-state function $\mathcal{J}(\mathbf{x})$ plays an important role in determining the probability of failure. However, almost for all practical cases, explicit form for $\mathcal{J}(\mathbf{x})$ is not known and one need to perform simulations to determine $\mathcal{J}(\mathbf{x})$. As a consequence, the procedure becomes computationally intensive and time consuming.

In surrogate modelling technique, an explicit expression $\hat{\mathcal{J}}(\mathbf{x})$ for the limit-state function $\mathcal{J}(\mathbf{x})$ is formulated. To formulate the surrogate, firstly suitable training points in the design space need to be selected. Next, deterministic analyses are performed to determine the response at the training points. Finally, the surrogate is formulated by utilizing the training points and responses obtained at the training points. Once the surrogate model is obtained, probability of failure is obtained by performing the classical reliability analysis procedure on the generated surrogate.

It is evident from the above discussion that the accuracy of surrogate plays a vital role in evaluating the probability of failure accurately. Furthermore, efficiency of a surrogate model is inversely proportional with the number of training points required. In the next section, an accurate as well as efficient surrogate modelling technique for structural reliability analysis has been presented.

3. Generalised ANOVA

In the first part of this section, the basic concepts of ANOVA decomposition and PCE have been discussed. Towards, the end of this section the proposed generalised ANOVA, derived by coupling the ANOVA decomposition with PCE, has been presented.

3.1. ANOVA

Suppose $\mathbf{i} = (i_1, i_2, \dots, i_N) \in \mathbb{N}_0^N$ be a multi-index with $|\mathbf{i}| = i_1 + i_2 + \dots + i_N$. Now considering $\mathbf{x} = (x_1, x_2, \dots, x_N)$ to be the random inputs, the unknown response $g(\mathbf{x})$ can be expressed as [42]:

$$\begin{aligned} g(\mathbf{x}) &= \sum_{|\mathbf{i}|=0}^N g_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}}) \\ &= g_0 + \sum_{k=1}^N \sum_{i_1 < i_2 < \dots < i_k} g_{i_1 i_2 \dots i_k}(x_{i_1}, x_{i_2}, \dots, x_{i_k}) \end{aligned} \quad (4)$$

where $g_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}})$ are the component functions.

Definition 1. The univariate terms (i.e., the terms corresponding to $k = 1$) in Eq. (4) are termed as first order component functions. Similarly, the bivariate terms (terms corresponding to $k = 2$) are termed as second order component functions. g_0 is the zeroth order component function.

Remark 1. First order component function does not indicate linear variation and consists of terms having higher order. Order in ANOVA indicates the maximum order of cooperativity considered.

Remark 2. The component functions defined in Eq. (4) must be orthogonal to each other. This criteria is known as the *hierarchical orthogonality criteria*. Note that this is an essential condition to ensure uniqueness of the solution.

Now if, x_1, x_2, \dots, x_N are assumed to be independent, the component function over the problem space can be determined by imposing a vanishing condition [51]:

$$\int_{p_i}^{q_i} \varpi_k(x_k) g_{i_1 i_2 \dots i_m}(x_{i_1}, x_{i_2}, \dots, x_{i_m}) dx_k = 0, \quad \forall k \in \{i_1, i_2, \dots, i_m\} \quad (5)$$

where ϖ_k denoted the PDF of x_k and p_i, q_i denote the bounds of the variable. Using Eq. (5), the component function of ANOVA can be written as [52]:

$$\begin{aligned} g_0 &= E(g(\mathbf{x})) \\ g_i(x_i) &= E(g(\mathbf{x}) | x_i) - g_0 \\ g_{ij}(x_i, x_j) &= E(g(\mathbf{x}) | x_i, x_j) - g_i(x_i) - g_j(x_j) - g_0 \\ &\vdots \end{aligned} \quad (6)$$

where $E(\bullet)$ denotes expectation. Classical ANOVA utilizes the equations specified in Eq. (6) to determine the component functions. However, practical determination of the component function using classical ANOVA is computationally cumbersome and requires large number of training points. In order to address this issue, other variants of ANOVA has emerged. While the anchored ANOVA [53,54] utilizes interpolation functions to represent the component function, the polynomial based ANOVA [51,52,55] expresses the component functions in term of some suitable basis. However, all the above mentioned ANOVA are only suitable for system involving only independent variables.

Of late, the concept of polynomial based ANOVA has been extended for systems involving both correlated and independent random variables [45,46]. This method, referred to as G-ANOVA express the component functions in term of extended bases. The unknown coefficients are determined by enforcing the orthogonality of the component functions. As already demonstrated, this method yields excellent result for high dimensional system [47–50].

Next the concepts of PCE, which is utilized as a tool in the proposed framework, has been briefly discussed.

3.2. PCE

Considering the multi-index notation defined in preceding section and assuming $n \geq 0$ to be an integer, the n^{th} order PCE of a function $g(\mathbf{x})$ can be expressed as:

$$g_n(\mathbf{x}) = \sum_{|\mathbf{i}|=0}^n \alpha_{\mathbf{i}} \psi_{\mathbf{i}} \quad (7)$$

where $\alpha_{\mathbf{i}}$ in Eq. (7) are the unknown coefficients and $\psi_{\mathbf{i}}$ represents the orthogonal polynomial of degree i . The orthogonal polynomials satisfies:

$$E(\psi_{\mathbf{i}} \psi_{\mathbf{j}}) = m \delta_{\mathbf{ij}}, \quad 0 \leq |\mathbf{i}|, |\mathbf{j}| \leq n, \quad m \geq 0 \quad (8)$$

where

$$\begin{aligned} \delta_{\mathbf{ij}} &= \prod_{k=1}^n \delta_{i_k j_k} = 1 \text{ if } \mathbf{i} = \mathbf{j} \\ &= 0, \text{ elsewhere} \end{aligned} \quad (9)$$

Based on the orthogonality criteria specified in Eq. (8), appropriate orthogonal polynomial needs to be selected. The detailed description of orthogonal polynomials can be found in [30].

Different variants of PCE have been suggested by various researchers. The variants of PCE, primarily differ in the way the unknown coefficients associated with the bases are determined. Xiu and Karniadakis [56] suggested the use of Galerkin projection for determining the unknown coefficients. Other methods for determining the unknown coefficients include least-squared regression [57] and least-angle regression [31]. Note that, in this work, the interest resides in the functional form of PCE described in Eq. (7). The unknown coefficients will be determined after PCE has been coupled into the framework of ANOVA.

3.3. Proposed G-ANOVA

In this section, a new variant of G-ANOVA has been proposed. The basic idea is to represent the unknown component functions by using the PCE. Utilizing the functional form of PCE shown in Eq. (7), Eq. (4) can be rewritten as:

$$g(\mathbf{x}) = g_0 + \sum_{|i|=1}^N \sum_{|j_i|=1}^{\infty} \alpha_{j_i}^i \psi_{j_i} \quad (10)$$

Considering an M^{th} order ANOVA and r^{th} order PCE, Eq. (10) reduces to:

$$\hat{g}(\mathbf{x}) = g_0 + \sum_{|i|=1}^M \sum_{|j_i|=1}^r \alpha_{j_i}^i \psi_{j_i} \quad (11)$$

Eq. (11) can be rewritten in matrix form as:

$$\Psi \alpha = \mathbf{d} \quad (12)$$

where Ψ is a matrix consisting of the orthogonal basis, vector α consists of the unknown coefficient vector and $\mathbf{d} = \mathbf{g} - \bar{\mathbf{g}}$ where, $\mathbf{g} = (g_1, g_2, \dots, g_{N_S})^T$ is a vector consisting of the observed responses at N_S training points and $\bar{\mathbf{g}} = (g_0, g_0, \dots, g_0)^T$ is the mean response vector. Premultiplying Eq. (12) by Ψ^T yields:

$$\mathbf{B} \alpha = \mathbf{C} \quad (13)$$

where $\mathbf{B} = \Psi^T \Psi$ and $\mathbf{C} = \Psi^T \mathbf{d}$. Assuming that the system under consideration is having q unknown coefficients, it is obvious that \mathbf{B} is a matrix of dimension $q \times q$. Similarly, α and \mathbf{C} are vectors having dimension $q \times 1$.

