1 **Effects of categorization method, regression type, and variable distribution on the**

2 **inflation of Type-I error rate when categorizing a confounding variable**

3 **Running head:** Categorized confounders and Type-I error

4

5 Jean-Louis Barnwell-Ménard[1], Qing Li[2], and Alan A. Cohen[2]*

6

7 [1]Department of Economics, University of Sherbrooke, Sherbrooke, QC, Canada

8

9 [2]Department of Family Medicine, University of Sherbrooke, Sherbrooke, QC, Canada

10

11 *Corresponding author: Groupe de recherche PRIMUS, Faculté de Médecine et Sciences de la Santé,

12 CHUS-Fleurimont, 3001 12[e] Ave Nord, Sherbrooke, QC J1H 5N5 Canada Alan.Cohen@USherbrooke.ca,

13 (819) 346-1110 x12589, (819) 820-6419 (fax)

14

15

16

17    **Abstract**

18    The loss of signal associated with categorizing a continuous variable is well known, and

19    previous studies have demonstrated that this can lead to an inflation of Type-I error when

20    the categorized variable is a confounder in a regression analysis estimating the effect of

21    an exposure on an outcome. However, it is not known how the Type-I error may vary

22    under different circumstances, including logistic versus linear regression, different

23    distributions of the confounder, and different categorization methods. Here we

24    analytically quantified the effect of categorization, and then performed a series of 9600

25    Monte Carlo simulations to estimate the Type-I error inflation associated with

26    categorization of a confounder under different regression scenarios. We show that Type-I

27    error is unacceptably high (>10% in most scenarios, and often 100%). The only exception

28    was when the variable categorized was a continuous mixture proxy for a genuinely

29    dichotomous latent variable, where both the continuous proxy and the categorized

30    variable are error-ridden proxies for the dichotomous latent variable. As expected, error

31    inflation was also higher with larger sample size, fewer categories, and stronger

32    associations between the confounder and the exposure or outcome. We provide online

33    tools that can help researchers estimate the potential error inflation and understand how

34    serious a problem this is.

35

36    **Keywords:** Type-I error, confounding, categorization, dichotomization, simulation,

37    distribution

## Introduction

39      Researchers and clinicians in epidemiologic and medical studies often categorize

40   continuous variables for purposes of facilitating the interpretability of results [1]

41   (common examples include age, body-mass index, socio-economic status, and levels of

42   blood biomarkers). The unnecessary use of categorical variables has been criticized by

43   many for the potential increase in statistical bias and the loss of information [2-17], but

44   use of categorized continuous variables is still standard practice in the epidemiologic and

45   medical literature [18]. There is a consensus among statisticians that statistical tools

46   treating variables as continuous (e.g. with non-parametric or spline regressions) are

47   preferred and more robust when the latent trend is not easily captured by classical

48   parametric models [2, 7, 17]. Such tools are, however, more complex to apply and

49   interpret for clinicians, which might be a reason for the continued abundant use of

50   categorized data in epidemiological publications.

51      A specific situation prominent in epidemiologic and medical research where

52   categorized continuous variables are regularly used is for control variables (confounding

53   variables) in regression models when assessing the potential impact of an exposure (risk

54   factor, independent variable of interest) on an outcome (dependent variable).

55   Confounding variables are defined here as variables that are associated with both the

56   exposure of interest and the outcome of interest, but which are not affected by either

57   variable [19]. Unlike a categorization of the exposure or outcome variables, which can

58   lead to an inflation of Type-II error [20], categorization of a confounding variable can

59   lead to increased residual confounding, i.e., effects of confounding variables that are

60   unmeasured and thus not accounted for in the model. Such residual confounding

61    generally results in the detection of spurious relationships between the exposure and the

62    outcome, and thus false rejection of null hypotheses (inflated Type-I error) because the

63    model does not replicate perfectly the statistical relationship between the confounding

64    variable and the concerned variables in models [17]. Austin & Brunner [7] assessed the

65    influence such methodology has on the statistical performance of models under the

66    hypothesis of normal variable distributions and logistic regression. They demonstrated

67    important residual confounding sufficient to suggest that researchers may often falsely

68    detect a potential association between an exposure and an outcome.

69          Quantiles and clinical cut-offs are the most common methods for categorizing

70    continuous confounding variables [18]. Clinicians and epidemiologists frequently study

71    variables with various distribution shapes and select their cut-offs (i.e., through a

72    categorization method) in order to minimize the loss of information or to group similar

73    observations. In spite of the common categorization methodologies, little is known about

74    how cut-off selection, variable distributions, or type of regression model (linear, logistic)

75    might affect the statistical bias and robustness of the results induced by the categorization

76    of confounding variables.

77          Because unnecessary categorization is such a rampant problem, it is important to

78    understand what factors contribute to greater error inflation when categorizing, and to

79    quantify error inflation under different scenarios. The ability to quantify error inflation

80    could become a tool to force researchers to consider more carefully the consequences of

81    categorization on their conclusions. In this paper, we assessed how generalizable the

82    conclusions of Austin & Brunner [7] were across a wide range of realistic data analysis

83    scenarios, and whether there might be some cases where the implications of

84   categorization were particularly severe. We simulated the rate of falsely rejecting the true

85   value of the coefficient relating an exposure to an outcome (the Type-I error) under

86   different scenarios where a confounding variable is categorized. In addition, we

87   mathematically show the effect of categorization for the case of linear regression. We

88   have also developed a statistical application available on the web allowing easy

89   estimation of the Type-I error rate under different categorization algorithms for varying

90   statistical hypotheses.

# Mathematical Derivation

92        The categorization of a confounding variable generates measurement error with

93   respect to the original variable. We recapitulate this effect with the following

94   mathematical derivation in the case of linear regression because it is possible to get a

95   closed-form expression of the asymptotic bias which allows seeing immediately the

96   determinants. The literature origin of the effect is well exposed in [21], as well as the risk

97   for measurement error in general for different error sources and regression scenarios.

98   Under these circumstances the estimators are asymptotically biased, affecting the

99   estimated values, the confidence intervals and consequently the Type-1 error rate.

100       For individual $i$ the model is

101
$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$$

102   The confounding variable $x_{2i}$ is categorized into

$$x_{2i}^c = x_{2i} + u_i$$

103   where the superscript "c" denotes "categorical". Under the assumption that $\mathrm{E}(x_{1i} x_{2i}) \neq$

104   $0$ , and since the value of $x_{2i}$ decides which category the individual $i$ goes into, we know

105    that $cov(x_{2i}, u_i) \neq 0$ and hence $cov(x_{1i}, u_i) \neq 0$. The term $u_i$ is the difference between

106    $x_{2i}^c$ and $x_{2i}$ for individual $i$, i.e. the measurement error introduced by categorization. Note

107    that $E(u_i) \neq 0$, and in addition, the measurement error $u_i$ is correlated with the true

108    value $x_{2i}$ which is different from classical measurement error models; this case was

109    discussed in [22] and the correlation between $u_i$ and $x_{2i}$ has an influence on the

110    analytical expression of the bias, making the bias more unpredictable. We make the

111    classical assumptions of orthogonality for linear regression, i.e. $E(x_{1i}\varepsilon_i) = 0$ and

112    $E(x_{2i}\varepsilon_i) = 0$, which leads to $E(u_i\varepsilon_i) = 0$ and $E(x_{2i}^c\varepsilon_i) = 0$. Plugging $x_{2i}^c$ into the

113    regression gives

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}^c + \varepsilon_i - \beta_2 u_i$$

114    Letting $v_i = \varepsilon_i - \beta_2 u_i$, we get $cov(x_{1i}, v_i) = -\beta_1\beta_2 cov(x_{1i}, u_i) \neq 0$ and

