# A UNIFIED APPROACH TO POWER CALCULATION AND SAMPLE SIZE DETERMINATION FOR RANDOM REGRESSION MODELS

## GWOWEN SHIEH

### NATIONAL CHIAO TUNG UNIVERSITY

The underlying statistical models for multiple regression analysis are typically attributed to two types of modeling: fixed and random. The procedures for calculating power and sample size under the fixed regression models are well known. However, the literature on random regression models is limited and has been confined to the case of all variables having a joint multivariate normal distribution. This paper presents a unified approach to determining power and sample size for random regression models with arbitrary distribution configurations for explanatory variables. Numerical examples are provided to illustrate the usefulness of the proposed method and Monte Carlo simulation studies are also conducted to assess the accuracy. The results show that the proposed method performs well for various model specifications and explanatory variable distributions.

Key words: asymptotic distribution, effect size, noncentral $F$ distribution.

## 1. Introduction

Multiple regression analysis is one of the widely used statistical methods. Conventionally, there are two approaches to the statistical modeling of these regression applications. They are referred to as fixed (conditional) and random (unconditional) models. In the context of regression analysis, it is quite common in the behavioral and social sciences to have studies in which not only the values of response variables for each experimental unit are just available after the observations are made, but also the levels of explanatory variables cannot be fixed in advance. Therefore, the explanatory variables are also outcomes of the study under such circumstances. In order to take account of this extra variability, the appropriate strategy is to consider the random regression setting. On the other hand, the fixed regression model is suitable for studies in which the configurations of the explanatory variables are preset by the researcher.

Sample size calculations and power analyses are often critical for researchers to address specific scientific hypotheses and confirm credible treatment effects. Thus, they should be an integral part of the whole study. Accordingly, it is of practical importance to be able to perform these tasks in a multiple regression setup. For fixed regression models, the procedures are well documented in the literature. Regarding random regression models, Gatsonis and Sampson (1989) gave an excellent and thorough description of exact power and sample size calculations when the response and explanatory variables have a joint multivariate normal distribution. Traditionally, the problem is referred to as multiple correlation analysis and the parameter of interest is the squared multiple correlation coefficient. In contrast, Algina and Olejnik (2000) presented results for determining the sample size required for adequate estimation accuracy of the squared multiple correlation coefficient. Furthermore, related treatments can be found in Kelley and Maxwell (2003), Mendoza and Stafford (2001), Shieh (2006) and Steiger and Fouladi (1992). It is important to note that

the studies cited above for random regression models and multiple correlation analysis are applicable in the circumstance that the response and explanatory variables have a joint multivariate normal distribution. However, there are many situations in which assuming normal distribution for explanatory variables is inappropriate. For instance, consider the simple interaction model in the formulation of $Y = \beta_0 + X\beta_1 + Z\beta_2 + XZ\beta_3 + \varepsilon$. It is commonly assumed that the continuous measurements $X$ and $Z$ are normally distributed. However, the product of two normally distributed variables ($XZ$) does not have a normal distribution. Therefore, the existing results for power analysis of a multinormal situation do not apply in this application. In fact, Gatsonis and Sampson (1989) noted that when the joint distribution of $Y$ and all explanatory variables is nonnormal, it is doubtful that their power calculations give results that are accurate and reasonable. Therefore, neglecting other configurations of explanatory variables is an obvious limitation of available methods. A natural generalization to incorporate both normal and nonnormal explanatory variables should be essential to the existing approaches for performing power and sample size calculations in practice. It should be observed that the prescribed simple interaction model is directly connected to the moderated multiple regression formulation which has been pervasively used for testing moderator effects in all areas of social sciences, see Aguinis, Beaty, Boik, and Pierce (2005) for further details.

This paper aims to provide a unified approach to the determinations of power and sample size for random regression models. The distinct feature of the proposed method is the accommodation of arbitrary discrete and/or continuous distribution formulations for explanatory variables. Therefore, the aforementioned multivariate normal setting can be viewed as a special case. For related results and various extensions in hierarchical linear models and multivariate linear models, the interested reader is referred to Raudenbush and Liu (2000, 2001), Shieh (2003, 2005), and the references therein. In fact, the suggested two-stage methodology can be viewed as an extension of the results for the univariate case in Shieh (2005). The rest of the paper is organized as follows. In Section 2, the important analytical details of the proposed method are described. Numerical examples are provided in Section 3 to demonstrate the proposed power and sample size calculations for several random regression formulations. Since the approach considered here uses large sample approximations, simulation studies are conducted to assess its adequacy for finite sample and robustness under various model specifications and distributions of explanatory variables. Finally, Section 4 contains some final remarks.

## 2. The Proposed Method

To facilitate the illustration of the proposed method for random regression models, it is instructive to review first the situation under the fixed regression models where the results would be specific to the particular values of the explanatory variables that are observed or predetermined by the researcher.

### 2.1. Review of Fixed Regression Models

Consider the standard multiple linear regression model with response variable $Y$ and all the levels of $p$ explanatory variables $X(1), \ldots, X(p)$ fixed a priori:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \tag{1}$$

where $\mathbf{Y} = (Y_1, \ldots, Y_N)^{\mathrm{T}}$, $Y_i$ is the value of the response variable $Y$; $\mathbf{X} = (\mathbf{1}_N, \mathbf{X}_D)$ where $\mathbf{1}_N$ is the $N \times 1$ vector of all 1's, $\mathbf{X}_D = (\mathbf{X}_1, \ldots, \mathbf{X}_N)^{\mathrm{T}}$ is often called the design matrix, $\mathbf{X}_i = (x_{i1}, \ldots, x_{ip})^{\mathrm{T}}$, $x_{i1}, \ldots, x_{ip}$ are the known constants of the $p$ explanatory variables for $i = 1, \ldots, N$; $\boldsymbol{\beta} = (\beta_0, \beta_1, \ldots, \beta_p)^{\mathrm{T}}$ where $\beta_0, \beta_1, \ldots, \beta_p$ are unknown parameters; and $\boldsymbol{\varepsilon} =$