Close inspection of Ψ reveals identical columns. Thus, \mathbf{B} has identical rows. These rows are redundants and can be removed. Removing identical rows of \mathbf{B} and corresponding rows of \mathbf{C} , one obtains:

$$\mathbf{B}' \alpha = \mathbf{C}' \quad (14)$$

Eq. (14) represents a set of underdetermined equations and naturally there exists an infinite number of solutions. Assume \mathbf{B}' to be a $p \times q$ matrix, α to be a $q \times 1$ vector and \mathbf{C}' to be a $p \times 1$ vector. Then all the solutions of Eq. (14) can be represented as:

$$\alpha(s) = (\mathbf{B}')^{-1} \mathbf{C}' + [\mathbf{I} - (\mathbf{B}')^{-1} \mathbf{B}'] v(s) \quad (15)$$

where $(\mathbf{B}')^{-1}$ denotes the generalised inverse of \mathbf{B}' , $v(s)$ is an arbitrary vector in \mathbb{R}^q and \mathbf{I} represents an identity matrix. One choice of $(\mathbf{B}')^{-1}$ in Eq. (14) is $(\mathbf{B}')^\dagger$, where $\alpha_0 = (\mathbf{B}')^\dagger \mathbf{C}'$ is the solution of Eq. (14) obtained using the least-squared regression. Replacing $(\mathbf{B}')^{-1}$ by $(\mathbf{B}')^\dagger$ in Eq. (15) yields:

$$\alpha(s) = (\mathbf{B}')^\dagger \mathbf{C}' + \mathbf{P} v(s) \quad (16)$$

where

$$\mathbf{P} = \mathbf{I} - (\mathbf{B}')^\dagger \mathbf{B}' \quad (17)$$

Definition 2. Out of all the possible solution defined by Eq. (16), the solution that minimizes the least squared error and satisfies the hierarchical orthogonality criteria of G-ANOVA defined in Remark 2, is termed as the 'best solution'.

In the proposed G-ANOVA, the best solution, from all the available solutions defined by Eq. (16) is obtained by employing homotopy algorithm (HA) [58–60]. HA determines the unknown coefficient by minimizing the least-squared error and an objective function. The solution using HA is given as:

$$\alpha_{HA} = \left[\mathbf{V}_{q-r} \left(\mathbf{U}_{q-r}^T \mathbf{V}_{q-r} \right)^{-1} \mathbf{U}_{q-r}^T \right] \alpha_0 \quad (18)$$

where \mathbf{U} and \mathbf{V} are matrices obtained by singular value decomposition of \mathbf{PW} matrix:

$$\mathbf{PW} = \mathbf{U} \begin{pmatrix} \mathbf{A}_r & 0 \\ 0 & 0 \end{pmatrix} \mathbf{V}^T \quad (19)$$

For detailed derivation of Eq. (18), interested readers may refer to [58,59]. \mathbf{P} in \mathbf{PW} is the matrix defined in Eq. (17) and \mathbf{W} is the weight matrix used to formulate the objective function in HA. For details regarding the weight matrix \mathbf{W} , interested readers may refer to [46,48].

Once the unknown coefficient vector α is determined, Eq. (10) provides an explicit mapping of the input and output variables. In this paper, the proposed G-ANOVA is utilized to generate an explicit expression for the limit-state function. Once determined, probability of failure P_f is calculated as [47]:

$$P_f = \frac{1}{2} \frac{\sqrt{\pi} \exp\left(\frac{1}{4} \frac{\lambda_1^2}{\lambda_2}\right) \left[\operatorname{erf}\left(\frac{1}{2} \frac{\lambda_1}{\sqrt{\lambda_2}}\right) - \operatorname{erf}\left(\frac{1}{2} \frac{2\lambda_2 y_l + \lambda_1}{\sqrt{\lambda_2}}\right) \right] \exp(-\lambda_0)}{\sqrt{\lambda_2}} \quad (20)$$

where $\operatorname{erf}(\bullet)$ denotes error function. y_l in above equation represents the lower limit of response. The λ_i 's, $i = 0, 1, 2$ in Eq. (20) are functions of the first two statistical moments of the output response [47].

Remark 3. As already demonstrated, G-ANOVA yields highly accurate estimate of the first two statistical moments [48]. Now using Eq. (20), it is possible obtain highly accurate estimate for the probability of failure based on the first two statistical moments estimated using G-ANOVA.

Remark 4. One advantage of the proposed G-ANOVA resides in the fact that it is possible to derive analytical formulas for the first two statistical moments. As a consequence, the moments obtained are free from the sampling error introduced when MCS is utilized for determining the statistical moments.

3.3.1. Statistical moments

Lemma 1. The first moment of all but the zeroth order component function is zero.

Proof. An arbitrary m^{th} order component function $g_{i_1 i_2 \dots i_m}(x_{i_1}, x_{i_2}, \dots, x_{i_m})$ can be represented as:

$$g_{i_1 i_2 \dots i_m}(x_{i_1}, x_{i_2}, \dots, x_{i_m}) = \sum_{|\mathbf{i}|=1}^r \alpha_{\mathbf{i}} \psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}}) \quad (21)$$

Applying the expectation operator $E(\bullet)$ on both sides:

$$E(g_{i_1 i_2 \dots i_m}) = \sum_{|\mathbf{i}|=1}^r \alpha_{\mathbf{i}} E(\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}})) \quad (22)$$

Now as already mentioned, $\psi_{\mathbf{i}}$ is a orthogonal polynomial and thus

$$E(\psi_{\mathbf{i}}) = 0, \quad \mathbf{i} = 1, 2, \dots, r \quad (23)$$

Thus,

$$E(g_{i_1 i_2 \dots i_m}) = 0 \quad (24)$$

From Eq. (24), it can be concluded that first moment of all but the zeroth order component function is zero. \square

Corollary 1. The mean of the proposed G-ANOVA is g_0 .

Proof. Eq. (4) is rewritten as:

$$g(\mathbf{x}) = \underbrace{g_0}_{0^{\text{th}} \text{ order}} + \sum_i \underbrace{g_i(x_i)}_{1^{\text{st}} \text{ order}} + \sum_{1 \leq i < j \leq N} \underbrace{g_{ij}(x_i, x_j)}_{2^{\text{nd}} \text{ order}} + \dots + \underbrace{g_{12 \dots N}(x_1, x_2, \dots, x_N)}_{N^{\text{th}} \text{ order}} \quad (25)$$

Applying expectation operator on both sides:

$$\begin{aligned} E(g(\mathbf{x})) &= E(g_0) + \sum_i E(g_i(x_i)) + \sum_{1 \leq i < j \leq N} E(g_{ij}(x_i, x_j)) + \dots \\ &\quad + E(g_{12 \dots N}(x_1, x_2, \dots, x_N)) \end{aligned} \quad (26)$$

Now g_0 is constant and hence $E(g_0) = g_0$. Furthermore from Lemma 1, expectation of all the other component function is zero. Thus,

$$E(g(\mathbf{x})) = g_0 \quad (27)$$

This completes the proof of Corollary 1. \square

Theorem 1. All the component functions are mutually uncorrelated.

Proof. Consider \mathcal{K}_1 and \mathcal{K}_2 to be two functions. From definition, \mathcal{K}_1 and \mathcal{K}_2 are orthogonal if

$$E(\mathcal{K}_1 \mathcal{K}_2) = \frac{\mathbf{K}_1^T \mathbf{K}_2}{N_r} = 0 \quad (28)$$

where \mathbf{K}_1 and \mathbf{K}_2 are vectors consisting of N_r realization of the function \mathcal{K}_1 and \mathcal{K}_2 . Similarly, \mathcal{K}_1 and \mathcal{K}_2 are uncorrelated if

$$E((\mathcal{K}_1 - \bar{\mathcal{K}}_1)(\mathcal{K}_2 - \bar{\mathcal{K}}_2)) = \frac{(\mathbf{K}_1 - \bar{\mathbf{K}}_1)^T (\mathbf{K}_2 - \bar{\mathbf{K}}_2)}{N_r} = 0 \quad (29)$$

From Lemma 1, the component functions (except the zeroth order) are having zero mean. Under this circumstances, Eq. (28) and Eq. (29) are identical. Thus, one may conclude that all the component functions are mutually uncorrelated. This completes the proof of Theorem 1. \square