115    $cov(x_{2i}^c, v_i) = -\beta_2^2 cov(x_{2i}^c, u_i) \neq 0$. In matrix form, defining

116
$$\beta := \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} \quad y := \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \quad v := \begin{bmatrix} \varepsilon_1 - \beta_2 u_1 \\ \vdots \\ \varepsilon_N - \beta_2 u_N \end{bmatrix} \quad X := \begin{bmatrix} 1 & x_{11} & x_{21}^c \\ \vdots & \vdots & \vdots \\ 1 & x_{1N} & x_{2N}^c \end{bmatrix}$$

117    for a sample with size N, we can write the regression as $y = X\beta + v$. Hence, the classical

118    least squares estimator $\hat{\beta}$ converges asymptotically to

$$\plim_{N\to\infty} \hat{\beta} = \plim_{N\to\infty}(X'X)^{-1} X'y = \beta + \plim_{N\to\infty}(X'X)^{-1} X'v$$

119    and the asymptotic bias generated by categorization (or by introducing measurement

120    error, in a broader sense) is

$$bias(\hat{\beta}) = \plim_{N\to\infty}(X'X)^{-1} X'v = \plim_{N\to\infty}\left(\frac{1}{N}X'X\right)^{-1} \plim_{N\to\infty}\frac{1}{N}X'v$$

121    Since,

$$\plim_{N\to\infty}\left(\frac{1}{N}X'X\right)^{-1} = \plim_{N\to\infty}\begin{bmatrix} 1 & \frac{\sum_{i=1}^{N}x_{1i}}{N} & \frac{\sum_{i=1}^{N}x_{2i}^{c}}{N} \\ \frac{\sum_{i=1}^{N}x_{1i}}{N} & \frac{\sum_{i=1}^{N}x_{1i}^{2}}{N} & \frac{\sum_{i=1}^{N}x_{1i}x_{2i}^{c}}{N} \\ \frac{\sum_{i=1}^{N}x_{2i}^{c}}{N} & \frac{\sum_{i=1}^{N}x_{1i}x_{2i}^{c}}{N} & \frac{\sum_{i=1}^{N}x_{2i}^{c2}}{N} \end{bmatrix}^{-1}$$

122    and

$$\plim_{N\to\infty}\frac{1}{N}X'v = \plim_{N\to\infty}\begin{bmatrix} \frac{\sum_{i=1}^{N}(\varepsilon_i - \beta_2 u_i)}{N} \\ \frac{\sum_{i=1}^{N}[x_{1i}(\varepsilon_i - \beta_2 u_i)]}{N} \\ \frac{\sum_{i=1}^{N}[x_{2i}^{c}(\varepsilon_i - \beta_2 u_i)]}{N} \end{bmatrix}$$

123    using Slutsky theorem and the property $\plim\limits_{N\to\infty}\frac{\sum_{i=1}^{N}x_i}{N} = E(x_i)$, we get

$$bias(\hat{\beta}) = \plim_{N\to\infty}\left(\frac{1}{N}X'X\right)^{-1}\plim_{N\to\infty}\frac{1}{N}X'v$$

$$= \begin{bmatrix} 1 & E(x_{1i}) & E(x_{2i}^{c}) \\ E(x_{1i}) & E(x_{1i}^{2}) & E(x_{1i}x_{2i}^{c}) \\ E(x_{2i}^{c}) & E(x_{1i}x_{2i}^{c}) & E(x_{2i}^{c2}) \end{bmatrix}^{-1}\begin{bmatrix} E(\varepsilon_i) - \beta_2 E(u_i) \\ E(x_{1i}\varepsilon_i) - \beta_2 E(x_{1i}u_i) \\ E(x_{2i}^{c}\varepsilon_i) - \beta_2 E(x_{2i}^{c}u_i) \end{bmatrix}$$

124    The last matrix product leads to a $3 \times 1$ matrix where the three elements correspond to

125    the asymptotic biases of $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$, respectively. With the assumptions $E(\varepsilon_i) = $

126    $E(x_{1i}\varepsilon_i) = E(x_{2i}\varepsilon_i) = 0$ and some basic calculations we get the following expression for

127    the second element of the matrix, namely $bias(\hat{\beta}_1)$ which is equal to

128    $\frac{\beta_2[E(u_i)E(x_{1i}x_{2i}^{c})E(x_{2i}^{c}) - E(x_{1i}u_i)E^2(x_{2i}^{c}) + E(x_{1i}u_i)E(x_{2i}^{c2}) - E(u_i)E(x_{1i})E(x_{2i}^{c2}) + E(x_{1i})E(x_{2i}^{c})E(x_{2i}^{c}u_i) - E(x_{1i}x_{2i}^{c})E(x_{2i}^{c}u_i)]}{E^2(x_{1i}x_{2i}^{c}) - 2E(x_{1i})E(x_{2i}^{c})E(x_{1i}x_{2i}^{c}) + E(x_{1i}^{2})\left(E^2(x_{2i}^{c}) - E(x_{2i}^{c2})\right) + E^2(x_{1i})E(x_{2i}^{c2})}$

129         The last expression, which shares similarities to the bias expression found by [21],

130    finds that the asymptotic bias of $\hat{\beta}_1$ depends on the value of $\beta_2$, but does not depend on

131    the value of $\beta_1$ itself. Also, the bias is affected by the first and second order moments

132    related to $x_{2i}^{c}$ and $u_i$ which depend on the method of categorization as well as the

7

133    distributions of the original variables. The analytical expression is non-linear in the

134    relevant moments and so it is not easy to characterize the effect of a single determinant

135    (e.g. method of categorization, data distribution, number of categories, etc.); in practice,

136    the expression will become even more complex when adding additional regressors, but in

137    a word, it is the introduction of measurement error that creates the bias, whatever the

138    nature of the original variables is. Importantly, the complexity of this expression shows

139    that the precise magnitude of the bias is not easily predictable. Simulations in the

140    following sections give intuitive results in different cases.


141    ## Methods

142    Our simulations were modeled on the approach of Austin & Brunner [7]. We

143    simulated data under the general scenario of the following regression model:

$$Y = \beta_0 + X_1\beta_1 + X_2\beta_2 + \nu$$

144    where $Y$ is an outcome of interest, $X_1$ is an exposure whose relationship to $Y$ we would

145    like to assess, and $X_2$ is a potential confounding variable which is available in continuous

146    format but which is categorized for analysis. The true values of $\beta_0$ and $\beta_1$ are assumed to

147    be zero (i.e., $X_1$ has no direct effect on $Y$, since we wish to evaluate the Type-I error), and

148    $\beta_2$ has a specified positive value. Parameters which were allowed to vary included (a)

149    type of regression model (linear versus logistic), (b) distribution of the underlying

150    confounder ($X_2$), (c) the covariance between $X_1$ and $X_2$, (d) $\beta_2$, (e) the method for

151    categorizing $X_2$ when continuous, (f) the number of categories into which $X_2$ is divided,

152    and (g) the sample size of the simulated data set.

153

154 ## Data generation

155 *Continuous confounding variable*

156    An exposure ($X_1$, assumed to be independent of the outcome) and a continuous

157 confounding variable ($X_2$) were generated with three different processes in order to assess

158 the confounding variable under (1) normal, (2) log-normal or (3) bimodal distributions.

159

160    (1)    The first process ("normal") for the generation of a normal exposure $X_{1,n}$ and a

161          normal confounding variable $X_{2,n}$ used a bivariate normal distribution of size $N$

162          with mean $\mu = (0,0)$ and covariance matrix $\Omega = \begin{pmatrix} 1 & \sigma_{1,2} \\ \sigma_{1,2} & 1 \end{pmatrix}$ where $\sigma_{1,2}$, the

163          confounder-exposure covariance, ranged from 0 to 0.9 in increments of 0.1.

