

DRAFT

VVS2020-13417

FULL FUNCTION SAMPLING OF UNCERTAIN CORRELATIONS

Kevin W. Irick, Jeff Engerer, Blake Lance, Scott A. Roberts, and Ben Schroeder

Sandia National Laboratories¹

P.O. Box 5800

Albuquerque, NM 87185-0828

ABSTRACT

Empirically-based correlations are commonly used in modeling and simulation but rarely have rigorous uncertainty quantification that captures the nature of the underlying data. In many applications, a mathematical description for a parameter response to some input stimulus is often either unknown, unable to be measured, or both. Likewise, the data used to observe a parameter response is often noisy, and correlations are derived to approximate the bulk response. Practitioners frequently treat the chosen correlation—sometimes referred to as the “surrogate” or “reduced-order” model of the response—as a constant mathematical description of the relationship between input and output. This assumption, as with any model, is incorrect to some degree, and the uncertainty in the correlation can potentially have significant impacts on system responses. Thus, proper treatment of correlation uncertainty is necessary. In this paper, a method is proposed for high-level abstract sampling of uncertain data correlations. Whereas uncertainty characterization is often assigned to scalar values for direct sampling, functional uncertainty is not always straightforward. A systematic approach for sampling univariable uncertain correlations was developed to perform more rigorous uncertainty analyses and more reliably sample the correlation space. This procedure implements pseudo-random sampling of a correlation with a bounded input range to maintain the correlation form, to respect variable uncertainty across the range, and to ensure function continuity with respect to the input variable.

NOMENCLATURE

a = coefficient
 \mathbf{a} = vector of coefficients
 α = confidence level
 dF = difference in F-statistic
 δx = independent variable perturbation
 Δx = independent variable spacing
 f = function

\mathbf{f} = vector of functions
 F = F-statistic
 i = index
 m = deviation rate
 N = parameter quantity
 PI = prediction interval
 r = random variable
 σ = standard deviation
 t = t-statistic
 x = independent variable
 \mathbf{x} = vector of independent variables
 y = dependent variable
 \mathbf{y} = vector of dependent variables

Subscripts

a = coefficients
 avg = average
 cmp = computed from samples
 $corr$ = correlation
 crt = critical
 $ctrl$ = control
 fit = fit
 i = parameter index
 max = maximum
 min = minimum
 nom = nominal
 o = objective
 PI = prediction interval
 $smpl$ = sample
 tar = target
 trn = training
 U = uncertainty

1. INTRODUCTION

In many types of critical systems—such as high-consequence, high-value, and/or high-production-volume

¹ This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

systems—various forms of modeling are often employed to provide stakeholders with an assurance that the system of interest will perform in a certain way. For example, government agencies may want to be assured that weapon systems will engage properly and not cause unintentional harm [1-3], energy corporations may want to realize increased safety and longevity in power generation systems [4-6], or medical companies may want to reduce equipment failure rates [7-9]. Models can support system design goals by taking forms ranging from simple analytical models to complex physics finite element models. Empirically-based correlations—sometimes referred to as “surrogate” or “reduced-order” models of a response—are commonly used in modeling and simulation to describe the behavior of systems in an analytical fashion. For example, in thermal engineering, correlations are typically used to describe material thermophysical properties such as thermal conductivity, mass density, and specific heat as functions of temperature. Likewise, a myriad of correlations exists to describe convective heat transfer based on thermodynamic conditions [10]. The underlying correlations can be coupled into simple or complex models.

Practitioners frequently treat chosen correlations as “true” or constant mathematical descriptions of a relationship between input and output. Rarely is rigorous uncertainty quantification (UQ) applied that captures the nature of the underlying data. In this sense, the uncertainty of the correlation is ignored. Oftentimes, if UQ is pursued at all, the analyst may apply an uncertainty scaling factor to the result of the nominal correlation, or the analyst may try to characterize uncertainty bounds in correlation coefficients [11-13].

Two examples of traditional approaches to UQ sampling of correlations are given here to illustrate common weaknesses of typical approaches. The first traditional method involves simply multiplying the chosen correlation by a scaling factor. For example, a practitioner might obtain a nominal correlation, y_{nom} , from a set of training data, y_{tm} . Then, it may be given that uncertainty on y_{nom} is simply $\pm 20\%$ of y_{nom} with a 95% confidence level. A random functional sample, y_{smp} , could be generated by

$$y_{smp} = y_{nom}r, \quad (1)$$

where r is a random number sampled from a normal distribution with a mean of 1.0 and standard deviation of 0.1. Figure 1 shows an example of what such an approach may yield, where $y_{smp,avg}$ is the average of the random samples and $y_{smp,U}$ is the 95% uncertainty bounds based on the aggregate sample variation. In the figure, y_{nom} and $y_{smp,avg}$ overlap across the entire data range, making them difficult to distinguish. This multiplier method employs a single random variable, r , but does not systematically and fully represent random uncertainty in the data, relying on practitioner judgement and the correlation value. When this approach is used, each y_{smp} realization is biased to being either greater than or less than y_{nom} across the entire range, not allowing for a single realization to span the \pm uncertainty space. This method also has issues with dependency on the

correlation value, where smaller y_{nom} values yield smaller absolute uncertainty bands, not generally respecting absolute local uncertainty.