From Theorem 1, it is evident that covariance between any two component function is zero. Thus applying variance operator $\text{var}(\bullet)$ on both sides of Eq. (25)

$$\begin{aligned} \text{var}(g(\mathbf{x})) &= \text{var}(g_0) + \sum_i \text{var}(g_i(x_i)) + \sum_{1 \leq i < j \leq N} \text{var}(g_{ij}(x_i, x_j)) + \cdots \\ &\quad + \text{var}(g_{12\dots N}(x_1, x_2, \dots, x_N)) \end{aligned} \quad (30)$$

Now g_0 is constant and hence, $\text{var}(g_0) = 0$. Thus, Eq. (30) reduces to

$$\text{var}(g(\mathbf{x})) = \sum_i \text{var}(g_i(x_i)) + \sum_{1 \leq i < j \leq N} \text{var}(g_{ij}(x_i, x_j)) + \cdots + \text{var}(g_{12\dots N}(x_1, x_2, \dots, x_N)) \quad (31)$$

Now considering an arbitrary m^{th} order component function $g_{i_1 i_2 \dots i_m}(x_{i_1}, x_{i_2}, \dots, x_{i_m})$ as before (Eq. (21)) and applying the variance operator on both sides:

$$\text{var}(g_{i_1 i_2 \dots i_m}(x_{i_1}, x_{i_2}, \dots, x_{i_m})) = \sum_{|\mathbf{i}|=1}^r (\alpha_{\mathbf{i}})^2 \text{var}(\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}})) \quad (32)$$

Using the property of orthogonal basis (Eq. (23)),

$$\begin{aligned} \text{var}(\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}})) &= E((\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}}))^2) - (E(\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}})))^2 \\ &= E((\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}}))^2) \end{aligned} \quad (33)$$

Substituting Eq. (33) into Eq. (32),

$$\text{var}(g_{i_1 i_2 \dots i_m}(x_{i_1}, x_{i_2}, \dots, x_{i_m})) = \sum_{|\mathbf{i}|=1}^r (\alpha_{\mathbf{i}})^2 E((\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}}))^2) \quad (34)$$

Furthermore if special category of orthogonal polynomial, known as the orthonormal polynomial is used,

$$E((\psi_{\mathbf{i}}(\mathbf{x}_{\mathbf{i}}))^2) = 1 \quad (35)$$

Thus,

$$\text{var}(g_{i_1 i_2 \dots i_m}(x_{i_1}, x_{i_2}, \dots, x_{i_m})) = \sum_{|\mathbf{i}|=1}^r (\alpha_{\mathbf{i}})^2 \quad (36)$$

Substituting Eq. (36) into Eq. (31) and writing using the multi-index notation

$$\text{var}(g(\mathbf{x})) = \sum_{|\mathbf{i}|=1}^N \sum_{|\mathbf{j}|=1}^r (\alpha_{\mathbf{j}\mathbf{i}}^{\mathbf{i}})^2 \quad (37)$$

and thus,

$$E((g(\mathbf{x}))^2) = (g_0)^2 + \sum_{|\mathbf{i}|=1}^N \sum_{|\mathbf{j}|=1}^r (\alpha_{\mathbf{j}\mathbf{i}}^{\mathbf{i}})^2 \quad (38)$$

Eq. (38) is the basic formula for calculating the second moment of the proposed G-ANOVA. Once the first two moments are determined using Eq. (37) and Eq. (38) respectively, the procedure described in [47] is employed to determine the λ 's. Finally, Eq. (20) is employed to determine the probability of failure.

4. Sequential experimental design (SED)

Design of experiments (DOE) is a systematic procedure for selecting appropriate training points. In practice, all the surrogate models depend on some DOE for selecting the training points. Inappropriate selection of training points may introduce unwanted errors into the system and appropriate selection of the training points can increase the efficiency of a surrogate model. Popular DOE available in literature includes, but are not limited to factorial design [61], central composite design [23], optimal design [62] and Latin Hypercube sampling [63]. However, most of the popular experimental designs require the design points (or training points) to be selected beforehand. As a consequence, these DOEs are prone to either oversampling or undersampling.

Recently, sequential experimental design (SED) [64] has emerged as a possible alternative. The basic idea of SED is to generate the design points in a sequential manner until convergence, defined by some appropriate criteria, is reached. In the following section, a novel distribution adaptive SED (DA-SED) has been proposed.

4.1. DA-SED

DA-SED is a novel sequential experimental design developed to be used in conjunction with the proposed G-ANOVA. This scheme generates distribution adaptive design points in a sequential manner. The algorithm stops when convergence, defined by some appropriate criteria, is achieved.

Let, $k^s = [0, 1]^s$ be a N dimensional unit hypercube and \mathcal{F} to be a real function that is integrable in k^s . The basic goal is to generate a sequence x_n , such that:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_i f(x_i) = \int_{k^N} f(x) \varpi(x) dx \quad (39)$$

with the highest possible convergence. $\varpi(x)$ is the PDF of x . It is evident from Eq. (39) that a sequence satisfying the following three criteria, will have high convergence rate:

- The first criteria, known as the 'distribution adaptivity criteria', ensures that the training points follow the distribution pattern of the variable. For example, if variable is lognormally distributed, the 'distribution adaptivity' forces the training points to be lognormally distributed as well. Note that this is an essential criteria and should be satisfied under all circumstances.
- The second criteria is known as the 'minimum intersite criteria'. According to this criteria, the intersite distance between two consecutive training points should be minimum.
- The final criteria is the 'projected distance criteria'. This criteria is especially important for cases where relative importance of the design variables are unknown. For example, if one of the design parameters does not influence the output behaviour, two design points which only differ in this (irrelevant) parameter have the same behaviour, and can be seen as the 'same' design point. Thus, the projected distance must also be maximized.

Remark 5. It is to be noted that the above mentioned three criteria are conflicting in nature. Thus, a trade-off between the above mentioned criteria is necessary.

A sequence that satisfy the second and the third criteria is known as the (t, m, s) -nets and (t, s) -sequence in base b [65,66].

Definition 3. An elementary n -interval in base b of k^s , having the form $\prod_{i=1}^s \left[\frac{a_i}{b^{d_i}}, \frac{a_i+1}{b^{d_i}} \right]$, is considered, where a_i and d_i are non-negative integers and $b^{d_i} > a_i$, $\forall i = \{1, 2, \dots, N\}$. Now considering two integers t and m , such that $0 \leq t \leq m$, (t, m, s) -net in base b is a sequence such that

$$\# P \cap \{x_1, x_2, \dots, x_{b^m}\} = b^t \quad (40)$$

for all elementary interval P in base b of hypervolume $\lambda(P) = b^{t-m}$, where $\#$ denotes cardinality of a set.

Definition 4. Considering an integer t such that $t \geq 0$, a (t, s) sequence in base b is an infinite sequence of points x_n such that for all the integers $k \geq 0$, $m \geq t$ and the sequence $\{x_{kb^m}, \dots, x_{(k+1)b^m-1}\}$ is a (t, m, s) -net in base b .

Remark 6. A (t, m, s) -net and (t, s) -sequence in base b is known as the Sobol's sequence [65,66]. Sobol's sequence is one of the most popular sequence available in literature. It has already been demonstrated that Sobol's sequence has higher convergence rate [67] as compared to other sequences such as Faure sequence [68] and Halton sequence [69].

The procedure for generating a Sobol sequence is quite straightforward. To generate the j^{th} component, a primitive polynomial of degree s_j in base 2

$$P_j = x^{s_j} + a_1^j x^{s_j-1} + \dots + a_{s_j-1}^j x + 1 \quad (41)$$

where the coefficient $a_i^j, \forall i$ are either 0 or 1, is considered. Now if a sequence of positive integers, denoted as $\{m_1^j, m_2^j, \dots\}$, is defined as:

$$m_k^j := 2a_1^j m_{k-1}^j \oplus 2^2 a_2^j m_{k-2}^j \oplus \dots \oplus 2^{s_j-1} a_{s_j-1}^j m_{k-s_j+1}^j \oplus 2^{s_j} a_{s_j}^j m_{k-s_j}^j, \quad (42)$$

where \oplus is the bit by bit exclusive-or operator, it is possible to determine all the points in a sequential manner. Also note that the starting point $\{m_1^j, m_2^j, \dots\}$ can be freely chosen provided m_k^j is odd and

$$m_k^j < 2^k, \forall k \text{ and } 1 \leq k \leq s_j \quad (43)$$

Using m_k^j , the direction numbers $\{v_1^j, v_2^j, \dots\}$ are defined as:

$$v_k^j := \frac{m_k^j}{2^k} \quad (44)$$

Finally the i^{th} component of the j^{th} point in Sobol's sequence, is obtained as:

$$x_i^j := i_1 v_1^j \oplus i_2 v_2^j \oplus \dots \quad (45)$$

where i_k is the k^{th} digit from right when i is written in binary:

$$i = (\dots i_2 i_1)_2 \quad (46)$$

where $(\bullet)_2$ denotes binary operator. Eq. (45) is the basic formula for generating the Sobol's sequence. It is evident that the x_i^j 's are generated sequentially and thus can be utilized as design points in SED.

