# Consequences of Ignoring Clustering in Linear Regression

**Georgia Ntani** ( ✉ gn@mrc.soton.ac.uk )

MRC Lifecourse Epidemiology Unit   https://orcid.org/0000-0001-7481-6860

**Hazel Inskip**

MRC Lifecourse Epidemiology Unit, University of Southampton

**Clive Osmond**

MRC Lifecourse Epidemiology Unit, University of Southampton

**David Coggon**

MRC Lifecourse Epidemiology Unit, University of Southampton

1    **Consequences of ignoring clustering in linear regression**

2    Georgia Ntani[1,2], Hazel Inskip[1], Clive Osmond[1], David Coggon[1,2]

3    [1] Medical Research Council Lifecourse Epidemiology Unit, University of Southampton, United

4    Kingdom

5    [2] Medical Research Council Versus Arthritis Centre for Musculoskeletal Health and Work, Medical

6    Research Council Lifecourse Epidemiology Unit, University of Southampton, United Kingdom

7    **Abstract**

8    *Background*

9    Clustering of observations is a common phenomenon in epidemiological and clinical research.

10   Previous studies have highlighted the importance of using multilevel analysis to account for such

11   clustering, but in practice, methods ignoring clustering are often used. We used simulated data to

12   explore the circumstances in which failure to account for clustering in linear regression analysis could

13   lead to importantly erroneous conclusions.

14   *Methods*

15   We simulated data following the random-intercept model specification under different scenarios of

16   clustering of a continuous outcome and a single continuous or binary explanatory variable. We fitted

17   random-intercept (RI) and cluster-unadjusted ordinary least squares (OLS) models and compared the

18   derived estimates of effect, as quantified by regression coefficients, and their estimated precision. We

19   also assessed the extent to which coverage by 95% confidence intervals and rates of Type I error were

20   appropriate.

21   *Results*

22   We found that effects estimated from OLS linear regression models that ignored clustering were on

23   average unbiased. The precision of effect estimates from the OLS model was overestimated when

24   both the outcome and explanatory variable were continuous. By contrast, in linear regression with a

25  binary explanatory variable, in most circumstances, the precision of effects was somewhat

26  underestimated by the OLS model. The magnitude of bias, both in point estimates and their precision,

27  increased with greater clustering of the outcome variable, and was influenced also by the amount of

28  clustering in the explanatory variable. The cluster-unadjusted model resulted in poor coverage rates

29  by 95% confidence intervals and high rates of Type I error especially when the explanatory variable

30  was continuous.

31  *Conclusions*

32  In this study we identified situations in which an OLS regression model is more likely to affect

33  statistical inference, namely when the explanatory variable is continuous, and its intraclass correlation

34  coefficient is higher than 0.01. Situations in which statistical inference is less likely to be affected

35  have also been identified.

36  **Keywords:** Clustering, linear regression, random intercept model, consequences, simulation,

37  comparison, bias

**Introduction**

Clinical and epidemiological research often uses some form of regression analysis to explore the relationship of an outcome variable to one or more explanatory variables. In many cases, the study design is such that participants can be grouped into discrete, non-overlapping subsets (clusters), such that the outcome and/or explanatory variables vary less within than between clusters. This might occur, for example, in cluster-randomised controlled trials (with the units of randomisation defining clusters), or in a multi-centre observational study (the participants from each centre constituting a cluster). The extent to which a variable is "clustered" can be quantified by the intra-class correlation coefficient (ICC), which is defined as the ratio of its variance between clusters to its total variance (both between and within clusters) (1).

Clustering has implications for statistical inference from regression analysis if the outcome variable is clustered after the effects of all measured explanatory variables are taken into account. If allowance is not made for such clustering as part of the analysis, parameter estimates and/or their precision may be biased. This possibility can be demonstrated by a hypothetical study of hearing impairment and noise exposure, in which observations are made in four different cities (clusters), as illustrated in Figure 1. In this example, the effect of cumulative noise exposure on hearing impairment is the same within each city (i.e. the regression coefficient for hearing impairment on noise exposure is the same in each cluster) (Figure 1a). However, after allowance for noise exposure, hearing impairment differs by city, such that it varies more between the clusters than within them. An analysis that ignored this clustering would give a misleading estimate for the regression coefficient of hearing loss on noise exposure (Figure 1b). Moreover, even if the distribution of noise exposures in each city was similar, so that the regression coefficient was unbiased, its precision would be underestimated as it would have made no allowance for the differences between clusters (at the intercept) (Figure 1c).

Where, as in the example above, the number of clusters is small relative to the total number of participants in the study sample, a categorical variable that distinguishes clusters can be treated as an additional explanatory variable in the regression model (2). However, when the number of clusters is

64  larger, use of the cluster variable as an additional explanatory variable in the regression model can

65  seriously reduce the precision with which effects are estimated.  In such circumstances, an alternative

66  approach is to assume that cluster effects are randomly distributed with a mean and variance that can

67  be estimated from the data in the study sample. Random intercept models assume that the effects of

68  explanatory variables are the same across all clusters, but that the intercepts of regression lines differ

69  with a mean and variance which can be estimated from the study data, along with the effect estimates

70  of primary interest.  Random slope models assume that the effects of explanatory variables also differ

71  between clusters, with a mean and variance that can be estimated.

72  In recognition of the potential implications of clustering for statistical inference, there has been a

73  growth over recent years in the use of statistical techniques that allow for clustering (3). Nevertheless,

74  many studies still ignore clustering of observations (4-8). Recent systematic reviews have reported

75  that clustering was taken into account in only 21.5% of multicentre trials (9) and 47% of cluster

76  randomised trials (10). This may in part reflect computational challenges and statistical complexities

77  (11), but, perhaps because of a lack of clarity about the effects of ignoring clustering, authors have

78  omitted to discuss the limitations of their chosen analytical techniques.

79  Several studies have investigated implications of ignoring clustering in statistical inference, most

80  being based on analysis of real data (1, 12-19). To date, no study has systematically investigated the

81  extent to which bias can occur in effect estimates when clustering is ignored, the determinants of that

82  bias, or the exact consequences for the precision of estimates according to different distributions of

83  the explanatory variable and, in particular, the extent to which the explanatory variable varies within

84  as compared with between clusters.

85  The first aim of the research described in this paper was to assess in detail the implications for effect

86  estimates (regression coefficients), and their precision (characterised by standard errors (SEs)), when

87  a linear regression analysis exploring the relation of a continuous outcome variable to an explanatory

88  variable fails to account for clustering. The second aim was to describe rates of Type I error and

4

89    coverage by 95% confidence intervals in the same setting. These research questions were explored

90    through simulation studies.

91                                        INSERT FIGURE 1 HERE

92    **Figure 1**. Hypothetical relationship of hearing impairment to cumulative noise exposure in four cities.  Units for

93    noise exposure and hearing impairment have been specified arbitrarily for ease of presentation.  Data for each

94    city are distinguished by the shading of data points.  Cluster-specific regression lines are indicated, along with

95    the regression line for the full dataset when clustering is ignored (dotted red line), and that when adjustment is

96    made for cluster (solid blue line)

97    **Methods**

98    In the simplest case, in which there is a single explanatory variable, the ordinary least squares (OLS)

99    linear regression is specified by a model of the form:

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

-1-

100

101   For a continuous outcome and a single explanatory variable, the random intercept (RI) multi-level

102   model can be viewed as an extension of the OLS model, and is specified as:

$$y_{ij} = \beta_{0j} + \beta_1 x_{ij} + e_{ij}$$
$$= \beta_0 + \beta_1 x_{ij} + e_{ij} + u_j$$

-2-

103   where the index $i$ refers to the individual and the index $j$ to the cluster, and $\beta_{0j} = \beta_0 + u_j$, the

104   estimate of the intercept for cluster $j$. The term $u_j$ represents the error for cluster $j$ around the fixed

105   intercept value of $\beta_0$, and is assumed to be normally distributed with $u_j | x_{ij} \sim N(0, SD_u^2)$. The term $e_{ij}$

106   represents the additional error within the cluster, also referred to as the individual level error term,

107   with $e_{ij} | x_{ij}, u_j \sim N(0, SD_e^2)$.

108   As described in the introduction, ICC is a measure which characterises the extent to which the

109   outcome variable $y_{ij}$ is similar within clusters, given the distribution of the explanatory variable $x_{ij}$

5

110 (20). For a continuous outcome variable, and with the nomenclature used above, the ICC is defined as

111 $\text{ICC} = \frac{SD_u^2}{SD_u^2 + SD_e^2}$ (21).

112 To explore the study questions, simulated datasets were generated according to the assumptions of the

113 RI model. For each Monte Carlo simulation, both the number of clusters and the number of

114 observations per cluster were set to 100. For simplicity, the size of the effect of $x_{ij}$ on $y_{ij}$ was

115 arbitrarily set to 1 ($\beta_1 = 1$), and the average value of $y_{ij}$ when $x_{ij} = 0$ was arbitrarily set to 0 ($\beta_0 =$

116 0).