The second traditional method involves sampling from coefficient distributions. If confidence intervals can be placed on the fit coefficients, then coefficient values can be sampled at random to create a sample realization of the uncertain correlation. This coefficient sampling method has sound statistical basis and better captures the underlying uncertainty in the training data, as seen in Figure 2. Again, y_{nom} and $y_{smp,avg}$ overlap and are not visually distinguishable. The set of random numbers required to sample the correlation is equal to the number of uncertain coefficients in the correlation, where the number of coefficients might be large for complicated physics-based or surrogate models. If coefficient parameters are correlated in unknown fashion, sampling the coefficients independently may yield spurious results. Also, the uncertainty generated by the aggregate samples only samples the uncertain space of coefficients themselves to describe the mean response of the correlation. This method does not necessarily accurately represent the prediction interval, the uncertainty bounds where future data may lie. A brief discussion on this is given later.



Figure 1. Uncertain correlation multiplier sampling example

Both methods mentioned above are viable and convenient options for implementing uncertainty into a model; however, the methods may not accurately represent the uncertain nature of the correlation itself with respect to the data the correlation represents. The work presented in this paper proposes a method for sampling uncertain correlations in a more abstract manner than is typically performed. The method implements pseudo-random sampling of a correlation with a bounded input range to maintain the nominal correlation form, to respect variable uncertainty across the range, and to ensure function continuity with respect to the input variable.

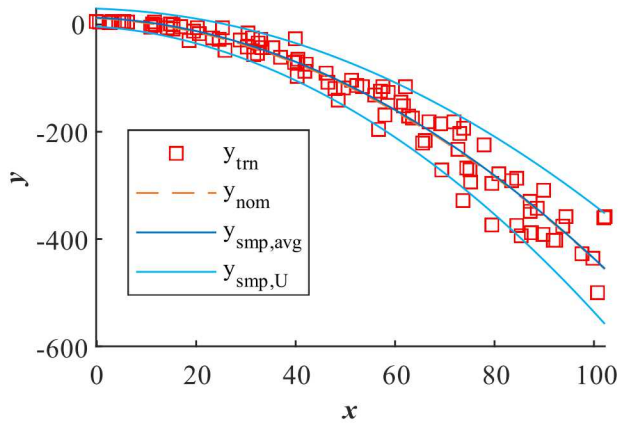


Figure 2. Uncertain coefficient sampling example

2. PROCEDURE

This section describes in detail the technical process for the proposed full function sampling method. The steps are listed here, and descriptions follow in subsequent subsections:

1. Gather training data
2. Define correlation form
3. Obtain nominal correlation
4. Determine prediction interval on nominal correlation
5. Select random numbers
6. Determine input values
7. Determine control response values
8. Determine target response
9. Find sample coefficients
10. Repeat steps 5 through 9 for multiple samples
11. Perform F-test

2.1 Step 1: Gather Training Data

The data set from which a correlation is built can be referred to as the training data. For this univariate application, the vectors of N_{trn} training input and output—or independent and dependent variable—values are denoted by \mathbf{x}_{trn} and \mathbf{y}_{trn} , respectively.

2.2 Step 2: Define Correlation Form

A form for the functional correlation, $y=f_{corr}(\mathbf{a},x)$, must be defined to map input x values to output y values using a set of N_a coefficients, \mathbf{a} . The f_{corr} form is carried throughout this sampling technique.

2.3 Step 3: Obtain Nominal Correlation

Using the f_{corr} form, a nominal set of correlation coefficients, \mathbf{a}_{nom} , must be determined to define a nominal correlation, $f_{nom}(x)=f_{corr}(\mathbf{a}_{nom},x)$, based on the training data. This can be done through any preferred fitting function method, $f_{fit}(\mathbf{a}_{nom},\mathbf{x}_{trn},\mathbf{y}_{trn})$,

but in this work, a simple least squares approach is used. The least squares approach simply seeks to minimize the object function, f_o , defined as

$$f_o = \sum_{i=1}^{N_{trn}} (f_{nom}(x_{trn,i}) - y_{trn,i})^2, \quad (2)$$

where $x_{trn,i}$ and $y_{trn,i}$ are the i^{th} values from \mathbf{x}_{trn} and \mathbf{y}_{trn} , respectively. The fitting problem is thus formulated as

$$f_{fit}(\mathbf{a}_{nom}, \mathbf{x}_{trn}, \mathbf{y}_{trn}) = \min_{\mathbf{a}_{nom}} f_o = \min_{\mathbf{a}_{nom}} \sum_{i=1}^{N_{trn}} (f_{corr}(\mathbf{a}_{nom}, x_{trn,i}) - y_{trn,i})^2. \quad (3)$$