$(\varepsilon_1, \ldots, \varepsilon_N)^{\mathrm{T}}$ where $\varepsilon_i$ are iid $N(0, \sigma^2)$ random variables. We are concerned with the general linear hypothesis $H_0 : \mathbf{L}\boldsymbol{\beta} = \boldsymbol{\theta}$ versus $H_1 : \mathbf{L}\boldsymbol{\beta} \neq \boldsymbol{\theta}$, where $\mathbf{L}$ is an $l \times (p+1)$ coefficient matrix of rank $l \leq p + 1$ and $\boldsymbol{\theta}$ is an $l \times 1$ vector of constants. It is well known that under the assumption given in (1), the likelihood ratio test for $H_0$ is based on

$$F = \frac{SSH/l}{SSE/(N - p - 1)}, \tag{2}$$

where $SSH = (\mathbf{L}\hat{\boldsymbol{\beta}} - \boldsymbol{\theta})^{\mathrm{T}}[\mathbf{L}(\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{L}^{\mathrm{T}}]^{-1}(\mathbf{L}\hat{\boldsymbol{\beta}} - \boldsymbol{\theta})$, $SSE = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^{\mathrm{T}}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ and $\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{Y}$ is an unbiased estimator of $\boldsymbol{\beta}$, see Rencher (2000, Chaps. 7–8) for further details. Under the alternative hypothesis, $F$ is distributed as $F(l, N - p - 1, \lambda)$, the noncentral $F$ distribution with $l$ and $N - p - 1$ degrees of freedom and noncentrality parameter

$$\lambda = (\mathbf{L}\boldsymbol{\beta} - \boldsymbol{\theta})^{\mathrm{T}}\left[\mathbf{L}(\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{L}^{\mathrm{T}}\right]^{-1}(\mathbf{L}\boldsymbol{\beta} - \boldsymbol{\theta})/\sigma^2. \tag{3}$$

If the null hypothesis is true, then $\lambda = 0$ and $F$ is distributed as $F(l, N - p - 1)$, a central or regular $F$ distribution with $l$ and $N - p - 1$ degrees of freedom. The test is carried out by rejecting $H_0$ if $F > F_{l,N-p-1,\alpha}$, where $F_{l,N-p-1,\alpha}$ is the upper $100\alpha$ percentage point of the central $F$ distribution $F(l, N - p - 1)$.

To calculate power and sample size, it is assumed that there are $m$ distinct configurations of $\mathbf{X}_i$ for $i = 1, \ldots, N$, and they are denoted by $\mathbf{Z}_j$ with the proportions $w_j$, $j = 1, \ldots, m$ ($\leq N$). Then, $\mathbf{X}^{\mathrm{T}}\mathbf{X}$ can be expressed as $\mathbf{X}^{\mathrm{T}}\mathbf{X} = N \cdot \boldsymbol{\Gamma}$, where $\boldsymbol{\Gamma} = \sum_{j=1}^{m} w_j \mathbf{Z}_j \mathbf{Z}_j^{\mathrm{T}}$. Accordingly, the noncentrality parameter $\lambda$ in (3) is rewritten as

$$\lambda = N\delta, \tag{4}$$

where $\delta = (\mathbf{L}\boldsymbol{\beta} - \boldsymbol{\theta})^{\mathrm{T}}[\mathbf{L}\boldsymbol{\Gamma}^{-1}\mathbf{L}^{\mathrm{T}}]^{-1}(\mathbf{L}\boldsymbol{\beta} - \boldsymbol{\theta})/\sigma^2$ is the so-called effect size, see Cohen (1988). Hence, given all model configurations and sample size $N$, the statistical power achieved for testing hypothesis $H_0 : \mathbf{L}\boldsymbol{\beta} = \boldsymbol{\theta}$ with specified significance level $\alpha$ against the alternative $H_1 : \mathbf{L}\boldsymbol{\beta} \neq \boldsymbol{\theta}$ is the probability

$$P\big\{F(l, N - p - 1, N\delta) > F_{l,N-p-1,\alpha}\big\}, \tag{5}$$

where $\delta$ is defined in (4). Furthermore, this power function can be utilized to calculate the sample size needed in order to attain the specified power. However, it usually involves an iterative process to find the solution because both $F(l, N - p - 1, N\delta)$ and $F_{l,N-p-1,\alpha}$ depend on the sample size $N$.

## 2.2. Random Regression Models

To extend the concept and interpretation of the aforementioned results, we assume that the explanatory variables $\{\mathbf{X}_i^* = \mathbf{x}_i^*, i = 1, \ldots, N\}$ in (1) have a probability function $f(\mathbf{X}_i^*)$ with finite moments. The form of $f(\mathbf{X}_i^*)$ is assumed to be dependent on none of the unknown parameters $\boldsymbol{\beta}$ and $\sigma^2$. Thus, the notations of $\mathbf{X}_i$ and $\mathbf{X}_D$ as observed values in the fixed regression model are replaced by $\mathbf{X}_i^*$ and $\mathbf{X}_D^* = (\mathbf{X}_1^*, \ldots, \mathbf{X}_N^*)^{\mathrm{T}}$ as random variables hereinafter. Frequently, the inferences are concerned mainly with the regression coefficients $\boldsymbol{\beta}_1 = (\beta_1, \ldots, \beta_p)^{\mathrm{T}}$ and the corresponding coefficient matrix is written in the form of $\mathbf{L} = \mathbf{L}_1$, where $\mathbf{L}_1 = (\mathbf{0}_c, \mathbf{C})$, $\mathbf{0}_c$ is the $c \times 1$ null vector of all 0's and $\mathbf{C}$ is a $c \times p$ coefficient matrix of rank $c \leq p$. It follows from the overall estimator $\hat{\boldsymbol{\beta}}$ given above that the prescribed estimator for $\boldsymbol{\beta}_1$ can be expressed as $\hat{\boldsymbol{\beta}}_1^* = (\mathbf{X}_C^{*\mathrm{T}}\mathbf{X}_C^*)^{-1}\mathbf{X}_C^{*\mathrm{T}}\mathbf{Y}$, where $\mathbf{X}_C^* = (\mathbf{I}_N - \mathbf{J}/N)\mathbf{X}_D^*$ is the centered form of $\mathbf{X}_D^*$, $\mathbf{I}_N$ is the

identity matrix of dimension $N$ and $\mathbf{J}$ is the $N \times N$ square matrix of all 1's. In view of the extra random nature of $\mathbf{X}_D^*$, it is easily seen that

$$\mathbf{C}\hat{\boldsymbol{\beta}}_1^* | \mathbf{X}_D^* \sim N_c\big(\mathbf{C}\boldsymbol{\beta}_1, \sigma^2 \mathbf{C}\big(\mathbf{X}_C^{*\mathrm{T}}\mathbf{X}_C^*\big)^{-1}\mathbf{C}^{\mathrm{T}}\big).$$