117 Separate simulation studies were generated for a continuous and a binary explanatory variable $x_{ij}$. To

118 set values $x_{ij}$ for the continuous explanatory variable in a cluster $j$, an individual level variable was

119 generated as $x_{0ij} \sim N(0,1)$, and a cluster-specific variable as $shift_j \sim N(0, SD_{shift}^2)$. The individual

120 level variable was then added to the cluster-specific shift, so that $x_{ij} = x_{0ij} + shift_j$. For a binary

121 explanatory variable $x_{ij}$, we set the prevalence in each cluster to be the sum of a constant (the same in

122 all clusters) set to 0.05, 0.1, 0.2 and 0.4 and a cluster-specific variable $shift_j \sim N(0, SD_{shift}^2)$. In both

123 cases, the corresponding values for the outcome variable $y_{ij}$ were generated according to equation -2-.

124 For this purpose, the individual-level error terms were drawn from a random standard normal

125 distribution $(N(0,1))$, and the cluster-level error terms were drawn from a random normal distribution

126 with mean zero and variance $SD_{u_j}^2$. Simulated data were generated for various different values for

127 $SD_{u_j}$ (0.0316, 0.05485, 0.1005, 0.1759, 0.3333 and 0.6547) chosen to give expected values for the

128 $ICC$ of 0.001, 0.003, 0.01, 0.03, 0.1 and 0.3 respectively, while $SD_{shift} \sim U[a, b]$, with the parameters

129 $a$ and $b$ being arbitrarily chosen to be 0 and 15, in the case of a continuous $x_{ij}$, and 0 and 0.05 in the

130 case of a binary $x_{ij}$.

131 For each simulated dataset, two linear regression models were fitted; an OLS model which ignored

132 the clustering (equation -1-), and a RI multi-level model which allowed for clustering effects

133 (equation -2-). For each of the models, the regression coefficient and its standard error (SE) were

134  estimated. To compare results from the two models, the difference between the estimated regression

135  coefficients ($\beta_1^{RI} - \beta_1^{OLS}$), and the ratio of their SEs ($SE^{RI}/SE^{OLS}$) were calculated.

136  To assess how the comparison between the two models was affected by the distribution of $x_{ij}$ within

137  and between clusters, these two measures were plotted against the dispersion (expressed as standard

138  deviation) of the mean values of $x_{ij}$ ($\bar{x}_j$) between clusters (dispersion of $shift_j$), for the case of

139  continuous $x_{ij}$, and dispersion of prevalence of $x_{ij}$, for the case of binary $x_{ij}$. In addition, descriptive

140  statistics were produced for the distributions of the two measures across simulated samples, according

141  to values for expected ICC and overall prevalence of $x_{ij}$, in the case of a binary explanatory variable.

142  The accuracy of the 95% confidence intervals for the regression coefficient $\beta_1$ from the two methods

143  was assessed by calculating the proportion of the estimated confidence intervals that included the true

144  value that had been used in the simulations. A method was considered to have appropriate coverage if

145  95% of the 95% confidence intervals included the value of the effect $\beta_1$ (i.e. the value 1) used in the

146  simulations. Deviations from this ideal could reflect bias in the estimates of effect, unsatisfactory

147  standard errors (22), or both.

148  To assess impacts on type I error, the simulations were repeated assuming no association between $x_{ij}$

149  and $y_{ij}$ (i.e. $\beta_1 = 0$), and the proportions of datasets for which the null hypothesis was rejected at a

150  5% significance level in OLS and RI modelling were compared according to ICC.

151  For each expected ICC, and each value of $shift_j$, 100 simulated datasets were produced with a

152  continuous $x_{ij}$, and another 100 for each of the four overall prevalence rates of a binary $x_{ij}$.

153  Due to random sampling variation the estimated ICC values were within given ranges of the target

154  levels of ICC. For target levels of 0.001, 0.003, 0.01, 0.03, 0.1 and 0.3, these ranges were 0.0005-

155  0.0014, 0.0025-0.0034, 0.005-0.014, 0.025-0.034, 0.05-0.14, and 0.25-0.34 respectively. Simulations

156  resulting in estimated ICC values outside of these ranges were discarded and not used further. In the

157  description of the results that follows ICC values are labelled according to the target levels.

158  All simulations and analysis were conducted using Stata software v12.1.

159     **Results**

160     *Difference in regression coefficients*

161     Differences in regression coefficients $(\beta_1^{RI} - \beta_1^{OLS})$ estimated from the two linear models are

162     illustrated in Figure 2. The two different subplots of the figure (A and B) correspond to the two

163     different distributions of the explanatory variable (continuous and binary respectively), and the

164     different shades of grey correspond to different ICC levels with darker shades corresponding to

165     simulated results for higher ICCs.

166                                     INSERT FIGURE 2 HERE

167     **Figure 2**. Difference between regression coefficients estimated from RI and OLS models ($\beta_1^{RI} - \beta_1^{OLS}$) plotted

168     against dispersion (expressed as SD) of mean value/prevalence of $x_{ij}$, for different levels of intraclass

169     correlation (shades of grey as indicated in the legend). Figure A: Continuous $x_{ij}$. Figure B: Binary $x_{ij}$

170

171     In all cases, differences in regression coefficients were on average zero, with $\beta_1^{RI}$ and $\beta_1^{OLS}$ being on

172     average $\cong 1$. For both continuous and binary distributions of $x_{ij}$, differences were on average more

173     narrowly spread for small ICCs and more widely spread for large ICCs. For a continuous explanatory

174     variable $x_{ij}$ (Figure 2A), and for each value of ICC, increasing the dispersion of $\bar{x}_j$ across clusters

175     resulted in larger differences in regression coefficients up to a dispersion of $\bar{x}_j = 1$ (i.e. same

176     dispersion of $x_{ij}$ between and within clusters). Beyond that point, further increase in the dispersion of

177     $\bar{x}_j$ resulted in smaller differences in regression coefficients from the two methods, approaching a

178     difference of zero.

179     For a binary explanatory variable $x_{ij}$, and for each value of ICC, small dispersion of cluster-specific

180     prevalence of $x_{ij}$ resulted in small differences between the regression coefficients. However,

181     increasing the dispersion of cluster-specific prevalence of $x_{ij}$, resulted in larger differences between

182     the regression coefficients from the two methods. Comparing the different subplots of **Error!**

183    **Reference source not found.** (note the different scales on the y-axes), higher overall prevalence of

184    $x_{ij}$ resulted in regression coefficients from the two models being more similar even for large

185    dispersion of the prevalence of $x_{ij}$ across clusters; for ICC=0.3, differences ranged from -0.2 to 0.2,

186    corresponding to a 20% difference in the regression coefficients from the two methods, when the

187    overall prevalence of $x_{ij}$ was 0.05, and this range decreased to approximately -0.05 to 0.05 for an

188    overall prevalence of $x_{ij}$ of 0.4.

189    *Ratio of standard errors*

190    The ratios of SEs derived from the RI and OLS models ($SE_{\beta_1^{RI}}/SE_{\beta_1^{OLS}}$) were examined in relation to

191    the dispersion across clusters of the mean value/prevalence of the continuous/binary explanatory

192    variable $x_{ij}$, and are presented in Figure 3. As in Figure 2, the different levels of ICCs are represented

193    by different shades of grey, with lighter shades corresponding to lower ICCs and darker shades to

194    higher ICCs. Subplots A and B correspond to the ratios of SEs when $x_{ij}$ was continuous and binary,

195    respectively.

196    For a continuous variable $x_{ij}$, the ratio took its minimum value for the smallest dispersion of $\bar{x}_j$ and

197    increased as dispersion of $\bar{x}_j$ increased, tending asymptotically to a maximum value. The minimum

198    and maximum values of the ratio of the SEs (the latter also corresponding to its asymptote) were ICC-

199    dependent, higher ICCs resulting in lower minimum and higher maximum values for the ratio. The

200    dispersion of $\bar{x}_j$ at which the ratio of SEs approached its asymptote was also ICC-dependent, being

201    higher for larger ICCs. For very small values of dispersion of $\bar{x}_j$, the minimum value of the ratio of

202    the SEs was approximately one for small levels of ICC and was less than one for higher ICCs.

203    Particularly for small values of the dispersion of $\bar{x}_j$ and ICC $\cong 0.10$ or $0.30$, the ratio of SEs was <1,

204    meaning that SEs from RI models were smaller than from OLS models.

205                                    INSERT FIGURE 3 HERE

9

206 **Figure 3**. Ratios of standard errors estimated from RI and OLS models ($SE_{\beta_1^{RI}}/SE_{\beta_1^{OLS}}$) plotted against relative

207 between- to within-clusters dispersion (expressed as SD) of explanatory variable $x_{ij}$. Figure A: Continuous $x_{ij}$.

208 Figure B: Binary $x_{ij}$

209 When $x_{ij}$ was binary, the ratios of the SEs were below one for most of the situations examined,

210 indicating that the SEs of the regression coefficients estimated from the RI model were smaller than

211 those under the OLS model in most circumstances. The ratio of the SEs achieved its minimum value

212 for the smallest dispersion of the prevalence of $x_{ij}$ across the clusters, and increased progressively

213 with increasing dispersion of $x_{ij}$ across clusters. For small ICCs (<0.1), the SEs from the two models

214 were very similar. However, increasing the ICC to 0.1 or higher led to the ratio of the SEs decreasing

215 to values much lower than 1. For constant ICC, comparison of subplots of Figure 3B, shows that the

216 rate of increase of the ratio of the SEs was higher for lower underlying prevalence rates of the $x_{ij}$.