164    (2)    The second process ("log-normal") for the generation of a normal exposure $X_{1,l}$

165          and a log-normal confounding variable $X_{2,l}$ was obtained with the exponential

166          transformation of a normal confounding variable $X_{2,n}$ generated with process

167          (1). The average sampled kurtosis of $X_{2,l}$, out of 1000 samples with arbitrary

168          covariance specification and sample size of 2000, was 62.04, with a 95%

169          bootstrap confidence interval for the sample kurtosis average ranging from

170          56.88 to 67.18, and the average skewness was 5.22, with a 95% bootstrap

171          confidence interval for the sample skewness average ranging from 5.10 to 5.35.

172    (3)    The third process ("bimodal") for the generation of a normal exposure $X_{1,b}$ and

173          a potentially correlated bimodal confounding variable $X_{2,b}$ was based on the

174          separate simulation of two groups of data, *I* and *II*, representing each of the

175          modes in $X_{2,b}$ (i.e., $X_{2,b}^1$ and $X_{2,b}^2$) along with their paired values in $X_{1,b}$ (i.e.,

176    $X_{1,b}^1$ and $X_{1,b}^2$). $X_{1,b}^1$ and $X_{2,b}^1$ ($I$) were simulated from a bivariate normal

177    distribution of size $N_1$ with mean $\mu_1 = (0,0)$ and covariance matrix $\Omega_1 =$

178    $\begin{pmatrix} 1 & \sigma_{1,2} \\ \sigma_{1,2} & 1 \end{pmatrix}$. $X_{1,b}^2$ and $X_{2,b}^2$ ($II$) were simulated from a bivariate normal

179    distribution of size $N_2$ with mean $\mu_2 = (0, U(3,4))$ and covariance matrix

180    $\Omega_2 = \begin{pmatrix} 1 & \sigma_{1,2} \\ \sigma_{1,2} & U(4,9) \end{pmatrix}$. Once the four variables were simulated, $X_{1,b}$ was

181    generated as the union of $X_{1,b}^1$ and $X_{1,b}^2$, and $X_{2,b}$ was generated as the union of

182    $X_{2,b}^1$ and $X_{2,b}^2$, keeping their relative orders so as to maintain the pairing of

183    values and thus the correlation. $\sigma_{1,2}$ ranged from 0 to 0.9 by increments of 0.1.

184    $U$ represents the uniform distribution (e.g., min = 3 and max = 4). Total sample

185    size $N = N_1 + N_2$, but $N_1 \neq N_2$. This simulation method allowed $X_{1,b}$ and $X_{2,b}$ to

186    covary at level $\sigma_{1,2}$ even while $X_{1,b}$ represents a unimodal normal distribution

187    and $X_{2,b}$ represents a bimodal distribution generated as a mixture of two normal

188    distributions with different means and variances. The average sampled kurtosis

189    of $X_{2,b}$, out of 1000 samples with arbitrary covariance specification and sample

190    size of 2000, was 5.29, with a 95% bootstrap confidence interval for the sample

191    kurtosis average ranging from 5.25 to 5.33, and the average skewness was 1.45,

192    with a 95% bootstrap confidence interval for the sample skewness average

193    ranging from 1.44 to 1.46.

194

195    *Proxy variable for a dichotomous underlying confounder*

196        In addition to the three above scenarios featuring continuous confounding

197    variables with different distributions, we simulated a fourth scenario in which the true

198 confounding variable is dichotomous but researchers only observe a continuous proxy.

199 This corresponds in reality to using blood glucose level as a continuous proxy for

200 underlying diabetes state, or to using a sex steroid level to assign sex when true sex is

201 unknown. If the true confounder is the underlying dichotomous variable, we might ask to

202 what extent we can categorize the proxy in order to better approach the true confounder

203 (supposing that it is known that proxy is not the true confounder). The exposure ($X_1$), the

204 proxy confounding variable ($X_2$), and the underlying dichotomous confounding variable

205 ($X_3$) were generated with the following fourth process ("mixture distribution"):

206    (4)   $X_{1,d}$ (the normal exposure) and $X_{2,d}$ (the bimodal proxy confounding variable)

207    were generated identically as in process (3), the mixture of two multivariate

208    normal distributions (*I*) and (*II*) of size $N = N_1 + N_2$. $X_{3,d}$ (the underlying

209    dichotomous confounding variable) is a dummy variable taking the following

210    values:

$$\begin{cases} if\ X_{2i,d} \in (I):\ X_{3i,d} = 0 \\ if\ X_{2i,d} \in (II):\ X_{3i,d} = 1 \end{cases}$$

211

212 *Outcome variable (continuous confounder)*

213    Once the unrelated exposure $X_{1,(n,l\ or\ b)}$ and the confounding variable $X_{2,(n,l\ or\ b)}$

214 were generated, the outcome (independent) variable $Y_{(n,l\ or\ b)}$ could be obtained using (a)

215 a logistic model or (b) a linear model for its generation in the following procedure:

216    (a) Logistic model:

217 $$logit(p_i) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}, \text{ for } i = 1, 2, \dots, N$$

218 $$\text{where, } p_i = \frac{\exp(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i})}{\exp(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}) + 1}, \text{ for } i = 1, 2, \dots, N$$

219 $$y_i \sim Binomial(p_i), \text{ for } i = 1, 2, \ldots, N$$

220     $X_1$ denotes the *de facto* unrelated exposure and $X_2$ is the confounding variable

221 correlated with $X_1$ and the outcome $Y$. The logistic model was assessed for five

222 confounder-outcome association scenarios:

223 $$\beta_0 = 0, \beta_1 = 0 \text{ and } \beta_2 = (0.2, 0.5, 1, 2, 3)$$

224 where the range for the predetermined values of $\beta_2$ was based upon Austin & Brunner [7]

225 modeling scenarios, with the addition of 0.2 and 2 for generality purposes.

226     (b) The linear model:

227 $$y_i = x_{2i} + \varepsilon_i, \text{ for } i = 1, 2, \ldots, N$$

228     where $x_{2i}$ is treated as a constant and $\varepsilon \sim N(0, \sigma^2)$. Therefore, as $\sigma^2$ increases, we

229 expect a lower predictive power of the outcome variable ($y$) by the confounding variable

230 ($x_2$), which correspond to the idea of a decreasing value of $\beta_2$ in the logistic model. The

231 linear model was assessed for five confounder-outcome association scenarios:

232 $$\sigma^2 = (9.95, 3.17, 1.73, 1.02, 0.48)$$

233     The values for $\sigma^2$ were chosen empirically via simulations to correspond as

234 closely as possible to values of $\beta_2$ for a residual confounding effect equivalent to those

235 used in (a) for the logistic model.

236

237 *Outcome variable (dichotomous underlying confounder)*

238     Once the unrelated exposure $X_{1,d}$, the bimodally distributed proxy representing

239 the dichotomous confounder $X_{2,d}$ and the underlying dichotomous confounder $X_{3,d}$ were

240 generated, the outcome (independent) variable $Y_d$ could be obtained using both models

241 (a) and (b), with the sole difference here that $X_2$ is replaced by $X_3$ in the generating

242 procedure. Therefore, the logistic and linear model become, respectively:

243     (a)

244
$$logit(p_i) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{3i}, \text{ for } i = 1, 2, \dots, N$$

245
$$where, p_i = \frac{\exp(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{3i})}{\exp(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{3i}) + 1}, \text{ for } i = 1, 2, \dots, N$$

246
$$y_i \sim Binomial(p_i), \text{ for } i = 1, 2, \dots, N$$

247     (b)

248
$$y_i = x_{3i} + \varepsilon_i, \text{ for } i = 1, 2, \dots, N$$

249 Both models use the same confounder-outcome association scenarios as with the

250 continuous confounding variable modeling. $X_{3,d}$ is used only to generate $Y_d$; once $Y_d$ is

251 generated, the dichotomous variable $X_{3,d}$ is represented by its proxy variable $X_{2,d}$

252 (bimodally distributed) in the model estimating the Type-I error rates. The mixture and

253 bimodal distributions thus differ only in that the outcome is determined directly by the

254 continuous bimodal confounder in the bimodal distribution, but is determined by the

255 underlying dichotomous variable in the mixture distribution.

