Statistical correlation in stratified sampling

M. Vořechovský & D. Novák

Institute of Structural Mechanics, Faculty of Civil Engineering, Brno University of Technology, Brno, Czech Republic

Keywords: correlation, Monte Carlo simulation, Latin hypercube sampling, statistical analysis, sensitivity, reliability, stochastic optimization, Simulated Annealing

ABSTRACT: A new efficient technique to impose the statistical correlation when using the Monte Carlo type method for the statistical analysis of computational problems is proposed. The technique is based on the stochastic optimization method called Simulated Annealing. The comparison with other techniques presently used and intensive numerical testing showed the superiority and robustness of the method. No significant obstacles have been found when also working with large problems (large number of random variables). The advantages and limitations of the approach will be discussed. Remarks on the positive definiteness of target correlation matrix are made. Numerical examples show the efficiency of the method.

1 INTRODUCTION

The aim of the statistical and reliability analyses of any computational problem which can be numerically simulated is mainly the estimation of the statistical parameters of the response variable and/or the theoretical failure probability. The pure Monte Carlo simulation cannot be applied for time-consuming problems, as it requires a large number of simulations (repetitive calculation of response). A small number of simulations can be used for the acceptable accuracy of statistical characteristics of response using the stratified sampling technique Latin Hypercube Sampling (LHS), e.g. McKey et al. (1979), Iman & Conover (1980, 1982) or Novák et al. (1998, 2000). Briefly, it is a special type of the Monte Carlo numerical simulation which uses the stratification of the theoretical probability distribution functions of the input random variables. It requires relatively a small number of simulations (from tens to hundreds) - repetitive calculations of the response resulting from the analyzed computational model. The LHS strategy has been used by many authors in different fields of engineering and with both a simple and a very complicated computational model, e.g. Novák et al. (1998). The classical reliability theory introduced the basic concept using formally the response variable $Z = g(\mathbf{X})$, where g (computational model) represents the functional relationship between the elements of vector X. The elements of vector X are generally uncertainties

(random variables). These quantities can also be naturally statistically correlated. The paper is focused on the problem of the efficient imposition of the statistical correlation within the framework of the Monte Carlo type simulation (preferably LHS), Vořechovský & Novák (2002). The techniques presently available are discussed first.

2 SAMPLING AND STATISTICAL CORRELATION

There are two stages to the Latin Hypercube Sampling. First, samples for each (marginal) variable are chosen strategically to represent the variable's probability density function PDF. The N_{Sim} samples (where N_{Sim} is the number of the planned simulations) for each random variable X_i are often chosen from the cumulative distribution function (CDF) by the inverse transformation of CDF.

Table 1 Sampling scheme for N_{Sim} deterministic calculations of $g(\mathbf{X})$

Simulation	Var. 1	Var. 2	Var. 3	 Var. _{Nv}	
1	$x_{I, I}$	$x_{I, 2}$	<i>x</i> _{1,3}	 $x_{I, Nv}$	
2	$x_{2, l}$	$x_{2, 2}$	$x_{2, 3}$	 $x_{2, Nv}$	
•••				 	
N_{Sim}	X _{NSim, 1}	X _{NSim, 2}	X _{NSim, 3}	 $x_{NSim,Nv}$	

Then the samples for the variables are ordered to match the target correlation among themselves. Further we presume using the LHS methodology for the sampling. Table 1 represents the sampling scheme, where the simulations represented by rows and columns are related to the random variables (NV is the number of the input variables).

As mentioned previously, the first stage of LHS is to generate representing samples for each random variable. The domain of each variable is divided into equiprobable disjunct intervals of the probability N_{Sim} . One sample is chosen from each interval. The current practice is to choose samples directly by inverse transformation of CDF, in the middle of the *k*-th strata:

$$x_{i,k} = F_i^{-1} \left(\frac{k - 0.5}{N_{Sim}} \right)$$
(1)

where $x_{i,k}$ is *k*-th sample of *i*-th variable X_i , and F_i^{-1} is the inverse CDF for the variable X_i . The objection against the approach deals mainly with the tails of PDF, which mostly influences the variance, skewness and kurtosis of the sample set. However this elementary approach was overcome by sampling the mean values related to the intervals (e.g. Huntington & Lyrintzis, 1998), see Fig. 1:

$$x_{i,k} = \frac{\int_{y_{i,k-1}}^{y_{i,k}} x \cdot f_i(x) \, dx}{\int_{y_{i,k-1}}^{y_{i,k}} f_i(x) \, dx} = N_{Sim} \cdot \int_{y_{i,k-1}}^{y_{i,k}} x \cdot f_i(x) \, dx \tag{2}$$

where f_i is PDF of the variable X_i , and the integration limits are:

(3)



Fig. 1 Illustration of sampling from marginals.

