

# Accepted Manuscript

Measurement error and timing of predictor values for multivariable risk prediction models are poorly reported

Rebecca Whittle, George Peat, John Belcher, Gary S. Collins, Richard D. Riley

PII: S0895-4356(18)30037-4

DOI: [10.1016/j.jclinepi.2018.05.008](https://doi.org/10.1016/j.jclinepi.2018.05.008)

Reference: JCE 9655

To appear in: *Journal of Clinical Epidemiology*

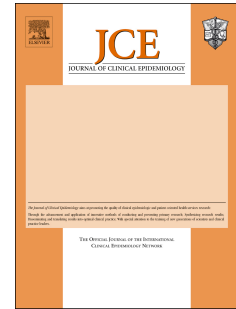
Received Date: 17 January 2018

Revised Date: 26 April 2018

Accepted Date: 14 May 2018

Please cite this article as: Whittle R, Peat G, Belcher J, Collins GS, Riley RD, Measurement error and timing of predictor values for multivariable risk prediction models are poorly reported, *Journal of Clinical Epidemiology* (2018), doi: 10.1016/j.jclinepi.2018.05.008.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



1 **Measurement error and timing of predictor values for multivariable risk**  
2 **prediction models are poorly reported**

3

4 Rebecca Whittle<sup>1</sup>, George Peat<sup>1</sup>, John Belcher<sup>1</sup>, Gary S. Collins<sup>2</sup>, Richard D. Riley<sup>1</sup>

5

6 <sup>1</sup> Centre for Prognosis Research, Arthritis Research UK Primary Care Centre, Research Institute for  
7 Primary Care & Health Sciences, Keele University, Keele, Staffordshire, UK (RW:

8 [r.l.whittle@keele.ac.uk](mailto:r.l.whittle@keele.ac.uk); GP: [g.m.peat@keele.ac.uk](mailto:g.m.peat@keele.ac.uk); JB: [j.belcher@keele.ac.uk](mailto:j.belcher@keele.ac.uk); RR:

9 [r.riley@keele.ac.uk](mailto:r.riley@keele.ac.uk))

10 <sup>2</sup> Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology and  
11 Musculoskeletal Sciences, Botnar Research Centre, University of Oxford, Oxford, UK (GC:

12 [gary.collins@csm.ox.ac.uk](mailto:gary.collins@csm.ox.ac.uk))

13

14 **Corresponding author:** Rebecca Whittle, Centre for Prognosis Research, Arthritis Research UK

15 Primary Care Centre, Research Institute for Primary Care & Health Sciences, Keele University, Keele,  
16 Staffordshire, UK. Tel: 00 44 1782 734848; Email: [r.l.whittle@keele.ac.uk](mailto:r.l.whittle@keele.ac.uk)

17

18

19 **ABSTRACT**

20 **Objective:** Measurement error in predictor variables may threaten the validity of clinical prediction  
21 models. We sought to evaluate the possible extent of the problem. A secondary objective was to  
22 examine whether predictors are measured at the intended moment of model use.

23 **Methods:** A systematic search of Medline was used to identify a sample of articles reporting the  
24 development of a clinical prediction model published in 2015. After screening according to a  
25 predefined inclusion criteria, information on predictors, strategies to control for measurement error  
26 and intended moment of model use were extracted. Susceptibility to measurement error for each  
27 predictor was classified into low and high risk.

28 **Results:** Thirty-three studies were reviewed, including 151 different predictors in the final prediction  
29 models. Fifty-one (33.7%) predictors were categorised as high risk of error, however this was not  
30 accounted for in the model development. Only 8 (24.2%) studies explicitly stated the intended  
31 moment of model use and when the predictors were measured.

32 **Conclusion:** Reporting of measurement error and intended moment of model use is poor in  
33 prediction model studies. There is a need to identify circumstances where ignoring measurement  
34 error in prediction models is consequential and whether accounting for the error will improve the  
35 predictions.

36 **Keywords:** Prediction models, prediction, prognosis, diagnosis, measurement error, error

37 **Word count:** 5416

38 **WHAT IS NEW?**39 **Key findings**

- 40 • Many published prediction models include predictors that are susceptible to measurement  
41 error and this measurement error is not being acknowledged or accounted for in the  
42 development of the models.
- 43 • Most prediction model articles do not explicitly state the intended moment of model use, or  
44 exactly when the predictors used in the model development were measured.

