





RESEARCH ARTICLE

# Minimum sample size for developing a multivariable prediction model: Part I – Continuous outcomes

Richard D. Riley<sup>1</sup>  | Kym I.E. Snell<sup>1</sup> | Joie Ensor<sup>1</sup>  | Danielle L. Burke<sup>1</sup>  |  
Frank E. Harrell Jr<sup>2</sup> | Karel G.M. Moons<sup>3</sup> | Gary S. Collins<sup>4</sup> 

<sup>1</sup>Centre for Prognosis Research, Research Institute for Primary Care and Health Sciences, Keele University, Keele, UK

<sup>2</sup>Department of Biostatistics, Vanderbilt University School of Medicine, Nashville, TN

<sup>3</sup>Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht, The Netherlands

<sup>4</sup>Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, University of Oxford, Oxford, UK

## Correspondence

Richard D Riley, Centre for Prognosis Research, Research Institute for Primary Care and Health Sciences, Keele University, Keele, Staffordshire ST5 5BG, UK.  
Email: r.riley@keele.ac.uk

## Funding information

National Institute for Health Research School for Primary Care Research (NIHR SPCR); Netherlands Organisation for Scientific Research, Grant/Award Number: project 9120.8004 and 918.10.615; CTSA, Grant/Award Number: UL1 TR002243; National Center for Advancing Translational Sciences; US National Institutes of Health; NIHR Biomedical Research Centre

In the medical literature, hundreds of prediction models are being developed to predict health outcomes in individuals. For continuous outcomes, typically a linear regression model is developed to predict an individual's outcome value conditional on values of multiple predictors (covariates). To improve model development and reduce the potential for overfitting, a suitable sample size is required in terms of the number of subjects ( $n$ ) relative to the number of predictor parameters ( $p$ ) for potential inclusion. We propose that the minimum value of  $n$  should meet the following four key criteria: (i) small optimism in predictor effect estimates as defined by a global shrinkage factor of  $\geq 0.9$ ; (ii) small absolute difference of  $\leq 0.05$  in the apparent and adjusted  $R^2$ ; (iii) precise estimation (a margin of error  $\leq 10\%$  of the true value) of the model's residual standard deviation; and similarly, (iv) precise estimation of the mean predicted outcome value (model intercept). The criteria require prespecification of the user's chosen  $p$  and the model's anticipated  $R^2$  as informed by previous studies. The value of  $n$  that meets all four criteria provides the minimum sample size required for model development. In an applied example, a new model to predict lung function in African-American women using 25 predictor parameters requires at least 918 subjects to meet all criteria, corresponding to at least 36.7 subjects per predictor parameter. Even larger sample sizes may be needed to additionally ensure precise estimates of key predictor effects, especially when important categorical predictors have low prevalence in certain categories.

## KEYWORDS

continuous outcome, linear regression, minimum sample size, multivariable prediction model, R-squared

## 1 | INTRODUCTION

Each year in the medical literature, hundreds of prediction models are developed to predict health outcomes in individuals.<sup>1–3</sup> Such models estimate an individual's predicted risk (for binary or categorical outcomes) or their expected outcome value (for a continuous outcome), conditional on the individual's observed value of multiple predictors. In this article, we focus on multivariable prediction models for continuous outcomes (eg, blood pressure, birth weight, depression score), which are typically developed using linear regression. This provides an equation containing an intercept term and multiple predictor effects (corresponding to mean differences), which is then used in new individuals to predict their

expected outcome value. Predictors (also known as variables, covariates, or prognostic factors<sup>4</sup>) typically include standard characteristics, such as age and stage of disease, or increasingly, biomarkers and genetic information.

Prediction models for continuous outcomes can potentially inform healthcare decisions and patient management, for example, to help decide on treatment and monitoring strategies.<sup>1</sup> Therefore, when developing their model, researchers should strive to use high quality datasets that allow a reliable model to be produced. This includes ensuring that the dataset has a suitable sample size. In particular, the number of subjects should be large enough relative to the number of predictor parameters to be estimated; otherwise, overfitting may be a serious problem. Overfitting refers to when a model is capturing idiosyncrasies in the development data; this leads to optimism in predictive performance such that the apparent performance is too high for the underlying population from which the development sample is drawn.<sup>5</sup> For example, in the development dataset, the developed model's apparent proportion of variation explained ( $R^2$ ) will often be too high, and the model's predicted outcome values will often be too extreme (ie, pushed too far from the mean). Therefore, it is good practice to ensure that sample sizes are large enough to minimize this problem.<sup>3</sup>

In this article, we build on the previous work of Harrell et al<sup>3,6</sup> to propose how to calculate a suitable sample size for development of a prediction model using linear regression. Specifically, we suggest that the minimum sample size required should minimize the potential for overfitting (and therefore optimism) and ensure precise estimates of key model parameters. We propose four criteria, ie, (i) small optimism in predictor effect estimates; (ii) small absolute difference in the apparent and adjusted  $R^2$ ; (iii) precise estimation of the residual standard deviation; and (iv) precise estimation of the mean predicted outcome value (model intercept when predictors are mean-centered). The number of subjects that meets all four criteria provides the minimum sample size required for model development.

The paper outline is as follows. In Section 2, we provide formulae to calculate the sample size required to meet criterion (i) and (ii), conditional on the user prespecifying the number of predictor parameters ( $p$ ) and the model's anticipated proportion of variation explained ( $R^2$ ), as informed by previous studies. Criteria (iii) and (iv) are then described in Section 3, and we show how to calculate sample sizes that ensure a small margin of error in the estimates, such as within 10% of their true values. Section 4 then provides an example to illustrate the approach. Section 5 briefly mentions that additional criteria may be important, such as precise estimation of predictor effects, and Section 6 concludes with discussion. A subsequent paper extends the ideas to binary and time-to-event models, such as logistic and Cox regression.<sup>7</sup>

## 2 | SAMPLE SIZE REQUIRED TO MINIMIZE OVERFITTING AND OPTIMISM

To adjust for overfitting during model development, statistical methods for penalization of predictor effect estimates are available, where regression coefficients are shrunk toward zero from their usual estimated value (eg, from traditional maximum likelihood estimation). There are many options for shrinkage,<sup>8</sup> including a global shrinkage factor (sometimes referred to as a uniform shrinkage factor) that is derived and applied postestimation,<sup>9,10</sup> or more holistic options such as ridge regression, elastic net, and the Lasso, which operate during the estimation process.<sup>11,12</sup> However, the penalization factors used within these methods are often estimated with large uncertainty, which increases as the magnitude of overfitting increases. Van Houwelingen notes that, "... shrinkage works on the average but may fail in the particular unique problem on which the statistician is working."<sup>8</sup> Therefore, it is important to minimize the potential for overfitting. In this section, we outline how researchers can target a sample size ( $n$ ) to minimize the potential for overfitting in advance of model development. Our formula is motivated by the concept of a global shrinkage factor, and so we begin by introducing this.