Than the samples represent the marginal PDF better. The estimated mean value is achieved accurately (analytical determination, definition) and the variance of the sample set is much closer to the target one. For some probability density functions (inclusive e.g. Gaussian, Exponential, Laplace, Rayleigh, Logistic, Pareto, etc.) the integral (2) can be solved analytically. For others, the extra effort of doing the numerical integration is definitely worthwhile. Samples determined by both the approaches are nearly identical excluding the tail samples. Therefore the second approach is recommended especially around the tails of the distributions.

Once samples are generated, the correlation structure according to the target correlation matrix must be taken into account. There are generally two problems related to the statistical correlation: First, during sampling an undesired correlation can be introduced between the random variables (columns in Table 1). For example, instead of the correlation coefficient zero for the uncorrelated random variables, i.e. an undesired correlation, e.g. 0.6 can be generated. It can happen especially in the case of a very small number of simulations (tens), where the number of interval combination is rather limited. The second task is to introduce the prescribed statistical correlation between the random variables defined by the correlation matrix. The columns in Table 1 should be rearranged in such a way that they may fulfill the following two requirements: to diminish the undesired random correlation and to introduce the prescribed correlation. The efficiency of the LHS technique was showed for the first time by McKay and Conover W.J (1979), but only for the uncorrelated random variables. The first technique for the generation of the correlated random variables has been proposed by Iman & Conover (1982). The authors also showed the alternative of diminishing the undesired random correlation. The technique is based on the iterative updating of the sampling matrix; the Cholesky decomposition of the correlation matrix has to be applied. As a measure of the statistical correlation, the Spearman correlation coefficient is used:

$$T_{ij} = 1 - \frac{6\sum_{k} (R_{ki} - R_{kj})^2}{N(N-1)(N+1)}$$
(4)

where *R* is the $(N_{Sim}:N_V)$ matrix containing a permutation of the rank numbers in each column and coefficients T_{ij} represents the Spearman's correlation coefficient between the variables *i* and *j*, $T_{ij} \in \langle -1; 1 \rangle$. The correlation matrix *T* is symmetric, positive definite (unless some columns have identical ordering). Therefore the Cholesky decomposition may be performed:

$$T = Q^T \cdot Q \tag{5}$$

and the new ordering matrix can be generated as follows:

$$R_B = R \cdot Q^{-1} \tag{6}$$

Then the rank numbers in each column of the ordering matrix R are then arranged to have the same ordering as the numbers in each column of R_B . The technique can be applied iteratively and can result in a very low correlation coefficient if generating uncorrelated random variables. But Huntington & Lyrintzis (1998) have found that the approach tends to converge to an ordering which still gives significant correlation errors between some variables.

The scheme has more difficulties when simulating correlated variables. The correlation procedure can be performed only once, there is no way to iterate it and to improve the result.

The described scheme is linked to the Spearman correlation measure which is very robust in the cases of the nonGaussian (and different) marginal densities. It uses the ranks only instead of the sample values but the limitation is that the number of simulations have to be higher than the number of random variables to achieve the positive definite correlation matrix. It can be understood as a serious drawback in the cases of the utilization of the LHS technique for the cases of a very high number of the variables and a limited number of simulations executable, e.g. in the random field simulation for the stochastic finite element calculations. method is needed for the simulation of both the uncorrelated and the correlated random variables. Such a method should be efficient enough: reliable, robust and fast.