2.4 Step 4: Determine Prediction Interval

A prediction interval for f_{nom} describes the possible range of y -values within some statistical confidence level, α , about f_{nom} , assuming a normal distribution [14]. At a given x -value, the prediction interval deviation rate, $m_{PI}(x)$, can be computed by

$$m_{PI}(x)^2 = \left(\sum_{i=1}^{N_{trn}} (f_{nom}(x_{trn,i}) - y_{trn,i})^2 \right) * \dots \left(1 + \frac{1}{N_{trn}} + \frac{(x-x_{trn,avg})^2}{\sum_{i=1}^{N_{trn}} (x_{trn,i}-x_{trn,avg})^2} \right), \quad (4)$$

where $x_{trn,avg}$ is the average of \mathbf{x}_{trn} . The total prediction interval amplitude for a given α , $PI(x, \alpha)$, can be found as

$$PI(x, \alpha) = t(0.5(1 + \alpha), N_{trn} - 2) m_{PI}(x), \quad (5)$$

with $t(0.5(1+\alpha), N_{trn}-2)$ being the Student's t-statistic for probability $0.5(1-\alpha)$ and degrees of freedom $N_{trn}-2$. Thus, for a given confidence level, a predicted value based on the correlation would be in the range $[f_{nom}(x) - PI(x, \alpha), f_{nom}(x) + PI(x, \alpha)]$.

2.5 Step 5: Select Random Numbers

The full function sampling method requires sampling two random values to determine a single correlation realization, $f_{smp,i}$. A convenient approach is to sample two random variables, r_1 and r_2 , each uniformly distributed between 0 and 1. The first random variable, r_1 , is used to define control and target x -values, and r_2 is used to determine a target y -value. Control and target points are used to find the correlation realization.

2.6 Step 6: Determine Input Values

In order to sample the full function range, N_a control points and 1 target point are chosen from the training data range, $[x_{min}, x_{max}]$, where x_{min} and x_{max} are the minimum and maximum values from \mathbf{x}_{trn} . The target point input value, x_{tar} , and the set of control point input values, \mathbf{x}_{ctrl} , are evenly distributed in the range in a cyclic fashion. The cyclic range implies that any x -value greater than x_{max} by some δx amount is interpreted as

$$x_{max} + \delta x \rightarrow x_{min} + \delta x. \quad (6)$$

By the same token, any x -value lower than x_{min} by some δx amount is interpreted as

$$x_{min} - \delta x \rightarrow x_{max} - \delta x. \quad (7)$$

It follows that the target and control input values are evenly spaced by Δx , where

$$\Delta x = \frac{x_{max} - x_{min}}{N_a + 1}. \quad (8)$$

The x_{tar} value is chosen as

$$x_{tar} = x_{min} + r_1(x_{max} - x_{min}), \quad (9)$$

and the remaining i^{th} $x_{ctrl,i}$ values are chosen by

$$x_{ctrl,i} = x_{tar} + i\Delta x \quad (10)$$

for $i=1$ to N_a , following the cyclic range rules stipulated above. Although spacing between the x -values is fixed, locations of the sampled x -values are randomized by r_1 .

2.7 Step 7: Determine Control Response Values

The control point values are used to guide the shape of $f_{smp,i}$ to be in the neighborhood of f_{nom} . Thus, control point response values, y_{ctrl} , are simply made to fall directly on f_{nom} such that

$$y_{ctrl} = f_{nom}(x_{ctrl}). \quad (11)$$

2.8 Step 8: Determine Target Response

The target response point, (x_{tar}, y_{tar}) , is used as a constraint for determining the correlation realization coefficients, a_{smp} . The y_{tar} value is randomly selected from the prediction interval of f_{nom} at x_{tar} , where,

$$y_{tar} = f_{nom}(x_{tar}) + t(r_2, N_{trn} - 2)m_{PI}(x_{tar}). \quad (12)$$

The equality constraint forces $f_{smp,i}$ to pass through (x_{tar}, y_{tar}) and is given to be

$$f_{smp,i}(x_{tar}) - y_{tar} = 0. \quad (13)$$

2.9 Step 9: Find Sample Coefficients

Determining a_{smp} for a realization of the uncertain correlation is accomplished by solving the optimization problem

$$f_{fit}(a_{smp}, x_{trn}, y_{trn}) \quad (14)$$

subject to the constraint given in Eq. (11). By performing this optimization problem, the realization of the uncertain correlation

is sampled from the prediction interval distribution respecting the underlying training data while still maintaining the form and shape expected, as suggested by f_{nom} .

2.10 Step 10: Iterate Sampling

Full function sampling—Step 5 through Step 9—can be performed N_{smp} times to generate some set of realizations, f_{smp} .