217 *Coverage of 95% confidence intervals*

218 Table 1 shows the extent to which 95% confidence intervals covered the simulated effect of the

219 explanatory continuous variable on the outcome ($\beta_1$=1), when derived from the two statistical models,

220 for different levels of ICC, and for fifths of the distribution of the dispersion of $\bar{x}_j$.

221 Irrespective of ICC and type of explanatory variable, coverage with the RI model was approximately

222 95%. For a continuous $x_{ij}$, coverage for the OLS model was close to 95% for very low ICC and

223 decreased for increasing levels of ICC. For the highest ICC level examined (ICC=0.3), OLS showed a

224 notably poor coverage of 30%. For a given ICC, coverage of 95% confidence intervals did not vary

225 much by dispersion of $\bar{x}_j$, although it was somewhat higher in the bottom fifth as compared to the 2nd,

226 3rd, 4th, and 5th fifth of the distribution of dispersion of $\bar{x}_j$.

227

228

229　**Table 1**. Coverage (%) by 95% confidence intervals of simulated effect $\beta_1=1$ under the RI and OLS models

230　according to fifths of the distribution of dispersion (expressed as SD) of the continuous $\bar{x}_j$

| ICC | Bottom fifth=1 | | 2 | | 3 | | 4 | | Top fifth=5 | | Total | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | RI | OLS | RI | OLS | RI | OLS | RI | OLS | RI | OLS | RI | OLS |
| **0.001** | 95.04 | 94.37 | 95.00 | 94.09 | 95.33 | 94.10 | 95.15 | 93.99 | 94.92 | 93.84 | 95.08 | 94.08 |
| **0.003** | 95.14 | 93.14 | 95.37 | 92.33 | 95.20 | 91.83 | 95.53 | 92.25 | 95.42 | 91.99 | 95.33 | 92.30 |
| **0.01** | 94.99 | 88.47 | 94.64 | 83.95 | 94.75 | 83.65 | 94.74 | 83.72 | 94.90 | 84.07 | 94.80 | 84.75 |
| **0.03** | 94.59 | 76.21 | 95.11 | 68.79 | 94.80 | 67.62 | 95.06 | 67.57 | 94.74 | 67.15 | 94.87 | 69.39 |
| **0.1** | 94.68 | 59.58 | 94.80 | 45.45 | 94.86 | 44.83 | 94.39 | 44.73 | 95.04 | 44.37 | 94.76 | 47.80 |
| **0.3** | 94.84 | 41.32 | 94.53 | 28.06 | 94.95 | 28.24 | 94.64 | 26.98 | 94.79 | 27.41 | 94.75 | 30.36 |

231

232　For a binary $x_{ij}$, coverage for the OLS model was close to 95% but only for ICC≤0.03. As ICC

233　increased, coverage from the OLS model deviated from the nominal value of 95%. As shown in

234　Figure 4, when ICC was 0.1 or 0.3, coverage was on average lower for lower prevalence of $x_{ij}$; it fell

235　below the nominal value of 95% for 0.05 prevalence of $x_{ij}$ and it increased to values higher than 95%

236　for 0.40 prevalence of $x_{ij}$ (comparison of the four sub-plots of the figure). Also, for any given

237　prevalence of $x_{ij}$, coverage was lower for increasing dispersion of prevalence of $x_{ij}$ across clusters.

238　Variation of the average coverage by categories of prevalence rates of $x_{ij}$ and overall prevalence of

239　$x_{ij}$ was higher when ICC was higher (ICC=0.3) than when it was lower (ICC=0.1). The smallest and

240　the largest values of coverage were 87% and 98% and they were observed when overall prevalence of

241　$x_{ij}$ was 0.05, ICC=0.3, and in the bottom and top thirds respectively of the distribution of dispersion

242　of prevalence of $x_{ij}$ across clusters. Coverage as high as 98% was also seen in the bottom third of the

243　distribution of dispersion of prevalence of $x_{ij}$ across clusters for the other prevalence rates (0.10,

244　0.20, and 0.40) explored when ICC was high (ICC=0.3).

246   **Figure 4**. Coverage (%) by 95% confidence intervals from the OLS model for ICC=0.1 and 0.3, by overall

247   prevalence rates of $x$ (A) 0.05, B) 0.10, C) 0.20, and D) 0.40), and thirds of the distribution of the dispersion

248   (expressed as SD) of prevalence of across clusters

249   *Type I error*

250   To assess the frequency of type I error, defined as incorrect rejection of a true null hypothesis, under

251   the OLS and the RI multi-level models, simulations were repeated assuming no association between

252   the explanatory variable $x_{ij}$ and the outcome variable $y_{ij}$ ($\beta_1^{RI} = \beta_1^{OLS} = 0$).

253   Figure 5 shows the proportion of datasets for which the null hypothesis was rejected at a 5%

254   significance level for varying levels of ICC, when $x_{ij}$ was continuous. Using the RI multi-level

255   model, the association between $x_{ij}$ and $y_{ij}$ was statistically significant in approximately 5% of the

256   datasets for all ICCs. However, using the OLS models, type I error varied with ICC. For a very small

257   ICC, type I error was very close to that under the RI model (~6%) but increased rapidly as the ICC

258   increased, reaching ~70% for ICC$\cong$0.30. Type I error did not vary by dispersion of mean value of $x_{ij}$

259   (data not shown).

261   **Figure 5**. Proportion (%) of datasets for which the null hypothesis was rejected according to level of ICC when

262   $\beta_1^{RI} = 0$ and $x_{ij}$ was continuous

263   When the explanatory variable $x_{ij}$ was binary, type I error rates varied very little around the nominal

264   level of 5% when an OLS model was fitted instead of the RI model, when ICC values were less than

265   0.1; the average value was 5% and varied from 4.8% to 5.3% for different ICC values (<0.1), overall

266   prevalence rates of $x_{ij}$, and dispersion of prevalence of $x_{ij}$ across clusters. However, for ICC values

267   of 0.1 and 0.3, type I error rates diverged from 5%. The variation of rates in those cases is illustrated

268   in Figure 6 for the four prevalence rates of $x_{ij}$ (subplots A, B, C, and D of the figure), and for thirds

269    of the distribution of dispersion of prevalence of $x_{ij}$ across clusters. For small dispersion of

270    prevalence rates of $x_{ij}$ (bottom third of the distribution), type I error was lower than 5%, and it

271    increased as dispersion increased. This trend was more prominent for lower values of overall

272    prevalence of $x_{ij}$, and for ICC=0.3 compared to ICC=0.1. The smallest and the largest values of type I

273    error were 2% and 13% and they were observed when overall prevalence of $x_{ij}$ was 0.05 and in the

274    bottom and top thirds respectively of the distribution of dispersion of prevalence of $x_{ij}$ across clusters.


275                                          INSERT FIGURE 6 HERE


276    **Figure 6**. Type I error rates (%) from the OLS model for ICC=0.1 and 0.3, by overall prevalence rates of $x_{ij}$ (A)

277    0.05, B) 0.10, C) 0.20, and D) 0.40), and thirds of the distribution of the dispersion (expressed as SD) of

278    prevalence of $x$ across clusters


279    **Discussion**

280    In this paper we focused on the implications of ignoring clustering in statistical inference regarding

281    the relationship between a continuous outcome and a single explanatory variable $x_{ij}$. Two different

282    types of $x_{ij}$ were considered – continuous and binary. For each of the two categories of $x_{ij}$, the

283    implications for statistical inference of failing to account for clustering were explored by comparison

284    of effect estimates and their precision, assessment of the coverage by 95% confidence intervals, and

285    estimation of the frequency of type I error. In the cases of both a continuous and a binary $x_{ij}$, where

286    the true slope of the regression line was non-zero, we found that the cluster-unadjusted OLS and RI

287    models gave on average very similar estimates of effect for any level of ICC. However, despite the

288    average value of difference in point estimates from the two methods being zero, differences occurred

289    in both directions and varied more when the level of ICC increased. The largest differences in

290    estimates of effect between OLS and multi-level RI regression modelling were only about 20% of the

291    true value and they occurred when the ICC was high (0.3). For a continuous $x_{ij}$, the largest errors in

292    the differences of estimated effects occurred when the dispersion of the $x_{ij}$ within clusters was

293    approximately the same as that between clusters, while, for a binary explanatory variable, differences

294    increased with increasing dispersion of prevalence of $x_{ij}$ across clusters.

13

295    Conclusions drawn from comparison of SEs estimated from cluster-unadjusted OLS and RI models

296    are somewhat different for continuous as compared with binary $x_{ij}$. When $x_{ij}$ was continuous, the

297    SEs of regression coefficients were generally larger for the multi-level RI model than for the cluster-

298    unadjusted OLS model, their ratio being highest (>4) for a high ICC (0.3) and where the dispersion of

299    the mean value of $x_{ij}$ was large.  However, contrary to what is widely stated, the spuriously greater

300    precision of OLS method was not universal.  When dispersion of mean values of $x_{ij}$<1, OLS

301    regression gave larger SEs than multi-level modelling. When $x_{ij}$ was binary, SEs estimated from the

302    RI model, were higher than those from the cluster-unadjusted OLS model for lower ICCs (<0.03) and

303    larger dispersion of prevalence of $x_{ij}$ across clusters, and lower than those from the cluster-unadjusted

304    OLS model for smaller dispersion of prevalence of $x_{ij}$ across clusters. The SEs differed by up to 15%

305    for the highest ICC value (ICC=0.3).

