# Using Simulation to Test the Reliability of Regression Models

Fred J. Rispoli[1] & Vishal Shah[2]

[1] Department of Mathematics and Computer Sciences, Dowling College, Oakdale, The United States

[2] Department of Biology, Dowling College, Oakdale, The United States

Correspondence: Professor Fred J. Rispoli, Department of Mathematics, Dowling College, Oakdale, NY 11769, The United States. Tel: 1-631-244-3179. E-mail: rispolif@dowling.edu

## Abstract

In many sciences, it is standard laboratory practice to use a statistical design of experiment and a regression model to study the influence of multiple parameters under a wide range of conditions. The current study aims at investigating the reliability of regression models by examining recently published models. Of particular interest are the assumptions that are not robust to violation such as the reliability of measurements, constant variation of residuals, and sample size. To test regression models simulation is used to model potential measurement error and the importance of sample sizes on parameter estimation. The randomly perturbed designs are then used together with associated mathematical models obtained from the original designs to simulate experiments and obtain new regression models. A comparison of the original model to the new model, and various statistical tests are performed to determine how accurate the original parameters have been predicted when exposed to simulated measurement error.

**Keywords:** design of experiments, environmental science, multiple regression analysis, simulation

## 1. Introduction

Scientists perform experiments in virtually all areas of study, often to determine a relationship between numerous input factors and either one or multiple output factors. The area known as the Design of Experiments (DOE) is concerned with planning and conducting of experiments, as well as analyzing the resulting data so that valid and objective conclusions are obtained (Montgomery, 2009). Factorial experiments are often used as an experimental strategy in which inputs are varied together. These experiments are an important class of experiments because they may be used to accomplish a variety of different goals, such as perform factor screening, or to determine optimum factor levels. In this study we focus on experiments in environmental sciences, where factorial experiments are primarily used to study the influence of multiple parameters (physical, chemical or biological.) Upon performing experiments dictated by a factorial design, a multiple regression model is created to make predictions and inferences. This is a widely used method, in fact more than 4,000 hits were obtained with the keywords "multiple regression analysis" in the Science Citation Index just within the areas of Environmental Sciences.

Typically, when evaluating a regression model one uses the coefficient of determination ($R^2$) to confirm the goodness of fit of the model to that of experimental data. Coefficients of each variable and its associated *p* value are also used to help assess the influence of the variable on the process under study. But $R^2$ values can be made artificially large by including an excessive number of terms, and p-values only indicate if a term is statistically significant and do not assess the accuracy of parameter estimation. Statisticians have studied the reliability of regression models and have identified a "reliability matrix" (Gleser, 1992) to help assess the model. It is widely known that measurement errors influence regression models (Pagano & Anoke, 2013). However, a reliability matrix is rarely used in an environmental science study.

For a multiple regression model to be reliable, a necessary condition is that it does not violate the regression model assumptions (Kahane, 2008). However, a literature search reveals that in environmental sciences, multiple regression models are often not tested for their robustness to regression assumptions prior to interpretation and drawing conclusions. Two of the key model assumptions that are made are:

1.    The independent and dependent variables do not contain measurement errors.

2.  Residuals of the model are independent over time, normally distributed, have mean zero, and exhibit constant variation.

To further exacerbate the problem, scientists have known for a long time that when carrying out experimental studies two types of errors are introduced during input and output analysis: precision errors and accuracy errors (Taylor, 1982). Precision errors are related to the random errors associated with an experiment (e.g. measurement error), whereas accuracy errors are related to the systematic differences observed between laboratories (e.g. different calibration of the instruments). A regression model developed based on experimental data should also be robust to these experimental errors. Compounded with a limited data set available due to practical considerations, a regression model developed could provide a false sense of parameter efficiency and output predictability.

The current study is aimed at developing a method to evaluate the reliability of regression models when the input and output parameters are subjected to small random perturbations created using simulation. It is our assumption that the coefficients of the variables and their *p* values should not significantly change in reliable models, even when the variables are subjected to simulated random perturbations. Also addressed is how residual analysis of the models can help to evaluate the reliability of the model prediction.