45 **What this adds to what is known**

- 46 • Reporting of measurement error and intended moment of model use is poor in prediction  
47 model studies.

48 **What is the implication, what should change now?**

- 49 • There is a need to identify circumstances where ignoring measurement error in prediction  
50 models is consequential and whether accounting for the error will improve the predictions.
- 51 • Future prediction model research studies must clearly report the intended moment of use of  
52 the prediction model, and be explicit about when the predictors were measured.

53

54

55

56 **BACKGROUND**

57 Predicting a patient's future outcome risk is an important part of medical research as it guides  
58 treatment, informs clinical decision making and helps patients understand their risk. Prognosis  
59 research can be used to help predict future outcomes in patients with a particular disease or health  
60 condition by developing a prediction model [1]. The number of articles reporting clinical prediction  
61 models has been increasing steadily over time, with approximately 500 articles published in 2011 [2],  
62 and these models utilise values of multiple predictors to enable individualised risk prediction [3].  
63 Such models are intended "to assist clinicians with their prediction of a patient's future outcome and  
64 to enhance informed decision making with the patient" [4]. Therefore, the predictions from these  
65 models should have optimal performance when being practically implemented at the "intended  
66 moment of using the model" [5].

67 However, when developing such models, measurement error may affect the observed predictor  
68 values, which could potentially lead to biased or incorrect estimates of predictor-outcome  
69 associations [6-9]. Measurement error is a difference between the measured values of a predictor  
70 and the true values of the predictor, or if the predictor is categorical, it is the classification to an  
71 incorrect category (misclassification). The term measurement error will be used throughout this  
72 article to refer generally to measurement error in continuous predictors and misclassification of  
73 categorical predictors. Measurement error is common within clinical studies, particularly  
74 observational studies [10], and has been found to be commonly neglected within the medical  
75 literature [11] Measurement error can occur for many different reasons such as biological variability,  
76 inaccuracy of measurement instruments, imperfect recall, cost or resource limitations, the  
77 subjective nature of measures, laboratory or measurer error and timing error. For example,  
78 measurement error in blood pressure commonly occurs due to biological variability [12]. Body mass  
79 index (BMI) is also commonly measured with error either due to the inaccuracy of measurement

80 instruments (i.e. the scales not being calibrated correctly), or due to imperfect recall by the patient,  
81 and this measurement error could then cause misclassification into an incorrect category.

82 Prognosis research is becoming increasingly more important [1], but there has been little research  
83 into the impact that measurement error in the predictors used to develop a prediction model may  
84 have, both in terms of the predictions made and model performance. It is also unclear how  
85 accounting for measurement error within the statistical modelling may improve this. A recent study  
86 demonstrated that measurement error in the predictors can dramatically reduce the c-statistic and  
87 increase the Brier score [13], and another study found that both random and systematic error in self-  
88 reported health data influences the calibration, discrimination and predicted risks [14], but in  
89 general the extent and impact of measurement error in prediction model research is often  
90 overlooked. However, the STRATOS (STRengthening Analytical Thinking for Observational Studies)  
91 initiative ([www.stratos-initiative.org](http://www.stratos-initiative.org)) have identified measurement error as a common issue in  
92 observational studies which is often ignored and for which guidance is needed. There is a vast  
93 amount of literature on the statistical effect of measurement error in general, but whether  
94 investigators consider measurement error when developing a prediction model, has not previously  
95 been evaluated. Models developed with predictors containing measurement error could therefore  
96 provide inaccurate estimates of patient risk and the model may not perform as well as expected in  
97 practice. A summary of the most commonly used methods to correct for measurement error is given  
98 by Brakenhoff et al [11] with more detailed reviews of these (and other) methods given by Carroll et  
99 al. [8] and Gustafson [9]. Several other methods that can be used to account for measurement error  
100 in the particular context of prediction research have been developed, including methods in a  
101 Bayesian framework, using an item response theory model to handle the measurement error [15]  
102 and bootstrap regression calibration [16], based on resampling techniques.