Remark 7. The sequence obtained using Eq. (45) does not incorporate available information regarding the probability distribution of the variables. Naturally, the first criteria is not satisfied.

In order to ensure the 'distribution adaptivity criteria', the authors propose to transform the obtained x_i^j as per the probability distribution of the variables.

Suppose \mathbf{X}_{samp} to be the design matrix obtained using Eq. (45). In the first step, the design matrix \mathbf{X}_{samp} is transformed into standard normal variate as:

$$(\mathbf{X}_{\text{samp}})_{\text{normal}} = \Phi^{-1}(\mathbf{X}_{\text{samp}}) \quad (47)$$

where $\Phi(\bullet)$ denotes the cumulative distribution of standard normal variable. In the next step, the information regarding the dependencies of the variables, defined in term of correlation matrix is incorporated as:

$$(\mathbf{X}_{\text{samp}})_{\text{dep_normal}} = (\mathbf{X}_{\text{samp}})_{\text{normal}} \times \text{chol}(\rho) \quad (48)$$

where ρ denotes the correlation matrix and $\text{chol}(\bullet)$ denotes the Cholesky decomposition operator. The third step involves transforming the dependent standard normally distributed training points, obtained in Eq. (48), to dependent uniformly distributed training points as:

$$\mathbf{X}_{\text{dep_uniform}} = \Phi\left((\mathbf{X}_{\text{samp}})_{\text{dep_normal}}\right) \quad (49)$$

If the variables under considering follows uniform distribution in $[0, 1]$, the design matrix obtained using Eq. (49) is utilized as the training points. However, if one or more variables are non-uniformly distributed, the design matrix is transformed. Some suitable transformations are provided in the Appendix A. An algorithm depicting the steps involved for generating the training points is shown in Algorithm 1. Representative training points generated using the proposed scheme are shown in Fig. 1.

Algorithm 1 Algorithm for generating distribution adaptive training points.**Initialize:** Input number of random variables and their type

- 1) Generate uniformly distributed training point. → Eqs. (41)–(46)
- 2) Input correlation matrix
- 3) Obtain independent standard normally distributed training points → Eq. (47)
- 4) Obtain dependent standard normally distributed training points → Eq. (48)
- 5) Obtain dependent uniformly distributed training points → Eq. (49)
- 6) Transform to obtain distribution adaptive training points

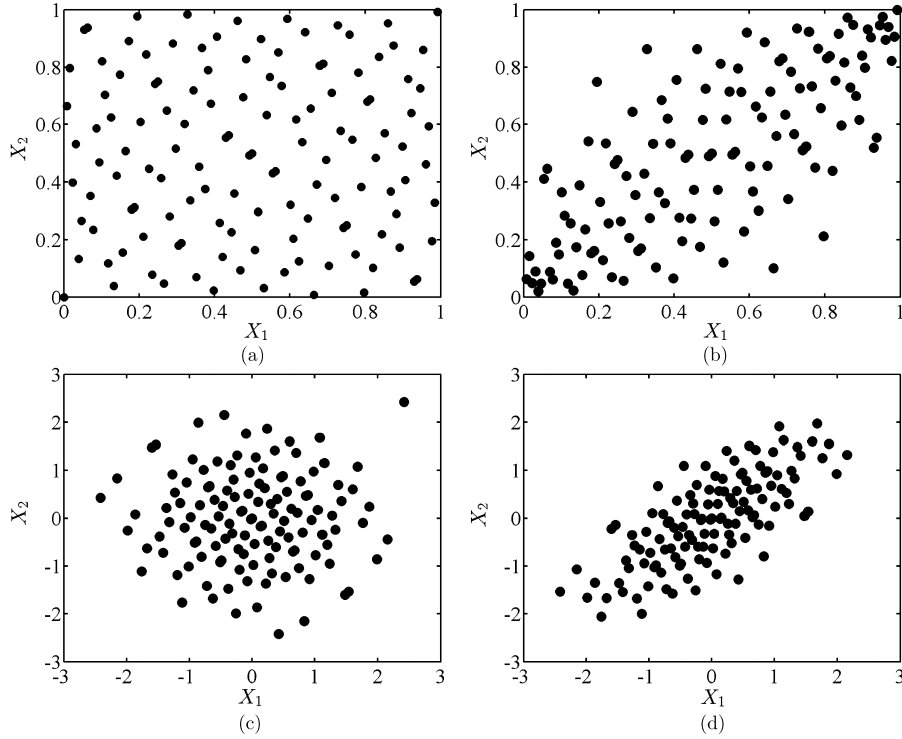


Fig. 1. Distribution adaptive training points for (a) two uniformly distributed independent random variables in $[0, 1]$, (b) two uniformly distributed dependent random variables ($\rho = 0.8$) in $[0, 1]$, (c) two independent standard normal variates and (d) two dependent ($\rho = 0.8$) standard normal variates.

Remark 8. The number of rows in $\mathbf{X}_{\text{dep_uniform}}$ corresponds to the number of training points and the number of columns corresponds to the number of variables. The final transformation is performed for each column (i.e., each variable) at a time.

Remark 9. In present work, the proposed distribution adaptive experimental design has been utilized to generate design points sequentially and hence the name ‘distribution adaptive sequential experimental design’.

5. Sequential experimental design based G-ANOVA

In this section, a unified framework that couples the proposed G-ANOVA with the DA-SED has been presented. The proposed framework consists of the two steps:

- Select initial number of training points. Generate training points. Obtain responses at selected training points.
- Formulate the proposed G-ANOVA and check for accuracy. If the surrogate is accurate enough proceed with the computation of structural reliability. Otherwise, update the design matrix by generating more training points and repeat this step until convergence

However, there are some issues associated with the above mentioned steps. For example, there is no guideline regarding the initial number of training points. Furthermore, checking accuracy of any surrogate model is a tricky job. Almost all the available statistical tests require actual responses at some points that has not been utilized for formulating the surrogate. In this work, the initial training points is considered to be equal to the number of variables and leave one out (LOO) statistical test has been utilized. The reason being, LOO is highly efficient and requires only one additional training point. The third issue associated with above procedure resides in selecting the increment size. It is essential to determine appropriate increment size for satisfactory performance of the proposed algorithm. In this paper, a scheme for determining the step size