All approaches discussed above and further presume the imposition of the target correlation structure only by matrix manipulations. The task can be understood as the simulation from the multivariate distribution model consistent with the prescribed marginals and covariances. Most existing models for random vectors, however, are restricted to the bivariate case and/or can only describe the small correlation between variables. Two models based on the earlier works of Nataf and Morgenstern are recommended by Liu and Kiureghian (1986).

Due to some limitations of both the models, this paper tries to find a solution of the problem by changing ranks of the samples instead of their values, while the marginal probability density functions remain intact.

The representation of 2D marginals for two correlation coefficients is illustrated in Figures 2 and 3, the sample values for each variable are used as coordinates for the samples representing the joint PDF.



Fig. 2 Example of negative statistical dependence between samples representing a random vector.

These obstacles stimulated the work of Huntington & Lyrintzis (1998), they proposed the so called single-switch-optimized ordering scheme. The approach is based on the iterative switching of the pair of samples of Table 1 which gives the greatest reduction in error. The authors showed that their technique performs clearly well enough but it may still converge to a non-optimum ordering. A different



Fig. 3 Example of statistically independent samples.

Note that the accurate best result is obtained if all possible combinations of the ranks for each column (variable) itself in Table 1 are treated. It would be necessary to try an extremely large number of the rank combinations $(N_{Sim}!)^{Nv-1}$. It is clear that this rough approach is hardly applicable in spite of the fast development of computer hardware.

3 STOCHASTIC OPTIMIZATION SIMULATED ANNEALING

The imposition of the prescribed correlation matrix into the sampling scheme can be understood as an optimization problem: The difference between the prescribed **K** and the generated **S** correlation matrices should be as small as possible. A suitable measure of quality of the overall statistical properties can be introduced, e.g. the maximal difference of the correlation coefficients between matrices:

$$E_{\max} = \max_{1 \le i < j \le N_V} \left| S_{i,j} - K_{i,j} \right|$$
(7)

or a norm which takes into account the deviations of all correlation coefficients:

$$E_{oveall} = \sqrt{\sum_{i=1}^{N_{v}-1} \sum_{j=i+1}^{N_{v}} (S_{i,j} - K_{i,j})^{2}}$$
(8)

The norm *E* has to be minimized from the point of view of the definition of the optimization problem: the objective function is E and the design variables are related to the ordering in the sampling scheme (Table 1). It is well known that the deterministic optimization techniques and the simple stochastic optimization approaches can very often fail to find the global minimum. Such a technique fails in some local minimum and then there is no chance to escape from it and to find the global minimum. It can be intuitively predicted that in our problem we are definitely facing the problem with a multiple local minima. Therefore we need to use the stochastic optimization method which works with some probability of escaping from the local minimum. The simplest form is the two-member evolution strategy which works in two steps: Mutation and selection.

1. Step 1 (mutation): In generation a new arrangement of the random permutations matrix \mathbf{X} is obtained using random changes of the ranks, one change is applied for one random variable. The generation should be performed randomly. Then the objective function (norm E) can be calculated using the newly obtained correlation matrix - it is usually called "offspring". The norm E calculated by using former arrangement is called "parent".

2. Step 2 (selection): The selection chooses the best norm between the "parent" and "offspring" to survive: For the new generation (permutation table arrangement) the best individual (table arrangement) has to give a value of the objective function (norm E) that is smaller than before.

Such an approach has been tested intensively using numbers of examples. It was observed that the method in most cases could not capture the global minimum. It failed in the local minimum and there was no chance to escape from it, as the only improvement of the norm E resulted in acceptance of "offspring".

The step "Selection" can be improved by Simulated Annealing approach (SA), a technique which is very robust concerning the starting point (initial arrangement of random permutations table). The SA is an optimization algorithm based on the randomization techniques and the incorporates aspects of the iterative improvement algorithms. Basically it is based on the Boltzmann probability distribution:

$$P_r(E) \approx e^{\left(\frac{-\Delta E}{k_b \cdot T}\right)} \tag{9}$$

where ΔE is the difference of norms E before and after the random change. This probability distribution expresses the concept of a system in thermal equilibrium at temperature T having its energy distributed probabilistically among all different energy states ΔE . The Boltzmann constant k_b relates to the temperature and energy of the system. Even at low temperatures there is a chance (although very small) of a system being locally in a high energy state. Therefore, there is a corresponding possibility for the system to move from the local energy minimum in favor of finding a better minimum. In other words, there is some probability of escaping from the local minimum. There are two alternatives in step 2 (mutation).