### 2.1 | Global shrinkage factor

Consider a continuous outcome ( $Y_i$ ), for  $i = 1$  to  $n$  subjects (participants) in a study, to which we want to fit a linear regression model of the form

$$Y_i = \mu_i + e_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \cdots + e_i \quad (1)$$

$$e_i \sim N(0, \sigma^2).$$

Assume that the unknown parameters of the equation (ie, the  $\beta$ s and  $\sigma^2$ ) are estimated using the data, usually, via ordinary least squares or maximum likelihood estimation. The intercept term, ie,  $\alpha$ , is the true mean outcome value for individuals whose  $X$  values are all zero, and each  $X$  term denotes values of included predictors. For example,  $X_{1i}$  could be

the age of the subject in years,  $X_{2i}$  could be 1 for males and 0 for females, and so on. Each  $\beta$  denotes the change in mean outcome value (ie, the mean difference) for each 1-unit increase in the corresponding predictor, after adjusting for other predictors. The error term, ie,  $e_i$ , represents the residuals, and these are assumed to follow a normal distribution with a mean of zero and variance of  $\sigma^2$ .

After fitting this regression model using traditional methods (eg, ordinary least squares), to adjust for overfitting a global (uniform) shrinkage factor ( $S$ ) can be applied to all estimated predictor effects ( $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$ , etc). That is, when making predictions in new individuals, we can use the modified equation of

$$E(Y_i) = \alpha^* + S(\hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots), \quad (2)$$

where  $\alpha^*$  is the revised intercept, which is re-estimated to ensure the overall predicted mean agrees with the observed mean in the development dataset (for details on how to do this, see Table 1 and the work of Harrell<sup>3</sup>). Compared to the original (nonpenalized) model, this will shrink predicted values in new individuals away from the extremes and move them toward the mean.

Implementation of this global shrinkage approach requires  $S$  to be estimated. A popular approach postestimation is bootstrapping.<sup>13</sup> An alternative is to utilize the closed-form solution of Copas<sup>10,14</sup>

$$S_C = 1 - \frac{(p - 2)}{LR}, \quad (3)$$

which, for linear regression (see equation (8.5) in the Copas paper<sup>10</sup>), provides an unbiased estimate of the shrinkage factor. Here,  $p$  is the total number of predictor parameters (assumed  $\geq 2$ ) and LR is the likelihood ratio (chi-squared) statistic for the model, which can be defined as

$$LR = -2(\ln L_{\text{null}} - \ln L_{\text{model}}), \quad (4)$$

where  $\ln L_{\text{null}}$  is the log-likelihood of a model with no predictors (ie, intercept-only null model), and  $\ln L_{\text{model}}$  is the log-likelihood of the final developed model. In Section 2.2, we also show how LR can be expressed in terms of  $R^2$  (see Equation (5)).

This shrinkage estimate of  $S_C$  relates to a model developed without any variable selection procedure, and thus  $p$  represents the entire set of predictor parameters in the model. When selection procedures are used during model development,  $p$  will be closer to the number of parameters based on the entire set of candidate predictors (ie, all those considered for inclusion, regardless of whether they were all included in the final model).<sup>3</sup> Therefore, we generally define  $p$  as the total number of predictor parameters *considered* within the model development. Note that, if a predictor is categorical with three or more categories, or continuous and modeled as a nonlinear trend, then it will contribute two or more parameters. For these reasons, we refer to subjects per predictor parameter, rather than subjects per variable, in this article.

**TABLE 1** Example of global shrinkage applied to a traditional (nonpenalized) linear regression model for predicting systolic blood pressure at the end of treatment in hypertension patients

|                  | Developed Model | Final Model After Adjustment for Overfitting <sup>s</sup> |
|------------------|-----------------|---|
| <b>Intercept</b> | $\hat{\alpha}$  | $\alpha^*$  |
|                  | 28.096          | 39.057  |
| <b>Predictor</b> | $\hat{\beta}$   | $S\hat{\beta} = 0.928\hat{\beta}$                         |
| SBP at baseline  | 0.462           | 0.429   |
| DBP at baseline  | 0.411           | 0.381   |
| BMI              | 0.013           | 0.012   |
| Age              | 0.450           | 0.418   |
| Sex              | -2.050          | -1.902  |
| Treatment        | -17.807         | -16.525   |
| Smoker           | -2.082          | -1.932  |

<sup>s</sup>The revised intercept is obtained by  $\alpha^* = (1 - S_C)\bar{Y} + S_C\hat{\alpha}$ , where  $\bar{Y}$  is the mean outcome value in the development dataset and  $\hat{\alpha}$  is the estimated intercept from a traditional (nonpenalized) model. For further details, see the work of Harrell.<sup>3</sup> DBP, diastolic blood pressure; BMI, body mass index; SBP, systolic blood pressure.

**Example of a global shrinkage factor.**

For illustration, we used data from a randomized trial of 262 hypertension patients to develop a linear regression model for predicting systolic blood pressure (SBP) at the end of treatment.<sup>15</sup> We forced inclusion of seven predictors, ie, age, sex, treatment group (treatment/control), smoker (yes/no), body mass index (BMI), baseline SBP, and baseline diastolic blood pressure (DBP). These correspond to seven predictor parameters (ie,  $p = 7$ ). The model parameter estimates are shown in Table 1, and the LR statistic was 69.295. The corresponding global shrinkage factor estimate from Equation (3) is

$$S_C = 1 - \frac{(p-2)}{LR} = 1 - \frac{(7-2)}{69.295} = 0.928.$$

We also used 5000 bootstrap samples to estimate the global shrinkage factor (as described elsewhere<sup>5,16</sup>), and this gave a very similar value of 0.94. Furthermore, the adjusted  $R^2$  was 0.21, which is about 0.91 times the apparent  $R^2$  value of 0.23. Therefore, even though there was no automated predictor selection based on  $p$ -values, there is still some evidence of overfitting (and thus optimism in apparent predictive performance). For a more robust prediction of SBP in new individuals, Table 1 also shows the original beta coefficients multiplied by 0.928, which shrink the model's predictions toward the overall mean.

**2.2 | Shrinkage expressed in terms of sample size and  $R^2$** 

We now propose utilizing the Copas shrinkage factor, ie,  $S_C$ , to inform sample size calculations at the start of a study, ie, before individual participant data have been obtained. Specifically, we derive an expression that allows the researcher to identify the sample size and number of predictor parameters that gives an expected value of  $S_C$  close to 1 (eg, 0.9). Our approach specifically builds on the work of Harrell et al,<sup>3,6</sup> who shows how (after the development dataset is obtained and a model fitted including all predictors) the shrinkage estimate can inform whether to reduce the number of predictors (using so-called data reduction techniques). Our premise is the same, except we focus on calculating the *expected* shrinkage before data collection to inform sample size calculations for a new study.