2.11 Step 11: Perform F-Test

An F-test can be used to compare the statistical equivalence of two sample populations [15,16]. In this case, the test compares the variance of f_{smp} with the predicted variance of f_{nom} , assuming the two represent sample sets are drawn from the same population. First, assume the standard deviation of f_{nom} is σ_{nom} and to vary with x such that

$$\sigma_{nom}(x) = m_{PI}(x). \quad (15)$$

This assumption states that the standard error computed for f_{nom} is used to represent the standard deviation of f_{nom} because the standard deviation is not known.

Then, the standard deviation of f_{smp} , σ_{smp} , is easily computed at any x -value since all a_{smp} values are known. For the following, let index 1 represent either nom or smp , where 1 represents the set that has the higher standard deviation between σ_{nom} and σ_{smp} . Index 2 represents the set with the lower standard deviation. The F-statistic for the samples, F_{cmp} , is computed as

$$F_{cmp} = \sigma_1^2 / \sigma_2^2, \quad (16)$$

and the degrees of freedom for the sets, ν_1 and ν_2 , are either $N_{trn} - 1$ or $N_{smp} - 1$, according to the same index references used for the standard deviations.

The critical F-statistic, F_{crt} , should be determined from α , ν_1 , and ν_2 . If F_{crt} is greater than F_{cmp} , then the null hypothesis that the f_{smp} and the f_{nom} variations are the same cannot be rejected. This F-test can be computed across the independent variable range of the correlation to show a relative evaluation of the statistical comparison between the two function statistics, $y_{nom,PI}$ and $y_{smp,U}$. However, the test does not need to be used to accept the sample set as viable.

3. RESULTS

A handful of correlation forms were investigated as specimens for the full function sampling approach, including linear, quadratic, cubic, exponential, and logarithmic function forms. Noisy data were generated against which to compute fits and statistics. The full function sampling method was implemented in MATLAB®, using the *fmincon* function with the Sequential Quadratic Programming algorithm [17]. For each of the example problems, 10, 100, and 1,000 sample correlations were produced for a given f_{nom} . The results showed no drastic

difference between 100 and 1,000 sample sizes, so only the results from using 100 samples are shown here.

The linear f_{corr} takes the form

$$f_{corr} = a_0 + a_1x. \quad (17)$$

Figure 3a shows the training data with all sample correlations overlaid for the linear problem. Figure 3b shows the comparison of where the linear f_{nom} and its corresponding 95% prediction interval falls with respect to the 95% confidence level of the f_{smp} set. Lastly, Figure 3c shows the relative F-test result, where

$$dF_{rel} = \frac{F_{crt} - F_{cmp}}{F_{crt}}. \quad (18)$$

Positive values of dF_{rel} suggest that the null hypothesis that the variances of f_{nom} and f_{smp} are the same cannot be rejected.

The quadratic f_{corr} takes the form

$$f_{corr} = a_0 + a_1x + a_2x^2. \quad (19)$$

Figure 4a shows the training data with all sample correlations overlaid for the quadratic problem. Figure 4b shows the comparison of where the linear f_{nom} and its corresponding 95% prediction interval fall with respect to the 95% confidence level of the f_{smp} set. Lastly, Figure 4c shows the relative F-test result.

The cubic f_{corr} takes the form

$$f_{corr} = a_0 + a_1x + a_2x^2 + a_3x^3. \quad (20)$$

Figure 5a shows the training data with all sample correlations overlaid for the cubic problem. Figure 5b shows the comparison of where the linear f_{nom} and its corresponding 95% prediction interval fall with respect to the 95% confidence level of the f_{smp} set. Lastly, Figure 5c shows the relative F-test result.

The exponential f_{corr} takes the form

$$f_{corr} = a_0 + a_1e^{x/a_2}. \quad (21)$$

Figure 6a shows the training data with all sample correlations overlaid for the exponential problem. Figure 6b shows the comparison of where the linear f_{nom} and its corresponding 95% prediction interval fall with respect to the 95% confidence level of the f_{smp} set. Lastly, Figure 6c shows the relative F-test result.

The logarithmic f_{corr} takes the form

$$f_{corr} = a_0 + a_1\ln(a_2x + a_3). \quad (22)$$

Figure 7a shows the training data with all sample correlations overlaid for the logarithmic problem. Figure 7b shows the comparison of where the linear f_{nom} and its corresponding 95% prediction interval fall with respect to the 95% confidence level of the f_{smp} set. Lastly, Figure 7c shows the relative F-test result.