## 2. Methodology

Twenty studies within the field of environmental sciences were selected to determine if the published models were robust enough to withstand simulated random perturbations of the input and output values uniformly distributed between ±5%. The studies examined along with errors measures (as described in this section later) are given in Table 1. The different types of experimental designs considered are: full factorial, fractional factorial, mixture design, Box-Behnken and central composite. The perturbations are used to assess how sensitive the model is to small changes in the input and output data. The changes could be a result of measurement error or possibly due to other types of process variation. To introduce a random perturbation of some value between ±5%, every coordinate of the design points and the output values as well, are multiplied by a simulated random number between 0.95 and 1.05. Table 2a provides an example of design points used in a $2^3$ full factorial design that have been modified by multiplying by a simulated random number between 0.95 and 1.05 to obtain the design matrix given in Table 2b.

Table 1. The list of 20 studies examined along with error measures of the regression models published

| Reference | Type | MAPE | $APE_{90}$ |
|---|---|---|---|
| (Bhunia & Ghangrekar, 2007) | Full Factorial | 11% | 24% |
| (Prasad & Srivasta, 2009) | Full Factorial | 54% | 125% |
| (Lima, et. Al., 2007) | Full Factorial | 5% | 9% |
| Saadat & Karimi-Jashni, 2010) | Full Factorial | 135% | 444% |
| (Srinivasan & Viraraghavan, 2010) | Fractional Factorial | 18% | 39% |
| (Mobilia, Scipioni, Veglio & Sciano, 2010) | Fractional Factorial | 27% | 67% |
| (Dobrev, Pishtiyski, Stanchev & Mircheva, 2007) | Fractional Factorial | 5% | 9% |
| (Chen, Lin, Jones, Fu & Zhan, 2009) | Fractional Factorial | 37% | 109% |
| (Rispoli, et. al, 2010) | Mixture Design | 3% | 6% |
| (Abdullah & Chin, 2010) | Mixture Design | 3% | 6% |
| (Chen, Huang, Hsiao & Tsai, 2010) | Mixture Design | 7% | 14% |
| (Santafe-Moros, et. Al., 2005) | Mixture Design | 38% | 113% |
| (Gurkok, Cekmecelioglu & Ogel, 2011) | Box-Behnken | 5% | 8% |
| (Anunziata & Cussa, 2008) | Box-Behnken | 674% | 1679% |
| (Baskan & Pala, 2010) | Box-Behnken | 33% | 53% |
| (Dopar, Kusic & Koprivanac, 2011) | Box-Behnken | 7% | 19% |
| (Djoudi, Aissani-Benissad & Bourouina-Bacha, 2007) | Central Composite | 21% | 42% |
| (Landaburu-Aguirre, Pongracz, Peramaki & Keiski,2010) | Central Composite | 57% | 151% |
| (Mohajeri, Aziz, Isa & Zahed, 2010) | Central Composite | 10% | 17% |
| (Imandi, Bandaru, Somalanka, Bandaru & Garapati, 2008) | Central Composite | 14% | 27% |

Table 2. (a) Full factorial $2^3$ experimental design.    (b) The factorial design with 5% random perturbations

| $x_1$ | $x_2$ | $x_3$ | $x_1$ | $x_2$ | $x_3$ |
|---|---|---|---|---|---|
| +1 | +1 | −1 | +1.04 | +1.04 | −1.02 |
| −1 | −1 | +1 | −1.03 | −0.99 | +1.00 |
| −1 | +1 | +1 | −1.01 | +1.00 | +0.99 |
| +1 | +1 | +1 | +1.00 | +1.03 | +1.05 |
| +1 | −1 | +1 | +0.98 | −0.97 | +0.96 |
| +1 | −1 | −1 | +1.03 | −0.96 | −0.97 |
| −1 | −1 | −1 | −0.97 | −1.02 | −1.03 |
| −1 | +1 | −1 | −0.96 | +1.05 | −1.04 |
| | | **(a)** | | | **(b)** |