We start by re-expressing  $S_C$  in terms of sample size ( $n$ ), number of predictor parameters ( $p$ ), and  $R^2$ , the proportion of variability explained. Let  $R_{app}^2$  be the apparent estimate of a prediction model's  $R^2$  in the dataset used to develop the model. That is,  $R_{app}^2 = 1 - (\hat{\sigma}_{model}^2 / \hat{\sigma}_{null}^2)$ , and thus  $0 \leq R_{app}^2 \leq 1$ . As shown elsewhere,<sup>17,18</sup> the LR statistic can be expressed in terms of the sample size and  $R_{app}^2$  as follows:

$$LR = -n \ln(1 - R_{app}^2). \quad (5)$$

Applying Equation (5) within Equation (3), the Copas shrinkage formula becomes

$$S_C = 1 + \frac{p-2}{n \ln(1 - R_{app}^2)}. \quad (6)$$

Equation (6) cannot be used to directly inform the sample size ( $n$ ) in advance of model fitting because  $R_{app}^2$  is a *postestimation* measure of model fit. However, an approximately unbiased (optimism-adjusted) estimate of the proportion of variation explained is  $R_{adj}^2$ ,<sup>10,19</sup> ie,

$$R_{adj}^2 = 1 - \left( (1 - R_{app}^2) \frac{(n-1)}{(n-p-1)} \right) = \frac{(n-1)R_{app}^2 - p}{(n-p-1)}, \quad (7)$$

and rearranging gives

$$R_{app}^2 = \frac{R_{adj}^2(n-p-1) + p}{(n-1)}. \quad (8)$$

Therefore, applying Equation (8) within Equation (6) provides

$$S_C = 1 + \frac{p-2}{n \ln \left( 1 - \left( \frac{R_{adj}^2(n-p-1) + p}{(n-1)} \right) \right)}. \quad (9)$$

Hence, we now have an expression for the expected shrinkage factor conditional on a particular  $R_{adj}^2$  and, crucially, the sample size ( $n$ ), and number of predictor parameters ( $p$ ). When studying Equation (9), we observe that the expected

shrinkage will decrease (ie,  $S_C$  will move closer to 1) as  $n$  increases, as  $p$  decreases, as  $n/p$  increases, and as  $R_{adj}^2$  increases. Therefore, shrinkage (overfitting) will be a larger concern in development datasets with a small number of subjects, a large number of predictor parameters (relative to the number of subjects), and when the proportion of variance explained by the model is low.

### 2.3 | Criterion (i): calculating sample size to ensure a shrinkage factor $\geq 0.9$

Recall, at the start of Section 2, we explained that it is important to minimize the potential for overfitting. Therefore, when designing a new model development study, we propose researchers should utilize Equation (9) to reveal the sample size ( $n$ ) needed to obtain a targeted value of  $S_C$ . We suggest using a value of  $S_C \geq 0.9$ , such that predictor effects would shrink by  $\leq 10\%$ , which represents small overfitting. This is in accordance with the work of Harrell,<sup>3</sup> who suggests that, if the shrinkage estimate “falls below 0.9, for example, we may be concerned with the lack of calibration the model may experience on new data.”

Although there is no closed-form solution for  $n$  based on Equation (9), an iterative process can be used to identify the value of  $n$  that gives the desired  $S_C$  conditional on a chosen  $p$  and  $R_{adj}^2$ . For example, to obtain an expected  $S_C$  of 0.9 for a hypothetical model with up to 30 predictor parameters and an anticipated  $R_{adj}^2$  of 0.7, a sample size of 206 subjects is required to meet criterion (i) as follows:

$$S_C = 1 + \frac{30 - 2}{206 \ln \left( 1 - \left( \frac{0.70(206-30-1)+30}{(206-1)} \right) \right)} = 0.90.$$

This equates to  $206/30 = 6.87$  subjects per predictor parameter; that is, about 6.87 subjects are required for each predictor parameter considered. If it was rather considered that up to 50 predictor parameters are needed for an anticipated  $R_{adj}^2$  of 0.7, then a sample size of 355 subjects is required to obtain an expected  $S_C$  of 0.9, corresponding to 7.10 subjects per predictor parameter. Hence, the number of subjects per predictor parameter changes depending on the number of predictor parameters considered.

In situations where the calculated sample size is considered unrealistic (eg, due to time and cost constraints),  $p$  could be lowered by reducing the number of candidate predictor parameters. For example, those predictors known from previous studies (or systematic reviews) to have predictive value could be prioritized, or two or more predictors could be combined into one, such as BMI instead of weight and height. Alternatively, after data collection, unsupervised learning techniques such as principal component analysis could be used, which are blind to the outcome values. In the aforementioned example, reducing the number of predictor parameters to 25 leads to a sample size of 169 subjects to meet criterion (i)

$$S_C = 1 + \frac{25 - 2}{169 \ln \left( 1 - \left( \frac{0.70(169-25-1)+25}{(169-1)} \right) \right)} = 0.90.$$

Thus, by removing five predictor parameters, the required sample size to meet criterion (i) is reduced by 37 subjects.

When using a more stringent shrinkage factor, say of 0.95, then the necessary sample size to meet criterion (i) will be increased substantially. For instance, in the previous example with 25 predictor parameters and an anticipated  $R_{adj}^2$  of 0.7, the sample size required is increased from 169 to 361 when increasing  $S_C$  from 0.90 to 0.95. Hence, over twice the sample size is needed to reduce the expected shrinkage from 10% to 5%. For this reason, we anticipate that an  $S_C$  of 0.90 will often be a pragmatic choice for criterion (i).<sup>3</sup>

### 2.4 | Criterion (ii): calculating sample size to ensure a small absolute difference in $R_{adj}^2$ and $R_{app}^2$

Criterion (i) focuses on shrinkage of predictor effects, which is a multiplicative measure of overfitting (ie, on the relative scale), and therefore Harrell also suggested to evaluate overfitting on the absolute scale.<sup>3</sup> To address this, our second criterion for minimum sample size is to ensure that the difference ( $\delta$ ) between  $R_{app}^2$  and  $R_{adj}^2$  is small, say  $\leq 0.05$ , such that the optimism in the developed model's apparent proportion of variance explained is small.

Utilizing Equation (8), the difference in  $R_{\text{app}}^2$  and  $R_{\text{adj}}^2$  can be written as follows:

$$\begin{aligned}\delta &= R_{\text{app}}^2 - R_{\text{adj}}^2 \\ &= \frac{R_{\text{adj}}^2 (n - p - 1) + p}{(n - 1)} - R_{\text{adj}}^2 \\ &= \frac{R_{\text{adj}}^2 (n - p - 1) + p - (n - 1) R_{\text{adj}}^2}{(n - 1)} \\ &= \frac{-p R_{\text{adj}}^2 + p}{(n - 1)} \\ &= \frac{p (1 - R_{\text{adj}}^2)}{(n - 1)}.\end{aligned}$$

After rearranging this solution, we find that, to meet criterion (ii), we require the number of subjects to be

$$n \geq 1 + \frac{p (1 - R_{\text{adj}}^2)}{\delta}, \quad (10)$$

where  $\delta$  is a small value, such as  $\leq 0.05$ . For example, returning to our hypothetical model with an anticipated  $R_{\text{adj}}^2$  of 0.7 and up to 30 potential predictor parameters, the sample size required to meet criterion (ii) is as follows:

$$n \geq 1 + \frac{p (1 - R_{\text{adj}}^2)}{\delta} = 1 + \frac{30 (1 - 0.7)}{0.05} = 181.$$

This is slightly lower than the sample size of 206 identified for criterion (i) in Section 2.3.

Equation (10) reveals that, similar to the shrinkage approach for criterion (i), the required sample size for criterion (ii) will increase as  $p$  increases and  $R_{\text{adj}}^2$  decreases. For example, if our hypothetical model had an anticipated  $R_{\text{adj}}^2$  of 0.3 rather than 0.7, then the sample size to meet criterion (ii) increases substantially to 421.