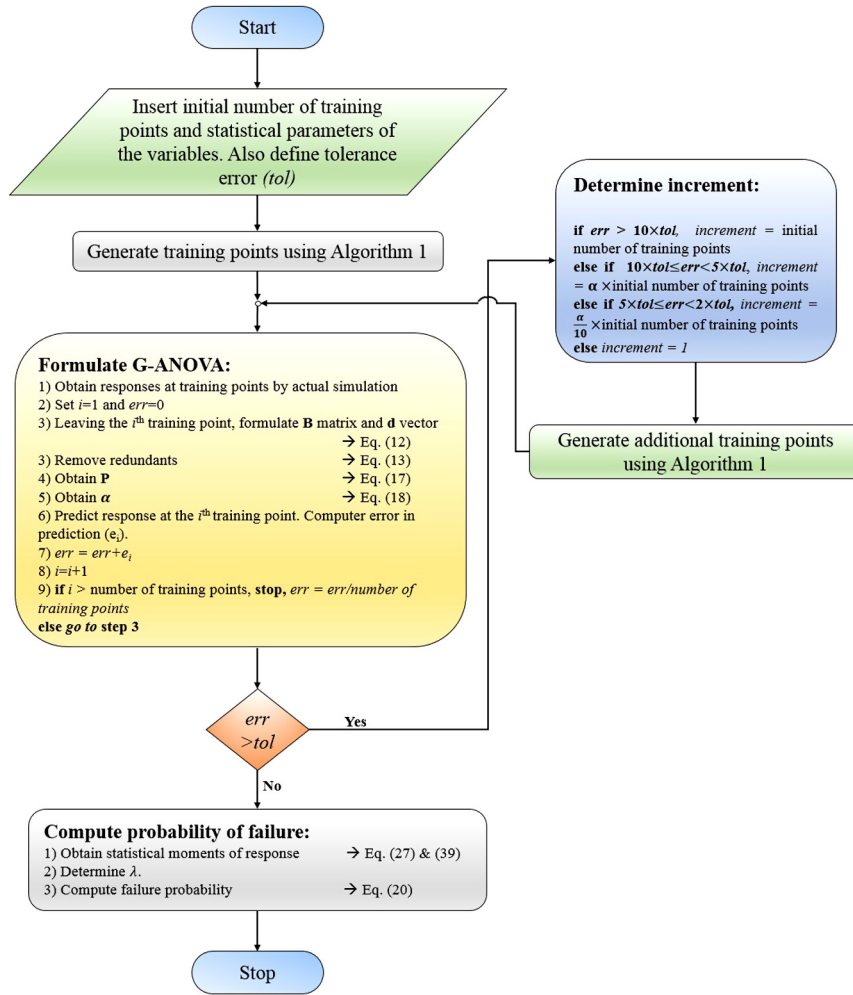


Fig. 2. Flowchart depicting the proposed approach.

based on the error obtained in previous step has been prescribed. The proposed scheme considers four intervals, namely (i) error is greater than ten times the tolerance, (ii) error is within ten times the tolerance but greater than five times the tolerance, (iii) error resides between twice the tolerance and five times the tolerance and (iv) error is greater than the tolerance but less than twice the tolerance. Among the above mentioned four cases, case (i) is considered to be the 'worst case scenario' and hence, it can be argued that *significant* increase in number of training points is required to achieve convergence. Thus, the initial number of training points is considered to be the size of increment. Case (ii) is comparatively better than case (i). Naturally, the increment size ($\alpha \times$ initial training points, where $0 < \alpha < 1$) considered is comparatively smaller than case (a). In case (iii), the increment size is reduced even further (a factor of $\alpha/10$ is considered). In case (iv), the error is very close to the tolerance and hence, the training points are increased one at a time. A detailed flowchart depicting the proposed SED based G-ANOVA is shown in Fig. 2.

6. Numerical examples

In this section, structural reliability analysis of three structural engineering problems has been presented. Results obtained has been benchmarked against full-scale MCS results. The coefficient of variation (COV) δ of the estimated probability of failure P_f obtained using full-scale MCS for the sample size N_s is computed as:

$$\delta = \sqrt{\frac{(1 - P_f)}{N_s P_f}} \quad (50)$$

Furthermore, comparison with FORM, HO-SRSM and conventional G-ANOVA have also been presented in order to demonstrate the elegance of the proposed approach.

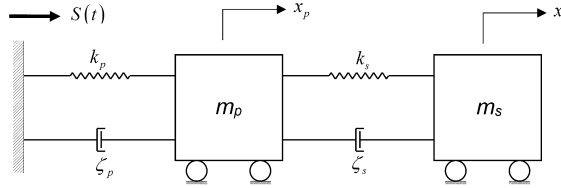


Fig. 3. 2-DOF damped oscillator considered in Example 1. The system involves 8 random variables.

Table 1
Statistical properties of random variables for Example 1.

Variable	Distribution	Mean	COV
m_p	Lognormal	1.5	0.1
m_s	Lognormal	0.01	0.1
k_p	Lognormal	1	0.2
k_s	Lognormal	0.01	0.2
ζ_p	Lognormal	0.05	0.4
ζ_s	Lognormal	0.02	0.5
F_s	Lognormal	15	0.1
S_0	Lognormal	100	0.1

All the finite element analysis in this paper has been carried out using *ABAQUS*TM version 6.8 [70]. For performing structural reliability analysis, a MATLAB-PYTHON-ABAQUS interface has been developed. Results obtained have been compared on the basis of accuracy and computational cost. When comparing computational cost in estimating failure probability, the number of actual FE analyses is chosen as the basic comparison tool because it indirectly represents the CPU time usage.

6.1. Example 1: 2-DOF damped oscillator [50]

In this example, a 2-DOF damped oscillator has been considered. Fig. 3 shows a schematic diagram of the system under consideration. The system is characterised by masses m_p and m_s , damping ratios ζ_p and ζ_s and stiffnesses k_p and k_s . The governing differential equation of the system is given as:

$$\mathbf{M}\ddot{\mathbf{x}} + \mathbf{C}\dot{\mathbf{x}} + \mathbf{K}\mathbf{x} = \mathbf{F} \quad (51)$$

where

$$\begin{aligned} \mathbf{M} &= \begin{bmatrix} m_p & 0 \\ 0 & m_s \end{bmatrix} \\ \mathbf{C} &= \begin{bmatrix} \zeta_p \omega_p m_p + \zeta_s \omega_s m_s & -\zeta_s \omega_s m_s \\ -\zeta_s \omega_s m_s & \zeta_s \omega_s m_s \end{bmatrix} \\ \mathbf{K} &= \begin{bmatrix} \omega_p^2 m_p + \omega_s^2 m_s & -\omega_s^2 m_s \\ -\omega_s^2 m_s & \omega_s^2 m_s \end{bmatrix} \end{aligned} \quad (52)$$

The external force \mathbf{F} in Eq. (51) arises due to a base acceleration $S(t)$. In this work, $S(t)$ is modelled as a white noise.

The limit-state equation of the system is considered to be:

$$g(\mathbf{x}) = F_s - f_p k_s \cdot \max_{t \in [0; \tau]} |x_s(t)| \quad (53)$$

where τ denotes the duration of excitation, F_s is the force capacity of the secondary spring and f_p is the peak factor. In this work, $f_p = 3$ has been considered. All the parameters, namely $[m_p, m_s, \zeta_p, \zeta_s, k_p, k_s, S_0, F_s]$ are considered to be random. The statistical properties of the random variables are shown in Table 1.

Proposed approach has been utilized for structural reliability analysis of the above mentioned problem. Moreover, full-scale MCS with target COV of 1.5%, FORM and HO-SRSM have also been performed. Results obtained are summarized in Table 2. While FORM ($P_f = 0.0219$, $N_s = 2862$) overestimates the failure probability by 356.25%, HO-SRSM ($P_f = 0.00126$, $N_s = 1594$) underestimate the failure probability by 73.75%. It is observed that the proposed DA-SED based G-ANOVA and ordinary G-ANOVA yields almost identical results. However, the number of training points required using the proposed approach ($N_s = 248$) is less as compared to the ordinary G-ANOVA ($N_s = 300$) [50].

6.2. Example 2: 10 storey building with mass damper [50]

In this example, a 10 storey building with a mass-damper installed at roof has been considered. Fig. 4 shows a schematic diagram of the building. The system is subjected to an external excitation generated due to ground acceleration \ddot{a}_g . The i^{th}

Table 2
Results obtained in Example 1.

Method	Reliability index	Failure probability	N_s^a
FORM	2.0160	0.0219	2862
HO-SRSM	3.0209	0.00126	1594
direct MCS	2.5899	0.0048	10^6
G-ANOVA	2.5828	0.0049	300
PA ^b	2.5899	0.0048	248

^a N_s indicates total number of actual function evaluations.

^b PA = proposed approach, i.e., DA-SED based G-ANOVA.