Let **X** denote the entire set of original experimental design points stored as a matrix, and let **Y** denote the vector of experimental outputs. To introduce a random perturbation of some value between ±5%, every coordinate of the design points in **X**, and the output values as well, are multiplied by a simulated random number uniformly distributed between 0.95 and 1.05. This procedure yields the modified design matrix denoted by $\widehat{\mathbf{X}}$. Once $\widehat{\mathbf{X}}$ has been obtained we use the initial polynomial model together with $\widehat{\mathbf{X}}$ to compute predicted output values. The predicted output values are then multiplied by a random number between 0.95 and 1.05 to yield $\widehat{\mathbf{Y}}$. Next, perform a regression analysis using $\widehat{\mathbf{X}}$ and $\widehat{\mathbf{Y}}$, which generates a new polynomial model. Let $\beta_i$ denote the coefficients of the original model. Observe that these are the parameters we are trying to estimate. We let $\hat{\beta}_i$ denote coefficients of the model obtained from $\widehat{\mathbf{X}}$ and $\widehat{\mathbf{Y}}$. For every parameter $\beta_i$, we compute the absolute percentage error (APE) given by $\left| \frac{\beta_i - \hat{\beta}_i}{\beta_i} \right|$. Using these values we then compute the mean absolute percentage error (MAPE). However, since the MAPE can be greatly influenced by one or two extreme terms, we also calculate the $90^{th}$ percentile absolute percentage error ($APE_{90}$). Observe that 90% of the absolute percentage errors will be below $APE_{90}$ and 10% will be larger. Additional nonparametric statistical tests such as the Sign Test can also be performed to test the reliability of the coefficients. A summary of the complete simulation process is given in Figure 1.
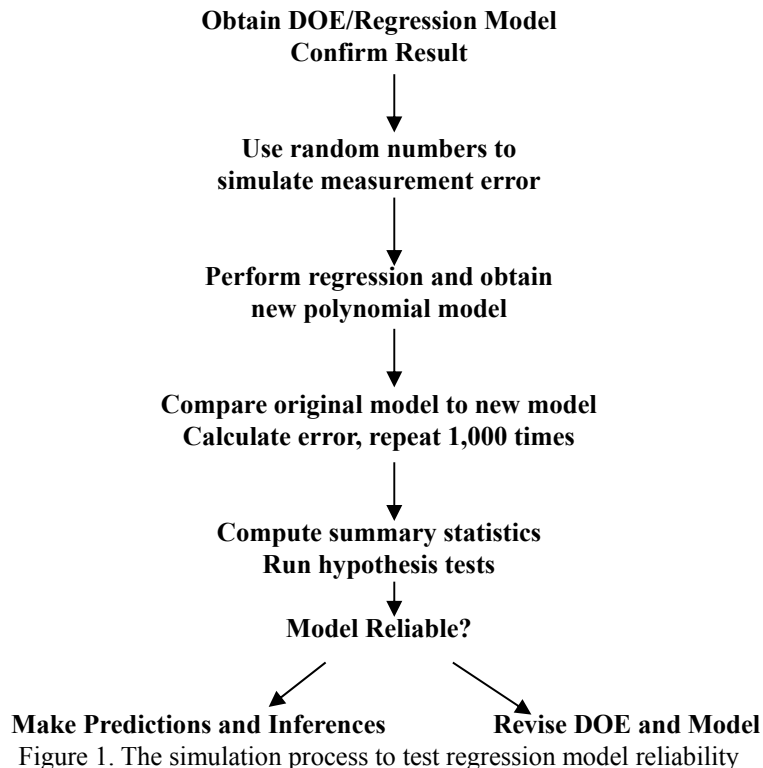
<div align="center">

**Obtain DOE/Regression Model**
**Confirm Result**

↓

**Use random numbers to**
**simulate measurement error**

↓

**Perform regression and obtain**
**new polynomial model**

↓

**Compare original model to new model**
**Calculate error, repeat 1,000 times**

↓

**Compute summary statistics**
**Run hypothesis tests**

↓

**Model Reliable?**

↙          ↘

**Make Predictions and Inferences**          **Revise DOE and Model**

</div>

Figure 1. The simulation process to test regression model reliability

## 3. Results and Discussion

We propose through this study that once a regression model has been developed using a DOE, a random perturbation of ±5% should be introduced into the design points and the output values to help assess the model performance. The perturbation is intended to represent the measurement and systematic errors introduced when performing experiments with limited measurement resolution.