## 2.5 | How to prespecify $R_{\text{adj}}^2$

To identify a sample size to meet our criteria (i) and (ii), researchers have to prespecify a value for the model's anticipated  $R_{\text{adj}}^2$ . How should this be done? We recommend identifying previous prediction model studies for the same or similar populations and outcomes of interest and extracting their  $R_{\text{adj}}^2$  values, which are usually well reported for linear regression models. Helpful for this purpose are systematic reviews of existing models<sup>20</sup> and registries that record the prediction models available in a particular field.<sup>21</sup> If only an  $R_{\text{app}}^2$  value is reported in a model development study, then its  $R_{\text{adj}}^2$  can be derived using Equation (7) as long as the study's  $n$  and  $p$  can also be obtained. Note that, if  $R_{\text{app}}^2$  is reported from an external validation study of an existing model, there is no need for adjustment (ie,  $R_{\text{app}}^2 = R_{\text{adj}}^2$ ), as the validation dataset provides a direct estimate of the model's performance in new individuals (free from overfitting concerns as there is no model development therein). In other words, deriving  $R_{\text{adj}}^2$  based on a reported  $R_{\text{app}}^2$  is necessary when the latter is from a model development study, but not when it is from an appropriate external validation study of an existing model.<sup>5</sup> Guidance for choosing an  $R_{\text{adj}}^2$  value in the absence of any prior information is given in the Discussion section.

## 3 | SAMPLE SIZE REQUIRED FOR PRECISE ESTIMATION OF THE RESIDUAL STANDARD DEVIATION AND MEAN PREDICTED OUTCOME VALUE

In addition to reducing the potential for overfitting, Harrell noted that sample sizes should be large enough to precisely estimate key model parameters such as the intercept or residual variance.<sup>3</sup> We now address this.



### 3.1 | Criterion (iii): precise estimate of the residual standard deviation

A precise estimate of the residual standard deviation is essential, as it is subsequently used to estimate  $R^2$ , and also to derive the standard errors and confidence intervals for the intercept and predictor effects (betas). For simplicity, Harrell suggested focusing on the standard deviation, ie,  $\sigma_{\text{null}}$ , say, in a null model (ie, intercept-only model),<sup>3</sup> and ensuring the lower and upper bounds of a 95% confidence interval for  $\sigma_{\text{null}}$  have a small multiplicative margin of error (MMOE) around the true value of  $\sigma_{\text{null}}$ . Assuming residuals are approximately normally distributed, this approach can be extended to consider the MMOE for estimating  $\sigma_{\text{model}}$ , the residual standard deviation in the developed prediction model, by

$$\text{MMOE} = \sqrt{\max \left( \frac{\chi^2_{1-\frac{\alpha}{2}, n-p-1}}{n-p-1}, \frac{n-p-1}{\chi^2_{\frac{\alpha}{2}, n-p-1}} \right)}, \quad (11)$$

where  $\chi^2_{1-\frac{\alpha}{2}, n-p-1}$  and  $\chi^2_{\frac{\alpha}{2}, n-p-1}$  are the critical values of a  $\chi^2$  distribution with  $n-p-1$  degrees of freedom for which there is, respectively, a probability of  $1 - \frac{\alpha}{2}$  and  $\frac{\alpha}{2}$  of being less than the critical value. The second term within the bracket of Equation (11) will typically give the largest MMOE.

For example, consider that we wish to ensure (with 95% confidence) that the margin of error is within 20% of the true value, ie,  $1.0 \leq \text{MMOE} \leq 1.2$ . For a null model (ie, one containing no predictors), then Equation (11) reveals that a sample size of at least 70 subjects is needed to meet this criterion,<sup>3</sup> as this gives an MMOE of 1.2. Therefore, in a multivariable model with  $p$  predictor parameters, the minimum sample required to meet an  $\text{MMOE} \leq 1.2$  for criterion (iii) is simply  $70 + p$ .

However, we recommend a more stringent margin of error of within 10% of the true value, ie,  $1.0 \leq \text{MMOE} \leq 1.1$ . In a null model, Equation (11) reveals that a sample size of at least 234 subjects is needed to ensure an  $\text{MMOE} \leq 1.1$ . Therefore, in a multivariable model with  $p$  predictor parameters, the minimum sample required to meet an  $\text{MMOE}$  of  $\leq 1.1$  for criterion (iii) is simply  $234 + p$ .

To illustrate this, let us return to the hypothetical example initially described in Section 2.3, where the sample size to meet criterion (i) was 206 and for criterion (ii) was 181. However, these values correspond to an MMOE of greater than 1.1 for  $\sigma_{\text{model}}$ . For example, taking the sample size of 206 subjects and 30 predictor parameters, this corresponds to a  $\chi^2_{0.975, 206-30-1} = 213.52$  and  $\chi^2_{0.025, 206-30-1} = 140.26$ , and thus (using Equation (11)) an MMOE of  $\sqrt{(206 - 30 - 1) / 140.26} = 1.12$ , and 12% potential margin of error in the estimate of  $\sigma_{\text{model}}$ .

Rather, with 30 predictor parameters, to achieve an expected MMOE of 1.1, we require  $234 + p = 234 + 30 = 264$  subjects because the maximum value of Equation (11) is then exactly 1.10

$$\sqrt{(264 - 30 - 1) / \chi^2_{0.025, 264-30-1}} = \sqrt{(264 - 30 - 1) / 192.615} = 1.10.$$

Hence, in this example, the minimum sample size of 264 subjects for criterion (iii) (ie, to ensure an  $\text{MMOE} \leq 1.1$ ) is more stringent than those identified for criteria (i) and (ii).

### 3.2 | Criterion (iv): precise estimate of the mean predicted outcome value (model intercept)

It is also important for model predictions to be precise; in particular, it is fundamental that the mean predicted outcome value is precisely estimated. If we assume our model will include predictors centered at their mean values in the developed dataset, then the fitted model's intercept ( $\hat{\alpha}_{\text{model}}$ ) will correspond to the predicted outcome value for an individual with mean predictor values. This estimate will be similar (though not identical) to the overall mean outcome in the population of interest; such a population mean estimate has variance of  $\hat{\sigma}_{\text{null}}^2 / n$ . However, in a linear regression model with multiple predictors the residual variance is  $\hat{\sigma}_{\text{model}}^2$ , and so the fitted model's intercept will have an approximate variance of\*

$$\text{var}(\hat{\alpha}_{\text{model}}) = \hat{\sigma}_{\text{model}}^2 / n \approx \sigma_{\text{null}}^2 (1 - R_{\text{adj}}^2) / n.$$

\*We use  $R_{\text{adj}}^2$  in this equation rather than  $R_{\text{app}}^2$  to be conservative.

Then, a 95% confidence interval for the model intercept is

$$\hat{\alpha}_{\text{model}} \pm \left( t_{1-\frac{0.05}{2}, n-p-1} \sqrt{\frac{\sigma_{\text{null}}^2(1-R_{\text{adj}}^2)}{n}} \right), \quad (12)$$

where  $t_{1-\frac{0.05}{2}, n-p-1}$  is the critical value of a  $t$ -distribution with  $n-p-1$  degrees of freedom for which there is a probability of  $1 - \frac{0.05}{2}$  below the critical value. Therefore, to derive this confidence interval in advance of model development, the researcher needs to prespecify (eg, from previous studies) sensible values for the anticipated mean outcome value ( $\hat{\alpha}_{\text{model}}$ ), the population (null model) variance ( $\sigma_{\text{null}}^2$ ), and  $R_{\text{adj}}^2$ . Then, the researcher can identify the sample size that ensures a sufficiently narrow confidence interval for  $\alpha_{\text{model}}$  to satisfy criterion (iv). For example, they might ensure the lower and upper bounds are within a small MMOE of the anticipated prediction mean (ie,  $1.0 \leq \text{MMOE} \leq 1.1$ ).