It therefore follows that the general linear hypothesis reduces to $H_0 : \mathbf{C}\boldsymbol{\beta}_1 = \boldsymbol{\theta}$ versus $H_1 : \mathbf{C}\boldsymbol{\beta}_1 \neq \boldsymbol{\theta}$ and the test statistic is of the form

$$F^* = \frac{SSH^*/c}{SSE^*/(N - p - 1)}, \tag{6}$$

$SSH^* = (\mathbf{C}\hat{\boldsymbol{\beta}}_1^* - \boldsymbol{\theta})^{\mathrm{T}}[\mathbf{C}(\mathbf{X}_C^{*\mathrm{T}}\mathbf{X}_C^*)^{-1}\mathbf{C}^{\mathrm{T}}]^{-1}(\mathbf{C}\hat{\boldsymbol{\beta}}_1^* - \boldsymbol{\theta})$, $SSE^* = (\mathbf{Y} - \mathbf{X}^*\hat{\boldsymbol{\beta}}^*)^{\mathrm{T}}(\mathbf{Y} - \mathbf{X}^*\hat{\boldsymbol{\beta}}^*)$, $\mathbf{X}^* = (\mathbf{1}_N, \mathbf{X}_D^*)$ and $\hat{\boldsymbol{\beta}}^* = (\mathbf{X}^{*\mathrm{T}}\mathbf{X}^*)^{-1}\mathbf{X}^{*\mathrm{T}}\mathbf{Y}$. Under random formulation, $F^*$ has the conditional distribution

$$F^* | \Lambda \sim F(c, N - p - 1, \Lambda), \tag{7}$$

where the noncentrality parameter $\Lambda = (\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})^{\mathrm{T}}[\mathbf{C}(\mathbf{X}_C^{*\mathrm{T}}\mathbf{X}_C^*)^{-1}\mathbf{C}^{\mathrm{T}}]^{-1}(\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})/\sigma^2$ is a random variable and the exact distribution depends ultimately on the joint distribution of $\mathbf{X}_D^*$. In order to provide a generally useful and versatile solution without specifically confining to any particular $\mathbf{X}_D^*$, the asymptotic property of $\Lambda$ is studied next.

Let $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ denote the mean vector and covariance matrix of the random explanatory variables $\mathbf{X}_i^* = (X_{i1}^*, \ldots, X_{ip}^*)^{\mathrm{T}}$, respectively. It follows from the standard asymptotic result (Muirhead, 1982, Cor. 1.2.18) that $\mathbf{S}^* = (\mathbf{X}_C^{*\mathrm{T}}\mathbf{X}_C^*)/(N - 1)$ has asymptotic normal distribution

$$(N - 1)^{1/2}\big[\mathrm{vec}(\mathbf{S}^*) - \mathrm{vec}(\boldsymbol{\Sigma})\big] \overset{.}{\sim} N_{p^2}\big(\mathbf{0}_{p^2}, \boldsymbol{\Psi} - \mathrm{vec}(\boldsymbol{\Sigma}) \cdot \mathrm{vec}(\boldsymbol{\Sigma})^{\mathrm{T}}\big),$$

where $\mathrm{vec}(\cdot)$ is a matrix operator which arranges the columns of a matrix into one long column, $\boldsymbol{\Psi} = E[(\mathbf{X}_i^* - \boldsymbol{\mu})(\mathbf{X}_i^* - \boldsymbol{\mu})^{\mathrm{T}} \otimes (\mathbf{X}_i^* - \boldsymbol{\mu})(\mathbf{X}_i^* - \boldsymbol{\mu})^{\mathrm{T}}]$, $E[\cdot]$ denotes the expectation taken with respect to the distribution of $\mathbf{X}_i^*$, and $\otimes$ represents the Kronecker product. Using the identity $\mathrm{vec}(\mathbf{ABC}) = (\mathbf{C}^{\mathrm{T}} \otimes \mathbf{A}) \cdot \mathrm{vec}(\mathbf{B})$, the noncentrality parameter $\Lambda$ given in (7) can be expressed as $\Lambda = (N - 1)\Delta$, where

$$\Delta = \big[(\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})^{\mathrm{T}} \otimes (\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})^{\mathrm{T}}\big] \cdot \mathrm{vec}\big[\big(\mathbf{C}\mathbf{S}^{*-1}\mathbf{C}^{\mathrm{T}}\big)^{-1}\big]/\sigma^2.$$

Let $\partial \Delta / \partial \mathrm{vec}(\mathbf{S}^*)$ denote the $p^2$-dimensional column vector whose $i$th component is the derivative of $\Delta$ with respect to the $i$th element of $\mathrm{vec}(\mathbf{S}^*)$. It can be shown by applying the algebraic manipulation and matrix differentiation results that

$$\frac{\partial \Delta}{\partial \mathrm{vec}(\mathbf{S}^*)} = \big\{\big[\mathbf{S}^{*-1}\mathbf{C}^{\mathrm{T}}\big(\mathbf{C}\mathbf{S}^{*-1}\mathbf{C}^{\mathrm{T}}\big)^{-1}(\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})\big] \otimes \big[\mathbf{S}^{*-1}\mathbf{C}^{\mathrm{T}}\big(\mathbf{C}\mathbf{S}^{*-1}\mathbf{C}^{\mathrm{T}}\big)^{-1}(\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})\big]\big\}/\sigma^2.$$

For operational ease, the derivative is computed ignoring the symmetry of $\mathbf{S}^*$. Then, it can be readily derived from the Cramer delta method that $\Delta$ has the following large-sample distribution

$$\Delta \overset{.}{\sim} N(\mu_\Delta, \Sigma_\Delta), \tag{8}$$

where

$$\mu_\Delta = (\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})^{\mathrm{T}}\big(\mathbf{C}\boldsymbol{\Sigma}^{-1}\mathbf{C}^{\mathrm{T}}\big)^{-1}(\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})/\sigma^2$$

and

$$\Sigma_\Delta = (\mathbf{G} \otimes \mathbf{G})^{\mathrm{T}}\boldsymbol{\Psi}(\mathbf{G} \otimes \mathbf{G})/\big\{(N - 1)\sigma^4\big\} - \mu_\Delta^2/(N - 1)$$

with $\mathbf{G} = \boldsymbol{\Sigma}^{-1}\mathbf{C}^{\mathrm{T}}(\mathbf{C}\boldsymbol{\Sigma}^{-1}\mathbf{C}^{\mathrm{T}})^{-1}(\mathbf{C}\boldsymbol{\beta}_1 - \boldsymbol{\theta})$. For $\mathbf{C} = \mathbf{I}_p$ and $\boldsymbol{\theta} = \mathbf{0}_p$, it leads to the useful results that $\mu_\Delta = \boldsymbol{\beta}_1^{\mathrm{T}}\boldsymbol{\Sigma}\boldsymbol{\beta}_1/\sigma^2$ and $\Sigma_\Delta = (\boldsymbol{\beta}_1 \otimes \boldsymbol{\beta}_1)^{\mathrm{T}}\boldsymbol{\Psi}(\boldsymbol{\beta}_1 \otimes \boldsymbol{\beta}_1)/\{(N-1)\sigma^4\} - \mu_\Delta^2/(N-1)$. Consequently, the conditional distribution of $F^*|\Lambda$ with $\Lambda = (N-1)\Delta$, and asymptotic distribution of $\Delta$ described in the last two equations specify fully the proposed approximate distribution for test statistic $F^*$. It is clear under the null hypothesis that the distribution of $F^*$ remains as $F(c, N-p-1)$ under both fixed and random settings as in the special case of the multinormal distribution of Sampson (1974). Hence, the test is conducted by rejecting $H_0$ if $F^* > F_{c,N-p-1,\alpha}$. However, the power function associated with the general linear hypothesis $H_0: \mathbf{C}\boldsymbol{\beta}_1 = \boldsymbol{\theta}$ versus $H_1: \mathbf{C}\boldsymbol{\beta}_1 \neq \boldsymbol{\theta}$ can be well approximated by synthesizing the results in (7) and (8) as

$$P\{F^* > F_{c,N-p-1,\alpha}\} \doteq \int_{-\infty}^{\infty} P\{F(c, N-p-1, (N-1)\Delta) > F_{c,N-p-1,\alpha}\} \cdot g(\Delta)\, \mathrm{d}\Delta, \quad (9)$$

where $g(\Delta)$ is the normal pdf of $\Delta$ defined in (8). The numerical computation of approximate power requires the evaluations of central and noncentral $F$ cdfs and the one-dimensional integration with respect to a normal pdf. Since all related functions are readily embedded in modern statistical packages such as the SAS system, no substantial computing efforts are required.