We begin with an analysis of how the type of experimental design influences the robustness with respect to measurement errors. A summary of the 20 studies is given in Table 3. The data illustrates a significant potential of errors influencing the results with Box-Behnken designs which produced results with an average MAPE of 180% and an average $APE_{90}$ of over 400%. In contrast, mixture design seems to be more resilient to error with minimum % errors in both of the measures used. The percentages in Table 3 were obtained by finding the average MAPE and $APE_{90}$ of the four studies in each design group shown in Table 1.

Table 3. A summary of the simulation results for the twenty studies examined

| Design Type | MAPE | $APE_{90}$ |
|---|---|---|
| Basic Factorial | 51% | 151% |
| Fractional Factorial | 22% | 56% |
| Mixture Design | 13% | 35% |
| Box-Benken | 180% | 440% |
| Central Composite | 26% | 59% |

A second factor considered is the sample size measured by the ratio of number of design points divided by the number of predictor variables used in the final model. These ratios varied from a low of 1.25 to a maximum of 4.33. Surprisingly, the increase in this ratio does not necessarily decrease the MAPE. Sample size alone does not seem to be critical with respect to the robustness of the model. What is more important than sample size, is the location of the design points within the design space. In particular, the design types with points located throughout the design space, especially at or near the center, are the designs with the lowest MAPE. To test this rigorously we identified designs that have included design points in the interior of the design space vs. those that contain points only on the boundary. The results indicate a significant improvement in reducing the MAPE when interior points are included in the design.

## 4. Design Specific Analysis

The full factorial designs studies cited in Table 1 all had either 3 or 4 input variables, and between 8 and 16 predictor variables. The $R^2$ values for the models constructed were typically good and all above 0.95. When looking at the full factorial designs as a group we see that the MAPE's vary considerably with a minimum of 5% (Lima et al., 2007), to a maximum of 135% (Saadat & Karimi-Jashni., 2011). If one averages the MAPE's over the 4 studies we obtain a grand mean of 51% for the full factorial designs (Table 2). The design used in (Bhunia & Ghangrekar, 2007) is a full factorial $2^3$ design with 11 design points which includes 3 center points, and yielded the model given in equation (1). The model constructed in (Saadat & Karimi-Jashni., 2011) is based on a full factorial $2^4$ design, no additional interior or exterior points, but with repeated trials and is given below in equation (2).

$$Y = -26.1 + 190.78x_1 + 0.09x_2 + 159.65x_3 - 0.19x_1x_2 - 676.56x_1x_3 \qquad (1)$$

$$Y = 57.4 + 18.0x_1 + 22.7x_2 - 1.2x_3 + 0.7x_4 - 12.8x_1x_2 + 4.9x_1x_3 - 0.1x_1x_4 - 6.0x_2x_3 + 0.1x_2x_4 -$$

$$0.1x_3x_4 - 2.3x_1x_2x_3 - 0.03x_1x_2x_4 + 0.08x_1x_3x_4 - 0.05x_2x_3x_4 - 0.2x_1x_2x_3x_4 \qquad (2)$$

The residual plot for the model described in equation 1 suggests that despite low number of sample size, the assumptions are not violated. When introducing measurement error into the model constructed in (Bhunia & Ghangrekar, 2007), we obtained an MAPE of 11% and an $APE_{90}$ of 24%.