However, what constitutes a sufficiently narrow confidence interval will be context specific. A sensible start point is to examine the confidence interval width when using the sample sizes identified for criterion (i) to (iii). For example, let us return to our hypothetical model with an anticipated  $R_{\text{adj}}^2$  of 0.7 and  $p = 30$  predictor parameters, and now assume that the target population has an anticipated mean blood pressure ( $\hat{\alpha}_{\text{model}}$ ) of 165 and variance ( $\hat{\sigma}_{\text{null}}^2$ ) of  $18^2$ . Then, using Equation (12) and a sample size of 264 subjects identified by criterion (iii), the 95% confidence interval for the mean predicted outcome value is

$$165 \pm \left( t_{0.975, (264-30-1)} \sqrt{\frac{18^2(1-0.7)}{264}} \right) = 163.8 \text{ to } 166.2.$$

This is reassuringly precise, with the upper bound just 1.2 higher than the true mean of 165; this corresponds to a margin of error within 10% of the true mean (indeed, MMOE is  $166.2/165 = 1.007$ , and thus margin of error  $< 1\%$ ).

## 4 | WORKED EXAMPLE: PREDICTION OF LUNG FUNCTION IN AFRICAN-AMERICANS

A step-by-step summary of our sample size proposal is given Figure 1, and we now apply it to a worked example. Kumar et al use linear regression to identify predictors of lung function (ie, the forced expiratory volume in 1 second, FEV1) in African-American participants.<sup>22</sup> Let us assume that we want to build on this work by formally developing a linear regression model to predict FEV1 in African-American women. The aim could be to flag those individuals with low FEV1 values, as these are at risk of chronic obstructive pulmonary disease. We now go through the sample size calculation process.

### 4.1 | Step-by-step application

#### **Steps 1 and 2: Identifying a value (lower bound) for the model's anticipated $R_{\text{adj}}^2$ and choosing $p$ .**

Kumar et al (see their supplementary material<sup>22</sup>) report the performance of a model for prediction of FEV1 in women, with the model containing three predictors (age, height, African ancestry) and four predictor parameters. The  $R_{\text{adj}}^2$  was on average about 0.2 across three different datasets. Therefore, we could use this value as a lower bound for the anticipated  $R_{\text{adj}}^2$  in the new model. Furthermore, let us assume that there will be up to 25 predictor parameters in this new model (including the four used in the original model), and thus  $p = 25$ .

#### **Step 3: Criterion (i) – ensuring $S_C$ is close to 1.**

Based on the chosen  $R_{\text{adj}}^2 = 0.2$  and  $p = 25$ , to ensure an expected  $S_C$  of 0.9, a sample size of 918 subjects is needed because (using equation (9))

$$\begin{aligned} S_C &= 1 + \frac{p-2}{n \ln \left( 1 - \left( \frac{R_{\text{adj}}^2(n-p-1)+p}{(n-1)} \right) \right)} \\ &= 1 + \frac{25-2}{918 \ln \left( 1 - \left( \frac{0.2(918-25-1)+25}{(918-1)} \right) \right)} \\ &= 0.90. \end{aligned}$$

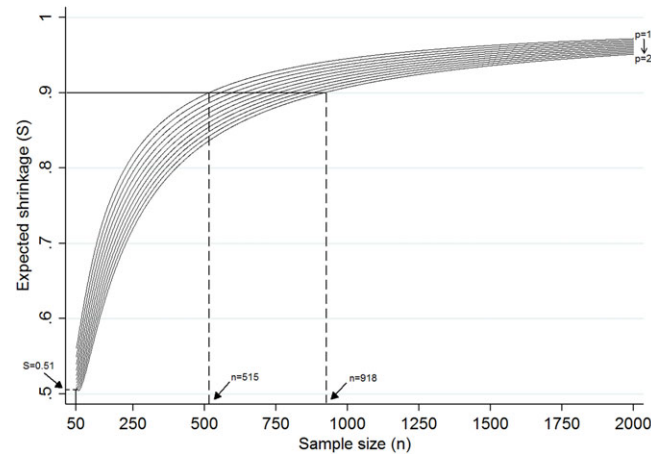


- **STEP 1: Choose the number of candidate predictors of interest for inclusion in the model, and calculate the corresponding number of predictor parameters ( $p$ ).** Recognise that one predictor may require two or more parameters; for example, a  $k$  category predictor requires  $k-1$  parameters, and a continuous predictor modelled with a non-linear trend requires  $> 1$  parameter to be estimated. Also include any potential interaction terms toward the total  $p$ .
  - **STEP 2: Choose a value for the anticipated proportion of variance explained ( $R^2$ ) for the new model** (referred to as  $R^2_{adj}$  in our article). For example, the  $R^2_{adj}$  value for a previously published model in the same setting and population could be used as a lower bound for the anticipated  $R^2$  of the new model.
  - **STEP 3: Criterion (i)** – use equation (9) to calculate the minimum sample size required to ensure Copas' global shrinkage factor ( $S_C$ ) is close to 1. We generally recommend a value of  $S_C \geq 0.90$ , which reflects a small amount of overfitting during model development.
  - **STEP 4: Criterion (ii)** – use equation (10) to calculate the minimum sample size required to ensure a small absolute difference of  $\leq 0.05$  in the developed model's  $R^2_{adj}$  and  $R^2_{app}$ .
  - **STEP 5: Criterion (iii)** – based on equation (11), calculate the minimum sample size required to ensure a precise estimate of the residual standard deviation ( $\sigma_{model}$ ). We generally recommend at least  $234 + p$  subjects, which ensures the  $\sigma_{model}$  estimate has no more than a 10% margin of error from the true value.
  - **STEP 6: Criterion (iv)** – based on equation (12), calculate the minimum sample size required to ensure a precise estimate of the mean predicted outcome value in the developed model (precise model intercept in a model with mean-centred predictors). This requires the researcher to pre-specify (e.g. from previous studies) sensible values for the anticipated mean outcome value ( $\hat{\alpha}_{model}$ ) and the population (null model) variance ( $\sigma^2_{null}$ ), in addition to  $R^2$ . What constitutes a precise estimate is context specific, but a broad suggestion is to at least ensure a confidence interval with lower and upper values within a 10% multiplicative margin of error from the true mean.
  - **STEP 7: Final sample size** - the required minimum sample size is the maximum value from steps 3 to 6, to ensure that each of criteria (i) to (iv) are met. Researchers might also examine whether the sample size would give precise estimates of key predictor effects (see Section 5).
- If the calculated sample size is not considered achievable due to criteria (i), (ii) or (iv), consider reducing the number of candidate predictors (and thus  $p$ ) to reduce the required sample size (whilst still meeting criterion (iii)). For example, prioritise those predictors identified as important from existing systematic reviews, or consider data reduction techniques such as principal component analysis (blinded to predictor-outcome associations in the development dataset). We do not recommend reducing the size of  $S_C$  or increasing  $R^2_{adj}$ .