1. New arrangement - "offspring" results in the decrease of the norm E. Naturally "offspring" is accepted for a new generation.

2. New arrangement - "offspring" does not decrease the norm *E*. Such "offspring" is accepted with some probability according to the probability distribution (9). This probability changes as the temperature changes. Particularly, the offspring is accepted if a realization of random variable $Z = e^{(-\Delta E/T)} - R$ is positive, otherwise the offspring leading to negative *Z* is rejected. In the formula, *R* is uniformly distributed (rectangular) random variable over the interval (0,1) $R \sim R(0,1)$. The result is that there is a much higher probability of the global minimum being found in the comparison with the deterministic methods and the simple evolution strategies.

The method represents the analogy with the annealing of crystals. In the treated case k_b can be considered to be one. In the classical application of the SA approach for optimization there is one problem: how to set the initial temperature? Usually it should be considered heuristically. However our problem is constrained: the acceptable elements of the correlation matrix are always from interval $\langle -1, 1 \rangle$. Based on this fact, the maximum of the norm (2) can be estimated using the prescribed and hypothetically "most remote" unit correlation coefficients, plus or minus. This approach represents a significant advantage: The heuristic estimation of the initial temperature is neglected; the estimation can be performed without the guess of the user and the "trial and error" procedure. The initial temperature has to be decreased step by step, e.g. using reduction factor f_T after the constant number of iteration (e.g. thousands):

$$T_{i+1} = T_i \cdot f_T \tag{10}$$

The simple case is to use e.g. $f_T = 0.95$, note that the more sophisticated cooling schedules are known in the Simulated Annealing theory, e.g. Otten & Ginneken (1989).

The process of imposition of the correlation structure should be monitored through the graph similar to the graph in Figure 4, where the decrease of the norm vs. the number of switches is plotted. Such a figure is typical of most solutions using the Simulated Annealing.



Fig. 4 The norm E (error) vs. number of random changes (rank switches).

4 NUMERICAL EXAMPLES

4.1 Correlated properties of concrete

In order to illustrate the efficiency of the proposed technique, let us consider an example of the correlation matrix which corresponds to the properties of concrete. They are described by 7 random variables; the prescribed correlation matrix is presented in the lower triangle. The upper triangle shows the imposed statistical correlation using the Simulated Annealing (SA), for two different numbers of the LHS-simulations (8, 64). The final values of the norms are included on the right side: the first line corresponds to the norm (1), the second line (bold) means the overall norm (2).

It can be seen that as the number of simulations increases the correlation matrix is closer to the target one.

Another example of utilization is given in (Bažant, Novák & Vořechovský, 2003 or Lehký & Novák, 2002), where the simulation of the uncorrelated random variables was needed to represent the material strength over a structure instead of the random field approach.

	(1	-0.017	0.7	0.863	-0.026	5 0.48	7 0.8'	75)	
	0	1	0.02	0.067	-0.039	0.10	4 -0.0	17	
	0.7	0	1	0.729	-0.016	5 0.82	3 0.7	7 _	0.0007
$S_8 =$	0.9	0.1	0.8	1	0.024	0.54	3 0.80	53	0.0990
	0	0	0	0	1	0.03	1 -0.0	26 L	0.180
	0.5	0.1	0.9	0.6	0	1	0.48	87	
	0.9	0	0.6	0.9	0	0.5	1		
	(1	-0.001	0.697	0.902	0	0.502	0.898		
$\mathbf{S}_{64} =$	0	1	0.004	0.099	0	0.099	0.001		
	0.7	0	1	0.793	0	0.895	0.605	Γορο	727
	0.9	0.1	0.8	1	0	0.604	0.894	0.00	13
	0	0	0	0	1	0	0	0.0	[4]
	0.5	0.1	0.9	0.6	0	1	0.497		
	0.9	0	0.6	0.9	0	0.5	1		