306    The rates of coverage of 95% confidence intervals for estimates of effect, whether of a continuous or

307    a binary $x_{ij}$, when derived from a RI model were at the nominal level of 95%, irrespective of other

308    parameters (i.e. ICC, dispersion across clusters of the mean value of a continuous $x_{ij}$, or dispersion of

309    the prevalence of the binary $x_{ij}$ across clusters). When $x_{ij}$ was binary, the cluster-unadjusted OLS

310    model also resulted in an appropriate coverage of the 95% confidence intervals when ICC was low ($\leq$

311    0.01). However, for higher values of ICC, coverage varied slightly (range: 87% - 98%) around the

312    nominal value of 95% depending on the overall prevalence and the dispersion of the cluster-specific

313    prevalence rates of $x_{ij}$. In contrast, when $x_{ij}$ was continuous, the model that failed to account for

314    clustering resulted in poor coverage rates, especially as ICC increased, reaching a rate as low as 30%

315    for ICC=0.3.

316    Setting the effect of $x_{ij}$ on the outcome variable to zero allowed exploration of the frequency of type I

317    error. With the RI model, in all of the scenarios explored, type I error was very close to 5%. When $x_{ij}$

318    was continuous, we found that failure to allow for clustering increased rates of Type I error, and that

319    the inflation of type I error was particularly pronounced (up to 70%) when the degree of clustering

320    was high (ICC=0.3). In contrast to this, when $x_{ij}$ was binary, type I error under the OLS model was

321   close to the expected value of 5% for low levels of clustering (ICC<0.1). However, when ICC was

322   high (0.1 or 0.3), type I error rates varied more widely around 5%, with values as low as 2% (for low

323   overall prevalence of $x_{ij}$ and small dispersion of its prevalence across clusters) and as high as 13%

324   (for low overall prevalence of $x_{ij}$ and large dispersion of its prevalence across clusters).

325   The analysis for each specification of parameters (expected ICC, dispersion of $x_{ij}$, overall prevalence

326   or dispersion of prevalence rates across clusters of a binary $x_{ij}$) was based on 1,000 simulated

327   samples of 10,000 observations grouped in 100 clusters, each of 100 individuals. By using such a

328   large sample size (larger than in most epidemiological investigations), we reduced random sampling

329   variation, making it easier to characterise any systematic differences between the two methods of

330   analysis. However, the approach may have led to underestimation of the maximum differences

331   between estimates of effect that could arise from OLS as compared with multi-level modelling.

332   Additionally, the number of observations per cluster was the same in all simulations, making it

333   impossible to draw conclusions about effects of ignoring clustering for varying cluster sizes. Also,

334   data were simulated following the specification of the RI regression model rather than that of the

335   random-effects model described in section -2-. That was done because the RI model is more

336   frequently used, especially when there is no a priori expectation of differential effects of the

337   explanatory on the outcome variables across the different clusters. Simulating data following the

338   specification of the random effects model would have added complexity to the algorithm used for

339   simulation, and the computational time required.

340   The effect of clustering when a cluster-unadjusted model is fitted could also have been assessed by

341   calculating bias as [(estimated effect – true effect)/true effect], as defined in earlier studies (23).

342   Instead, we defined bias by the difference in the effect estimates derived from the two analytical

343   models. The data were simulated following the model specification of RI linear regression, which is

344   one of the most well established and frequently chosen analytical approaches to account for

345   clustering. As such, given that all resulting effect estimates were positive, deviations of the difference

346   in regression coefficients from the value of zero can only represent deficiencies of the OLS model,

347   provided that the assumptions of the RI model are met. Therefore, there is no reason to expect that the

348  conclusions one would draw from an alternative definition of bias would be more reliable, provided

349  that the conditions under which data were simulated and the models fitted were the same.

350  When multilevel RI modelling was applied to the simulated clustered datasets with a continuous or a

351  binary explanatory variable, the rate of Type I error was 5%, and the coverage by 95% CIs was 95%,

352  as would be expected, given the method by which the simulated samples were generated. In

353  comparison, when cluster-unadjusted models were fitted to clustered data with a continuous $x_{ij}$, rates

354  of Type I error were higher, particularly when the ICC was high. For the highest level of ICC

355  examined (0.3), type I errors were as frequent as 70%. However, even with an ICC of only 0.01, rates

356  of Type 1 error were more than 10%. Consistent with this, coverage by 95% confidence intervals was

357  considerably lower than the nominal value for higher ICC levels. The lowest coverage of 30% was for

358  the highest ICC level. In contrast to these results Huang et al (24) have reported values of coverage

359  very close to 95% from the OLS model for a continuous explanatory variable. Differences between

360  findings presented in this study and those presented by Huang et al (24) can be explained by zero

361  clustering in the explanatory variable assumed in the latter. Sensitivity analysis restricting the

362  simulated datasets only to those in which clustering in the explanatory variable was not meaningful

363  showed that interval coverage rates were very close to 95% independent of clustering in the outcome

364  variable (results not shown). When $x_{ij}$ was binary, both the interval coverage and Type I error rates

365  varied little around the nominal values of 5% and 95%, and only for ICC values higher than 0.01.

366  Overall coverage rates were higher for higher ICCs and decreased for increasing dispersion of the

367  cluster-specific prevalence rates of $x_{ij}$ across clusters and for decreasing overall prevalence of the $x_{ij}$.

368  A similar observation of small variation of interval coverage around 95% for higher ICC values has

369  been made before (25). Type I error when $x_{ij}$ was binary and its prevalence was low, varied around

370  5% with values falling below 5% for small dispersion of prevalence of $x_{ij}$, and above 5% for large

371  dispersion. For larger overall prevalence of $x_{ij}$, Type I error rates fell below 5%. In accordance with

372  these findings, Galbraith et al (26) have shown that cluster-unadjusted models resulted in relatively

373  conservative Type I error. Also, in a context of individually randomised trials, Kahan et al (27) have

374     shown that Type I error increased with increasing ICC and increasing difference in the probability of

375     assignment of patients to treatment arms.

376     It has been widely stated that when data are clustered, effects estimated by OLS regression are

377     unbiased (23, 25, 27-30). Our results confirm that for data of the type simulated, coefficients from

378     OLS regression were on average very similar to those from RI multi-level modelling. Previous

379     studies based on simulation data have shown similar results (23-25, 31). However, for individual

380     simulated samples, the estimates may differ, and the potential magnitude of the differences depends

381     on the level of within-cluster similarity of the outcome variable. For an ICC of 0.3, the estimates of

382     effect from the two analytical methods could differ by up to 20%. In addition, when $x_{ij}$ is continuous,

383     the error in estimates of the regression coefficient is larger when the between-cluster dispersion of $x_{ij}$

384     is similar to that within-cluster. When $x_{ij}$ is binary, the error increases as the dispersion of the

385     prevalence rates across clusters increases, and when the overall prevalence rate across all clusters is

386     lower (<10%). These errors in the estimated effect indicate that in an individual study, failure of

387     regression analysis to account for clustering of observations could result in considerably higher or

388     lower estimates of effect than those derived from multilevel analysis. This has been illustrated in

389     numerous published papers of real data, which have shown that estimates from the two analytical

390     methods can differ to a lesser or greater extent (1, 8, 14, 17, 32). However, in those publications, no or

391     very limited information is provided to establish whether the error observed was due to dispersion of

392     the cluster-specific mean values of the continuous $x_{ij}$, or dispersion of prevalence rates for the binary

393     $x_{ij}$ across clusters.

394     Most often it is stated that regression coefficients are spuriously precise when clustering is not taken

395     into account in regression models. However, in several reports, authors have failed to specify the

396     conditions under which this applies (1, 31, 33-36). Other authors have pointed out that when $x_{ij}$ is

397     identical within each cluster, and a cluster-unadjusted approach is followed, SEs tend to be spuriously

398     low, and that the opposite occurs when $x_{ij}$ varies within clusters (24, 27, 37, 38). Bias in SEs for

399     effects of cluster-varying $x_{ij}$ has been shown in results from real data when both models were fitted

400 (17, 32, 39). However, others have reported contradictory results in which SEs of effects of

401 individual-level $x_{ij}$ from OL regression were very similar to, or lower than, those from a multi-level

402 model (14-16, 40). It should be noted that the dichotomy between cluster- and individual-level

403 variables is not clear-cut. There can be varying degrees of clustering in $x_{ij}$, with the extremes being

404 variables for which the values are completely unclustered (mean values are the same for all clusters),

405 and variables for which the values are the same within each cluster. However, in real data, an

406 explanatory variable can lie anywhere in between. An early report focused on this issue by

407 considering the level of clustering in $x_{ij}$ as the main driver for the expected bias of the precision of

408 the effect estimates (29), rather than the absolute distinction between cluster-constant and cluster-

409 varying $x_{ij}$. The authors reported that as clustering in $x_{ij}$ decreases, the bias in SEs from a cluster-

410 unadjusted model is expected to be upwards, and the opposite is expected when clustering in $x_{ij}$

411 increases. Taking into consideration clustering in $x_{ij}$ ($\rho_x$) as well as in the outcome variable ($\rho_y$), a

412 later study using simulated data showed that for a given level of $\rho_y$, increasing $\rho_x$ resulted in

413 increasing the ratio of estimated SEs ($SE_\beta^{RI}/SE_\beta^{OLS}$) from values <1 to values $\approx 1$ (41). Our results for

414 continuous explanatory variables differ slightly from this, with ratios of SEs ($SE_\beta^{RI}/SE_\beta^{OLS}$) moving

415 from values <1 to values >1, as clustering in the explanatory variable, expressed as dispersion of $\bar{x}_j$

416 across clusters, increased.