103 A particular aspect of measurement error in the predictors is timing error, so whether the predictors  
104 used in the model development were measured at the moment the model is intended to be used in

105 practice. When time-dependent predictors are not able to be measured at 'baseline' this creates  
106 time-dependent bias, which has been shown to often have an impact on the estimates of key  
107 predictors and study conclusions [17]. Additionally, the TRIPOD (Transparent Reporting of a  
108 multivariable prediction model for Individual Prognosis or Diagnosis) statement recommends to  
109 clearly define when the predictors used in the development of the model were measured [18] and  
110 states that "all predictors should be measured before or at the study time origin and known at the  
111 intended moment the model is intended to be used" [19]. Nevertheless, for a range of practical and  
112 ethical reasons, researchers may design prognosis studies that collect time-varying predictor  
113 information after the intended moment of use, which itself may lead to errors and misleading  
114 predictions [20].

115

116 The aim of this article is to present a systematic review of recent studies developing prediction  
117 models, to ascertain how susceptible to measurement error the predictors used in the final models  
118 are and how often the measurement error was acknowledged or accounted for within the  
119 development of the models. A secondary objective is to determine whether the predictors were  
120 measured at a different time point to the intended moment of using the prediction model.

121

## 122 **METHODS**

### 123 Data source and search

124 A systematic search was carried out in Medline on 27<sup>th</sup> November 2015 to identify the 30 most  
125 recent articles reporting the development of a multivariable prediction model for either  
126 individualised diagnosis or prognosis. It was decided a priori that approximately 30 articles would be  
127 sufficient in providing qualitative saturation of whether measurement error and incorrect timings  
128 was a general concern for the prediction model field.

129 The search strategy used was an adaptation of a published search string for finding prognostic and  
130 diagnostic prediction studies in Medline [21]. The search string was adapted by changing the term  
131 “OR ‘Multivariable’” to “AND (Multivariable OR Multivariate)” to refine the search further to studies  
132 developing multivariable prediction models for individualised prediction, which would hopefully  
133 remove other studies just examining associations between specific factors and an outcome but not  
134 developing a prediction model (see Supplementary Table 1 for the full search string).

135

### 136 Selection of articles

137 The titles and abstracts of the 1000 most recently published articles found using the search string  
138 were screened for inclusion, as we estimated this would return approximately 30 articles to be  
139 included. The full article was then obtained for any articles which were deemed to be potentially  
140 eligible or for any articles in which it was unclear from the title and abstract whether they met the  
141 eligibility criteria. These full articles were then screened for suitability and categorised into one of  
142 three groups: ‘include’, ‘exclude’ and ‘unsure’. The selection of articles until this stage was  
143 undertaken by a single reviewer (RW). Articles in the ‘include’ and the ‘unsure’ groups were sent to  
144 two additional reviewers (GP & RR). Both reviewers checked all ‘unsure’ articles, and the ‘include’  
145 articles were split between the two reviewers to check they met the eligibility criteria. Any ‘unsure’  
146 articles on which an agreement could not be reached were checked by a fourth reviewer (JB) and the  
147 decision to include or exclude was based on the verdict of the fourth reviewer.

148

### 149 Inclusion/exclusion criteria

150 Articles were included if they reported the development of a clinical prediction (prognostic or  
151 diagnostic) model for individualised prediction in human participants, based on a multivariable  
152 regression model or studies updating a previously developed prediction model by adding new



153 predictors. Articles were excluded if they developed a model using non-regression based techniques,  
154 validated a previously developed prediction model, created a risk score from an existing prediction  
155 model, used a multivariable model to examine whether a particular predictor is associated with the  
156 outcome when adjusting for other factors (prognostic factor research), estimated the prognostic  
157 effect (e.g. hazard ratio) of a previously developed score, updated a previously developed model  
158 without adding any new predictors to the model or investigated the optimal cut off value of a  
159 previously developed model. Any articles in excess of the required 30 that met the eligibility criteria  
160 were also retained for inclusion, to avoid any potential selection bias concerns when choosing which  
161 articles to remove.

162

### 163 Data extraction

164 Data were extracted from the selected articles by a single reviewer (RW). The list of items presented  
165 in Table 1, where available, were extracted from each article and were based on the CHARMS  
166 checklist [5], with the addition of information related to the intended moment of using the model  
167 and measurement error.