4.2 Non-positive definite prescribed correlation matrix?

In real applications of the simulation technique in engineering (e.g. LHS), the statistical correlation represents very often the weak part of the a priori assumptions. Because of this pure knowledge, the prescribed correlation matrix on input can be nonpositive definite. The user may have difficulties in updating the correlation coefficients in order to make the matrix positive definite. The example presented here demonstrates that if the non-positive definite matrix is on input, the Simulated Annealing can work with it and the resulting correlation matrix is always positive definite. It is as close as possible to the originally prescribed matrix but the dominant constraint (positive definiteness) is satisfied automatically. Let us consider a very unrealistic simple case of the statistical correlation for three random variables A, B and C according to the matrix K (the columns and rows correspond to the ranks of the variables A, B, C):

$$\mathbf{K} = \begin{pmatrix} \mathbf{1} & 0.9 & 0.9 \\ \mathbf{1} & -0.9 \\ \text{sym.} & \mathbf{1} \end{pmatrix}, \ \mathbf{S}_1 = \begin{pmatrix} \mathbf{1} & 0.499 & 0.499 \\ \mathbf{1} & -0.499 \\ \text{sym.} & \mathbf{1} \end{pmatrix} \begin{bmatrix} 0.401 \\ \mathbf{0.695} \end{bmatrix}$$

The correlation matrix is obviously not positive definite. Strong positive statistical correlation is required between variables (A, B) and variables (A, C), but strong negative correlation between variables (B, C). It is clear that only the compromise solution can be done. The method resulted in such a compromise solution without any problem, S_1 (number of simulations N_{Sim} was high enough to avoid limitation in a number of the rank combinations). This feature of the method can be accepted and interpreted as an advantage of the method. In practice, there are the reliability problems with the non-positive definite-

ness (lack of knowledge). It represents limitation when using some other methods (the Cholesky decomposition of the prescribed correlation matrix).

In real applications there may be a greater confidence in one correlation coefficient (good data) and a smaller confidence in another one (just estimation). The solution of the mentioned problems is weighted calculations of both the norms (7) and (8). For example the norm (8) can be modified in the following way:

$$E_{overall} = \sqrt{\sum_{i=1}^{N_{v}-I} \sum_{j=i+1}^{N_{v}} w_{i,j} \cdot \left(S_{i,j} - K_{i,j}\right)^{2}}$$
(11)

where $w_{i,j}$ is the linear weight of a certain correlation coefficient. Several examples of choices and resulting correlation matrices (with both the norms) follow. The resulting matrices S_2 and S_3 demonstrate the similarity of the resulting errors (equivalent weights), while S_4 and S_5 illustrate the significance of the proportions between the weights. The weights are in lower the triangle and the matrix **K** is targeted again. The weights of the accentuated members and the resulting values are underlined.

$$\mathbf{S}_{2} = \begin{pmatrix} \mathbf{1} & 0.311 & 0.311 \\ l & \mathbf{1} & \underline{-0.806} \\ l & \underline{l0} & \mathbf{1} \end{pmatrix} \begin{bmatrix} 0.589 \\ \mathbf{0.884} \end{bmatrix}, \ \mathbf{S}_{3} = \begin{pmatrix} \mathbf{1} & 0.311 & \underline{0.806} \\ l & \mathbf{1} & -0.311 \\ \underline{l0} & l & \mathbf{1} \end{pmatrix} \begin{bmatrix} 0.589 \\ \mathbf{0.884} \end{bmatrix}$$
$$\mathbf{S}_{3} = \begin{pmatrix} \mathbf{1} & 0.355 & 0.355 \\ l & \mathbf{1} & \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{1} & \mathbf{1} \end{pmatrix} \begin{bmatrix} 0.644 \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}$$

$$\mathbf{S}_{4} = \begin{bmatrix} I & \mathbf{1} & \frac{-0.747}{1} \\ I & \underline{5} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{0.843} \end{bmatrix}, \ \mathbf{S}_{5} = \begin{bmatrix} I & \mathbf{1} & \frac{-0.888}{1} \\ I & \underline{100} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{0.947} \end{bmatrix}$$

5 CONCLUSIONS

The new efficient technique of imposing the statistical correlation when using the Monte Carlo type simulation (e.g. LHS) is suggested. The technique is robust, efficient and very fast. The method is implemented in the *FREET* multipurpose software package based on LHS for the reliability analysis of computational problems, Novák et al, (2002, 2003), Pukl et al. (2003). The method has several advantages in comparison with the former techniques:

1. The technique uses only the random changes of the ranks in the sampling matrix. The number of the simulations does not increase CPU time in the practical cases but for the increasing number of the random variables more SA simulations are needed to achieve a good accuracy. The technique is robust, the Simulated Annealing can be terminated if the error (norm) is acceptable (users decision).