417 Bias in the precision of effect estimates for binary $x_{ij}$ when clustering is ignored has received very

418 limited attention in the published literature. Several of the reported studies have used real data to

419 compare standard and multi-level models, using both continuous and binary individual-level $x_{ij}$ (14,

420 17). For the majority of binary $x_{ij}$ used in the models fitted in these studies, SEs derived from the

421 OLS model were larger than those derived from the multi-level model. The same conclusion was

422 drawn from a study using simulated data (25). However, none of the studies using real data has

423 explored the level of bias in relation to variation in the prevalence of the binary $x_{ij}$, and the study of

424 simulated data assumed constant prevalence of $x_{ij}$ in all clusters. Simulation results presented here

425 suggest that, irrespective of the dispersion of prevalence of $x_{ij}$ across clusters and the overall

18

426    prevalence in all clusters, in most circumstances SEs from the multi-level model are lower than those

427    from the OLS model, and the bias is higher for higher ICC values.

428    The focus of this paper was on the association between a continuous outcome and an explanatory

429    variable that was defined at the individual level ($x_{ij}$ within cluster). We showed that when $x_{ij}$ was

430    continuous, and most of the variation occurred within rather than between clusters, the cluster-

431    unadjusted OLS model gave larger SEs for the regression coefficient than multi-level modelling. This

432    is consistent with reports in which ignoring clustering resulted in spuriously high SEs when $x_{ij}$ varied

433    within cluster. The reverse occurred when most of the dispersion of $x_{ij}$ was between rather than

434    within clusters. In this situation $x_{ij}$ approaches the characteristics of a cluster-specific variable. We

435    additionally showed that when $x_{ij}$ under investigation was binary, ignoring clustering in statistical

436    modelling in most cases resulted in higher SEs for the estimated effect than those derived from the

437    random-intercept model. The SEs differed more for higher ICCs but not with the overall prevalence of

438    $x_{ij}$, nor with the dispersion of its prevalence across clusters (Figure 3B). Unlike SEs, the point

439    estimates were unbiased for either continuous or binary $x_{ij}$ (Figure 2 A and B).

440    In conclusion, our results support the use of multi-level modelling to account for clustering effects in

441    linear regression analyses of data that are hierarchically structured, especially where ICCs might

442    exceed 0.01. Failure to do so is likely to result in incorrect estimates of effect (either too high or too

443    low) mostly with spurious precision in the case of continuous $x_{ij}$ or with underestimated precision in

444    the case of binary $x_{ij}$, and may lead to incorrect inferences. The errors in estimates of effect of a

445    continuous $x_{ij}$ will be smaller when most of its dispersion is between rather than within clusters – i.e.

446    the variable comes closer to being cluster-specific. Similarly, when $x_{ij}$ is binary, smaller differences

447    in the effect estimates occur when the dispersion of the prevalence of $x_{ij}$ across clusters is small, or

448    when its overall prevalence across clusters is high.

449    Additionally, we identified situations in which a standard analytical approach is more likely to

450    importantly affect statistical inference, i.e. when rates of Type I error and interval coverage deviate

19

451    more from the nominal values of 5% and 95% respectively. These occur when $x_{ij}$ is continuous, and

452    ICC levels are greater than 0.01. It is then that Type I error rates are higher than 10% and interval

453    coverage rates are lower than 80%. On the other hand, statistical inference when a standard regression

454    model is fitted is less likely to be of concern when $x_{ij}$ is binary, as the error and coverage rates

455    deviate very little from the nominal values. However, even for a binary $x_{ij}$, error rates can sometimes

456    be greater than 10%, and corresponding interval coverage rates lower than 90% (but possibly not

457    lower than 80%). This occurs when ICC is high, the overall prevalence of $x_{ij}$ is low (approximately

458    5%), and the dispersion of the cluster-specific rates is large. In all circumstances in which the ICC is

459    very small, clustering is minimal and there is little difference between RI and OLS regression.

460    **Abbreviations**

461    RI: random-intercept; OLS: ordinary least squares; ICC: intra-class correlation coefficient; SE:

462    standard error

463    **Acknowledgements**

464    Not applicable

465    **Authors' contribution**

466    GN, HI, and DC conceived the concept of this study. GN carried out the simulations, analysed the

467    data and drafted the manuscript. CO provided expert statistical advice on aspects of results presented.

468    DC and HI critically reviewed and made substantial contributions to the manuscript. All authors read

469    and approved the final manuscript.

470    **Funding**

473    **Availability of data and materials**

474 The simulated datasets used and analysis described in the current study are available from the

475 corresponding author on reasonable request.

476 **Ethics approval and consent to participate**

477 Not applicable

478 **Consent for publication**

479 Not applicable

480 **Competing interests**

481 The authors declare that they have no competing interests

482 **References**

483 1.      Park S, Lake ET. Multilevel modeling of a clustered continuous outcome: nurses' work hours

484 and burnout. Nursing research. 2005;54(6):406-13.

485 2.      Stimson JA. Regression in Space and Time: A Statistical Essay. American Journal of Political

486 Science. 1985;29(4):914-47.

487 3.      Bingenheimer JB, Raudenbush SW. Statistical and substantive inferences in public health:

488 issues in the application of multilevel models. Annu Rev Public Health. 2004;25:53-77.

489 4.      Bland JM. Cluster randomised trials in the medical literature: Two bibliometric surveys.

490 BMC Medical Research Methodology. 2004;4.

491 5.      Crits-Christoph P, Mintz J. Implications of therapist effects for the design and analysis of

492 comparative studies of psychotherapies. Journal of consulting and clinical psychology. 1991;59(1):20.

493 6.      Lee KJ, Thompson SG. Clustering by health professional in individually randomised trials.

494 Bmj. 2005;330(7483):142-4.

495 7.      Simpson JM, Klar N, Donner A. Accounting for cluster randomization: A review of primary

496 prevention trials, 1990 through 1993. American Journal of Public Health. 1995;85(10):1378-83.

497    8.    Biau DJ, Halm JA, Ahmadieh H, Capello WN, Jeekel J, Boutron I, et al. Provider and center

498    effect in multicenter randomized controlled trials of surgical specialties: an analysis on patient-level

499    data. Ann Surg. 2008;247(5):892-8.

500    9.    Oltean H, Gagnier JJ. Use of clustering analysis in randomized controlled trials in orthopaedic

501    surgery. BMC Medical Research Methodology. 2015;15(1).

502    10.    Diaz-Ordaz K, Froud R, Sheehan B, Eldridge S. A systematic review of cluster randomised

503    trials in residential facilities for older people suggests how to improve quality. BMC Medical

504    Research Methodology. 2013;13(1).

505    11.    Goldstein H. Multilevel Mixed Linear Model Analysis Using Iterative Generalized Least

506    Squares. Biometrika. 1986;73(1):43-56.

507    12.    Astin AW, Denson N. Multi-campus studies of college impact: Which statistical method is

508    appropriate? Research in Higher Education. 2009;50(4):354-67.

509    13.    Cheong YF, Fotiu RP, Raudenbush SW. Efficiency and robustness of alternative estimators

510    for two- and three-level models: The case of NAEP. Journal of Educational and Behavioral Statistics.

511    2001;26(4):411-29.

512    14.    Grieve R, Nixon R, Thompson SG, Normand C. Using multilevel models for assessing the

513    variability of multinational resource use and cost data. Health economics. 2005;14(2):185-96.

514    15.    Niehaus E, Campbell CM, Inkelas KK. HLM Behind the Curtain: Unveiling Decisions

515    Behind the Use and Interpretation of HLM in Higher Education Research. Research in Higher

516    Education. 2014;55(1):101-22.

517    16.    Steenbergen MR, Jones BS. Modeling multilevel data structures. american Journal of political

518    Science. 2002:218-37.

519    17.    Wendel-Vos GCW, Van Hooijdonk C, Uitenbroek D, Agyemang C, Lindeman EM,

520    Droomers M. Environmental attributes related to walking and bicycling at the individual and

521    contextual level. Journal of Epidemiology and Community Health. 2008;62(8):689-94.

522    18.    Walters SJ. Therapist effects in randomised controlled trials: what to do about them. Journal

523    of clinical nursing. 2010;19(7-8):1102-12.

524   19.   Newman D, Newman I, Salzman J. Comparing OLS and HLM models and the questions they

525   answer: Potential concerns for type VI errors. Multiple Linear Regression Viewpoints. 2010;36(1):1-

526   8.

527   20.   Goldstein H. Multilevel Statistical Models: Wiley; 2010.

528   21.   Rabe-Hesketh S, Skrondal A. Multilevel and Longitudinal Modeling Using Stata: Taylor &

529   Francis; 2005.

530   22.   Bradburn MJ, Deeks JJ, Berlin JA, Russell Localio A. Much ado about nothing: a comparison

531   of the performance of meta-analytical methods with rare events. Stat Med. 2007;26(1):53-77.

532   23.   Clarke P. When can group level clustering be ignored? Multilevel models versus single-level

533   models with sparse data. Journal of Epidemiology and Community Health. 2008;62(8):752-8.