For the purpose of sample size determination, the approximate power function defined in (9) can be employed to calculate the sample size needed to test hypothesis $H_0: \mathbf{C}\boldsymbol{\beta}_1 = \boldsymbol{\theta}$ versus $H_1: \mathbf{C}\boldsymbol{\beta}_1 \neq \boldsymbol{\theta}$ in order to attain the specified power for the chosen significance level $\alpha$, parameter values $\boldsymbol{\beta}$ and $\sigma^2$, and probability distribution $f(\mathbf{X}_i^*)$. The necessary sample size can be found through a simple iterative search. To reduce the computational effort in the search process, the starting sample size can be selected from the following simplified examination. Consider the even stronger asymptotic results for $\mathbf{S}^*$ and $SSE^*$ that $\mathbf{S}^*$ converges in probability to $\boldsymbol{\Sigma}$ and $SSE^*/\{(N-p-1)\sigma^2\}$ converges in probability to 1. It follows from the application of Slutsky's theorem that $SSH^*/\{(N-1)\sigma^2\}$ converges in distribution to the chi-square distribution $\chi^2(c, \mu_\Delta)$, the noncentral chi-square distribution with $c$ degrees of freedom and noncentrality parameter $\mu_\Delta$, where $\mu_\Delta$ is defined in (8). More importantly, the distribution of the $c \cdot F^*$ statistic can be alternatively approximated by the distribution $\chi^2(c, (N-1)\mu_\Delta)$. Therefore, the corresponding approximate power function is $P\{\chi^2(c, (N-1)\mu_\Delta) > \chi^2_{c,\alpha}\}$, where $\chi^2_{c,\alpha}$ is the upper $100\alpha$ percentage point of the central chi-square distribution $\chi^2(c)$. Hence, the sample size, say $N_{\mathrm{CS}}$, required to achieve the specified power level is a one-time direct inversion of a noncentral chi-square cdf. In general, the resulting sample size provides a close but smaller value than the desired outcome according to the proposed mixture of the noncentral $F$ cdf in (9). Note that the probability $P\{F^* > F_{c,N-p-1,\alpha}\}$, for fixed values of $c$, $p$, $\alpha$ and model parameters, is increasing in sample size $N$. Hence, by starting with sample size $N_{\mathrm{CS}}$ for $N$, it only requires a small number of incremental searches in order to find the minimum sample size that attains the nominal power.

It is noteworthy that the proposed approach avoids the need for a full specification of the joint distributional form of $\mathbf{X}_i^*$ by only assuming the second- and fourth-order mixed central moments of the underlying distribution. However, the calculations are fairly straightforward for some well-known distributions and it may require more involved mathematical manipulations (integration or summation) for complex and nonstandard situations. On the contrary, the mean of $\mathbf{X}_i^*$ is immaterial to the distribution of $\mathbf{S}^*$ and, more importantly, the suggested approximation. Additionally, it should be noted that the effect size $\delta$ given in (4) plays an important role in power and sample size determinations for fixed regression models. Owing to the proposed two-stage distribution approximation to the $F^*$ statistic described in (7) and (8), there is no simple closed-form expression for the effect size in (9). However, it can be comprehended from the mean value $\mu_\Delta$ of the random noncentrality parameter $\Delta$. Therefore, the computed value of $\mu_\Delta$ is viewed as a pseudo effect size.

### 3. Numerical Examples

For illustrative purposes, we present in this section the power and sample size calculations for a random simple regression model and the moderated multiple regression or interaction regression model with two continuous predictor variables and their cross-product term.

First, the random simple regression ($p = 1$) of the form $Y = \beta_0 + X\beta_1 + \varepsilon$ is investigated, where $\varepsilon$ has a normal distribution $N(0, \sigma^2)$. Without loss of generality, both the intercept parameter $\beta_0$ and the variance $\sigma^2$ are taken to be 1. As suggested by a referee, we consider two classes of distributions for the explanatory variable $X$, namely standardized gamma and standardized Poisson distributions. Therefore, the mean and variance of $X$ are identically $\mu = 0$ and $\Sigma = 1$, respectively. In order to investigate the finite-sample properties of the suggested procedure with respect to various shapes of distributions, the gamma distributions with shape parameter 9, 4 and 1 and scale parameter 1, denoted by gamma(9, 1), gamma(4, 1) and gamma(1, 1), and the Poisson distributions with mean 9, 4 and 1, denoted by Poisson(9), Poisson(4) and Poisson(1), are considered. Note that the skewness and kurtosis of the gamma($a$, 1) distribution are $2/a^{1/2}$ and $3 + 6/a$, respectively. Hence, the actual values of skewness and kurtosis for the three prescribed gamma distributions are $(0.67, 3.67)$, $(1, 4.5)$ and $(2, 9)$, respectively. In the case of Poisson distribution, the skewness and kurtosis of the Poisson($\lambda$) distribution are $1/\lambda^{1/2}$ and $3 + 1/\lambda$, respectively. Thus, for the three Poisson distributions, the corresponding skewness and kurtosis are $(0.33, 3.11)$, $(0.5, 3.25)$ and $(1, 4)$. For the test of slope coefficient ($c = p = 1$) $H_0 : \beta_1 = 0$ versus $H_1 : \beta_1 \neq 0$, it follows from (9) that the sample size needed to obtain the power $1 - \gamma$ at the significance level $\alpha = 0.05$ is the minimum number $N$ such that the approximate power function

$$P\{F^* > F_{1,N-2,\alpha}\} \doteq \int_{-\infty}^{\infty} P\{F(1, N-2, (N-1)\Delta) > F_{1,N-2,\alpha}\} \cdot g(\Delta)\, d\Delta \geq 1 - \gamma, \quad (10)$$