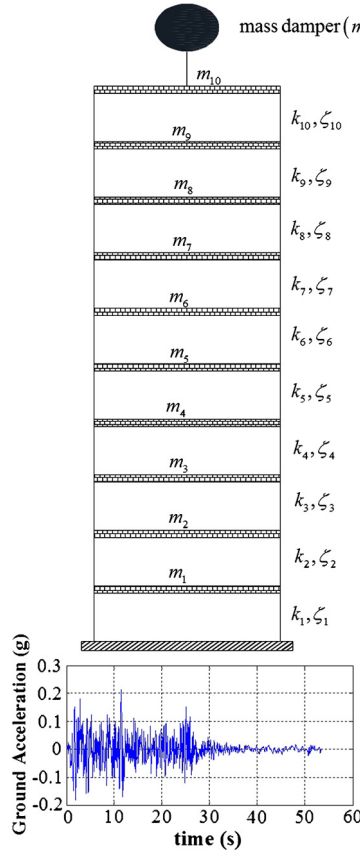


Fig. 4. Schematic diagram of 10-storey building with mass damper considered in Example 2. Time history of ground acceleration is also shown below. The system considered involves 33 random variables.

floor of the building is having a mass m_i , stiffness k_i and damping ratio ζ_i . The mass damper, installed at the roof, is having mass, stiffness and damping ratio of m_0 , k_0 and ζ_0 respectively. In this example, all the above mentioned parameters are considered to be random. Therefore, the system under consideration is having 33 random variables. Statistical properties of the random variables are shown in Table 3. The limit-state function of the above system is given as:

$$g(\mathbf{x}) = V_{\text{allowable}} - V_{\text{Base}}(\mathbf{x}) \quad (54)$$

where $V_{\text{Base}}(\mathbf{x})$ is the maximum base shear and $V_{\text{allowable}} = 7 \times 10^5$ kN. The structure has been modelled using FE package ABAQUSTM version 6.8 [70]. Proposed approach has been employed to determine probability of failure of the above mentioned system. For validation purpose, a full-scale MCS with target COV of 4.5% have been performed. Moreover in order to establish the superiority of the proposed approach, results obtained have been compared with that obtained using FORM, HO-SRSM and conventional/ordinary G-ANOVA. The probability of failure obtained are tabulated in Table 4. FORM did not converged even after 200 iterations and 7994 function evaluations. It is observed that the proposed approach yields the best result. Moreover, the number of training points required using the proposed approach is significantly less as compared to other approaches.

Table 3

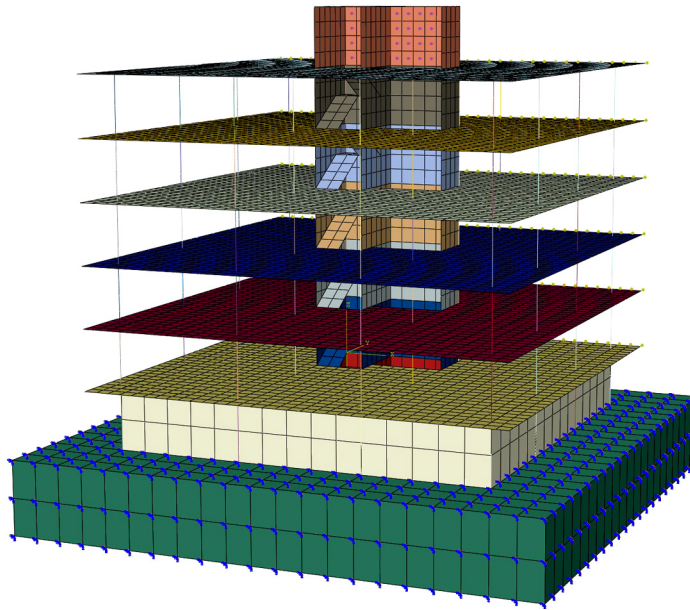
Statistical parameters of random variables for Example 2.

Random variables	m_1, m_2, \dots, m_{10}	k_1, k_2, \dots, k_{10}	m_0	k_0	$\zeta_0, \zeta_1, \dots, \zeta_{10}$
Distribution	Lognormal	Lognormal	Lognormal	Lognormal	Lognormal
Mean	8.75×10^4 kg	2.1×10^8 N/m	7.16×10^4 kg	3.85×10^6 N/m	0.05
COV	0.2	0.2	0.2	0.2	0.3

Table 4

Probability of failure of the system described in Example 2.

Method	Reliability index	Failure probability	N_s^a
FORM		not converged	
HO-SRSM	1.3330	0.09127	12193
Direct MCS	1.4010	0.08060	6000
G-ANOVA	1.3877	0.08261	1800
PA ^b	1.3951	0.08159	1167

^a N_s indicates total number of actual function evaluations.^b PA = proposed approach, i.e., DA-SED based G-ANOVA.**Fig. 5.** Six-storey building considered in Example 3. This system involves 39 random variables.

6.3. Example 3: a six storey buildings [47]

This example examines the performance of the proposed approach in predicting probability of failure of a large scale system. A six-storey building, as shown in Fig. 5, has been considered. The building is subjected to lateral and gravity loads. The plan area of the building is 18(m) \times 18(m). Each floor is having 16 columns at spacing of (5 m). Furthermore, a lift-wall is situated at the centre of the building. The floor to floor height is considered to be (3.6 m).

The building described above is modelled using ABAQUSTM [70]. While the floors and walls (lift wall and load bearing wall of ground floor) are modelled as 4-noded shell element (S4), the columns are modelled as 2-noded beam element (B31). The soil is modelled as 8-noded brick element (C3D8I).

The cross-section (both depth (d_c) and width (w_c)) of columns, thickness of slabs (both floor (t_{floor}) and staircase slab (t_{stair})), thickness of walls ($t_{cellular}$), density (ρ) and elastic modulus (E) of concrete, cellular material (used for load bearing wall) and soil, upper (Y_U) and lower yield point (Y_L) of concrete are considered to be independent random variables. The seven lateral loads, namely P_{f1} , P_{f2} , P_{f3} , P_{f4} , P_{f5} , P_{f6} and P_{surf} , acting at floor 1–6 and terrace lift wall respectively are also considered to be independent random variables. The system has 39 random variables. The details of random variables are shown in Table 5.

The limit-state for the system described above is given as:

$$g(\mathbf{X}) = \sigma_0 - (\sigma_{mises})_{\max} \quad (55)$$

Table 5

Details of random variables for six storey building defined in Example 3.

Variable no.	Variable	Mean	COV	Distribution
1	$E_{cellular}$ (N/m ²)	2.5×10^{10}	0.05	Lognormal
2–4	$E_{col}, E_{floor}, E_{stair}$ (N/m ²)	3.5×10^{10}	0.05	Lognormal
5	G_{col} (N/m ²)	1.4×10^{10}	0.05	Lognormal
6–8	$\rho_{col}, \rho_{floor}, \rho_{stairs}$ (kg/m ³)	2500	0.1	Normal
9	Y_L (N/m ²)	2.0×10^7	0.05	Lognormal
10	Y_U (N/m ²)	2.2×10^7	0.05	Lognormal
11	t_{stair} (m)	0.1	0.2	Uniform
12	$t_{cellular}$ (m)	0.5	0.2	Uniform
13–18	t_{floor1} (m) – t_{floor6} (m)	0.125	0.2	Uniform
19	E_{soil} (N/m ²)	5×10^9	0.1	Lognormal
20	ρ_{soil} (kg/m ³)	1800	0.15	Normal
21–24	$d_{cf1}, w_{cf1}, d_{cf2}, w_{cf2}$ (m)	0.8	0.2	Uniform
25–28	$d_{cf3}, w_{cf3}, d_{cf4}, w_{cf4}$ (m)	0.6	0.2	Uniform
29–32	$d_{cf5}, w_{cf5}, d_{cf6}, w_{cf6}$ (m)	0.4	0.2	Uniform
33	P_1 (N/m)	2722	0.15	Rayleigh
34	P_2 (N/m)	2528	0.15	Rayleigh
35	P_3 (N/m)	2917	0.15	Rayleigh
36	P_4 (N/m)	3111	0.15	Rayleigh
37	P_5 (N/m)	3306	0.15	Rayleigh
38	P_6 (N/m)	1750	0.15	Rayleigh
39	P_{surf} (N/m ²)	800	0.15	Uniform

Table 6

Probability of failure of the system described in Example 3.

Method	Reliability index	Failure probability	N_s^a
FORM	1.4473	0.0739	4594
HO-SRSM	1.2751	0.1011	13758
Direct MCS	1.0450	0.1480	20000
G-ANOVA	1.0572	0.1452	5000
PA ^b	1.0490	0.1471	2628

^a N_s indicates total number of actual function evaluations.^b PA = proposed approach, i.e., DA-SED based G-ANOVA.

where $(\sigma_{mises})_{max}$ is the maximum von Mises stress and σ_0 is the threshold value. In this example, $\sigma_0 = 1.1 \times 10^7$ N/m² is considered.