534   24.   Huang FL. Alternatives to multilevel modeling for the analysis of clustered data. The Journal

535   of Experimental Education. 2016;84(1):175-96.

536   25.   Chu R, Thabane L, Ma J, Holbrook A, Pullenayegum E, Devereaux PJ. Comparing methods

537   to estimate treatment effects on a continuous outcome in multicentre randomized controlled trials: a

538   simulation study. BMC medical research methodology. 2011;11(1):1.

539   26.   Galbraith S, Daniel JA, Vissel B. A study of clustered data and approaches to its analysis. The

540   journal of Neuroscience. 2010;30(32):10601-8.

541   27.   Kahan BC, Morris TP. Assessing potential sources of clustering in individually randomised

542   trials. BMC Medical Research Methodology. 2013;13(1).

543   28.   Arceneaux K, Nickerson DW. Modeling Certainty with Clustered Data: A Comparison of

544   Methods. Political Analysis. 2009;17(2):177-90.

545   29.   Scott AJ, Holt D. The effect of two-stage sampling on ordinary least squares methods. Journal

546   of the American Statistical Association. 1982;77(380):848-54.

547   30.   Barrios T, Diamond R, Imbens GW, Koleśar M. Clustering, spatial correlations, and

548   randomization inference. Journal of the American Statistical Association. 2012;107(498):578-91.

549   31.   Maas CJ, Hox JJ. The influence of violations of assumptions on multilevel parameter

550   estimates and their standard errors. Computational Statistics & Data Analysis. 2004;46(3):427-40.

551   32.     Dickinson LM, Basu A. Multilevel modeling and practice-based research. The Annals of

552   Family Medicine. 2005;3(suppl 1):S52-S60.

553   33.     Austin PC, Goel V, van Walraven C. An introduction to multilevel regression models.

554   Canadian Journal of Public Health. 2001;92(2):150.

555   34.     Lemeshow S, Letenneur L, Dartigues JF, Lafont S, Orgogozo JM, Commenges D. Illustration

556   of analysis taking into account complex survey considerations: The association between wine

557   consumption and dementia in the PAQUID study. American Journal of Epidemiology.

558   1998;148(3):298-306.

559   35.     Roberts C, Roberts SA. Design and analysis of clinical trials with clustering effects due to

560   treatment. Clinical Trials. 2005;2(2):152-62.

561   36.     Hox J. Multilevel Modeling: When and Why. In: Balderjahn I, Mathar R, Schader M, editors.

562   Classification, Data Analysis, and Data Highways: Proceedings of the 21st Annual Conference of the

563   Gesellschaft für Klassifikation eV, University of Potsdam, March 12–14, 1997. Berlin, Heidelberg:

564   Springer Berlin Heidelberg; 1998. p. 147-54.

565   37.     Chuang J-H, Hripcsak G, Heitjan DF. Design and Analysis of Controlled Trials in Naturally

566   Clustered Environments: Implications for Medical Informatics. Journal of the American Medical

567   Informatics Association : JAMIA. 2002;9(3):230-8.

568   38.     Sainani K. The importance of accounting for correlated observations. PM & R : the journal of

569   injury, function, and rehabilitation. 2010;2(9):858-61.

570   39.     Jones K. Do multilevel models ever give different results? 2009.

571   40.     Hedeker D, McMahon SD, Jason LA, Salina D. Analysis of clustered data in community

572   psychology: with an example from a worksite smoking cessation project. American journal of

573   community psychology. 1994;22(5):595-615.

574   41.     Bliese PD, Hanges PJ. Being Both Too Liberal and Too Conservative: The Perils of Treating

575   Grouped Data as though They Were Independent. Organizational Research Methods. 2004;7(4):400-
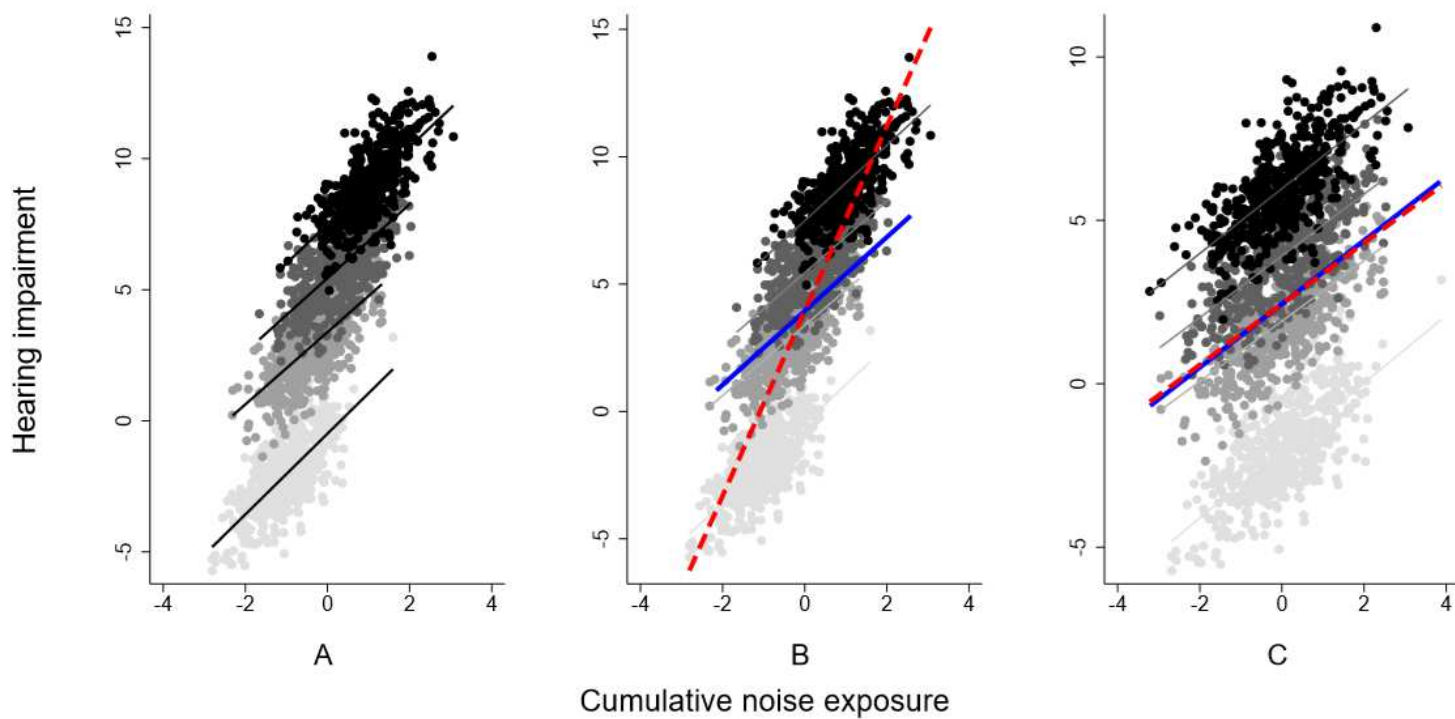
576   17.

577

# Figures



**Figure 1**

Hypothetical relationship of hearing impairment to cumulative noise exposure in four cities. Units for noise exposure and hearing impairment have been specified arbitrarily for ease of presentation. Data for each city are distinguished by the shading of data points. Cluster-specific regression lines are indicated, along with the regression line for the full dataset when clustering is ignored (dotted red line), and that when adjustment is made for cluster (solid blue line)
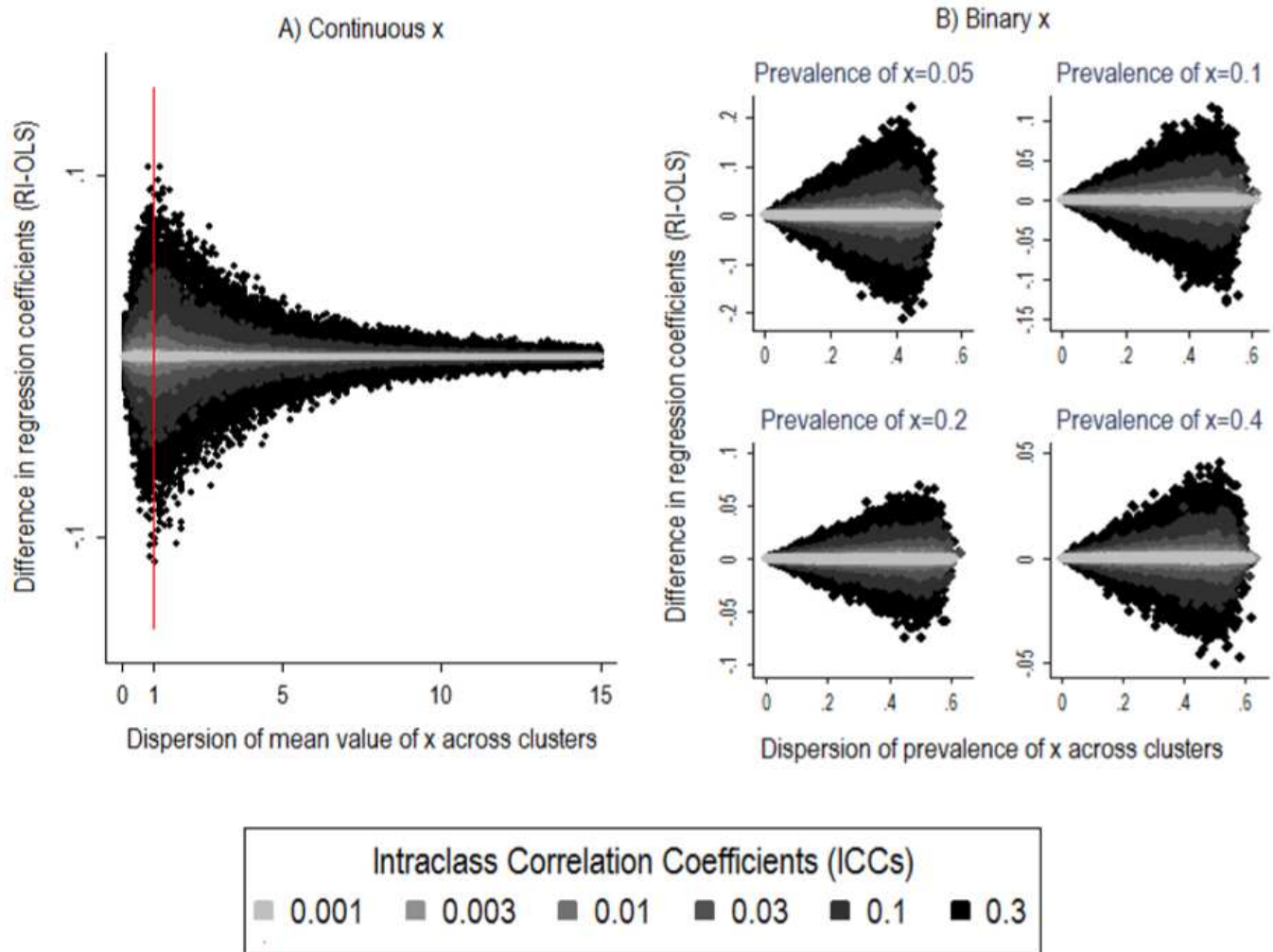
## Figure 2