where $\Delta \overset{\cdot}{\sim} N(\mu_\Delta, \Sigma_\Delta)$, $\mu_\Delta = \beta_1^2 \cdot \Sigma/\sigma^2$ and $\Sigma_\Delta = \beta_1^4 \cdot (\Psi - \Sigma^2)/[(N-1)\sigma^4]$. For regression coefficient $\beta_1 = 0.3$, 0.4 and 0.5 and power level $1 - \gamma = 0.80$, 0.90 and 0.95, the calculated sample sizes of the proposed method are presented in Tables 1 and 2 for the standardized gamma and standardized Poisson distributions of $X$, respectively. The results in the two tables reveal the general relation that sample sizes increase with increasing power and kurtosis, and decrease with increasing value of $\beta_1$. To demonstrate the power computation, the precise achieved powers associated with the derived sample sizes are recalculated with the proposed approximation given in (10) for all cases. As expected, the resulting approximate powers are slightly larger than their corresponding nominal power levels. Specifically, the calculated sample sizes of the proposed method for gamma(9, 1) distribution with $\beta_1 = 0.3$ are 93, 124 and 152 for power 0.80, 0.90 and 0.95, respectively. The corresponding approximate powers are 0.8027, 0.9020 and 0.9500, and almost identical to 0.80, 0.90 and 0.95, respectively.

The second model under consideration is the simple interaction model: $Y = \beta_0 + X\beta_1 + Z\beta_2 + XZ\beta_3 + \varepsilon$, where $\varepsilon \sim N(0, 1)$. The two predictors $(X, Z)$ are jointly normally distributed with mean $(0, 0)$, variance $(1, 1)$ and correlation $\rho$. It is important to note that, although both $X$ and $Z$ are normally distributed, the interaction term $XZ$ is obviously not a normal random variable. Therefore, the established methods for multinormal covariates are inappropriate for the power and sample size calculations of this interaction regression model. In this case, it can be shown that $E[XZ] = \rho$, $E[X^2Z] = E[XZ^2] = 0$ and $V[XZ] = 1 + \rho^2$, see Aiken and West (1991, Appendix A). Therefore, we have

$$\Sigma = E[\mathbf{H}] = \begin{bmatrix} 1 & \rho & 0 \\ \rho & 1 & 0 \\ 0 & 0 & 1+\rho^2 \end{bmatrix}, \quad \text{where } \mathbf{H} = \begin{bmatrix} X^2 & XZ & X(XZ-\rho) \\ XZ & Z^2 & Z(XZ-\rho) \\ X(XZ-\rho) & Z(XZ-\rho) & (XZ-\rho)^2 \end{bmatrix}.$$

TABLE 1.

Calculated sample sizes, approximate powers and simulated powers of the proposed method for random simple regression models $Y = \beta_0 + X\beta_1 + \varepsilon$ with standardized gamma predictor ($p = c = 1$ and $\alpha = 0.05$).

| | Gamma(9, 1) with kurtosis = 3.67 | | | | Gamma(4, 1) with kurtosis = 4.5 | | | | Gamma(1, 1) with kurtosis = 9 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| N | Approximate power | Simulated power | Error | N | Approximate power | Simulated power | Error | N | Approximate power | Simulated power | Error |
| | | | | | (i) $\beta_1 = 0.3$ | | | | | | | |
| 93 | 0.8027 | 0.8003 | −0.0024 | 94 | 0.8039 | 0.8060 | 0.0021 | 97 | 0.8001 | 0.7985 | −0.0016 |
| 124 | 0.9020 | 0.9065 | 0.0045 | 125 | 0.9017 | 0.9009 | −0.0008 | 131 | 0.9012 | 0.9083 | 0.0071 |
| 152 | 0.9500 | 0.9550 | 0.0050 | 154 | 0.9506 | 0.9522 | 0.0016 | 162 | 0.9500 | 0.9537 | 0.0037 |
| | | | | | (ii) $\beta_1 = 0.4$ | | | | | | | |
| 55 | 0.8058 | 0.8080 | 0.0022 | 55 | 0.8006 | 0.7942 | −0.0064 | 59 | 0.8012 | 0.8079 | 0.0067 |
| 73 | 0.9036 | 0.8987 | −0.0049 | 74 | 0.9031 | 0.9105 | 0.0074 | 80 | 0.9013 | 0.9122 | 0.0109 |
| 89 | 0.9503 | 0.9511 | 0.0008 | 91 | 0.9512 | 0.9481 | −0.0031 | 100 | 0.9510 | 0.9562 | 0.0052 |
| | | | | | (iii) $\beta_1 = 0.5$ | | | | | | | |
| 37 | 0.8050 | 0.8097 | 0.0047 | 38 | 0.8079 | 0.8036 | −0.0043 | 42 | 0.8083 | 0.8186 | 0.0103 |
| 49 | 0.9032 | 0.8996 | −0.0036 | 50 | 0.9022 | 0.9059 | 0.0037 | 57 | 0.9037 | 0.9181 | 0.0144 |
| 60 | 0.9510 | 0.9530 | 0.0020 | 62 | 0.9521 | 0.9563 | 0.0042 | 71 | 0.9505 | 0.9651 | 0.0146 |

The values of $\mu_\Delta$ associated with the three distributions are 0.09, 0.16 and 0.25 for the three coefficient values $\beta_1 = 0.3, 0.4$ and 0.5, respectively.

TABLE 2.
Calculated sample sizes, approximate powers and simulated powers of the proposed method for random simple regression models $Y = \beta_0 + X\beta_1 + \varepsilon$ with standardized Poisson predictor ($p = c = 1$ and $\alpha = 0.05$).

| | Poisson(9) with kurtosis = 3.11 | | | | Poisson(4) with kurtosis = 3.25 | | | | Poisson(1) with kurtosis = 4 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| N | Approximate power | Simulated power | Error | N | Approximate power | Simulated power | Error | N | Approximate power | Simulated power | Error | |
| | | | | | (i) $\beta_1 = 0.3$ | | | | | | | |
| 92 | 0.8004 | 0.7986 | −0.0018 | 93 | 0.8042 | 0.7978 | −0.0064 | 93 | 0.8015 | 0.7952 | −0.0063 | |
| 123 | 0.9014 | 0.9051 | 0.0037 | 123 | 0.9010 | 0.8987 | −0.0023 | 124 | 0.9010 | 0.9036 | 0.0026 | |
| 151 | 0.9500 | 0.9481 | −0.0019 | 152 | 0.9509 | 0.9494 | −0.0015 | 153 | 0.9505 | 0.9500 | −0.0005 | |
| | | | | | (ii) $\beta_1 = 0.4$ | | | | | | | |
| 54 | 0.8019 | 0.8001 | −0.0018 | 54 | 0.8010 | 0.8021 | 0.0011 | 55 | 0.8037 | 0.8032 | −0.0005 | |
| 72 | 0.9027 | 0.9038 | 0.0011 | 72 | 0.9019 | 0.9050 | 0.0031 | 73 | 0.9019 | 0.9000 | −0.0019 | |
| 88 | 0.9503 | 0.9506 | 0.0003 | 89 | 0.9519 | 0.9526 | 0.0007 | 90 | 0.9511 | 0.9526 | 0.0015 | |
| | | | | | (iii) $\beta_1 = 0.5$ | | | | | | | |
| 37 | 0.8103 | 0.8051 | −0.0052 | 37 | 0.8090 | 0.8098 | 0.0008 | 37 | 0.8019 | 0.8040 | 0.0021 | |
| 48 | 0.9018 | 0.9033 | 0.0015 | 48 | 0.9006 | 0.9012 | 0.0006 | 49 | 0.9005 | 0.9080 | 0.0075 | |
| 59 | 0.9511 | 0.9510 | −0.0001 | 59 | 0.9503 | 0.9495 | −0.0008 | 61 | 0.9520 | 0.9547 | 0.0027 | |

The values of $\mu_\Delta$ associated with the three distributions are 0.09, 0.16 and 0.25 for the three coefficient values $\beta_1 = 0.3$, 0.4 and 0.5, respectively.