The model (2) was obtained with a sample of 32 design points for 16 predictor variables, so the sample size is only twice the number of predictor variables. A potential source of error that arises in the model given in equation (2) is the inclusion of 3-way and 4-way cross terms with relatively small coefficients. Terms such as these are very sensitive to small measurement errors. Unless there is extreme confidence in the measurement system, we believe that a high degree term should not be included. Indeed, this is consistent with the

"sparsity-of-effects-principle" which states that a system is usually dominated by main effects and low order interactions. The principle has been explained in depth by (Wu & Hamada, 2000). The residuals are essentially normally distributed with mean zero. However, the constant variation assumption is violated with a large variation in the vertical spread, and there is a nonrandom pattern with decreasing residuals illustrated in the residuals versus observation order plot. Hence time is related to decreasing residuals indicating a violation of time independent residuals. After introducing random perturbations we obtained an MAPE of 135% and an $APE_{90}$ of 444%. Hence, the model is neither reliable nor robust to the input data variations.   The study (Prasad & Srivasta, 2009) exhibits similar problems. There are too many variables artificially inflating $R^2$ and the constant variation of residuals requirement is violated.

Mixture designs have recently been developed and are included as an option in most computer aided experimental design software. For a comprehensive reference on mixture designs, see (Cornel, 1981). Unlike the previous two types of designs, mixture design have more than half of their points in the interior of the design space, including one point in the center.  In full factorial designs the design space is an n-dimensional hypercube where n is the number of input variables. However, in a mixture design each point gives the proportion of each input into the mixture. Hence, there is a constraint that requires the sum of the proportions to be one. This implies that the design space is now an n-dimensional simplex.

The study (Abdullah & Chin, 2010) utilizes a simplex-centroid design for optimizing the composting of kitchen waste. The response is the carbon to nitrogen ratio Y, and the model developed is shown in equation (3). Observe that the model contains 3 input variables and 3 interaction terms for a total of 6 predictor terms.

$$Y = 14.9x_1 + 8.2x_2 + 281.6x_3 + 17.7x_1x_2 - 273x_1x_3 - 509.3x_2x_3 \qquad (3)$$

The experimental design consists of 13 design points which yields a ratio of roughly 2 times as many sample points as predictor variables. The model shown in (3) had a high goodness of fit with a $R^2$ of 0.98. The residuals do not violate any of the residual assumptions. Perhaps a better assessment of accurate parameter estimation is indicated by the MAPE being only 2.7% and the $APE_{90}$ at 5.8% (Table 1).   The other studies in Table 1 based on a mixture design all produced remarkably similar results with respect to the error measures MAPE and $APE_{90}$, even when the residuals appeared to be violated. The mixture design is an example of a design that is very robust to small perturbations in the input data.

The next group considered are the Box-Behnken designs listed in Table 1, which are a special case of response surface designs. In (Gurkok et al., 2011) a three-level Box-Behnken design was used in an optimization study. The model developed is a second order model with three independent variables and eight predictor terms. An inspection of the residual plots confirms the residual assumptions. When perturbations are introduced the MAPE is 4.5% and the $APE_{90}$ is 8.3%.   Clearly this design is robust with respect to the regression assumptions.

In study (Anunziata & Cussa, 2008) a Box-Behnken was also used to develop a model with 11 predictor terms using 27 data points. The model is given in equation (4).

$$Y = 21.6404 + 1.52833x_1 - 2.52083x_2 - 12.5292x_3 + 8.875x_4 - 2.17x_1x_2 - 0.005x_1x_3 - 9.475x_1x_4 \qquad (4)$$
$$- 0.0125x_2x_3 + 0.825x_2x_4 - 1.525x_3x_4$$

Inspection of the residual plots does not reveal any violations.   However, when applying perturbations to the design we obtained an MAPE of 674%, with a $APE_{90}$ of 1,679%. We believe that the model is hyper-sensitive to measurement error because of low number of data points used in the study.

We observed an inconsistency among Box-Behnken designs with respect to the inclusion of interior points in the design. The design used in (Annuziata and Cussa, 2008) did not include any center points. Whereas, the designs used in (Gurkok et al., 2011), (Baskan & Pala, 2010), and (Dopar et al., 2011) included 6, 5, and 3 center points respectively. The results obtained in (Baskan & Pala, 2010) and (Dopar et al., 2011) are much more reliable when subjected to random perturbations with an MAPE of 33% and 8% respectively.