The proposed approach has been utilized for predicting the probability of failure of the above mentioned system (Table 6). Result obtained has been validated against direct MCS results with target COV of 2%. In order to demonstrate the superiority of the proposed DA-SED based G-ANOVA, comparison with FORM, HO-SRSM and conventional G-ANOVA has also been presented. The proposed approach is found to yield highly accurate results. Furthermore, the number of training points using the proposed approach is significantly less, as compared to other approaches.

7. Conclusion

In this paper, a novel approach for addressing a fundamental issue in structural reliability analysis of a system by using surrogate approach has been presented. It is argued that if all the training points are selected at once, the resulting surrogate may not be optimized (in terms of training point required). In order to address this issue, a sequential experimental design (SED) based generalised ANOVA has been proposed. The primary findings/contribution of this work are summarized below:

1. A novel distribution adaptive sequential experimental design (DA-SED) scheme has been proposed. The proposed DA-SED incorporates the knowledge regarding the statistical properties of the input variables in selecting the training points.
2. The framework of G-ANOVA has been improved by incorporating the polynomial chaos expansion into the framework of G-ANOVA. It is further argued that the proposed G-ANOVA is more efficient, as compared to conventional G-ANOVA.
3. Generalise analytical expressions for the first two statistical moments of response have been derived. The statistical moments obtained are further coupled with the formula presented in [47]. This coupling results in an analytical framework for structural reliability analysis.
4. A unified framework that couples the proposed DA-SED with the modified G-ANOVA has been presented. The framework utilizes leave one out test to determine the error in each step.

The proposed approach has been utilized for structural reliability analysis of three structural mechanics problems. It is demonstrated that the proposed approach is highly accurate in predicting the probability of failure. Furthermore, number of training points required using the proposed approach is significantly less as compared to conventional G-ANOVA.

Acknowledgements

SC acknowledges the support of MHRD, Government of India. RC acknowledges the support of Royal Society through Newton Alumni Funding.

Appendix A. Transformation to non-uniform variables

A.1. Uniform to Gaussian variable

$$z = \mu + \sigma \left(\Phi^{-1}(x) \right) \quad (\text{A.1})$$

where x is uniformly distributed in $[0, 1]$ and z is a normal variable with mean μ and standard deviation σ .

A.2. Uniform to lognormal distribution

$$z = \exp \left(\mu + \sigma \left(\Phi^{-1}(x) \right) \right) \quad (\text{A.2})$$

where x is uniformly distributed in $[0, 1]$ and z is a lognormal variable with parameters μ and σ .

A.3. Uniform to Gumbell distribution

$$z = a - b \log(-\log(x)) \quad (\text{A.3})$$

where x is uniformly distributed in $[0, 1]$ and z follows Gumbell distribution with parameters a and b .

A.4. Uniform to Rayleigh distribution

$$z = b + \sqrt{-2a^2 \log(x)} \quad (\text{A.4})$$

where x is uniformly distributed in $[0, 1]$ and z follows Rayleigh distribution with parameters a and b .

References

- [1] G. Muscolino, G. Ricciardi, P. Cacciola, Monte Carlo simulation in the stochastic analysis of non-linear systems under external stationary Poisson white noise input, *Int. J. Non-Linear Mech.* 38 (2003) 1269–1283.
- [2] R. Rubenstein, *Simulation and the Monte Carlo Method*, Wiley, New York, 1981.
- [3] R. Thakur, K. Misra, Monte Carlo simulation for reliability evaluation of complex systems, *Int. J. Syst. Sci.* 9 (1978) 1303–1308.
- [4] W.A. Ericson, Optimum stratified sampling using prior information, *J. Am. Stat. Assoc.* 60 (311) (1965) 750, <http://dx.doi.org/10.2307/2283243>.
- [5] Z. Feng, L. Zhenzhou, C. Lijie, S. Shufang, Reliability sensitivity algorithm based on stratified importance sampling method for multiple failure modes systems, *Chin. J. Aeronaut.* 23 (6) (2010) 660–669.
- [6] M. Hohenbichler, R. Rackwitz, Improvement of second-order reliability estimates by importance sampling, *J. Eng. Mech.* 114 (12) (1988) 2195–2199.
- [7] S.K. Au, J.L. Beck, A new adaptive importance sampling scheme for reliability calculations, *Struct. Saf.* 21 (2) (1999) 135–158.
- [8] W.L. Jin, E. Luz, Improving importance sampling method in structural reliability, *Nucl. Eng. Des.* 147 (3) (1994) 393–401.
- [9] O. Ditlevsen, P. Bjerager, R. Olesen, A.M. Hasofer, Directional simulation in Gaussian processes, *Probab. Eng. Mech.* 3 (4) (1988) 207–217.
- [10] O. Ditlevsen, R.E. Melchers, H. Gluwer, General multidimensional probability integration by directional simulation, *Comput. Struct.* 36 (2) (1990) 355–368.
- [11] J.S. Nie, B.R. Ellingwood, A new directional simulation method for system reliability. Part I: Application of deterministic point sets, *Probab. Eng. Mech.* 19 (4) (2004) 425–436.
- [12] S.K. Au, J.L. Beck, Estimation of small failure probabilities in high dimensions by subset simulation, *Probab. Eng. Mech.* 16 (4) (2001) 263–277.
- [13] K.M. Zuev, J.L. Beck, S.-K. Au, L.S. Katafygiotis, Bayesian post-processor and other enhancements of subset simulation for estimating failure probabilities in high dimensions, *Comput. Struct.* 92–93 (2012) 283–296.
- [14] A.D. Kiureghian, H. Lin, S. Hwang, Second order reliability approximations, *J. Eng. Mech.* 113 (8) (1987) 1208–1225.
- [15] A. Hasofer, N. Lind, An exact and invariant first order reliability format, *J. Eng. Mech.* 100 (1) (1974) 111–121.
- [16] G. Cai, I. Elishakoff, Refined second-order reliability analysis, *Struct. Saf.* 14 (1994) 267–276.
- [17] S. Adhikari, Reliability analysis using parabolic failure surface approximation, *J. Eng. Mech.* 130 (12) (2004) 1407–1427.
- [18] S. Adhikari, Asymptotic distribution method for structural reliability analysis in high dimensions, *Proc. R. Soc. A, Math. Phys. Eng. Sci.* 461 (2062) (2005) 3141–3158.
- [19] L. Faravelli, Response-surface approach for reliability-analysis, *J. Eng. Mech.* 115 (12) (1989) 2763–2781.
- [20] R. Myers, *Response Surface Methodology*, Allyn and Bacon, Inc., Boston, 1971.