Also, $\mathbf{\Psi}$ can be expressed as

$$\mathbf{\Psi} = \begin{bmatrix} \Psi_1 & \Psi_2 & \Psi_3 \\ \Psi_2 & \Psi_4 & \Psi_5 \\ \Psi_3 & \Psi_5 & \Psi_6 \end{bmatrix} = E[\mathbf{H} \otimes \mathbf{H}] = E \begin{bmatrix} X^2\mathbf{H} & XZ\mathbf{H} & X(XZ-\rho)\mathbf{H} \\ XZ\mathbf{H} & Z^2\mathbf{H} & Z(XZ-\rho)\mathbf{H} \\ X(XZ-\rho)\mathbf{H} & Z(XZ-\rho)\mathbf{H} & (XZ-\rho)^2\mathbf{H} \end{bmatrix},$$

where

$$\mathbf{\Psi}_1 = \begin{bmatrix} 3 & 3\rho & 0 \\ 3\rho & 1+2\rho^2 & 0 \\ 0 & 0 & 3+7\rho^2 \end{bmatrix}, \qquad \mathbf{\Psi}_2 = \begin{bmatrix} 3\rho & 1+2\rho^2 & 0 \\ 1+2\rho^2 & 3\rho & 0 \\ 0 & 0 & 7\rho+3\rho^3 \end{bmatrix},$$

$$\mathbf{\Psi}_3 = \begin{bmatrix} 0 & 0 & 3+7\rho^2 \\ 0 & 0 & 7\rho+3\rho^3 \\ 3+7\rho^2 & 7\rho+3\rho^3 & 0 \end{bmatrix}, \qquad \mathbf{\Psi}_4 = \begin{bmatrix} 1+2\rho^2 & 3\rho & 0 \\ 3\rho & 3 & 0 \\ 0 & 0 & 3+7\rho^2 \end{bmatrix},$$

$$\mathbf{\Psi}_5 = \begin{bmatrix} 0 & 0 & 7\rho+3\rho^3 \\ 0 & 0 & 3+7\rho^2 \\ 7\rho+3\rho^3 & 3+7\rho^2 & 0 \end{bmatrix}$$

and

$$\mathbf{\Psi}_6 = \begin{bmatrix} 3+7\rho^2 & 7\rho+3\rho^3 & 0 \\ 7\rho+3\rho^3 & 3+7\rho^2 & 0 \\ 0 & 0 & 9+42\rho^2+9\rho^4 \end{bmatrix}.$$

The evaluations of $\mathbf{\Psi}$ for this interaction model are more involved than those in the second model. The conditional distribution properties of $Z$ given $X$ and high-order moments of a standard normal distribution ($E[X^6] = 15$ and $E[X^8] = 105$) are required to carry out the calculations. For illustration, the coefficient parameters are set as $(\beta_0, \beta_1, \beta_2, \beta_3) = (1, 0.1, 0.3, 0.25)$. Two hypothesis tests are investigated in this numerical demonstration. They are the tests of overall effects ($c = 3$) and interaction effect ($c = 1$) with the null hypotheses $H_0 : \beta_1 = \beta_2 = \beta_3 = 0$ and $H_0 : \beta_3 = 0$, respectively. In a similar fashion, the proposed approach is employed to perform the sample size and the corresponding approximate power calculations for testing the specified hypothesis with significance level $\alpha = 0.05$ and nominal power (0.80, 0.90, 0.95). These numerical results are presented in Table 3 for three different values of $\rho = 0.3$, 0.5 and 0.7.

It is important to note that the major analytical justification considered here applies large-sample approximation to the distribution of the $F^*$ statistic. In order to assess the finite-sample accuracy of the proposed approach, simulation studies are conducted next. With given sample size and model configuration, an estimate of the true power or simulated power is then computed through simulation of 10,000 replicate data sets. For each replicate, $N$ sets of explanatory variables are generated from the selected distribution. These values in turn determine the mean responses for generating $N$ normal outcomes with the underlying regression model. Then the test statistic is computed and the simulated power is the proportion of the 10,000 replicates whose $F^*$ test statistic values exceed the critical value $F_{c,N-p-1,\alpha}$. The adequacy of the proposed sample size formula is determined by the difference (simulated power–approximate power) between the simulated power and approximate power specified above. All calculations are performed using programs written with SAS/IML (SAS Institute, 2003). Detailed numerical results of the simulation studies are reported in Tables 1–2 and 3 for the two models, respectively. For the simple regression models, under the standardized gamma(9, 1) predictor variable situation with $\beta_1 = 0.3$, the simulated powers are 0.8003, 0.9065 and 0.9550 for the three different power levels = 0.80, 0.90 and 0.95, respectively. Thus, the differences or errors between simulated

TABLE 3.
Calculated sample sizes, approximate powers and simulated powers of the proposed method for random multiple regression model $Y = \beta_0 + X\beta_1 + Z\beta_2 + XZ\beta_3 + \varepsilon$ ($p = 3$, $\alpha = 0.05$) with standard normal error.

| | Test of overall effects ($c = 3$) | | | | Test of interaction effect ($c = 1$) | | |
|---|---|---|---|---|---|---|---|
| $N$ | Approximate power | Simulated power | Error | $N$ | Approximate power | Simulated power | Error |
| | | | (i) $\rho = 0.3$ | | | | |
| 70 | 0.8029 | 0.8067 | 0.0038 | 127 | 0.8013 | 0.7970 | −0.0043 |
| 91 | 0.9002 | 0.9032 | 0.0030 | 171 | 0.9010 | 0.8985 | −0.0025 |
| 111 | 0.9502 | 0.9558 | 0.0056 | 212 | 0.9503 | 0.9494 | −0.0009 |
| | | | (ii) $\rho = 0.5$ | | | | |
| 65 | 0.8049 | 0.8015 | −0.0034 | 114 | 0.8012 | 0.7930 | −0.0082 |
| 85 | 0.9017 | 0.9103 | 0.0086 | 154 | 0.9007 | 0.8933 | −0.0074 |
| 104 | 0.9508 | 0.9557 | 0.0049 | 192 | 0.9505 | 0.9515 | 0.0010 |
| | | | (iii) $\rho = 0.7$ | | | | |
| 60 | 0.8067 | 0.8148 | 0.0081 | 99 | 0.8010 | 0.7982 | −0.0028 |
| 79 | 0.9028 | 0.9144 | 0.0116 | 135 | 0.9015 | 0.9010 | −0.0005 |
| 97 | 0.9511 | 0.9583 | 0.0072 | 169 | 0.9510 | 0.9535 | 0.0025 |