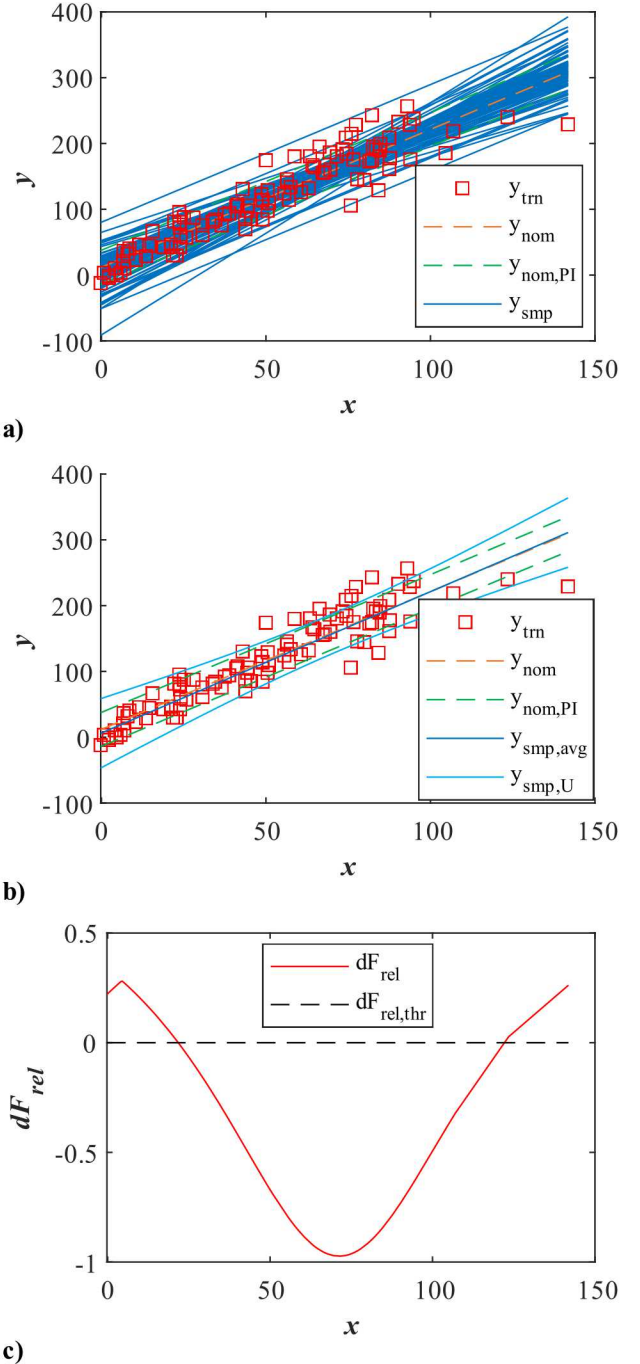


Figure 3. Linear example training data and sample correlations

Note that for the logarithmic problem, an additional inequality constraint was added to avoid a negative argument for the natural logarithm function, where

$$a_2(x_{min} - 1) + a_3 \geq 0. \quad (23)$$

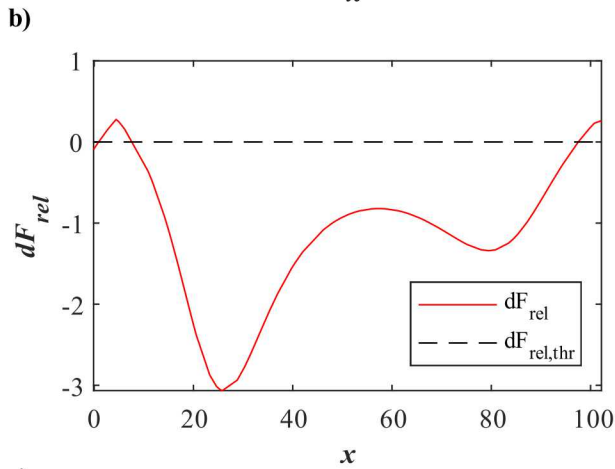
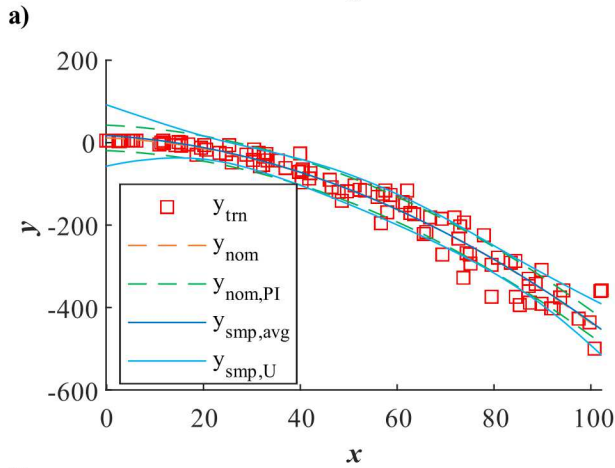
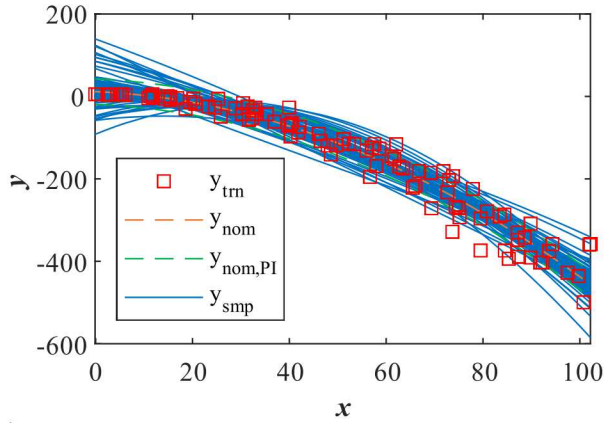