256 ## Categorization algorithms

257     The Type-I error for the true null hypothesis of the unrelated exposure was assessed

258 with the confounding or proxy variable categorized in two, three, four and five

259 categories, or kept continuous for comparison. The confounding variable was categorized

260 using two different methods: (A) quantile and (B) maximized $R^2$.

261     (A) The first method consists in dividing the confounding variable into quantiles, i.e.

262 separating the sorted $x_{2i}$, for $i = 1, 2, \dots, N$, in groups with an equal number of

263    observations. This method is arguably the most frequently used in practice, and was

264    explained in detail by [7].

265       (B) The second method finds category cut-offs that optimize the linear fit of a

266    continuous variable by the same categorized variable. The optimal cut-offs define the

267    categories that maximize the adjusted $R^2$ of the following preliminary linear model

268    (which differs from models (1) and (2)):

$$X_2 = \alpha_0 + X_2^c \alpha_1 + \mu$$

269    where $X_2$ corresponds to the continuous variable and $X_2^c$ to the same categorized variable.

270    The optimal cut-offs are found using a linear optimization function for a one cut-off

271    search and a non-linear optimization function for a 2-4 cut-off search (with the optimize

272    and optim functions in R). We applied this method with 1, 2, 3 and 4 cut-offs, giving a

273    categorized confounding variable ($X_2^c$) with two, three, four and five categories

274    respectively.

275

## Simulations of Type I error

277    Using the framework above, we had eight independent parameters that could be adjusted

278    in the simulations: (1) Underlying confounder distribution (4 levels: normal, log-normal,

279    bimodal, or dichotomous); (2) Regression type (2 levels: logistic or linear); (3)

280    Categorization method (2 levels: quantile or maximized $R^2$); (4) Category number (4

281    levels: 2-5); (5) Confounder-exposure covariance (10 levels: 0 - 0.9 in increments of 0.1);

282    (6) Confounder-outcome association (5 levels: $\beta_2$ or $\sigma^2$); (7) Sample size of the

283    simulated study (3 levels: 100, 500 or 2000); (8) Number of Monte Carlo iterations per

284    scenario (1 level used: 1000 iterations). Monte Carlo simulations were performed for all

14

285    9600 combinations of these parameters. For each parameter combination, we calculated

286    the Type-I error rate as the percentage of the 1000 Monte Carlo iterations in which the *p*-

287    value of the following parameter significance t-test:

$$H_0: \beta_1 = 0$$
$$H_1: \beta_1 \neq 0$$

288    was less than α=0.05, i.e. falsely rejecting the true null hypothesis of no relationship

289    between the exposure and the outcome with a confidence level of 95%.

290

## Summarizing results

291    

292    Because of the large number of results generated by these simulations, we used three

293    parallel methods to summarize our results. First, we conducted linear regression models

294    on the database of simulation results, modeling the Type-I error rate among the thousand

295    iterations as a function of the seven varying parameters included in the models. We also

296    stratified and included interactions as necessary. Presentation of results is stratified

297    between the normal, log-normal and bimodal confounder distributions on the one hand

298    and the mixture distribution on the other, given that the latter is a special case with

299    particular properties. In order to show the approximate magnitude of effects, we present

300    results of regression models as if effects were linear and additive (e.g., change in Type-I

301    error for each change of 0.1 in *σ*), though clearly this is not strictly true and should not be

302    taken overly literally. Second, we developed an online interactive interface that allows

303    users to choose parameters of interest and generate figures similar to those shown here in

304    order to graphically examine several parameters and their interactions,

305   https://usherbrookeprimus.shinyapps.io/resultsApp/. Third, we present a selection of

306   results from the online tool as figures to illustrate key points.


## Results

### Performance of categorization methods

309   For a normally distributed confounding variable, the quantile and maximized $R^2$ methods

310   provided essentially identical categories. For a log-normally distributed confounding

311   variable, the maximized $R^2$ method provided cut-offs that were substantially further

312   toward the tail of the distribution than those chosen by the quantile method. For the

313   bimodal distribution ($X_{2,b}$ or $X_{2,d}$), the maximized $R^2$ method was substantially better at

314   separating the two modes near the bottom of the trough (Figure 1), especially with only 2

315   categories (referred to hereafter as "optimal categorization").

316

### Type-I error: Normal, log-normal, and bimodal confounder distributions

318   The results from our simulations demonstrated a substantial inflation of the Type-I error

319   rate for detecting an effect of the unrelated exposure ($X_1$) on the outcome ($Y$) when the

320   confounder ($X_2$) was categorized, except when the confounder was very weakly

321   associated with either the exposure or the outcome (Table 1, top). As expected, Type-I

322   error rate always increased as the correlation between the exposure and the confounder ($\sigma$)

323   increased, with approximately 5.6% additional error for each increment of 0.1 in $\sigma$

324   (Figure 2). Type-I error rate decreased monotonically as the number of categories

325   increased, with approximately 7.9% fewer errors for each additional category added

326    (Figure 2). Accordingly, a confounder categorized in five categories obtained a lower

327    Type-I error rate compared to a confounder with two, three and four categories, with the

328    exception of the bimodal distribution (3) categorized with the quantile method (A) under

329    the linear model (a), where 3 categories minimized the type one error rate. Each

330    additional 100 added to sample size increased the Type-I error by about 0.96%, or about

331    14.4% higher rates with sample size =2000 than =500 (Figure 3). Additionally, there was

332    about a 8.8% increase in Type-I error for each additional increment of association

333    between the confounder ($X_2$) and the outcome ($Y$) (Figure 3). The quantile categorization

334    method obtained lower Type-one error rates (Figure 4) for the three distributions. The

335    distribution type did not express a clear pattern for minimizing the error rate.

336        In sum, under all scenarios, with the exception of a very weak confounder-outcome or

337    confounder-exposure association (where the addition of a confounding variable is not as

338    relevant), categorizing a continuous confounding variable substantially inflated the risk

339    for type-I error rate. Although it might seem intuitive to dichotomize a bimodal

340    confounder, we found that with the bimodal confounder distribution (3) and the

341    maximized $R^2$ categorization method (B) even an "optimal" categorization process

342    significantly inflated the type-I error rate, performing even worse than an arbitrary

343    categorization criterion such as with the quantile method.

344

345    **Type-I error: Dichotomous unmeasured confounder (mixture**

346    **distribution)**

347    With a dichotomous unmeasured confounder represented by a bimodal continuous proxy,

348    results were qualitatively similar to results under other distributions for sample size, the

349   strength of the confounder-exposure correlation ($\sigma$), and the strength of the confounder-

350   outcome association ($\beta_2$ or $\sigma^2$), and are not discussed further (Table 1, bottom). However,

351   inversed results were found for the number of categories and the categorization method

352   on the proxy variable (Figure 5). Two categories with the maximized $R^2$ method now

353   performed best, with worse results for three (4.5% more error), four-five categories (6%

354   more error), and the quantile method in general. The dichotomized proxy confounder

355   gave lower Type-I error rates than its continuous state, although its error rates were still

356   substantial. The maximized $R^2$ method performed better, with a 10% lower error rate,

357   though this effect was attenuated substantially with more than two categories: by 5% for

358   three categories and by 6% for four or five categories. In sum, the dichotomized proxy

359   confounder, representing a dichotomous underlying state, minimized the type-I error rate

360   and performed worst when left as continuous.