The values of $\mu_\Delta$ associated with the two tests are $(0.1861, 0.0681)$, $(0.2081, 0.0781)$ and $(0.2351, 0.0931)$ for the three correlation values $\rho = 0.3$, $0.5$, and $0.7$, respectively.

powers and approximate powers are $0.8003 - 0.8027 = -0.0024$, $0.9065 - 0.9020 = 0.0045$ and $0.9550 - 0.9500 = 0.0050$, respectively. Similarly, all other results in Tables 1–3 are obtained.

Examination of Table 1 shows that the absolute errors associated with the standardized gamma$(9, 1)$ and gamma$(4, 1)$ predictor distributions do not exceed 0.01 for all combinations of $\beta_1$ and power levels. However, the results for the standardized gamma$(1, 1)$ distribution vary with the value of $\beta_1$ and power levels. Specifically, the cases associated with $\beta_1 = 0.3$ are less sensitive to the influence of the outsized skewness 2 and kurtosis 9 of the standardized gamma$(1, 1)$ distribution than those for $\beta_1 = 0.4$ and $\beta_1 = 0.5$. Obviously, the errors 0.0103, 0.0144 and 0.0146 for the three power levels of $\beta_1 = 0.5$ are greater than 0.01. Hence, the accuracy of the large-sample power approximation is not as satisfactory as other circumstances for strongly skewed gamma distributions, especially for small samples. Nonetheless, the relative performance of the proposed method for Poisson distributions in Table 2 is excellent even for the most skewed Poisson$(1)$ distribution with comparatively small sample sizes.

With respect to the second model, the results in Table 3 suggest that there is a close agreement between the simulated power and the approximate power because the absolute errors are less than 0.01. The only exception is 0.0116 which is associated with the test of overall effects for correlation $\rho = 0.7$. Overall, the accuracy of the proposed approach increases slightly with the sample size, and varies marginally with the model configurations. According to these findings, the performance of the proposed method appears to be excellent for the range of random regression specifications considered here. It is important to note that, in the context of moderated multiple regression, the test for the existence of moderator effect is examined by the significance of regression coefficient $\beta_3$ of the cross-product term in the simple interaction model. As noted in the review by Aguinis et al. (2005), however, it has been widely recognized that moderated multiple regression analyses have suffered low statistical power in detecting moderator effects. Therefore, the proposed power formula can be employed to determine the minimum sample size required for testing the hypothesis $H_0 : \beta_3 = 0$ with specified model configurations, significance level, and nominal power.

TABLE 4.
Calculated sample sizes, approximate powers and simulated powers of the proposed method for random multiple regression model $Y = \beta_0 + X\beta_1 + Z\beta_2 + XZ\beta_3 + \varepsilon$ ($p = 3$, $\alpha = 0.05$) with standardized uniform error.

| | Test of overall effects ($c = 3$) | | | | | | Test of interaction effect ($c = 1$) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $N$ | Simulated $\alpha$ | Error | Approximate power | Simulated power | Error | $N$ | Simulated $\alpha$ | Error | Approximate power | Simulated power | Error |
| | | | | | (i) $\rho = 0.3$ | | | | | | |
| 70 | 0.0472 | −0.0028 | 0.8029 | 0.7999 | −0.0030 | 127 | 0.0490 | −0.0010 | 0.8013 | 0.7833 | −0.0180 |
| 91 | 0.0500 | 0.0000 | 0.9002 | 0.9047 | 0.0045 | 171 | 0.0484 | −0.0016 | 0.9010 | 0.8903 | −0.0107 |
| 111 | 0.0468 | −0.0032 | 0.9502 | 0.9554 | 0.0052 | 212 | 0.0472 | −0.0028 | 0.9503 | 0.9523 | 0.0020 |
| | | | | | (ii) $\rho = 0.5$ | | | | | | |
| 65 | 0.0457 | −0.0043 | 0.8049 | 0.8139 | 0.0090 | 114 | 0.0471 | −0.0029 | 0.8012 | 0.7887 | −0.0125 |
| 85 | 0.0508 | 0.0008 | 0.9017 | 0.9110 | 0.0093 | 154 | 0.0426 | −0.0064 | 0.9007 | 0.9044 | 0.0037 |
| 104 | 0.0459 | −0.0041 | 0.9508 | 0.9612 | 0.0104 | 192 | 0.0484 | −0.0016 | 0.9505 | 0.9469 | −0.0036 |
| | | | | | (iii) $\rho = 0.7$ | | | | | | |
| 60 | 0.0499 | −0.0001 | 0.8067 | 0.8123 | 0.0056 | 99 | 0.0504 | 0.0004 | 0.8010 | 0.7912 | −0.0098 |
| 79 | 0.0484 | −0.0016 | 0.9028 | 0.9155 | 0.0127 | 135 | 0.0501 | 0.0001 | 0.9015 | 0.8994 | −0.0021 |
| 97 | 0.0476 | −0.0024 | 0.9511 | 0.9638 | 0.0127 | 169 | 0.0510 | 0.0010 | 0.9510 | 0.9516 | 0.0006 |

TABLE 5.

Calculated sample sizes, approximate powers and simulated powers of the proposed method for random multiple regression model $Y = \beta_0 + X\beta_1 + Z\beta_2 + XZ\beta_3 + \varepsilon$ ($p = 3$, $\alpha = 0.05$) with standardized gamma error.