The last group of regression models examined is those based on a central composite design listed in Table 1. Central composite designs are often used to obtain a quadratic model that will facilitate optimization. As illustrated in Table 3, this group would be ranked second best with respect to the error measures. The model that performed best when subjected to input perturbations was (Mohajeri et al., 2010), which involves 4 input variables and 30 design points; where 16 of the points consist of a $2^4$ basic factorial design, 8 points are axial points, and 6 points are center points. The response Y for this model is the weathered crude oil removal percentage and is given in equation (5). Observe that there are 11 predictor terms which implies a run to predictor variable ratio of roughly 2.7. Moreover, there are no predictor terms with three factors.

$$Y = 68.43 - 10.55x_1 + 3.66x_2 + 10.14x_3 + 5.57x_4 + 4.79(x_2)^2 - 14.89(x_3)^2 - 7.08(x_4)^{\wedge 2} \qquad (5)$$
$$- 2.28x_1x_2 + 8.63x_1x_3 + 1.81x_3x_4$$

The axial points are identical to the center points except for one factor, which will take on values both below and above the median of the two factorial levels, and typically both outside their range. Interestingly, the number of center points used in the central composite designs of Table 1 is: 12, 3, 6, 6 respectively, and the number of axial points used in these designs is 8, 0, 8 ,8, respectively. The design with 3 center points and 0 axial points was the only "non-robust" central composite design. As in the case of Box-Behnken designs, we see an inconsistency in the number of center and axial points used.

## 5. Conclusions

While the goal of the study was not intended to be comprehensive in the sense of testing a majority of designs, we present a simulation tool that could be used by investigators to test their designs and developed models for coefficient reliability. The study results in the following conclusions:

- The reliability of a regression model is dependent on the type and parameters of the experimental design. Including numerous design points in the interior of the design space, such as center and axial points, increases the reliability and robustness of the model.

- Scientists should evaluate the robustness of the regression model assumptions before making inferences. Models with very good $R^2$ and *p*-values may not be very reliable when measurement errors are taken into consideration leading to false inferences. Testing the regression model using simulation would provide a greater degree of confidence in the scientific inferences made. It may also provide new insight into the sensitivity of the scientific process being studied.

- High order terms should be avoided as much as possible. Primary effect terms and two-way interactions and squared term are usually sufficient and lead to models that are much more stable than models with high order terms.

## Acknowledgement

## References

Abdullah, N. C. N. (2010). Simplex-centroid mixture formulation for optimised composting of kitchen waste. *Bioresource Technol, 101*, 8205-8210.

Anunziata, O. A, & Cussa, J. (2008). Applying response surface design to the optimization of methane activation with ethane over Zn-H_ZSM-11 zeolite. *Chemical Eng. J., 138*, 510-516.

Baskan, M. B., & Pala, A. (2010). A statistical experiment design approach for arsenic removal by coagulation process using aluminum sulfate. *Desalination, 254*, 42-48.

Bhunia, P., & Ghangrekar, M. M. (2007). Statistical modeling and optimization of biomass granulation and COD removal in UASB reactors treating low strength wastewaters. *Bioresource Technology, 99*, 4229-4238.

Chen, L. C, Huang, C. H, Hsiao, M. C, & Tsai, F. R. (2010). Mixture design optimization of the composition of S, C, $S_nO_2$-codoped $TiO_2$ for degradation of phenol under visible light. *Chemical Eng. J., 165*, 482-489.

Chen, Y, Lin, C, Jones, G, Fu, S., & Zhan, H. (2009). Enhancing biodegradation of wastewater by microbial consortia with fractional factorial design. *J. Hazardous Mat., 171*, 948-953.

Cornell, J. A. (1981). Experiments with mixtures: designs, models and the analysis of mixture data. New York: Wiley.

Djoudi, W., Aissani, B. F., & Bourouina, B. S. (2007). Optimization of copper cementation process by iron using central composite design experiments. *Chemical Eng, J., 133*, 1-6.