361

362   ## Online interactive results tool

363   For a further analysis of our results, we propose an interactive online application that

364   allows users to manipulate the different parameters used in this study to assess their

365   impact on the Type-I error rate, represented graphically. The application can be accessed

366   through: https://usherbrookeprimus.shinyapps.io/resultsApp/.

367   # Discussion

368       The results of these simulations confirm and expand the general conclusions of other

369   authors: categorizing a continuous confounding variable leads to a surprisingly large and

370   robust inflation of the Type-I error rate, nearly regardless of model parameters. Only with

371     a very weak association between the confounder and either the outcome or the exposure

372     (i.e., in the absence of a real confounding effect) did this inflation disappear; under many

373     realistic scenarios, the Type-I error was 100%. When applied across hundreds or

374     thousands of studies, even a small inflation of the Type-I error rate – from the expected

375     5% to, say, 10% – should have a large impact on our confidence in the results generated

376     by a body of literature, especially given the many other biases that tend to lead toward

377     false positive results [23]. The Type-I error rates observed here suggest the problem may

378     be much larger than this small inflation, given the pervasiveness of categorization of

379     important confounders such as age, socio-economic status, and many others.

380         We identified one highly specific case where categorization diminished the Type-I

381     error, and it is a case chosen specifically to be the exception that proves the rule. This

382     case is when the outcome (i.e., dependent) variable is determined not by the measured

383     confounding variable, but by an underlying dichotomous process for which the measured

384     confounder is a proxy. (Real-world examples might be using blood glucose level to

385     determine diabetes status, or identifying a patient's sex, when unknown, using levels of

386     steroid sex hormones, when it is diabetes status or sex rather than glucose level or

387     hormone level that affects the outcome.) Even in this case, categorization only reduced

388     Type-I error rate relative to the continuous proxy, and when the number of categories

389     corresponded to the number of underlying groups (i.e., 2). And even when all these

390     criteria were met, Type-I error was still substantially higher than the expected $\alpha=0.05$,

391     reaching error rates greater than 50% under some scenarios.

392         This counter-example is an example of the principle that all measurement error of a

393     confounding variable increases the risk of Type-I error [21]. In the case of the counter-

394   example, the true variable that should have been measured is the underlying (latent)

395   dichotomous variable, and using a continuous proxy introduces measurement error which

396   can be partially but not completely eliminated by dichotomizing the proxy. The

397   conditions for categorization are thus highly restrictive (and thus may never be met in

398   practice) – one would need to know *a priori* (a) that the continuous variable was a proxy

399   for a true categorical variable, (b) exactly how many underlying categories (sometimes

400   referred to as "latent classes") there were, and (c) that it was the underlying variable

401   rather than the proxy that was the true confounder. Because confounding variables are

402   generally measured with some measurement error to begin with, the effect of the

403   categorization is over and above the Type-I error inflation due to the original

404   measurement error [21].

405      The details of our results offer some guidance as to which situations present the

406   greatest Type-I error inflation due to categorization. Type-I error inflation is worse when

407   fewer categories are used.  Stronger associations between the confounder and either the

408   exposure or the outcome rapidly increase the Type-I error. Counter-intuitively, large

409   sample size also makes the problem worse, increasing the power to detect the residual

410   confounding present when a confounder is imperfectly measured. All of these effects are

411   quite large.

412      The effects of the confounder's distribution and the categorization method are more

413   nuanced. When the confounder has a normal or log-normal distribution, the maximized

414   $R^2$ method performs worse than the standard quantile method. However, maximized $R^2$

415   performed substantially better than quantile under the mixture distribution, a special case

416   when two categories also performs better than more categories. This case demonstrates

417    the limits of simulations for inferring the precise error rate in cases where particular but

418    unknown data generating processes are likely to underlie data structure. In theory, it

419    might be possible to use *a priori* clinical knowledge to slightly diminish the Type-I error

420    rate by choosing optimal cut-offs based on (a) the relationship between the confounder

421    and the outcome; (b) the relationship between the confounder and the exposure, and (c)

422    knowledge of underlying biological/ sociological/ psychological processes. In practice,

423    such *a priori* knowledge is unlikely to be sufficient. Our mathematical derivation of the

424    estimator bias shows substantial complexity in the interactions between such factors and

425    therefore how difficult the task of theoretically controlling for the introduction of

426    measurement error becomes. Traditional clinical cut-offs are unlikely to be valid, for

427    example, unless they approximate underlying biological thresholds, or unless there are

428    threshold effects in their relationships with the other variables. Also, we note that even

429    the best-case scenario for such dichotomization in our simulations still produced

430    substantial Type-I error; such error is unavoidable under the mixture distribution, where

431    the true confounder is unmeasured and an imperfect proxy is used. Even the use of a raw

432    continuous confounding variable in a regression model may sometimes be insufficient: if

433    the relationship of the confounder with the outcome is non-linear, there may still be

434    substantial residual confounding [24]. Quadratic regression, fractional polynomials [25],

435    non -parametric regression [26],  and splines are potential solutions to this problem.

436      All of which is to say that categorization is, in general, a conscious and unnecessary

437    introduction of measurement error. In lay terms, to drive the point home, categorizing a

438    continuous confounder is the equivalent of saying, "Hey, my study is pretty good, but

439    what it could really use is some measurement error. Why don't I categorize the

440 confounders? That way I will be essentially assured of detecting a positive result whether

441 or not one exists!" In order to help researchers understand the magnitude of the problem,

442 we propose a second interactive online application that allows the users, manually or with

443 the use of the quantile categorization method, to choose cut-offs and assess the probable

444 Type-I error rate of an unrelated exposure controlled for the given categorized

445 confounder. The user can also choose between the distributions and the models presented

446 in this study. The application can be accessed through:

447 https://usherbrookeprimus.shinyapps.io/simulationApp/. Our hope is that this tool will

448 allow many researchers to simulate a situation similar enough to their research question

449 that they get a sense of how bad the problem is likely to be.

450

451 **Acknowledgments**

456

457 **References**

458 1. Altman DG, and Royston P. The cost of dichotomising continuous variables.
459 *British Medical Journal* **2006**; **332**, pp. 1080, DOI: 10.1136/bmj.332.7549.1080.

460 2. Greenland S. Dose-response and trend analysis in epidemiology: alternatives to
461 categorical analysis. *Epidemiology* **1995**; **6(4)**, pp. 356–365, DOI:
462 10.1097/00001648-199507000-00005.

463　3.　Morabia A, Ten Have T. Controlling for Continuous Confounders with Non-
464　　　Gaussian Distribution in Epidemiologic Research. *Epidemiology* **2000**; **11(1)**, pp.
465　　　93–94, DOI: 10.2307/3703664.

466　4.　Greenland S. Avoiding power loss associated with categorization and ordinal
467　　　scores in dose-response and trend analysis. *Epidemiology* **1995**; **6(4)**, pp. 450–454,
468　　　DOI: 10.1097/00001648-199507000-00025.

469　5.　Bennette C, Vickers AJ. Against quantiles: categorization of continuous variables
470　　　in epidemiologic research, and its discontents. *BMC Medical Research*
471　　　*Methodology* **2012**; **12(1)**, pp. 21–25, DOI: 10.1186/1471-2288-12-21.