Difference between regression coefficients estimated from RI and OLS models ($\beta_1^{RI}-\beta_1^{OLS}$) plotted against dispersion (expressed as SD) of mean value/prevalence of $x_{ij}$, for different levels of intraclass correlation (shades of grey as indicated in the legend). Figure A: Continuous $x_{ij}$. Figure B: Binary $x_{ij}$
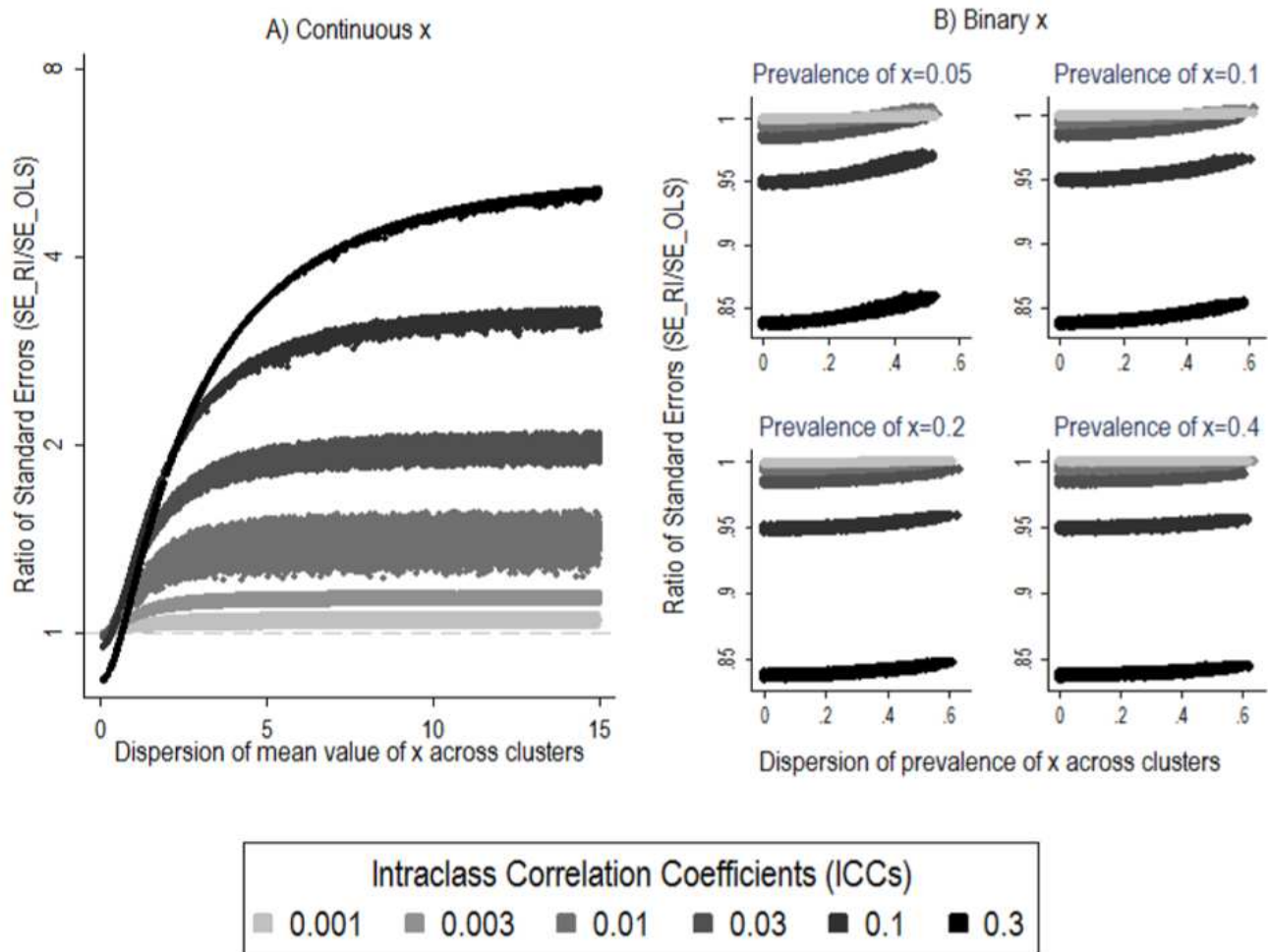
**Figure 3**

Ratios of standard errors estimated from RI and OLS models ($SE_{(\beta_1^{RI})}/SE_{(\beta_1^{OLS})}$) plotted against relative between- to within-clusters dispersion (expressed as SD) of explanatory variable $x_{ij}$. Figure A: Continuous $x_{ij}$. Figure B: Binary $x_{ij}$
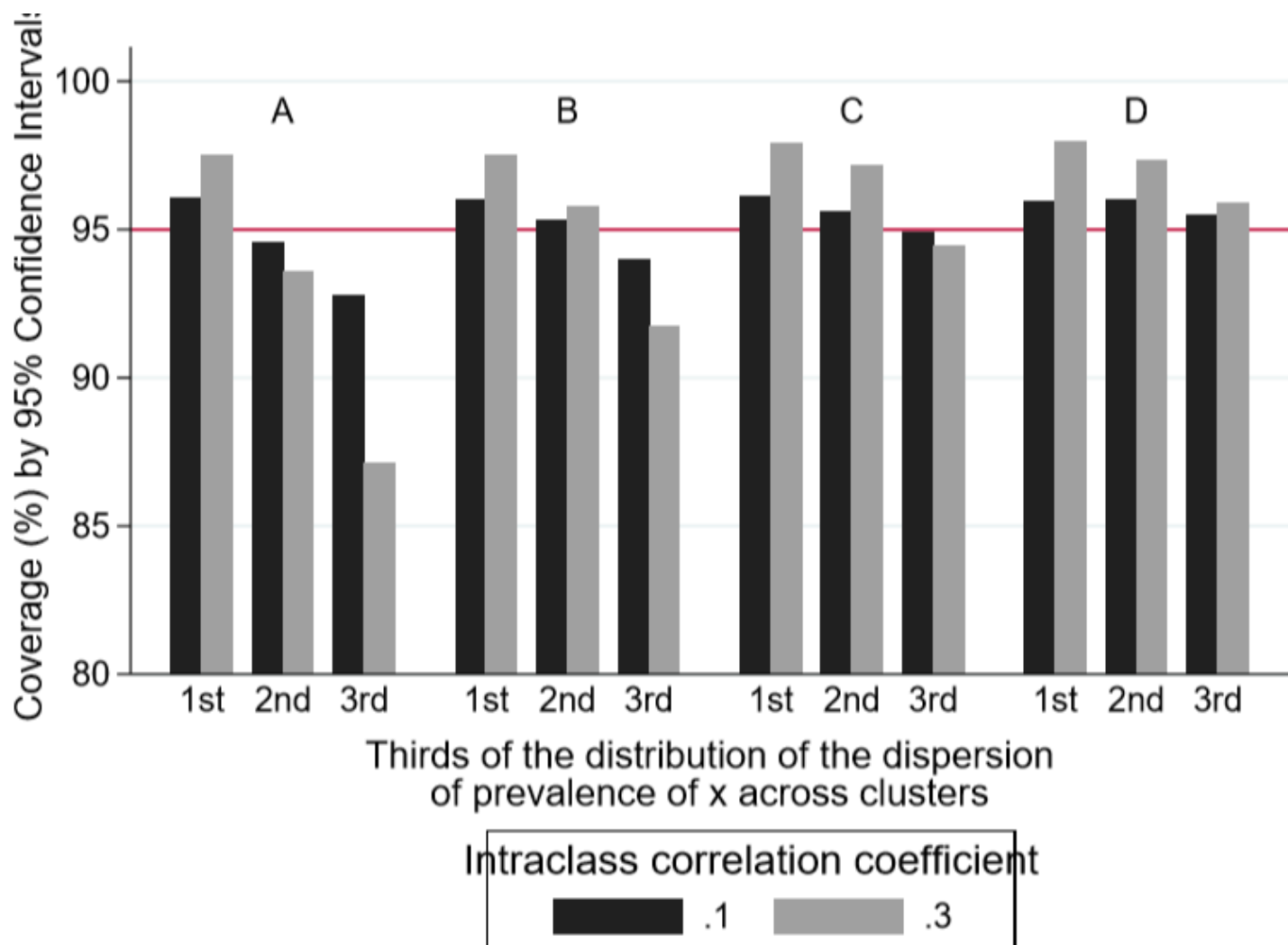
**Figure 4**

