



Functions of random variables

Abdel-Hamid Soubra, Emilio Bastidas-Arteaga

► To cite this version:

Abdel-Hamid Soubra, Emilio Bastidas-Arteaga. Functions of random variables. ALERT Doctoral School 2014 - Stochastic Analysis and Inverse Modelling, Michael A. Hicks; Cristina Jommi, pp.43-52, 2014, 978-2-9542517-5-2. .

HAL Id: in2p3-01082914

<http://hal.in2p3.fr/in2p3-01082914v2>

Submitted on 19 Feb 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Functions of random variables

A.-H. Soubra and E. Bastidas-Arteaga

*University of Nantes – GeM Laboratory – UMR CNRS 6183 – Nantes
– France*

In many engineering problems, the uncertainty associated with one random variable needs to be estimated indirectly from the information on uncertainty in another random variable. In most cases, functional relationships (linear or nonlinear) between the response and basic random variables are known; however, in some cases, the exact relationship may not be known explicitly. Since the response variable is a function of other random variables, it will also be random, whether the exact functional relationship between them is known or not. The subject of this chapter is the quantification of the uncertainty in the response variable when it is related to other random variables with a known or unknown relationship.

1 Introduction

This chapter deals with the study of functions of random variable(s). Engineering problems often involve the determination of a relationship between a dependent variable and one or more basic or independent variables. If any one of the independent variables is random, the dependent variable will likewise be random. The probability distribution (as well as its statistical moments) of the dependent variable will be functionally related to and may be derived from those of the basic random variables. As a simple example, the deflection D of a cantilever beam of length L subjected to a concentrated load P (applied at the end of the cantilever) is functionally related to the load P and the modulus of elasticity E of the beam material [$D = (PL^3)/(3EI)$] in which I is the moment of inertia of the beam cross section. Clearly, we can expect that if P and E are both random variables, with respective PDFs, f_P and f_E , the deflection D will also be a random variable with PDF, f_D , that can be derived from the PDFs of P and E . Moreover, the first two statistical moments (i.e. the mean and variance) of D can also be derived as a function of the respective moments of P and E .

In this chapter, we shall develop and illustrate the relevant concepts and procedures for determining the PDF of the response variable or the statistical moments of this

response variable. Both cases where the functional relationship between the response variable and the independent variables is known or unknown are considered.

2 Exact distributions of functions of random variable(s)

The exact distribution of a function of random variables is considered herein only in case where the response variable is a function of a single random variable. The case where this response variable is a function of several random variables can be found elsewhere (cf. [Hal00], [Ang07] and [Fen08] among others).

2.1 Function of a single random variable

Consider a general case in which the functional relationship between the response variable and the basic random variable is not linear. Assume that the response variable Y is functionally related to X as:

$$Y = g(X) \quad (1)$$

If Y is a monotonically increasing function of X , then

$$P(Y \leq y) = P(X \leq x) \quad (2)$$

Or:

$$F_Y(y) = F_X(x) = F_X[g^{-1}(y)] \quad (3)$$

The value $g^{-1}(y)$ can be evaluated by inverting equation (1). If both sides are differentiated with respect to y , the *PDF* of Y can be obtained as:

$$f_Y(y) = f_X[g^{-1}(y)] \frac{dg^{-1}(y)}{dy} \quad (4)$$

Thus, if the functional relationship g and the *PDF* of X are known, the uncertainty in Y in terms of its *PDF* can be obtained from equation (4).