| | Test of overall effects ($c = 3$) | | | | | | Test of interaction effect ($c = 1$) | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| $N$ | Simulated $\alpha$ | Error | Approximate power | Simulated power | Error | $N$ | Simulated $\alpha$ | Error | Approximate power | Simulated power | Error |
| | | | | | (i) $\rho = 0.3$ | | | | | | |
| 70 | 0.0507 | 0.0007 | 0.8029 | 0.8115 | 0.0086 | 127 | 0.0509 | 0.0009 | 0.8013 | 0.7936 | −0.0077 |
| 91 | 0.0510 | 0.0010 | 0.9002 | 0.9062 | 0.0060 | 171 | 0.0480 | −0.0020 | 0.9010 | 0.8958 | −0.0052 |
| 111 | 0.0501 | 0.0001 | 0.9502 | 0.9551 | 0.0049 | 212 | 0.0497 | −0.0003 | 0.9503 | 0.9512 | 0.0009 |
| | | | | | (ii) $\rho = 0.5$ | | | | | | |
| 65 | 0.0514 | 0.0014 | 0.8049 | 0.8051 | 0.0002 | 114 | 0.0511 | 0.0011 | 0.8012 | 0.7862 | −0.0150 |
| 85 | 0.0490 | −0.0010 | 0.9017 | 0.9051 | 0.0034 | 154 | 0.0478 | −0.0022 | 0.9007 | 0.9013 | 0.0006 |
| 104 | 0.0522 | 0.0022 | 0.9508 | 0.9573 | 0.0065 | 192 | 0.0489 | −0.0011 | 0.9505 | 0.9556 | 0.0051 |
| | | | | | (iii) $\rho = 0.7$ | | | | | | |
| 60 | 0.0503 | 0.0003 | 0.8067 | 0.8122 | 0.0055 | 99 | 0.0495 | −0.0005 | 0.8010 | 0.7900 | −0.0110 |
| 79 | 0.0511 | 0.0011 | 0.9028 | 0.9140 | 0.0112 | 135 | 0.0485 | −0.0015 | 0.9015 | 0.9048 | 0.0033 |
| 97 | 0.0502 | 0.0002 | 0.9511 | 0.9625 | 0.0114 | 169 | 0.0501 | 0.0001 | 0.9510 | 0.9559 | 0.0049 |

As pointed out by a referee, Anderson (1999) showed that the coefficients estimator within the multivariate multiple regression framework has an asymptotic multinormal distribution when the errors and predictors are mutually independently distributed, irrespective of whether they are normal. Since the multiple regression model considered here is a special case of the multivariate multiple regression, the asymptotic normality property of the multiple regression coefficient estimator can readily be established. However, Anderson (1999) did not explicitly discuss the distribution of the associated $F$ statistic and power calculation. It is important to note that the conditional distribution of the $F^*$ statistic given in (6) is no longer necessarily an exact $F$ distribution. Nonetheless, the proposed large-sample approximation for the conditional normal regression model given in (7–9) can be applied to the situation of nonnormal errors and predictors. To examine the robust issues of the proposed method against the extra complication of nonnormal errors, we have conducted a numerical evaluation for the simple interaction model: $Y = \beta_0 + X\beta_1 + Z\beta_2 + XZ\beta_3 + \varepsilon$, where $\varepsilon$ has a standardized uniform$(0, 1)$ or standardized gamma$(5, 1)$ distribution. For ease of exposition, the parameter settings are the same as those in Table 3, and the results are presented in Tables 4 and 5. The simulated Type I error rate and power are compared with the nominal $\alpha = 0.05$ and power level, respectively. Notably, all the absolute errors between simulated $\alpha$ and nominal value 0.05 are less than 0.01. In addition, the discrepancies between simulated powers and approximate powers calculated with the proposed power function (9) are slightly larger than those in Table 3 with standard normal errors. However, the performance seems completely acceptable, given the many unknowns in study planning. Therefore, the suggested procedures for conditional normal regression models with arbitrary distribution configurations for explanatory variables are not seriously affected by mild departures from the normality assumption of errors.

## 4. Conclusions

Procedures for power and sample size determinations in fixed regression models have been developed for years but none seems to have provided a comprehensive treatment or guideline for the calculations of power and sample sizes in the framework of random regression models. Within the context of random regression models, the current results are mainly under the situation of the normality assumption for explanatory variables. A natural generalization to incorporate other distributions of explanatory variables is essential to researchers for performing power and sample size calculations in practice. This paper discusses a feasible solution to this issue by providing both theoretical justification and numerical examination for the proposed unified approach. With this direct extension, one can perform power and sample size calculations in multiple regression models with any discrete and/or continuous distributions of explanatory variables. The remarkable performance for power and sample size calculations reveals that the proposed method may find useful applications in subsequent random regression analysis. Notably, the suggested methodology is applicable for the prominent moderated multiple regression models containing a continuous predictor and moderator.

### References

Aguinis, H., Beaty, J.C., Boik, R.J., & Pierce, C.A. (2005). Effect size and power in assessing moderating effects of categorical variables using multiple regression: A 30-year review. *Journal of Applied Psychology*, *90*, 94–107.

Aiken, L.S., & West, S.G. (1991). *Multiple regression: Testing and interpreting interactions*. Thousand Oaks: Sage.

Algina, J., & Olejnik, S. (2000). Determining sample size for accurate estimation of the squared multiple correlation coefficient. *Multivariate Behavioral Research*, *35*, 119–137.

Anderson, T.W. (1999). Asymptotic distribution of the reduced rank regression estimator under general conditions. *Annals of Statistics*, *27*, 1141–1154.

Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale: Lawrence Erlbaum.

Gatsonis, C., & Sampson, A.R. (1989). Multiple correlation: Exact power and sample size calculations. *Psychological Bulletin*, *106*, 516–524.

Kelley, K., & Maxwell, S.E. (2003). Sample size for multiple regression: Obtaining regression coefficients that are accurate, not simply significant. *Psychological Methods*, *8*, 305–321.

Mendoza, J.L., & Stafford, K.L. (2001). Confidence interval, power calculation, and sample size estimation for the squared multiple correlation coefficient under the fixed and random regression models: A computer program and useful standard tables. *Educational and Psychological Measurement*, *61*, 650–667.

Muirhead, R.J. (1982). *Aspects of multivariate statistical theory*. New York: Wiley.

Raudenbush, S.W., & Liu, X. (2000). Statistical power and optimal design for multisite randomized trials. *Psychological Methods*, *5*, 199–213.

Raudenbush, S.W., & Liu, X. (2001). Effects of study duration, frequency of observation, and sample size on power in studies of group differences in polynomial change. *Psychological Methods*, *6*, 387–401.

Rencher, A.C. (2000). *Linear models in statistics*. New York: Wiley.

Sampson, A.R. (1974). A tale of two regressions. *Journal of the American Statistical Association*, *69*, 682–689.

SAS Institute (2003). *SAS/IML user's guide, Version 8*. Cary, NC: SAS Institute Inc.

Shieh, G. (2003). A comparative study of power and sample size calculations for multivariate general linear models. *Multivariate Behavioral Research*, *38*, 285–307.

Shieh, G. (2005). Power and sample size calculations for multivariate linear models with random explanatory variables. *Psychometrika*, *70*, 347–358.

Shieh, G. (2006). Exact interval estimation, power calculation and sample size determination in normal correlation analysis. *Psychometrika*, *71*, 529–540.

Steiger, J.H., & Fouladi, R.T. (1992). R2: A computer program for interval estimation, power calculations, sample size estimation, and hypothesis testing in multiple regression. *Behavioral Research Methods, Instruments, and Computers*, *24*, 581–582.