Coverage (%) by 95% confidence intervals from the OLS model for ICC=0.1 and 0.3, by overall prevalence rates of x (A) 0.05, B) 0.10, C) 0.20, and D) 0.40), and thirds of the distribution of the dispersion (expressed as SD) of prevalence of across clusters
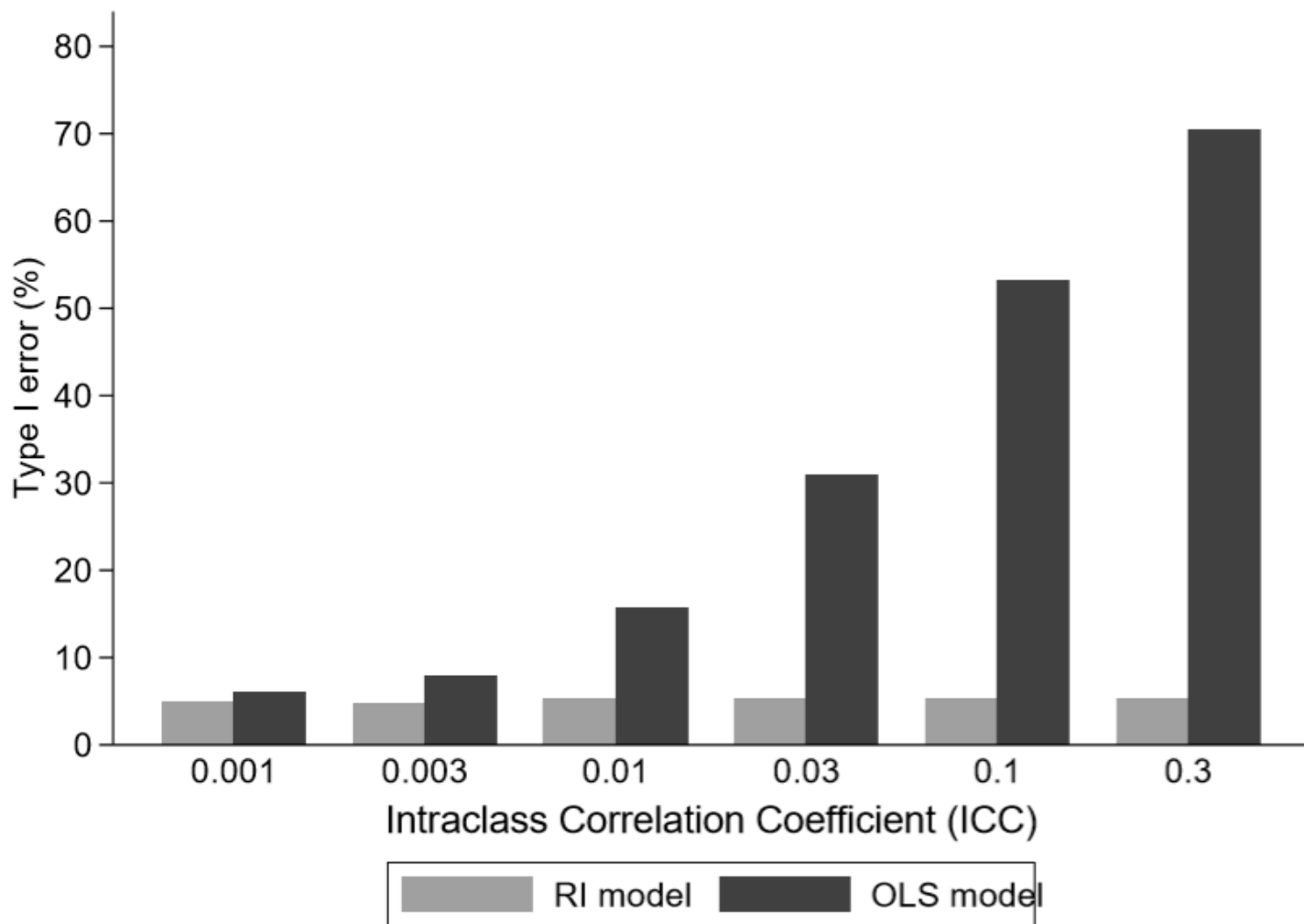
**Figure 5**

Proportion (%) of datasets for which the null hypothesis was rejected according to level of ICC when $\beta_1^{RI}=0$ and $x_{ij}$ was continuous
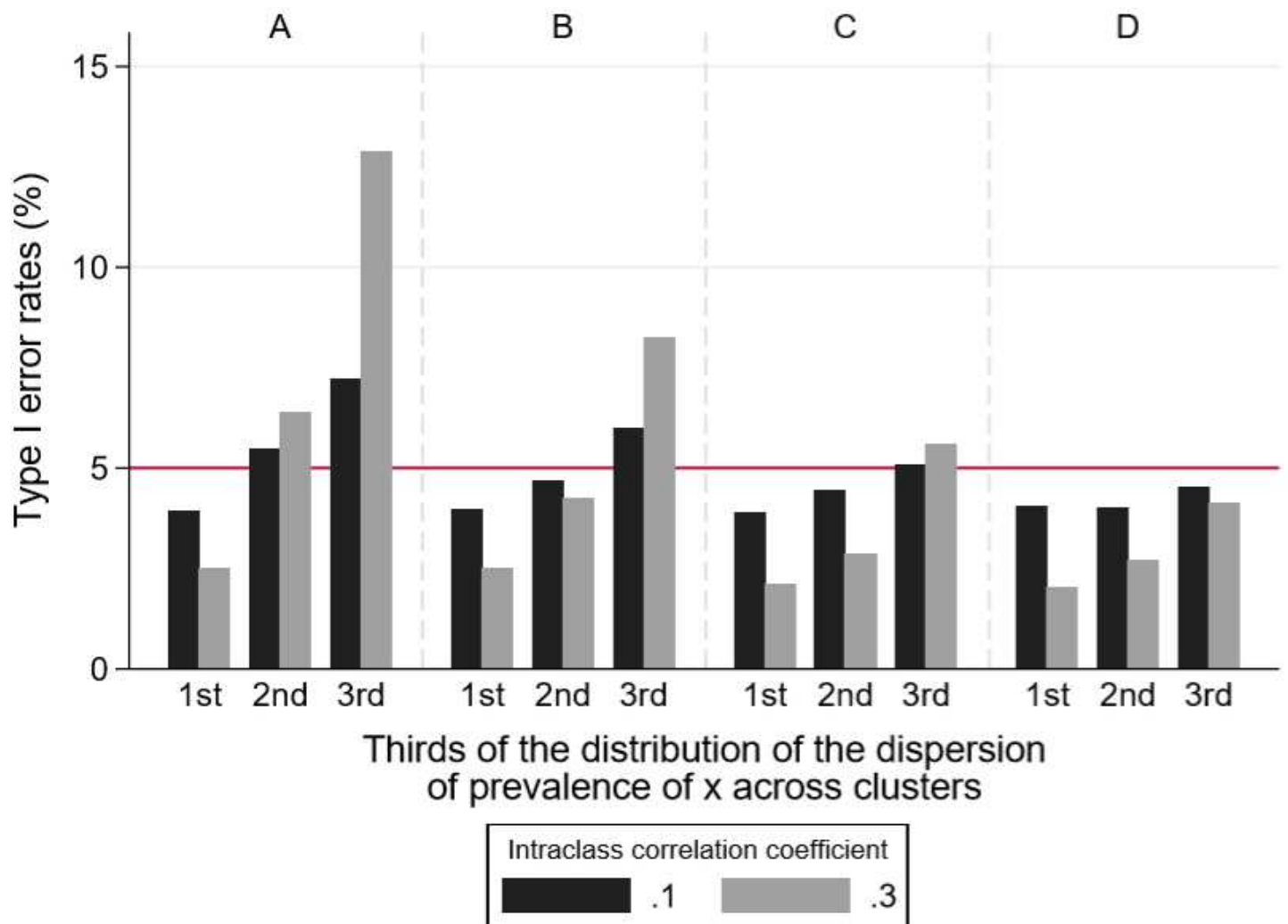
**Figure 6**

Type I error rates (%) from the OLS model for ICC=0.1 and 0.3, by overall prevalence rates of $x_{ij}$ (A) 0.05, B) 0.10, C) 0.20, and D) 0.40), and thirds of the distribution of the dispersion (expressed as SD) of prevalence of x across clusters