Dobrev, G. T., Pishtiyski, I. G., Stanchev, V. S., & Mircheva, R. (2007) Optimization of nutrient medium containing agricultural wastes for xylanase production by *Aspergillus niger BO3* using optimal composite experimental design. *Bioresource Technology, 98*, 2671-2678.

Dopar, M., Kusic, H., & Koprivanac, N. (2011). Treatment of simulated wastewater by photo-Fenton process. Part 1: The optimization of process parameters using design of experiments (DOE). *Chemical Eng. J., 173*, 267-279

Gleser, L. J. (1992). The importance of assessing measurement reliability in multivariate regression. *J. American*

*Statistical Association, 87*(419), 696–707.

Gurkok, S., Cekmecelioglu, D., & Ogel, Z. B. (2011). Optimization of culture conditions for *Aspergillus sojae* expressing an *Aspergillus fumigates* β-galactosidase. *Bioresource Technol, 102*, 4925-4929.

Imandi, S. B., Bandaru, V. V. R., Somalanka, S. R., Bandaru, S. R., & Garapati, H. R. (2008). Application of statistical experimental designs for the optimization of medium constituents for the production of citric acid from pineapple waste. *Bioresource Technol, 99*, 4445-4450.

Kahane, L. H. (2008). *Regression Basics* (2nd ed.). Sage.

Landaburu, A. J., Pongracz, E., Peramaki, P., & Keiski, R. L. (2010). Micellar-enhanced ultrafiltration for the removal of cadmium and zinc: Use of response surface methodology to improve understanding of process performance and optimization. *J. Hazardous Mat., 180*, 524-534.

Lima, E. C, Royer, B., Vaghetti, J. C. P., Brasil, J. L., & Simon, N. M., et al. (2007). Adsorption of Cu(II) on Araucaria angustifolia wastes: Determination of the optimal conditions by statistic design of experiments. *J. Hazardous Mat., 140*, 211-220.

Mabilia, R., Scipioni, C., Veglio, F., & Sciano, M. C. T. (2010). Fractional factorial experiments using a test atmosphere to assess the accuracy and precision of a new passive sampler for the determination of formaldehyde in the atmosphere. *Atmospheric Environment, 44*, 3942-3951.

Mohajeri, L., Aziz, H. A., Isa, M. H., & Zahed, M. A. (2010). A statistical experiment design approach for optimizing biodegradation of weathered crude oil in coastal sediments. *Bioresource Technol, 101*, 893-900.

Montgomery, D. C. (2009). *Design and Analysis of Experiments* (7th ed.). New York: Wiley

Pagano, M., & Anoke, S. (2013). Mommy's Baby, Daddy's Maybe: A Closer Look at Regression to the Mean. *Chance Magazine, 27*(1).

Prasad, R. K., & Srivasta, S. N. (2009). Electrochemical degradation of distillery spent wash using catalytic anode: Factorial Design of Experiments. *Chemical Eng. J., 146*, 22-29.

Rispoli, F., Angelov, A., Badia, D., Kumar, A., & Seal, S., et al. (2010). Understanding the toxicity of aggregated zero valent copper nanoparticles again*st Escherichia coli. J. Hazardous Mat., 185*, 212-216.

Saadat, S., & Karimi, J. A. (2011). Optimization of Pb(II) adsorption onto modified walnut shells iusing factorial design and simplex methodologies. *Chemical Eng. J., 173*, 743-749.

Santafe, M. A., Gozalvez, Z. J. M., Lora, G. L., & Garcia, D. J. C. (2005). Mixture design applied to describe the influence of ionic composition on the removal of nitrate ions using nanofiltration. *Desalination, 185*, 289-296.

Srinivasan, A., & Viraraghavan, T. (2010). Oil removal from water by fungal biomass: A factorial design analysis. *J. Hazardous Mat., 175*, 695-702.

Taylor, J. R. (1982). An introduction to error analysis. University Science Books.

Wu, C. F. J., & Hamada, M. (2000). *Experiments: Planning, analysis and parameter design optimization.* New York: Wiley.

**Copyrights**