Figure 4. Quadratic example training data and sample correlations

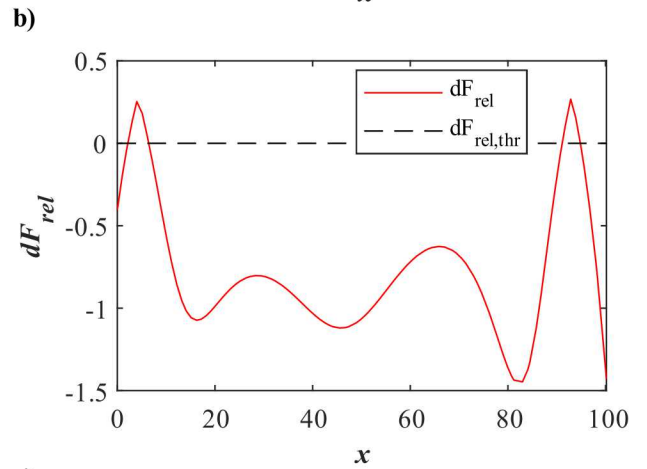
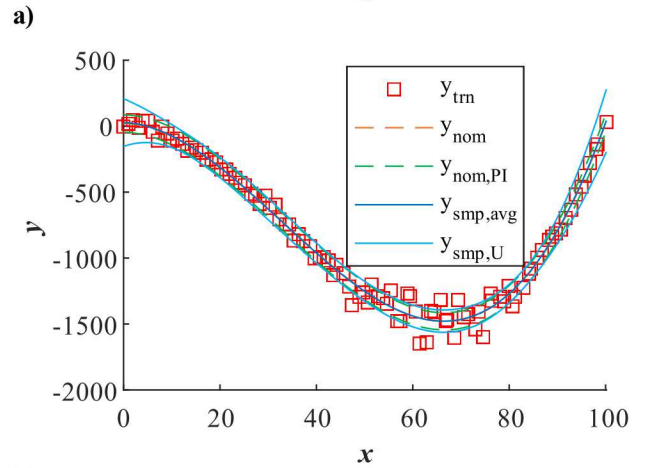
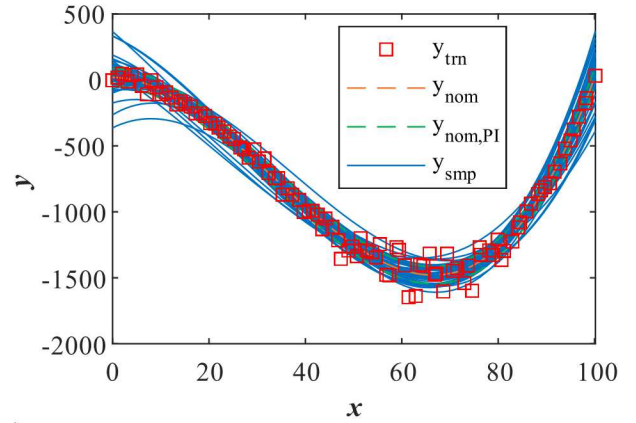
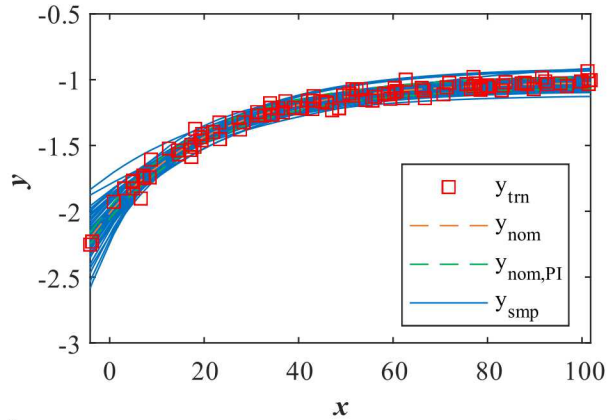
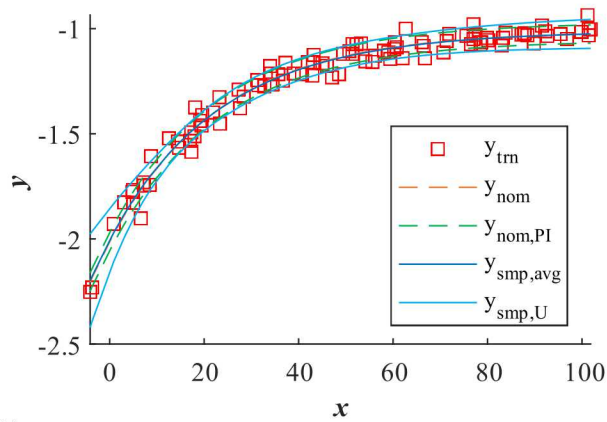