**FIGURE 1** Summary of the steps involved in our sample size calculation for developing a multivariable prediction model

This corresponds to requiring 36.7 subjects per predictor parameter to meet criterion (i).

Figure 2 shows how the expected shrinkage ( $S_C$ ) derived from Equation (9) changes according to  $n$  and  $p$ , conditional on an  $R^2_{adj}$  of 0.2. As  $n$  increases and  $p$  decreases, the expected  $S_C$  becomes closer to 1. Furthermore, for an  $S_C$  above 0.9, very large increases in the sample size are needed to improve the expected  $S_C$ . For example, if we wanted to use a more stringent criteria for low overfitting of  $S_C = 0.95$ , then a sample size of 1949 subjects is required (78 subjects per predictor parameter), which is over double the number when  $S_C$  is 0.9.



**FIGURE 2** Expected shrinkage ( $S_C$ ) from Equation (9) for a prediction model of lung function in African-Americans, conditional on a particular sample size ( $n$ ), number of predictor parameters ( $p$ ), and an assumed  $R^2_{adj}$  of 0.2 [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

**Step 4: Criterion (ii) – ensuring small absolute difference between  $R^2_{adj}$  and  $R^2_{app}$ .**

Using Equation (10), we can calculate the minimum sample size needed to ensure the absolute difference between  $R^2_{adj}$  and  $R^2_{app}$  is 0.05 or less

$$n \geq 1 + \frac{p(1 - R^2_{adj})}{\delta} = 1 + \frac{25(1 - 0.2)}{0.05} = 401.$$

Therefore, at least 401 subjects are required to ensure a small absolute magnitude of overfitting based on the difference between  $R^2_{adj}$  and  $R^2_{app}$ .

**Step 5: Criterion (iii) – ensuring a precise estimate of the residual standard deviation.**

When  $p = 25$ , then Equation (11) reveals that we need at least 259 subjects ( $= 234 + p$ ) to ensure a margin of error of  $\leq 10\%$  in the estimate of the model's residual standard deviation, ie,  $\sigma_{model}$ . This is because the maximum value of Equation (11) is as follows:

$$\sqrt{(259 - 25 - 1) / \chi^2_{0.025, 259 - 25 - 1}} = \sqrt{(259 - 25 - 1) / 192.615} = 1.10.$$

**Step 6: Criterion (iv) – ensure a precise estimate of the mean predicted outcome (model intercept).**

Based on the results in Kumar et al,<sup>22</sup> the population mean and variance of FEV1 are about 1.90 liters and 0.6,<sup>2</sup> respectively. Based on these values and using Equation (12) with an assumed sample size of 918 subjects (as needed for criterion (i)), we obtain a 95% confidence interval for the predicted mean outcome value (model intercept when predictors are mean-centered) of

$$1.90 \pm \left( t_{0.975, (918 - 25 - 1)} \sqrt{\frac{0.6^2 (1 - 0.2)}{918}} \right) = 1.87 \text{ to } 1.93.$$

This is reassuringly precise, and the upper bound indicates an MMOE  $< 1.1$  (ie, within 10%) of the true mean outcome value.

**Step 7: Identify sample size that ensures all criteria are met.**

Based on the largest sample size calculations identified in steps 3 to 6, the final minimum sample size required is 918 subjects. This is driven by criterion (i) to ensure an expected shrinkage factor  $\geq 0.9$ .

## 4.2 | What if the sample size is not considered achievable?

If the sample size of 918 subjects was not considered achievable (eg, due to time or cost constraints), then what should be done? For criterion (i), we do not recommend reducing  $S_C$  below 0.9, as our main premise is to minimize overfitting. That leaves two other potential options; either use a larger  $R^2_{adj}$  value or reduce  $p$ . We do not recommend the first option. The  $R^2_{adj}$  values reported in previous articles are themselves only estimates, and it is hard to judge how much  $R^2_{adj}$  will be

improved in a new model. It is far better to be conservative in the choice of  $R_{adj}^2$  and adopt larger sample sizes, rather than naively aiming for smaller sample sizes that ultimately may not be fit for purpose.

Therefore, the best approach is to reduce the number of candidate predictor parameters, ie,  $p$ . Returning to the hypothetical example where the assumed  $R_{adj}^2$  is 0.2, a reduced set of 15 predictor parameters would lower the sample size required for criterion (i) to 515 subjects (Figure 2). Thus, after sacrificing 10 parameters by removing some predictors, the researchers requires 403 fewer subjects to target an  $S_C$  of 0.9. The choice of which predictors to prioritize could be based on external evidence (eg, from systematic reviews) and, after data collection, data reduction techniques such as principal components analysis (which are based on observed correlation among predictors only, and not observed predictor-outcome associations). All those predictors within the existing Kumar et al model are best retained in the model to justify the assumption that  $R_{adj}^2$  is at least 0.2. A sample size of 515 subjects still ensures that criteria (ii) to (iv) are met. For example, for criterion (iv), the width of the 95% confidence interval for the prediction mean would still be very narrow (1.85 to 1.95).

### 4.3 | Comparison to other suggested sample size proposals

We now contrast our derived sample size of 918 subjects to those from two other suggested sample size approaches for linear regression models. Though not intended for informing prediction model development, a recent recommendation suggests two subjects per predictor parameter for adequate estimation of predictor effects in linear regression.<sup>23</sup> In our example, two subjects for each of the 25 predictor parameters leads to a substantially smaller sample size of 50 subjects. However, using Equation (9), this corresponds to an expected shrinkage of  $S_C = 0.51$ , which reflects substantial overfitting and does not meet criterion (i) (Figure 2).

Alternatively, Harrell suggests that there are at least 15 subjects per predictor parameter (see Chapter 4 in his book<sup>3</sup>), which in this example implies a sample size of at least 375. However, using Equation (9), a sample size of 375 subjects corresponds to an expected shrinkage of  $S_C = 0.79$ , which still suggests large overfitting.

## 5 | POTENTIAL ADDITIONAL CRITERIA

Criteria (i) to (iv) form our main proposal for the minimum sample size required when developing a prediction model for continuous outcomes. However, we now briefly mention two additional criteria that may also be important to consider.

### 5.1 | Ensuring precise estimation of $R_{adj}^2$ and the mean-square error

Criterion (ii) ensures that there is a small absolute difference between  $R_{adj}^2$  and  $R_{app}^2$  to reflect low overfitting. A related concept is to ensure a precise confidence interval for  $R_{adj}^2$ .<sup>24</sup> Tan gives an excellent overview of various exact and approximate approaches to calculate a confidence interval for  $R_{adj}^2$ ,<sup>25</sup> given a developed model's  $R_{app}^2$ ,  $n$ , and  $p$ . For example, Lee proposes a confidence interval based on a scaled noncentral F distribution approximation to the distribution of  $R^2$ .<sup>26</sup> This can be implemented in SAS,<sup>27</sup> or in R using the ci.R2 function of the MBESS package by Kelley.<sup>28-30</sup> Furthermore, the ss.aipe.R2 function within MBESS identifies the sample size required to ensure Lee's confidence interval is sufficiently narrow.