If Y decreases with X , $dg^{-1}(y)/dy$ is negative (since g^{-1} decreases with Y). Since the *PDF* of a random variable cannot be negative, its absolute value is of interest. Therefore, to account for both cases, the *PDF* of Y is written as

$$f_Y(y) = f_X[g^{-1}(y)] \left| \frac{dg^{-1}(y)}{dy} \right| \quad (5)$$

2.2 Example application for an exact distribution of a function of single random variable

Consider a normal variate X with parameters μ and σ ; i.e. $N(\mu, \sigma)$ with *PDF*

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right] \quad (6)$$

Let $Y = \frac{x-\mu}{\sigma}$. Using equation (5), we determine the *PDF* of Y as follows: First, we observe that the inverse function is $g^{-1}(y) = \sigma y + \mu$ and $\frac{dg^{-1}}{dy} = \sigma$. Then, according to equation (5), the *PDF* of Y is

$$f_Y(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \frac{(\sigma y + \mu - \mu)^2}{\sigma^2} \right] |\sigma| = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2} \quad (7)$$

which is the *PDF* of the standard normal distribution, $N(0,1)$.

3 Moments of functions of random variables

In the previous section, we derived the probability distribution of a function of one random variable. It was shown in literature that the linear function of normal variate remains normal. The product (or quotient) of lognormal variates remains also lognormal (see [Ang07] among others).

In general, the derived probability distributions of the function may be difficult (or even impossible) to derive analytically. Indeed, if the distributions of the X_i 's are not known, or if X_1 is normal, X_2 is lognormal, and so on, it is not possible to determine the exact distribution of the response variable Y ; however, its mean and variance can still be extracted from the information on the means and variances of the X_i 's, giving only limited information on its randomness.

If the functional relationship is linear, then the mean and variance of the response variable can be estimated without any approximation. For nonlinear functional relationships, the mean and variance of the response variable can only be estimated approximately. These are discussed next. Beforehand, remember that the expected value of a function $Z = g(X_1, X_2, \dots, X_n)$ of n random variables, called the mathematical expectation, is given by:

$$E(Z) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} g(x_1, x_2, \dots, x_n) f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n \quad (8)$$

Below, we shall use equation (8) to derive the moments of linear functions of random variables, as well as the first-order approximate moments of nonlinear functions.

3.1 Mean and variance of a linear function

Consider first the moments of the linear function

$$Y = aX + b \quad (9)$$

According to equation (8), the mean value of Y is:

$$\begin{aligned} E(Y) &= E(aX + b) = \int_{-\infty}^{\infty} (ax + b) f_X(x) dx \\ &= a \int_{-\infty}^{\infty} x f_X(x) dx + b \int_{-\infty}^{\infty} f_X(x) dx = aE(X) + b \end{aligned} \quad (10)$$

whereas the variance is:

$$\begin{aligned} VAR(Y) &= E[(Y - \mu_Y)^2] = E[(aX + b - a\mu_X - b)^2] \\ &= a^2 \int_{-\infty}^{\infty} (x - \mu_X)^2 f_X(x) dx = a^2 Var(X) \end{aligned} \quad (11)$$

For $Y = a_1 X_1 + a_2 X_2$,

where a_1 and a_2 are constants

$$E(Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (a_1 x_1 + a_2 x_2) f_{X_1, X_2}(x_1, x_2) dx_1 dx_2 \quad (12)$$

This equation may be written (in case where X_1 and X_2 are statistically independent) as follows:

$$E(Y) = a_1 \int_{-\infty}^{\infty} x_1 f_{X_1}(x_1) dx_1 + a_2 \int_{-\infty}^{\infty} x_2 f_{X_2}(x_2) dx_2 \quad (13)$$

We can recognize that the last two integrals above are, respectively, $E(X_1)$ and $E(X_2)$; hence, we have for the sum of two random variables

$$E(Y) = a_1 E(X_1) + a_2 E(X_2) \quad (14)$$

The variance of Y (for the general case of correlated random variables) is given by:

$$\begin{aligned} Var(Y) &= E[(a_1 X_1 + a_2 X_2) - (a_1 \mu_{X_1} + a_2 \mu_{X_2})]^2 \\ &= E[a_1(X_1 - \mu_{X_1}) + a_2(X_2 - \mu_{X_2})]^2 \\ &= E[a_1^2(X_1 - \mu_{X_1})^2 + 2a_1 a_2 (X_1 - \mu_{X_1})(X_2 - \mu_{X_2}) \\ &\quad + a_2^2(X_2 - \mu_{X_2})^2] \end{aligned} \quad (15)$$