Figure 5. Cubic example training data and sample correlations

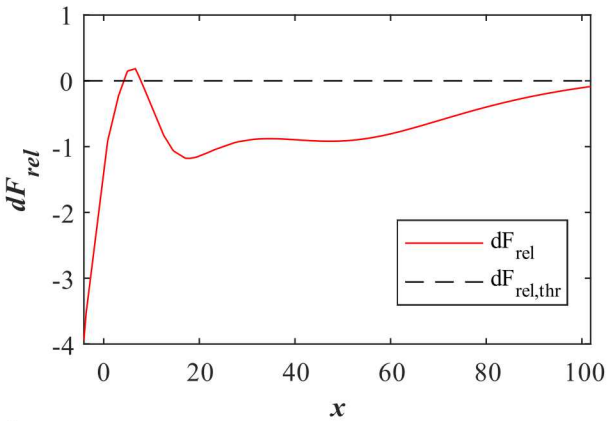
In each of the five scenarios presented above, only two random variables were required to generate a sampled realization from the respective uncertain correlation, regardless of the number of correlation coefficients.



a)

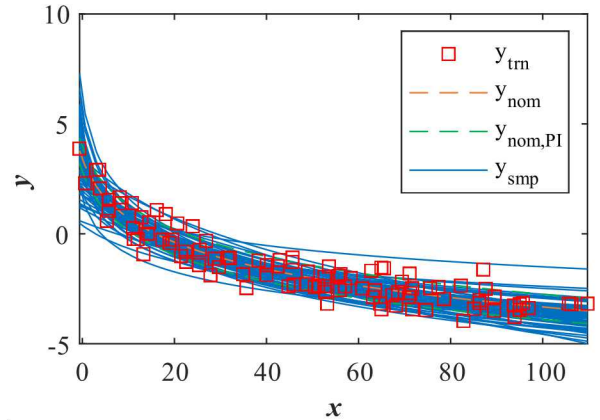


b)

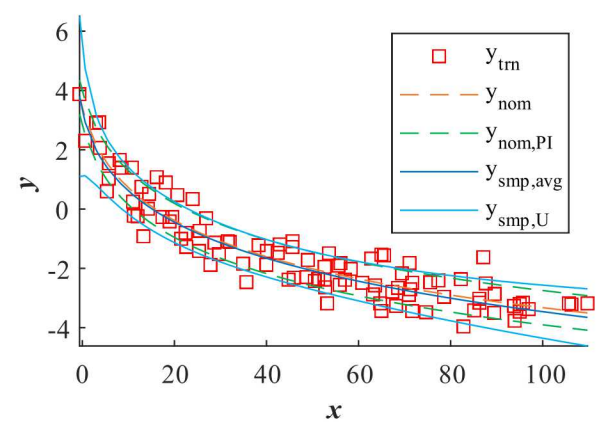


c)

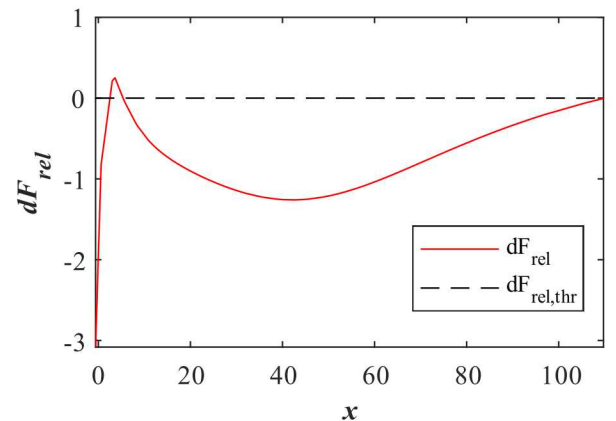
Figure 6. Exponential example training data and sample correlations



a)



b)



c)

Figure 7. Logarithmic example training data and sample correlations

As compared to the traditional sampling methods described previously, this full function sampling approach does not inherently bias a sample function realization as occurs with the multiplier method. Unlike the coefficient sampling approach, sampling the full function more abstractly—as is done here—

produces correlations that respect the uncertain function value space regardless of the uncertainty on the individual coefficients.

The proposed full function sampling method as presented here exhibits some behaviors that should be taken into consideration. It is clear from the presented results—through both qualitative analysis of the aggregate sample standard deviations and quantitative analysis using the F-test—that relatively wide uncertainties are generated near the upper and lower bounds of the range. This is due in most part to the unconstrained nature of the curve fitting procedure. If this is an issue, it is possible to add additional constraints to the optimization problem to restrict function values near x_{min} and x_{max} . Adding such tail constraints may change the overall statistical performance of the sampling method in other regions of the domain, but subsequent effects should be investigated on a case-by-case basis.