We applied the ss.aipe.R2 function to the lung function model described in Section 4. This identified that 835 subjects are required to ensure the expected width of the confidence interval for  $R_{adj}^2$  is exactly 0.10, assuming  $R_{adj}^2$  is 0.20 and  $p = 25$ . This sample size is lower than the 918 subjects required to meet criterion (i), and hence 918 subjects is still the minimum sample size required.

Ensuring precise estimates of  $R_{adj}^2$  and the residual standard deviation (criteria (iii)) also helps ensure a precise estimate of the mean-square error (MSE) of the model's predicted outcome values, as  $MSE = \sigma_{null}^2(1 - R_{adj}^2)$ .

### 5.2 | Ensuring precise estimation of key predictor effects

Criterion (iv) ensures a precise estimate of the mean predicted outcome value in the entire target population. Ideally, predictions should also be precise across the entire spectrum of predicted values, not just at the mean. This is challenging

but is helped by ensuring the effects of key predictors are estimated precisely. The precision of a particular predictor effect in a fitted linear regression model depends on the sample size, the estimated residual variance, the correlation of the predictor with other included predictors, and the variance of the predictor values.<sup>31</sup> For brevity, we do not consider this in detail here and refer the reader to other articles that focus on this.<sup>31–35</sup> In particular, the `ss.aipe.rc` function with the MBESS package identifies the sample size required to ensure the confidence interval around a predictor's effect is sufficiently narrow.<sup>30,33,35</sup>

Returning to the lung function example, let us consider that our new model will potentially include smoking as a predictor, defined as a binary variable (current/previous smokers versus nonsmokers). Furthermore, assume that the mean difference in FEV1 for smokers and nonsmokers is  $-0.5$ , and that (based on the work of Kumar et al<sup>22</sup>) 50% of subjects will be current/previous smokers. Moreover, assume (conservatively) that the final model will have  $R^2_{\text{adj}}$  of 0.2 and that the correlation is 0.5 between smoking and other included predictors. Using the `ss.aipe.rc` function in R, we identify that 619 subjects are required to ensure a confidence interval width of 0.2 (and thus the lower and upper bounds are within 0.1 of the true value of  $-0.5$ ). This is reassuring, and again, less than the 918 subjects are required to satisfy criterion (i).

Precise estimation of predictor effects will be especially difficult for those predictors with the smallest variance in their values, as their confidence intervals are likely to be the widest.<sup>31</sup> In particular, categorical predictors with low prevalence in certain categories have small variances.<sup>36</sup> In our example, had we assumed that the percentage of smokers was 10%, rather than 50%, then repeating the calculation identifies that 1668 subjects are needed for a confidence interval width of 0.2. In this situation, we would need to increase the sample size beyond 918 subjects previously identified to meet criteria (i) to (iv) or justify relaxing the magnitude of precision desired. For example if we were willing to widen the expected confidence interval width to 0.3, then this considerably reduces the number of the required subjects to 757, but the interval is still fairly precise and all well below zero ( $-0.65$  to  $-0.35$ ).

## 6 | DISCUSSION

Sample size calculations are a fundamental part of designing a study to develop a new prediction model. In this article, we proposed four criteria to identify the minimum sample size needed to minimize overfitting while ensuring precise estimates of key model parameters. Criterion (i) forms the most novel aspect of our sample size proposal, as it allows researchers to identify  $n$  and  $p$  that correspond to an expected shrinkage factor close to 1, such as 0.9, which reflects low overfitting. Furthermore, it allows the sample size to be tailored to each model of interest through the prespecification of the anticipated proportion of variation explained, ie,  $R^2_{\text{adj}}$ , which is a measure of overall model fit. The chosen value of  $R^2_{\text{adj}}$  strongly influences the amount of shrinkage required, with larger values requiring less shrinkage (with other things, such as  $p$ , being equal).<sup>37</sup> This issue is currently ignored when using blanket rules of thumb for sample size.

Researchers should use previous evidence from other prediction model's in the same setting to ascertain a (conservative) value for the new model's potential  $R^2_{\text{adj}}$  value. If no relevant prediction models exist, then information from predictor finding studies (ie, studies aiming to estimate the prognostic effect of a particular predictor adjusted for other existing factors<sup>4</sup>) might be relevant. Even though such studies are primarily focused on the estimation of the effect of a particular predictor, they typically involve multivariable modeling and therefore often also report  $R^2_{\text{app}}$  and  $R^2_{\text{adj}}$  values. Where truly no prior information exists about the potential  $R^2_{\text{adj}}$  value, researchers should recognize that medical diagnosis and prediction of health-related outcomes are, generally speaking, low signal:noise ratio situations. It is not uncommon in these situations to see  $R^2_{\text{adj}}$  values in the 0.1 to 0.2 range. Therefore, in the absence of other information, we suggest that sample sizes be derived assuming that  $R^2_{\text{adj}} = 0.15$ . An exception is when predictors include “direct” (mechanistic) measurements, such as the baseline version of the continuous outcome (eg, when predicting lung function one year after measuring baseline lung function). Then, in this special situation, an  $R^2_{\text{adj}} = 0.5$  may be a more appropriate default choice.

In practice, *after* a model development dataset is obtained, a better approach for estimating the shrinkage factor is to use a resampling approach such as bootstrapping.<sup>10,16</sup> However, as our sample size calculations are focused on situations before any data collection, it is not possible to incorporate such a resampling approach. In situations where a development dataset is already available, containing a specific number of subjects and predictors, our approach could be used to identify whether a reduction in the number of predictors is needed (prior to beginning the modeling). Indeed, Harrell previously illustrated this concept by using the shrinkage estimate from the full model (including all predictors) to gauge whether the number of predictors should be reduced.<sup>3</sup> This could then incorporate bootstrapping (rather than the Copas formula) to estimate the shrinkage. However, this should be done blind to the estimated predictor effects, as otherwise

the decisions about inclusion are already being made based on the full set of predictors. Similarly, when planning to use a predictor selection method (such as backwards selection) during model development, researchers should define  $p$  as the total number of parameters due to all predictors considered (screened) and not just the subset that are included in the final model.<sup>5</sup> As Harrell notes,<sup>3</sup> the value of  $p$  should be honest.

A potential limitation of our work is that multiple sample size calculations are required to address each of the criteria considered. However, this reflects the different elements that require consideration when developing a prediction model. Criteria (iii) and (iv) are needed to ensure that there will be a small margin of error in the estimates of the residual standard deviation and the mean predicted outcome value (model intercept). This is often overlooked when considering the sample size. In particular, at least  $234 + p$  subjects are always required to ensure an MMOE of  $\leq 1.1$  for estimating the model's residual standard deviation, ie,  $\sigma_{\text{model}}$ . Section 5 emphasized that further criteria may also be needed in some settings. In particular, ensuring precise estimates of predictor effects may be important, especially in settings where key predictors have low variance (eg, categorical predictors with few subjects in certain categories).

In summary, we have proposed how to ascertain the minimum sample size needed to develop a prediction model using linear regression. We hope this encourages researchers to move away from rules of thumb and to rather focus on attaining sample sizes that ensure precise estimates and reduce the potential for overfitting to develop more robust prediction models. We are currently writing software modules to implement the approach. Our accompanying paper extends the work to binary and time-to-event outcomes.<sup>7</sup>

## ACKNOWLEDGEMENTS

We wish to thank three reviewers and an Associate Editor for their constructive comments, which helped improve the article upon revision.