We may recognize that the expected values of the first and third terms within the brackets are variances of X_1 and X_2 , respectively, whereas the middle term is the covariance between X_1 and X_2 . Hence, we have:

$$Var(Y) = a_1^2 Var(X_1) + a_2^2 Var(X_2) + 2a_1 a_2 COV(X_1, X_2) \quad (16)$$

If the variables X_1 and X_2 are statistically independent, $COV(X_1, X_2) = 0$; thus, equation (16) becomes:

$$Var(Y) = a_1^2 Var(X_1) + a_2^2 Var(X_2) \quad (17)$$

The results we obtained above can be extended to a general linear function of n random variables, such as

$$Y = \sum_{i=1}^n a_i X_i \quad (18)$$

in which the a_i 's are constants. For this general case, we obtain the mean and variance of Y as follows:

$$E(Y) = \sum_{i=1}^n a_i E(X_i) = \sum_{i=1}^n a_i \mu_{X_i} \quad (19)$$

$$\begin{aligned}
Var(Y) &= \sum_{i=1}^n a_i^2 Var(X_i) + \sum_{\substack{i,j=1,\dots,n \\ i \neq j}} a_i a_j COV(X_i, X_j) \\
&= \sum_{i=1}^n a_i^2 \sigma_{X_i}^2 + \sum_{\substack{i,j=1,\dots,n \\ i \neq j}} a_i a_j \rho_{ij} \sigma_{X_i} \sigma_{X_j}
\end{aligned} \tag{20}$$

in which ρ_{ij} is the correlation coefficient between X_i and X_j .

3.2 Taylor's series and approximate moments of a general function

3.2.1. Function of a single random variable

For a general function of a single random variable X ,

$$Y = g(X) \tag{21}$$

the exact moments of Y may be obtained using

$$E(Y) = \int_{-\infty}^{\infty} g(x) f_X(x) dx \tag{22}$$

and

$$Var(Y) = \int_{-\infty}^{\infty} [g(x) - \mu_Y]^2 f_X(x) dx \tag{23}$$

Obviously, the determination of the mean and variance of the function Y with the above relations would require information on the *DF* $f_X(x)$. In many applications, however, the *PDF* of X may not be available. In such cases, we seek approximate mean and variance of the function Y as follows:

We may expand the function $g(X)$ in a Taylor series about the mean value of X , that is,

$$g(X) = g(\mu_X) + (X - \mu_X) \frac{dg}{dX} + \frac{1}{2} (X - \mu_X)^2 \frac{d^2 g}{dX^2} + \dots \tag{24}$$

where the derivatives are evaluated at μ_X .

Now, if we truncate the above series at the linear terms, i.e.,

$$g(X) \cong g(\mu_X) + (X - \mu_X) \frac{dg}{dX} \quad (25)$$

We obtain the first-order approximate mean and variance of Y as

$$E(Y) \cong g(\mu_X) \quad (26)$$

and

$$\text{Var}(Y) \cong \text{Var}(X) \left(\frac{dg}{dX} \right)^2 \quad (27)$$

We should observe that if the function $g(X)$ is approximately linear (i.e. not highly nonlinear) for the entire range of X , equations (26) and (27) should yield good approximations of the exact mean and variance of $g(X)$. Moreover, when $\text{Var}(X)$ is small relative to $g(\mu_X)$, the above approximations should be adequate even when the function $g(X)$ is nonlinear.