472　6.　Royston P, Altman DG, Sauerbrei W. Dichotomizing continuous predictors in
473　　　multiple regression: a bad idea. *Statistics in Medecine* **2006**; **25**, pp. 127–141,
474　　　DOI: 10.1002/sim.2331.

475　7.　Austin PC, Brunner J. Inflation of the type I error rate when a continuous
476　　　confounding variable is categorized in logistic regression analyses. *Statistics in*
477　　　*Medecine* **2004**; **23**, pp. 1159–1178, DOI: 10.1002/sim.1687.

478　8.　Altman DG. Categorising continuous variables. *British Journal of Cancer* **1991**;
479　　　**64**, p. 975, DOI: 10.1038/bjc.1991.441.

480　9.　Sigurdsson H, Baldetorp B, Borg A, Dalberg M, Ferno M, Killander D, Olsson H,
481　　　Ranstam J. Flow cytometry in primary breast cancer: improving the prognostic
482　　　value of the fraction of cells in the S-phase by optimal categorisation of cut-off
483　　　levels. *British Journal of Cancer* **1990B**; **62**, pp. 786–790.

484　10.　Taylor JMG, Yu M. Bias and efficiency loss due to categorising an explanatory
485　　　variable. *Journal of Multivariate Analysis* **2002**; **83**, pp. 248–263, DOI:
486　　　10.1006/jmva.2001.2045.

487　11.　Schellingerhout JM, Heymans MW, de Vet HCW, Koes BW, Verhagen AP.
488　　　Categorizing continuous variables resulted in different predictors in a prognostic
489　　　model for nonspecific neck pain. *Journal of Clininical Epidemiology* **2009**; **62(8)**,
490　　　pp. 868–874, DOI: 10.1016/j.jclinepi.2008.10.010.

491　12.　Lausen B, Schumacher M. Evaluating the effect of optimized cutoff values in the
492　　　assessment of prognostic factors. *Computational Statistics and Data Analysis*
493　　　**1996**; **21(3)**, pp. 307–326, DOI: 10.1016/0167-9473(95)00016-X.

494　13.　Froslie K, Roislien J, Laake P, Henriksen T, Qvigstad E, Veierod M.
495　　　Categorisation of continuous exposure variables revisited. A response to the
496　　　Hyperglycaemia and Adverse Pregnancy Outcome (HAPO) Study. *BMC Medical*
497　　　*Research Methodology* **2010**; **10(1)**, p. 103, DOI: 10.1186/1471-2288-10-103.

498  14.  Buettner P, Garbe C, Guggenmoos-Holzmann I. Problems in defining cutoff points
499      of continuous prognostic factors: Example of tumor thickness in primary
500      cutaneous melanoma. *Journal of Clininical Epidemiology* **1997**; **50(11)**, pp. 1201–
501      1210, DOI: 10.1016/S0895-4356(97)00155-8.

502  15.  Filardo G, Hamilton C, Hamman B, Ng HKT, Grayburn P. Categorizing BMI may
503      lead to biased results in studies investigating in-hospital mortality after isolated
504      CABG. *Journal of Clininical Epidemiology* **2007**; **60(11)**, pp. 1132–1139, DOI:
505      10.1016/j.jclinepi.2007.01.008.

506  16.  Groenwold RHH, Klungel OH, Altman DG, van der Graaf Y, Hoes AW, Moons
507      KGM. Adjustment for continuous confounders: an example of how to prevent
508      residual confounding. *Canadian Medical Association Journal* **2013**; **185(5)**, pp.
509      401–406, DOI: 10.1503/cmaj.120592.

510  17.  Brenner H, Blettner M. Controlling for Continuous Confounders in Epidemiologic
511      Research. *Epidemiology* **1997**; **8(4)**, pp. 429–434, DOI: 10.1097/00001648-
512      199707000-00014.

513  18.  Turner E, Dobson J, Pocock S. Categorisation of continuous risk factors in
514      epidemiological publications: a survey of current practice. *Epidemiologic
515      Perspectives and Innovation* **2010**; **7**, p. 9, DOI: 10.1186/1742-5573-7-9.

516  19.  Rothman KJ, Greenland S, Lash TL, *Modern Epidemiology* (3rd edn). Lippincott,
517      Williams & Wilkins: Philadelphia, 2008.

518  20.  Irwin J, McClelland G. Negative consequences of dichotomizing continuous
519      predictor variables. *Journal of Marketing Research* **2003**; **40(3)**, pp. 366-371, DOI:
520      10.1509/jmkr.40.3.366.19237.

521  21.  Brunner J, Austin PC. Inflation of Type I error rate in multiple regression when
522      independent variables are measured with error. *The Canadian Journal of Statistics*
523      **2009**; **37(1)**, pp. 33-46, DOI: 10.1002/cjs.10004.

524  22.  Li Q. Identifiability of mean-reverting measurement error with instrumental
525      variable. *Statistica Neerlandica* **2014**; **68(2)**, pp. 118–129, DOI:
526      10.1111/stan.12025.

527  23.  Ioannidis JPA. Why most published research findings are false. *PLoS Medecine*
528      **2005**; **2(8)**, p. e124, DOI: 10.1371/journal.pmed.0020124.

529  24.  Benedetti A, Abrahamowicz M. Using generalized additive models to reduce
530      residual confounding. *Statistics in Medecine* **2004**; **23**, pp. 3781–3801, DOI:
531      10.1002/sim.2073.

532    25.    Royston P, Ambler G, Sauerbrei W. The use of fractional polynomials to model
533           continuous risk variables in epidemiology. *International Journal of Epidemiology*
534           **1999**; **28(5)**, pp. 964–974, DOI: 10.1093/ije/28.5.964.

535    26.    Rosenberg PS, Katki H, Swanson CA, Brown LM, Wacholder S, Hoover RN.
536           Quantifying epidemiologic risk factors using non-parametric regression: model
537           selection remains the greatest challenge. *Statistics in Medecine* **2003**; **22**, pp. 3369-
538           3381, DOI: 10.1002/sim.1638.

539

540

541

542  **Table 1: Effects of model parameters on Type-1 error rate, modeled separately for (a) confounders**
543  **with normal, log-normal or bimodal continuous underlying distributions, or (b) confounders**
544  **with the dichotomous underlying distribution**

| | Beta | Std. Error | t-value | *p* |
|---|---|---|---|---|
| **Normal, log-normal or bimodal confounder** | | | | |
| Intercept | -0.67 | 0.018 | -38.1 | <0.0001 |
| # of categories (numeric variable) | -0.08 | 0.004 | -21.4 | <0.0001 |
| Confounder-exposure correlation | 0.56 | 0.008 | 67.1 | <0.0001 |
| Regression Type | | | | |
|    Logistic (ref) | 0 | - | - | - |
|    Linear | 0.09 | 0.005 | 19.3 | <0.0001 |
| Confounder-outcome association | 0.09 | 0.002 | 52.0 | <0.0001 |
| Sample size/100[a] | 0.10 | 0.002 | 49.1 | <0.0001 |
| Confounder distribution | | | | |
|    Normal (ref) | 0 | - | - | - |
|    Log-normal | 0.003 | 0.016 | 0.21 | 0.84 |
|    Bimodal | -0.16 | 0.016 | -10.0 | <0.0001 |
| Categorization method | | | | |
|    Quantile (ref) | 0 | - | - | - |
|    Max $R^2$ | -0.009 | 0.008 | -1.0 | 0.30 |
| Interaction: # Cat*Distribution[b] | | | | |
|    Log-normal | 0.009 | 0.005 | 1.7 | 0.10 |
|    Bimodal | 0.01 | 0.005 | 1.9 | 0.06 |
| Interaction: Distribution*Cat method[c] | | | | |
|    Log-normal* Max $R^2$ | 0.03 | 0.012 | 2.2 | 0.03 |
|    Bimodal* Max $R^2$ | 0.18 | 0.012 | 15.2 | <0.0001 |
| **Dichotomous confounder** | | | | |
| Intercept | -0.86 | 0.023 | -37.5 | <0.0001 |
| # of categories | | | | |
|    2 categories (ref) | 0 | - | - | - |
|    3 categories | 0.05 | 0.014 | 3.2 | 0.002 |
|    4 categories | 0.06 | 0.014 | 4.3 | <0.0001 |
|    5 categories | 0.07 | 0.014 | 4.6 | <0.0001 |
| Confounder-exposure correlation | 0.44 | 0.013 | 35.2 | <0.0001 |
| Regression Type | | | | |
|    Logistic (ref) | 0 | - | - | - |
|    Linear | 0.05 | 0.007 | 6.7 | <0.0001 |
| Confounder-outcome association | 0.10 | 0.003 | 38.5 | <0.0001 |
| Sample size/100[a] | 0.09 | 0.003 | 30.6 | <0.0001 |
| Categorization method | | | | |
|    Quantile (ref) | 0 | - | - | - |