Additional constraints may be added to the optimization problem by the user as needed. For example, if a cubic polynomial correlation function was being used, it may be necessary to avoid inflection points at certain locations, or if a linear correlation function was being used, a positive or negative slope might not be permitted by physics. The additional constraints added to the logarithmic function in Eq. (23) in this work is an example of such a constraint.

Lastly, it is important to note that this method may produce outlier realizations that are not acceptable for the application. Removal of these outliers should be considered to avoid unnecessary or unrealistic samples. Analyzing the function values at x_{min} and x_{max} may provide insight as to whether the sample correlation is an outlier—if the function value is too high or too low as compared to other samples. Likewise, one may filter samples based on the value of f_o as well as the fundamental $f_{smp,i}(x_{tar})$ constraint.

4. CONCLUSIONS AND FUTURE WORK

The results of the full function sampling approach present a more abstract method for sampling an uncertain correlation function than more traditional sampling methods. The presented method requires only two random variables for any function form and samples the uncertainty space of a correlation relatively well without any significant function bias. In future work, a comparison of F-test results should be made between traditional sampling methods and the full function sampling method presented in this work.

REFERENCES

[1] National Nuclear Security Administration, 2015, “Weapon Quality Policy,” *NAP-24A*.
[2] Hetreed, C. F., Carroll, M. D., Collard, J. E., and Snyder, R. C., 2018, “F-35 Weapons Separation Test and Verification,” *Proceedings of the 2018 Aviation Technology, Integration, and Operations Conference*, AIAA.

[3] Oberkampf, W. L. and Roy, C. J., 2010, *Verification and Validation in Scientific Computing*, Cambridge University Press, Cambridge, UK.
[4] He, Y.-N., Xiong W., Gu, P.-F., and Tang, J.-Z., 2019, “Research on the Verification and Validation Method of Safety Analysis Software in Nuclear Power Plants,” in *Nuclear Power Plants: Innovative Technologies for Instrumentation and Control Systems*, Lecture Notes in Electrical Engineering, **507**, Springer, Singapore.
[5] Irick, K. and Fathi, N., 2019, “Computational Evaluation of Thermal Barrier Coatings: Two-Phase Thermal Transport Analysis,” *Proceedings of the ASME 2019 Verification and Validation Symposium*.
[6] Stoots, C., Larson, T., Schultz, R., Gougar, H., McCarthy, K., Petti, D., Swiler, L., and Corradini, M., 2012, “Verification and Validation Strategy for LWRs Tools,” Idaho National Laboratory, *INL/EXT-12-27066*.
[7] Hargett, Z., Gutierrez, M., and Harman, M., 2019, “Verification of Manual Digitization Methods during Experimental Simulation of Knee Motion,” *Proceedings of the ASME 2019 Verification and Validation Symposium*.
[8] Keefe, D. F., Sotiropoulos, F., Interrante, V., Runesha, H. B., Coffey, D., Staker, M., Lin, C.-L., Sun, Y., Borazjani, I., Le, T., Rowe, N., and Erdman, A., 2010, “A Process for Design, Verification, Validation, and Manufacture of Medical Devices Using Immersive VR Environments,” *Journal of Medical Devices*, **4**(4): 045002.
[9] Jiang, Z., Pajic, M., Connolly, A., Dixit, S., and Mangharam, R., 2010, “Real-Time Heart Model for Implantable Cardiac Device Validation and Verification,” *2010 22nd Euromicro Conference on Real-Time Systems*, pp. 239-248.
[10] Incropera, F. P., DeWitt, D. P., Bergman, T. L., and Lavine, A. S., 2007, *Fundamentals of Heat and Mass Transfer*, 6th Ed., John Wiley & Sons, Inc., Hoboken, NJ.
[11] Rabe-Hesketh, S. and Skrondal, A., 2005, *Multilevel and Longitudinal Modeling Using Stata—Volume I: Continuous Responses*, Stata Press, College Station, TX.
[12] Bates, D. M. and Watts, D. G., 1988, *Nonlinear Regression Analysis and Its Applications*, John Wiley & Sons.
[13] Farrance, I. and Frenkel, R., 2012, “Uncertainty of Measurement: A Review of the Rules for Calculating Uncertainty Components through Functional Relationships,” *The Clinical Biochemist, Reviews*, **33**(2), pp. 49–75.
[11] Atchison, J. and Dunsmore, I., 1975, *Statistical Prediction Analysis*, Cambridge University Press, Cambridge, UK.
[12] Snedcor, G. W. and Cochran, W. G., 1989, *Statistical Methods*, 8th Ed., Iowa State University Press.
[13] “F-Test for Equality of Two Variances,” *NIST/SEMATECH e-Handbook of Statistical Methods*, <https://www.itl.nist.gov/div898/handbook/eda/section3/eda359.htm>, accessed 2 Jan 2020.
[14] MATLAB 2019b and Optimization Toolbox 2019b, The MathWorks, Inc., Natick, MA.