3.2.2. Function of multiple random variables

If Y is a function of several random variables,

$$Y = g(X_1, X_2, \dots, X_n) \quad (28)$$

We obtain the approximate mean and variance of Y as follows: Expand the function $g(X_1, X_2, \dots, X_n)$ in a Taylor series about the mean values $(\mu_{X_1}, \mu_{X_2}, \dots, \mu_{X_n})$, yielding

$$\begin{aligned} Y &= g(\mu_{X_1}, \mu_{X_2}, \dots, \mu_{X_n}) + \sum_{i=1}^n (X_i - \mu_{X_i}) \frac{\partial g}{\partial X_i} \\ &\quad + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (X_i - \mu_{X_i})(X_j - \mu_{X_j}) \frac{\partial^2 g}{\partial X_i \partial X_j} + \dots \end{aligned} \quad (29)$$

Where the derivatives are all evaluated at $\mu_{X_1}, \mu_{X_2}, \dots, \mu_{X_n}$. If we truncate the above series at the linear terms, i.e.,

$$Y = g(\mu_{X_1}, \mu_{X_2}, \dots, \mu_{X_n}) + \sum_{i=1}^n (X_i - \mu_{X_i}) \frac{\partial g}{\partial X_i} \quad (30)$$

We obtain the first-order mean and variance of Y , respectively as follows:

$$E(Y) \cong g(\mu_{X_1}, \mu_{X_2}, \dots, \mu_{X_n}) \quad (31)$$

and

$$Var(Y) \cong \sum_{i=1}^n (\sigma_{X_i})^2 \left(\frac{\partial g}{\partial X_i} \right)^2 + \sum_{\substack{i,j=1,\dots,n \\ i \neq j}} \rho_{ij} \sigma_{X_i} \sigma_{X_j} \frac{\partial g}{\partial X_i} \frac{\partial g}{\partial X_j} \quad (32)$$

We observe that if X_i and X_j are uncorrelated (or statically independent) for all i and j , i.e. $\rho_{ij} = 0$, then equation (32) becomes

$$Var(Y) \cong \sum_{i=1}^n (\sigma_{X_i})^2 \left(\frac{\partial g}{\partial X_i} \right)^2 \quad (33)$$

Equation (33) is a function of both the variances of the independent variables and of the sensitivity coefficients as represented by the partial derivatives.

3.2.3. Example application for the computation of the mean and variance of a general function of several variables

Assume that the random variable Y can be represented by the following relationship:

$$Y = X_1 X_2^2 X_3^{1/3} \quad (34)$$

Where X_1 , X_2 and X_3 are statistically independent random variables with means of 1.0, 1.5, and 0.8, respectively, and corresponding standard deviations of 0.10, 0.20, and 0.15, respectively. Using equations (31) and (33), we find the first-order mean and variance, respectively, to be:

$$E(Y) \approx 1.0 \times 1.5^2 \times (0.8)^{1/3} = 2.0887$$

and

$$\begin{aligned} Var(Y) &\approx Var(X_1)(\mu_{X_2}^{-2} \times \mu_{X_3}^{-1/3})^2 + Var(X_2)[\mu_{X_1} \times (2\mu_{X_2}) \times \mu_{X_3}^{-1/3}]^2 \\ &\quad + Var(X_3) \left[\mu_{X_1} \times \mu_{X_2}^{-2} \times \left(\frac{1}{3} \mu_{X_3}^{-2/3} \right) \right]^2 \\ &= (0.10)^2 (1.5^2 \times 0.8^{1/3})^2 + (0.20)^2 (1.0 \times 2 \times 1.5 \times 0.8^{1/3})^2 \\ &\quad + (0.15)^2 [1.0 \times 1.5^2 \times (1/3) \times 0.8^{-2/3}]^2 \\ &= (0.10)^2 (2.09)^2 + (0.20)^2 (2.78)^2 + (0.15)^2 (0.87)^2 \\ &= 0.04363 + 0.31024 + 0.01704 = 0.37091 \end{aligned}$$

and

$$\sigma_Y = 0.609$$

3.3 Mean and variance of an analytically-unknown functional relationship

In many cases, the exact form of g in equation (28) may not be known. In fact, the exact functional relationship is known in algorithmic form but not in any exact functional form. The implication is that the partial derivatives of the function with respect to the random variables cannot be calculated to approximate the first-order mean and the first-order variance of the response variable, as discussed before.