2. The problem of imposing the statistical correlation is constrained precisely; therefore the initial temperature for annealing can be estimated.

3. The technique can work also with the nonpositive definitive matrices defined unconsciously by the user as the input data. The important coefficients may be emphasized using weights while others may be suppressed.

ACKNOWLEDGEMENTS

The authors thank for support under the grant of Grant Agency of the Czech Republic No. 103/02/1030 and CEZ J22/98:261100007.

REFERENCES

- Bažant, Z.P., Novák, D. & Vořechovský M., 2003. Statistical size effect prediction in quasibrittle materials, *Proceedings* 2003, ICASP 9, San Francisco, USA. (in print)
- Huntington, D.E. & Lyrintzis, C.S. 1998. Improvements to and limitations of Latin hypercube sampling. *Probabilistic En*gineering Mechanics Vol. 13, No. 4: 245-253.
- Iman R.C. & Conover W.J. 1980. Small Sample Sensitivity Analysis Techniques for Computer Models with an Application to Risk Assessment. *Communications in Statistics: Theory and Methods*, Vol. A9 (No. 17): 1749-1842.
- Iman, R.C. & Conover, W.J. 1982. A Distribution Free Approach to Inducing Rank Correlation Among Input Variables. *Communications in Statistics* Vol. B11: 311-334.
- Lehký, D & Novák, D. 2002. Statistical size effect of concrete in uniaxial tension, 4th International Ph.D. Symposium in Civil Engineering, Proceedings 2002, Munich, Germany, September 19 – 21: 435-444
- Liu, P. and Der Kiureghian, A. 1986. Multivariate distribution models with prescribed marginals and covariances. *Probabilistic Engineering Mechanics* Vol. 1 (No. 2): 105-111.
- McKay, M.D., Conover, W.J. & Beckman, R.J. 1979. A Comparison of Three Methods for Selecting Values of Input Variables in the Analysis of Output from a Computer Code. *Technometrics* Vol. 21: 239-245.
- Novák, D., Vořechovský, M., Pukl, R. & Červenka, V. 2000. Statistical nonlinear analysis – size effect of concrete beams. Fracture Mechanics of Concrete and Concrete Structures, *Proceedings 2000, 4th Int. Conference FraMCoS, Cachan, France, Balkema.*
- Novák, D., Rusina, R., Vořechovský, M. 2002. FREET Feasible Reliability Engineering Efficient Tool, User's guide. Institute of Engineering Mechanics, Faculty of Civil Engineering, Brno University of Technology / Cervenka Consulting, Prague, Czech Republic.
- Novák, D., Teplý, B. & Keršner, Z. 1998. The role of Latin Hypercube Sampling method in reliability engineering. *ICOSSAR-97; Proceedings 1998, Kyoto, Japan 1997*: 403-409.
- Novák, D., Vořechovský, M., Rusina, R. 2003. Small-sample statistical analysis - software FREET. . Statistical size effect prediction in quasibrittle materials, *Proceedings 2003*, *ICASP 9, San Francisco, USA. (in print)*
- Otten, R. H. J. M. & Ginneken, L. P. P. P. 1989. *The Annealing* Algorithm. Kluwer Academic Publishers, USA.
- Pukl R., Strauss A., Novak D., Cervenka V., Bergmeister K. 2003. Probabilistic-based assessment of concrete structures using non-linear fracture mechanics, *Proceedings 2003*, *ICASP 9, San Francisco, USA.(in print)*
- Vořechovský, M., Novák, D. 2002. Correlated Random Variables in Probabilistic Simulation, 4th International Ph.D. Symposium in Civil Engineering, Proceedings 2002, Munich, Germany, September 19 – 21: 410-411