- [21] M.R. Rajashekhar, B.R. Ellingwood, Reliability of reinforced-concrete cylindrical Shells, *J. Struct. Eng.* 121 (2) (1995) 336–347.
- [22] H.P. Gavin, S.C. Yau, High-order limit state functions in the response surface method for structural reliability analysis, *Struct. Saf.* 30 (2) (2008) 162–179.
- [23] S. Goswami, S. Ghosh, S. Chakraborty, Reliability analysis of structures by iterative improved response surface method, *Struct. Saf.* 60 (2016) 56–66.
- [24] S. Chakraborty, A. Sen, Adaptive response surface based efficient finite element model updating, *Finite Elem. Anal. Des.* 80 (2014) 33–40.
- [25] S.-C. Kang, H.-M. Koh, J.F. Choo, An efficient response surface method using moving least squares approximation for structural reliability analysis, *Probab. Eng. Mech.* 25 (4) (2010) 365–371, <http://dx.doi.org/10.1016/j.pro bengmech.2010.04.002>.
- [26] V. Dubourg, B. Sudret, M. Bourinet, Reliability-based design optimization using kriging surrogates and subset simulation, *Struct. Multidiscip. Optim.* 44 (5) (2011) 673–690.
- [27] I. Kaymaz, Application of kriging method to structural reliability problems, *Struct. Saf.* 27 (2) (2005) 133–151, <http://dx.doi.org/10.1016/j.strusafe.2004.09.001>.
- [28] D. Lazzaro, L.B. Montefusco, Radial basis functions for the multivariate interpolation of large scattered data sets, *J. Comput. Appl. Math.* 140 (1–2) (2002) 521–536.
- [29] S.D. Marchi, G. Santin, A new stable basis for radial basis function interpolation, *J. Comput. Appl. Math.* 253 (2013) 1–13.
- [30] D. Xiu, G.E. Karniadakis, The Wiener–Askey polynomial chaos for stochastic differential equations, *SIAM J. Sci. Comput.* 24 (2) (2002) 619–644.
- [31] G. Blatman, B. Sudret, Adaptive sparse polynomial chaos expansion based on least angle regression, *J. Comput. Phys.* 230 (6) (2011) 2345–2367.
- [32] B. Pascual, S. Adhikari, A reduced polynomial chaos expansion method for the stochastic finite element analysis, *Sadhana. Acad. Proc. Eng. Sci.* 37 (3) (2012) 319–340.
- [33] S. Adhikari, A reduced spectral function approach for the stochastic finite element analysis, *Comput. Methods Appl. Mech. Eng.* 200 (21–22) (2011) 1804–1821.
- [34] A. Kundu, S. Adhikari, Dynamic analysis of stochastic structural systems using frequency adaptive spectral functions, *Probab. Eng. Mech.* 39 (2015) 23–38.
- [35] J.B. Cardoso, J.R. de Almeida, J.M. Dias, P.G. Coelho, Structural reliability analysis using Monte Carlo simulation and neural networks, *Adv. Eng. Softw.* 39 (6) (2008) 505–513.
- [36] A. Hosni Elhewy, E. Mesbahi, Y. Pu, Reliability analysis of structures using neural network method, *Probab. Eng. Mech.* 21 (1) (2006) 44–53.
- [37] J.-M. Bourinet, F. Deheeger, M. Lemaire, Assessing small failure probabilities by combined subset simulation and support vector machines, *Struct. Saf.* 33 (6) (2011) 343–353.
- [38] R. Lebrun, A. Dutfoy, An innovating analysis of the Nataf transformation from the copula viewpoint, *Probab. Eng. Mech.* 24 (3) (2009) 312–320.
- [39] L. Da-gang, First order reliability method based on linearized Nataf transformation, *Eng. Mech.* 24 (5) (2007) 0–086.
- [40] Ö.F. Aliş, H. Rabitz, Efficient implementation of high dimensional model representations, *J. Math. Chem.* 29 (2) (2001) 127–142.
- [41] H. Rabitz, Ö.F. Aliş, General foundations of high dimensional model representations, *J. Math. Chem.* 25 (2–3) (1999) 197–233.
- [42] I.M. Sobol, Sensitivity estimates for nonlinear mathematical models, *Math. Model. Comput. Exp.* 1 (4) (1993) 407–414.
- [43] S. Chakraborty, R. Chowdhury, Uncertainty propagation using hybrid HDMR for stochastic field problems, in: *International Conference on Structural Engineering and Mechanics*, 2013.
- [44] X. Ma, N. Zabarar, An adaptive high-dimensional stochastic model representation technique for the solution of stochastic partial differential equations, *J. Comput. Phys.* 229 (10) (2010) 3884–3915, <http://dx.doi.org/10.1016/j.jcp.2010.01.033>.
- [45] G. Hooker, Generalized functional ANOVA diagnostics for high-dimensional functions of dependent variables, *J. Comput. Graph. Stat.* 16 (3) (2007) 709–732.
- [46] G. Li, H. Rabitz, General formulation of HDMR component functions with independent and correlated variables, *J. Math. Chem.* 50 (1) (2012) 99–130.
- [47] S. Chakraborty, R. Chowdhury, A semi-analytical framework for structural reliability analysis, *Comput. Methods Appl. Mech. Eng.* 289 (1) (2015) 475–497.
- [48] S. Chakraborty, R. Chowdhury, Polynomial correlated function expansion for nonlinear stochastic dynamic analysis, *J. Eng. Mech.* 141 (3) (2014) 04014132.
- [49] S. Chakraborty, B. Mandal, R. Chowdhury, A. Chakrabarti, Stochastic free vibration analysis of laminated composite plates using polynomial correlated function expansion, *Compos. Struct.* 135 (2016) 236–249.
- [50] S. Chakraborty, R. Chowdhury, Assessment of polynomial correlated function expansion for high-fidelity structural reliability analysis, *Struct. Saf.* 59 (2016) 9–19.
- [51] S.W. Wang, P.G. Georgopoulos, G. Li, H. Rabitz, Random sampling-high dimensional model representation (RS-HDMR) with nonuniformly distributed variables: application to an integrated multimedia/multipathway exposure and dose model for trichloroethylene, *J. Phys. Chem. A* 107 (23) (2003) 4707–4718.
- [52] G. Li, J. Hu, S.-W. Wang, P.G. Georgopoulos, J. Schoendorf, H. Rabitz, Random sampling-high dimensional model representation (RS-HDMR) and orthogonality of its different order component functions, *J. Phys. Chem. A* 110 (7) (2006) 2474–2485.
- [53] R. Chowdhury, B.N. Rao, A.M. Prasad, High dimensional model representation for piece-wise continuous function approximation, *Commun. Numer. Methods Eng.* 24 (12) (2007) 1587–1609.
- [54] R. Chowdhury, B.N. Rao, A.M. Prasad, High-dimensional model representation for structural reliability analysis, *Commun. Numer. Methods Eng.* 25 (4) (2009) 301–337, <http://dx.doi.org/10.1002/cnm.1118>.
- [55] G. Li, H. Rabitz, J. Hu, Z. Chen, Y. Ju, Regularized random-sampling high dimensional model representation (RS-HDMR), *J. Math. Chem.* 43 (3) (2007) 1207–1232.
- [56] D. Xiu, G.E. Karniadakis, Modeling uncertainty in steady state diffusion problems via generalized polynomial chaos, *Comput. Methods Appl. Mech. Eng.* 191 (2002) 4927–4948.
- [57] B. Sudret, Global sensitivity analysis using polynomial chaos expansions, *Reliab. Eng. Syst. Saf.* 93 (7) (2008) 964–979.
- [58] S. Chakraborty, R. Chowdhury, Multivariate function approximations using the D-MORPH algorithm, *Appl. Math. Model.* 39 (23–24) (2015) 7155–7180, <http://dx.doi.org/10.1016/j.apm.2015.03.008>.
- [59] G. Li, H. Rabitz, D-MORPH regression: application to modeling with unknown parameters more than observation data, *J. Math. Chem.* 48 (4) (2010) 1010–1035.
- [60] G. Li, R. Rey-de Castro, H. Rabitz, D-MORPH regression for modeling with fewer unknown parameters than observation data, *J. Math. Chem.* 50 (7) (2012) 1747–1764.
- [61] R. Fisher, *The Design of Experiments*, first ed., Oliver and Boyd, 1935.
- [62] H. Dette, T. Holland-Letz, A geometric characterization of C-optimal designs for heteroscedastic regression, *Ann. Stat.* 37 (6B) (2009) 4088–4103.
- [63] M. Stein, Large sample properties of simulations using Latin hypercube sampling, *Technometrics* 29 (2) (1987) 143.
- [64] M. Schwaab, F.M. Silva, C.A. Queipo, A.G. Barreto, M. Nele, J.C. Pinto, A new approach for sequential experimental design for model discrimination, *Chem. Eng. Sci.* 61 (17) (2006) 5791–5806.
- [65] P. Bratley, B.L. Fox, Implementing Sobol's quasirandom sequence generator, *ACM Trans. Math. Softw.* 14 (1) (1988) 88–100.
- [66] I.M. Sobol, Uniformly distributed sequences with an additional uniform property, *USSR Comput. Math. Math. Phys.* 16 (1976) 236–242.

- [67] S. Galanti, A. Jung, Low-discrepancy sequences: Monte Carlo simulation of option prices, *J. Deriv.* 5 (1) (1997) 63–83.
- [68] H. Faure, Good permutations for extreme discrepancy, *J. Number Theory* 42 (1) (1992) 47–56.
- [69] J.H. Halton, On the efficiency of certain quasi-random sequences of points in evaluating multi-dimensional integrals, *Numer. Math.* 2 (1) (1960) 84–90.
- [70] Dassault Systèmes Simulia Corp., ABAQUS Documentation and Theory Manual, 2009.