In this case, the approximate (first-order) mean value of the response, represented by Y in equation (28), can be obtained by using the mean values of all the parameters in the problem, the same as in equation (31). Evaluating the variance of Y will be more involved since the functional form of g is unknown, and its partial derivatives with respect to the i^{th} random variable in equation (28) cannot be evaluated. The task is to calculate the variance of Y without information on the analytical partial derivatives. The Taylor series finite difference estimation procedure can be used to numerically evaluate the variance of Y , as discussed below:

To evaluate the variance, one needs to compute, for each random variable, the two following (intermediate) response variables:

$$Y_i^+ = g[\mu_{X_1}, \mu_{X_2}, \dots, (\mu_{X_i} + \sigma_{X_i}), \dots, \mu_{X_n}] \quad (35)$$

and

$$Y_i^- = g[\mu_{X_1}, \mu_{X_2}, \dots, (\mu_{X_i} - \sigma_{X_i}), \dots, \mu_{X_n}] \quad (36)$$

In simple terms, equation (35) states that the response variable Y_i^+ is calculated considering the mean of all the random variables except the i^{th} one, which is considered to be the mean plus one standard deviation value. Equation (36) indicates the same thing, except that for the i^{th} random variable, the mean minus one standard deviation value needs to be considered. Then, using the central difference approximation, we can show that

$$E_i = \frac{\partial g}{\partial X_i} = \frac{Y_i^+ - Y_i^-}{2\sigma_{X_i}} \quad (37)$$

Considering all the random variables, the first-order variance of Y is computed as

$$\text{Var}(Y) \approx \sum_{i=1}^n \left(\frac{Y_i^+ - Y_i^-}{2\sigma_{X_i}} \right)^2 \times \text{Var}(X_i) \approx \sum_{i=1}^n \left(\frac{Y_i^+ - Y_i^-}{2} \right)^2 \quad (38)$$

Thus, when the functional relationship among the random variables is not known explicitly, the mean and variance of the response variable can be approximated by

using equations (31) and (38). This requires the computation of the response variable several times. If there are n random variables present in a problem, the required total number of computations of the response variable is $(1 + 2n)$.

4. Conclusion

The probabilistic characteristics of a function of random variable(s) may be derived from those of the independent variates. These include, in particular, the probability distribution and the first two statistical moments (mean and variance) of the function. It was shown that for a function of a single random variable, the *PDF* of the function can be readily obtained analytically. However, it was shown in the literature that the derivation of the distribution of a function of multiple variables can be complicated mathematically, especially for nonlinear functions (see [Ang07] among others). Therefore, even though the required distribution of a function may theoretically be derived, it is often impractical to apply, except for special cases, such as a linear function of independent Gaussian variates or the strictly product/quotient of independent lognormal variates. In this light, it is often necessary, in many applications, to describe the probabilistic characteristics of a function approximately in terms only of its mean and variance. The mean and variance of linear functions can be estimated without any approximation; however, for a general nonlinear function, we must often resort to first-order (or second-order) approximations. In case of analytically-unknown functions, one may use the finite difference method for the (approximate) computation of the statistical moments. Finally, it should be mentioned that when the probability distribution of a general function is required, we may need to resort to Monte Carlo simulations or other numerical methods.

References

- [Ang07] A. Ang and W. Tang. Probability concepts in engineering, Emphasis on applications to civil and environmental engineering, John Wiley & Sons, 406 pages, 2007.
- [Fen08] G.A. Fenton and D.V. Griffiths. Risk assessment in geotechnical engineering, John Wiley & Sons, 461 pages, 2008.
- [Hal00] A. Haldar and S. Mahadevan. Probability, Reliability and Statistical Methods in Engineering desing, Jonh Wiley & Sons, Inc. 2000.