| | | | | |
|---|---|---|---|---|
| Max $R^2$ | -0.10 | 0.014 | -7.02 | <0.0001 |
| Interaction: # Cat*Cat method= Max $R^2$ | | | | |
| 2 categories (ref) | 0 | - | - | - |
| 3 categories | 0.06 | 0.020 | 2.8 | 0.006 |
| 4 categories | 0.07 | 0.020 | 3.3 | 0.0009 |
| 5 categories | 0.07 | 0.020 | 3.4 | 0.0006 |

[a]This is the effect of the natural logarithm of the continuous sample size on the Type-I error rate.
[b]This is the increase in Type-I error rate per additional category under a log-normal and bimodal distribution.
[c]This is the increase in Type-I error rate with the max $R^2$ method under a log-normal and bimodal distribution.

545

**Figure legends**

547  Figure 1. Thresholds/cut-offs found by the quantiles (A) and maximized $R^2$ (B) methods

548  for 2 categories in a sample size of 2000 under the bimodal (3) distribution.

549

550  Figure 2. Type-I error rates for logistic (a) models with the confounding variable

551  continuous and in 2-5 categories, using the quantiles (A) method, under normal (1), log-

552  normal (2) and bimodal (3) underlying distributions, $\beta_2 = 2$ and sample size=2000.

553  Vertical lines represent 95% confidence intervals for the simulated Type-1 error rates

554  based on an N of 1000 simulations.

555

556  Figure 3. Type-I error rates for linear (b) models with $\sigma^2 = \{0.48, 1.02, 1.73, 3.17, 9.95\}$

557  and sample size=$\{100, 500, 2000\}$, using the maximized $R^2$ (B) method for the

558  confounding variable in 2 categories. Vertical lines represent 95% confidence intervals

559  for the simulated Type-1 error rates based on an N of 1000 simulations.

560

561  Figure 4. Type-I error rates for logistic (b) models with the confounding variable in 3

562  categories using the quantiles (A) and the maximized $R^2$ (B) methods under normal (1),

563  log-normal (2) and bimodal (3) underlying distributions, $\beta_2 = 2$ and sample size=2000.

564     Vertical lines represent 95% confidence intervals for the simulated Type-1 error rates

565     based on an N of 1000 simulations.
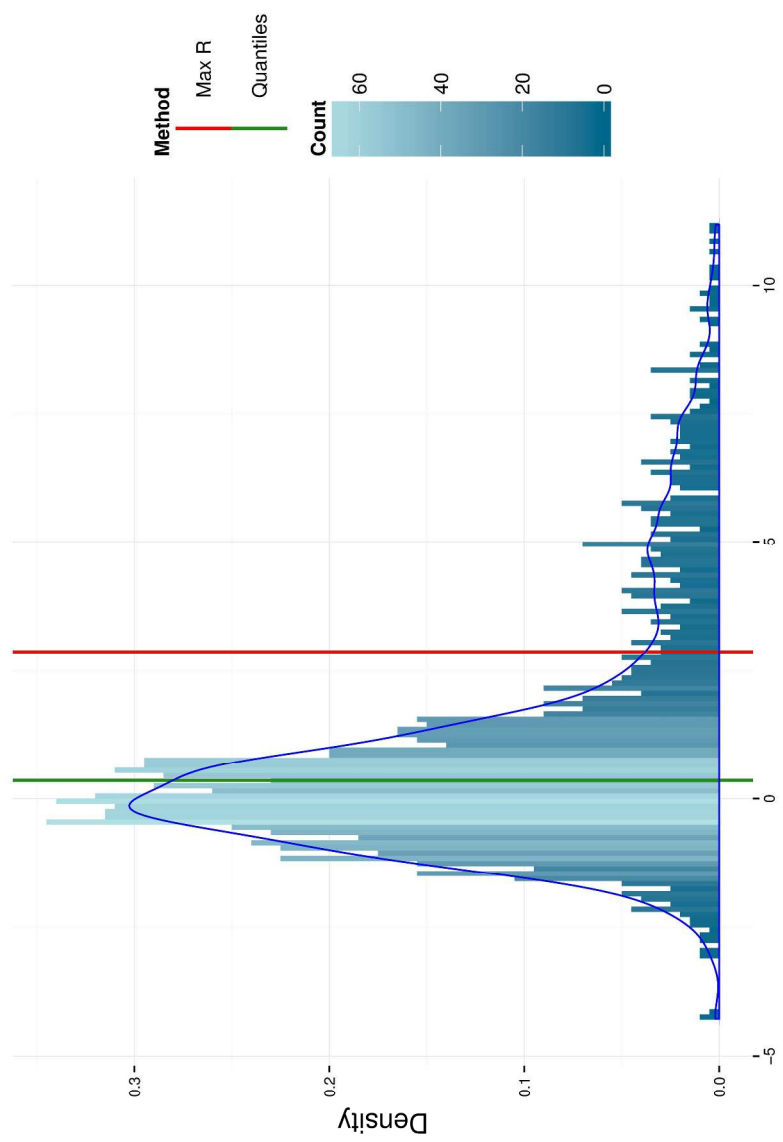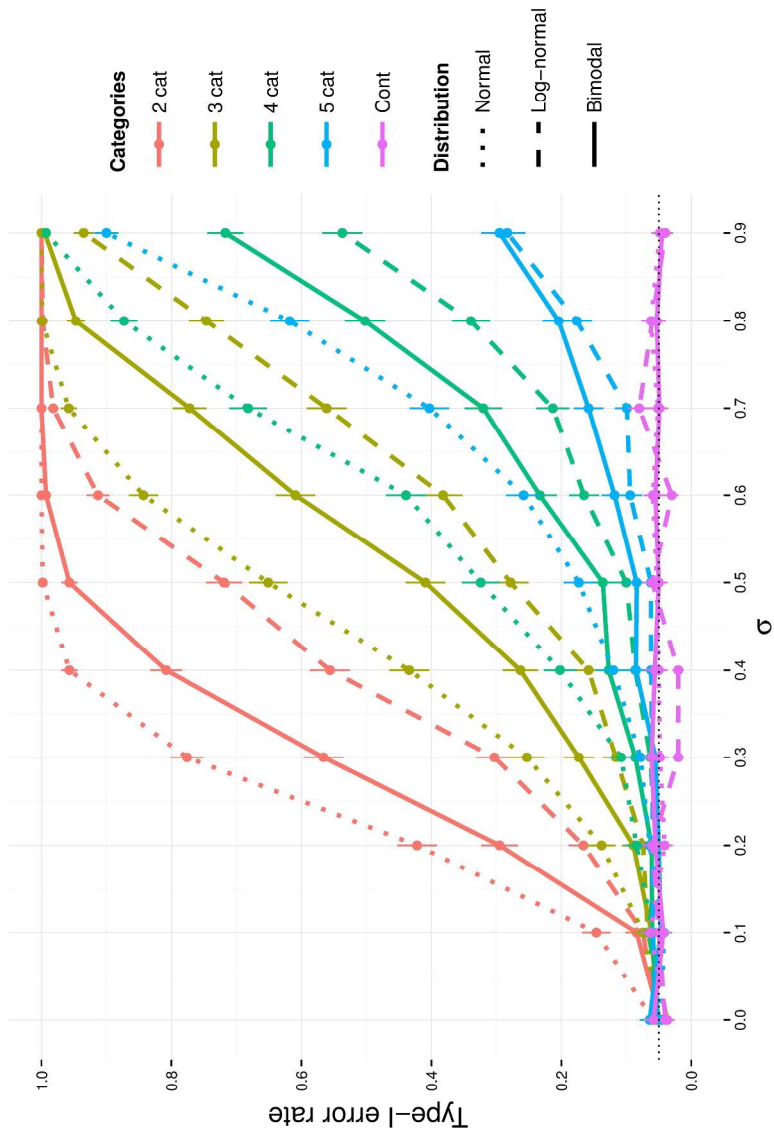
566