## FUNDING

Danielle Burke and Kym Snell are funded by the National Institute for Health Research School for Primary Care Research (NIHR SPCR). The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health. Karel G.M. Moons receives funding from the Netherlands Organisation for Scientific Research (project 9120.8004 and 918.10.615). Frank Harrell's work on this paper was supported by CTSA award No. UL1 TR002243 from the National Center for Advancing Translational Sciences. Its contents are solely the responsibility of the authors and do not necessarily represent official views of the National Center for Advancing Translational Sciences or the US National Institutes of Health. Gary Collins is supported by the NIHR Biomedical Research Centre, Oxford.

## ORCID

Richard D. Riley  <http://orcid.org/0000-0001-8699-0735>

Joie Ensor  <http://orcid.org/0000-0001-7481-0282>

Danielle L. Burke  <http://orcid.org/0000-0003-2803-1151>

Gary S. Collins  <http://orcid.org/0000-0002-2772-2316>

## REFERENCES

1. PROGRESS Group. Prognosis research strategy (PROGRESS) 3: prognostic model research. *PLoS Med.* 2013;10(2):e1001381.
2. Steyerberg EW. *Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating*. New York, NY: Springer; 2009.
3. Harrell FE. Ordinal Logistic Regression. In: *Regression Modeling Strategies*. 2nd Edition. Cham, Switzerland: Springer; 2015:311-325. Springer Series in Statistics (SSS).
4. Riley RD, Hayden JA, Steyerberg EW, et al. Prognosis research strategy (PROGRESS) 2: prognostic factor research. *PLoS Med.* 2013;10(2):e1001380.
5. Moons KG, Altman DG, Reitsma JB, et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): explanation and elaboration. *Ann Intern Med.* 2015;162(1):W1-W73.
6. Harrell FE Jr, Lee KL, Mark DB. Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Statist Med.* 1996;15(4):361-387.
7. Riley RD, Snell KIE, Ensor J, et al. Minimum sample size for developing a multivariable prediction model: PART II - binary and time-to-event outcomes. *Statist Med.* 2018. <https://doi.org/10.1002/sim.7992>
8. Van Houwelingen JC. Shrinkage and penalized likelihood as methods to improve predictive accuracy. *Stat Neerlandica.* 2001;55:17-34.



9. Van Houwelingen JC, Le Cessie S. Predictive value of statistical models. *Statist Med*. 1990;9(11):1303-1325.
10. Copas JB. Regression, prediction, and shrinkage. *J Royal Statist Soc B*. 1983;45(3):311-354.
11. Tibshirani R. Regression shrinkage and selection via the lasso. *J Royal Statist Soc B*. 1996;58:267-288.
12. Pavlou M, Ambler G, Seaman SR, et al. How to develop a more accurate risk prediction model when there are few events. *BMJ*. 2015;351:h3868.
13. Efron B. Bootstrap Methods: Another Look at the Jackknife. In: Kotz S, Johnson NL. (eds) *Breakthroughs in Statistics*. New York, NY: Springer; 1992:569-593. Springer Series in Statistics (SSS).
14. Copas JB. Using regression models for prediction: shrinkage and regression to the mean. *Stat Methods Med Res*. 1997;6(2):167-183.
15. Riley RD, Kauser I, Bland M, et al. Meta-analysis of randomised trials with a continuous outcome according to baseline imbalance and availability of individual participant data. *Statist Med*. 2013;32(16):2747-2766.
16. Steyerberg EW, Harrell FEJ, Borsboom GJ, Eijkemans MJC, Vergouwe Y, Habbema JDF. Internal validation of predictive models: efficiency of some procedures for logistic regression analysis. *J Clin Epidemiol*. 2001;54:774-781.
17. Magee L.  $R^2$  measures based on Wald and likelihood ratio joint significance tests. *Am Stat*. 1990;44(3):250-253.
18. Hendry DF, Nielsen B. *Econometric Modeling: A Likelihood Approach*. Princeton, NJ: Princeton University Press; 2012.
19. Goldberger AS. *Econometric Theory*. New York, NY: Wiley; 1964.
20. Debray TP, Damen JA, Snell KI, et al. A guide to systematic review and meta-analysis of prediction model performance. *BMJ*. 2017;356:i6460.
21. Wessler BS, Paulus J, Lundquist CM, et al. Tufts PACE clinical predictive model registry: update 1990 through 2015. *Diagn Progn Res*. 2017;1(1):20.
22. Kumar R, Seibold MA, Aldrich MC, et al. Genetic ancestry in lung-function predictions. *N Engl J Med*. 2010;363(4):321-330.
23. Austin PC, Steyerberg EW. The number of subjects per variable required in linear regression analyses. *J Clin Epidemiol*. 2015;68(6):627-636.
24. Algina J, Olejnik S. Determining sample size for accurate estimation of the squared multiple correlation coefficient. *Multivar Behav Res*. 2000;35(1):119-137.
25. Tan L. *Confidence Intervals for Comparison of the Squared Multiple Correlation Coefficients of non-Nested Models* [dissertation]. Ontario, Canada: School of Graduate and Postdoctoral Studies, The University of Western Ontario; 2012.
26. Lee YS. Tables of the upper percentage points of the multiple correlation. *Biometrika*. 1972;59:175-189.
27. Zou GY. Toward using confidence intervals to compare correlations. *Psychol Methods*. 2007;12(4):399-413.
28. Kelley K. Confidence intervals for standardized effect sizes: theory, application, and implementation. *J Stat Softw*. 2007;20(8):1-24.
29. Kelley K. Methods for the Behavioral, educational, and social sciences: an R package. *Behav Res Methods*. 2007;39(4):979-984.
30. Kelley K. MBESS (Version 4.0.0 and higher) [computer software and manual]. 2018. <https://CRAN.R-project.org/package=MBESS>
31. McClelland GH. Increasing statistical power without increasing sample size. *Am Psychol*. 2000;55(8):963-964.
32. Hsieh FY, Bloch DA, Larsen MD. A simple method of sample size calculation for linear and logistic regression. *Statist Med*. 1998;17(14):1623-1634.
33. Kelley K, Maxwell SE. Sample size for multiple regression: obtaining regression coefficients that are accurate, not simply significant. *Psychol Methods*. 2003;8(3):305-321.
34. Kelley K. Sample size planning for the coefficient of variation from the accuracy in parameter estimation approach. *Behav Res Methods*. 2007;39(4):755-766.
35. Maxwell SE, Kelley K, Rausch JR. Sample size planning for statistical power and accuracy in parameter estimation. *Annu Rev Psychol*. 2008;59:537-563.
36. Ogundimu EO, Altman DG, Collins GS. Adequate sample size for developing prediction models is not simply related to events per variable. *J Clin Epidemiol*. 2016;76:175-182.
37. Knofczynski GT, Mundfrom D. Sample sizes when using multiple linear regression for prediction. *Educ Psychol Meas*. 2008;68(3):431-442.

**How to cite this article:** Riley RD, Snell KIE, Ensor J, et al. Minimum sample size for developing a multivariable prediction model: Part I – Continuous outcomes. *Statistics in Medicine*. 2019;38:1262–1275. <https://doi.org/10.1002/sim.7993>