567     Figure 5. Type-I error rates for linear (b) models with the proxy variable continuous and

568     in 2-5 categories, using the quantiles (A) and the maximized $R^2$ (B) methods, under

569     dichotomous (4) underlying distribution, $\sigma^2 = 1.02$ and sample size=2000. Vertical lines

570     represent 95% confidence intervals for the simulated Type-1 error rates based on an N of

571     1000 simulations.

572

573

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
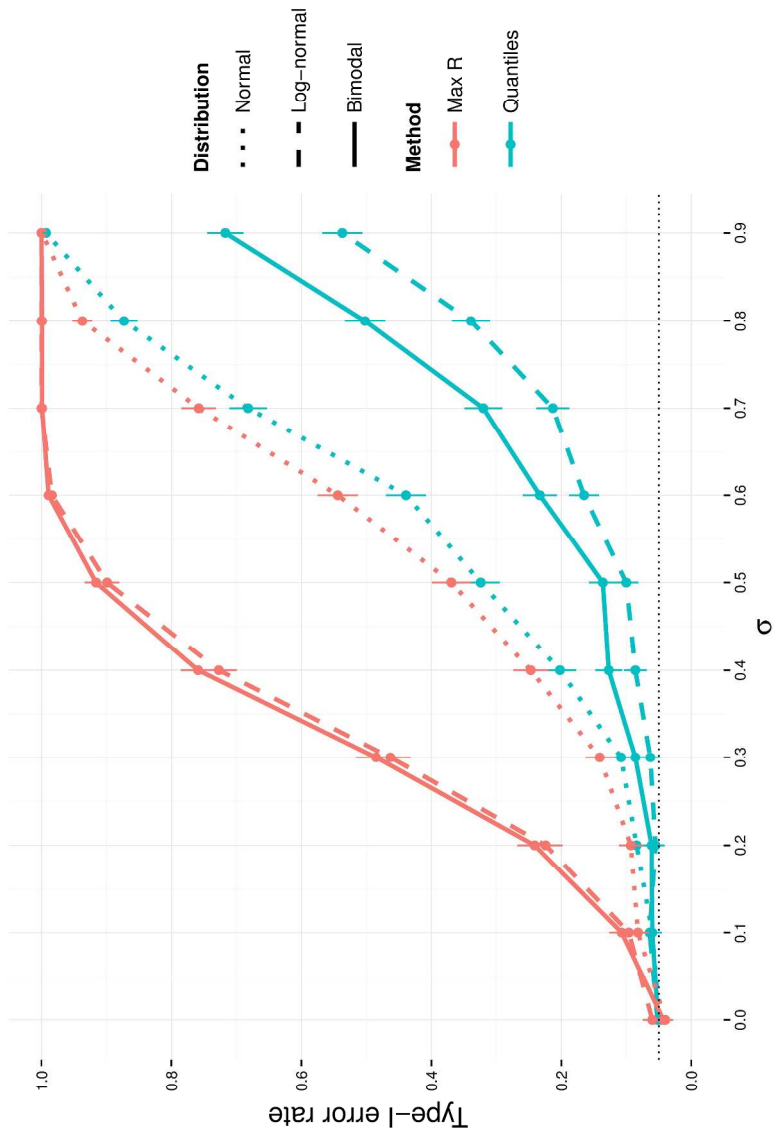48
49
50
51
52
53
54
55
56
57
58
59
60



279x361mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
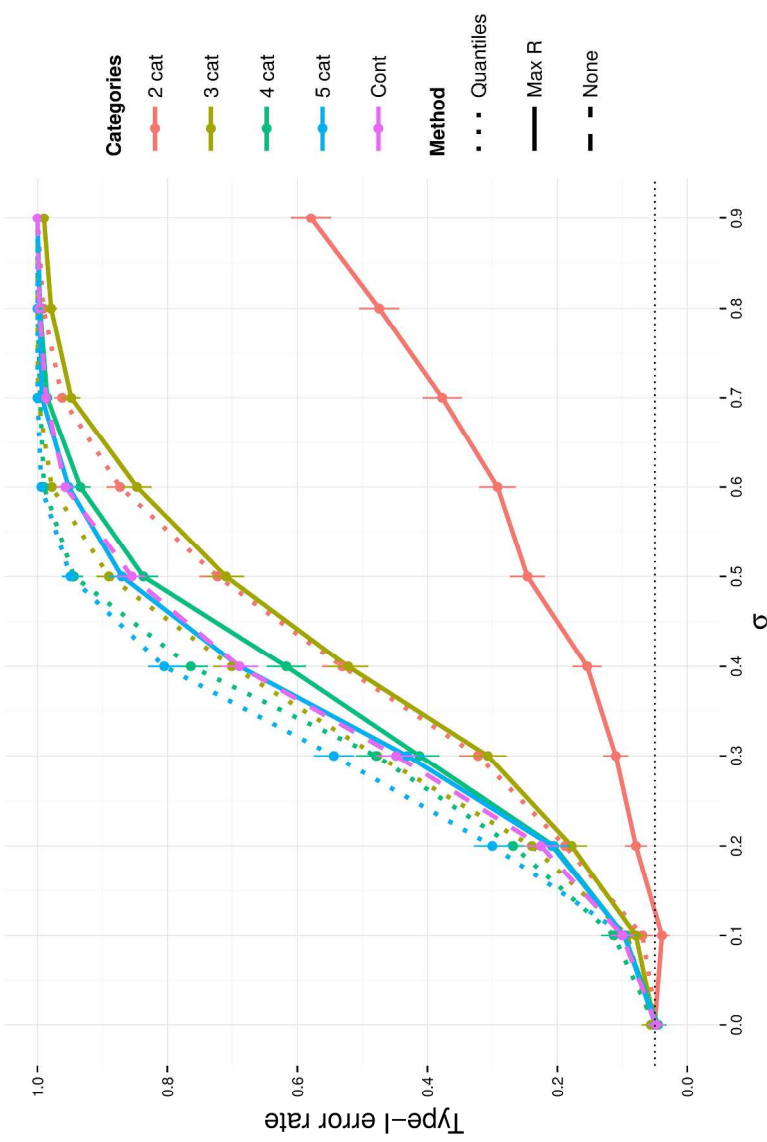44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



279x361mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



279x361mm (300 x 300 DPI)

279x361mm (300 x 300 DPI)

279x361mm (300 x 300 DPI)