168

169 >> insert Table 1: Data extracted from each article<<

170

171 Measurement error

172 The level of susceptibility to measurement error for each predictor used in the final models of the  
173 included articles were classified into two categories:

- 174 • Low risk: Unlikely to be measured with error, or possibly/likely to be measured with error  
175 but expected to be unimportant;
- 176 • High risk: Possibly/likely to be measured with error and may be important.

177 For example, age and gender are both extremely unlikely to be measured with error, and any error  
178 in age recorded would be expected to be negligible. Thus, age and gender would be classed as 'low  
179 risk' with regards to important measurement error. Whereas, blood pressure could be measured  
180 with error, as error in blood pressure measurement commonly occurs because of improper  
181 techniques such as talking during measurement or wrong cuff size [22] and blood pressure is also  
182 commonly measured with error due to biological variability [12]. This error could be large, and could  
183 be important when developing a prediction model for hypertension, for example, because blood  
184 pressure is an important component of the diagnostic evaluation for hypertension. Hence, blood  
185 pressure would be classed as 'high risk' of measurement error. Another high risk example would be  
186 body mass index (BMI), which the extent of the measurement error would depend on the way in  
187 which it was measured, but there would be a high chance it would be measured with some error.

188 To categorise the list of predictors into the two groups of susceptibility to measurement error, first  
189 the literature was searched for any publications discussing measurement error in any of the  
190 predictors of interest. For those where no evidence could be found, the categorisations were made  
191 based on the judgement of the reviewer (RW), which was corroborated by a postdoctoral academic  
192 General Practitioner.

193

194 Timing of measurements and intended moment of model use

195 If the timing of predictor measurements and intended moment of using the model was not explicitly  
196 stated in the article then, wherever possible, information on where the predictor information came  
197 from and the setting they were measured in were used to establish a likely time of measurement. If  
198 the intended moment of use of the model was not stated then, again where possible, information on  
199 what the model would be used for and the predictors that would be used within the model were  
200 considered to make a decision on the most probable intended moment of use.

201

202 **RESULTS**203 Included studies

204 A total of 1000 titles and abstracts were extracted and screened for inclusion. Of these, 876 were  
205 excluded based on not meeting the inclusion criteria from screening their title and abstract (Figure  
206 1). Study eligibility was then assessed for 124 full-text articles and 33 met eligibility criteria for  
207 inclusion [23-55], hence all 33 were retained for the review. The 33 included articles consisted of 27  
208 prognostic model studies and 6 diagnostic model studies published in 2015. The additional  
209 information extracted from the articles that is not related to measurement error or the timing of  
210 measurements is presented in Supplementary Table 3.

211 &gt;&gt;insert Figure 1: Flow chart of included studies&lt;&lt;

212

213

214

215 Measurement error

216 In the 33 articles reviewed, there was a total of 151 different predictors in the final prediction  
217 models. Many of the predictors were included in several different models, for example, age and  
218 gender were in many of the models (13 models and 5 models, respectively). Of the reported  
219 predictors included in the final models, we categorised 51 (33.8%) as high risk of being susceptible to  
220 measurement error (Table 2) and the remaining 100 (66.2%) as low risk (Supplementary Table 2).

221 >> insert Table 2: Predictors in final models at high risk of measurement error<<

222

223 Despite a third of the included predictors being at high risk of being susceptible to measurement  
224 error, only three studies acknowledged, or accounted for, measurement error within their model  
225 development. One study mentioned measurement error as a general limitation due to the study  
226 being from a single centre [23], but did not specify any particular predictor which may be at risk of  
227 error. Additionally, two studies used repeated measurements of a predictor within the modelling  
228 process [24, 27]. The first of the studies that used repeated measures [24] used generalised  
229 estimating equations (GEE's) to fit models accounting for the correlations among multiple biopsies  
230 that were performed on the same patients. The authors state that GEE's yield the same mean  
231 predictions as maximum likelihood, but result in inflated standard errors, wider confidence intervals  
232 and diminished statistical significance that more accurately reflect the amount of uncertainty in the  
233 data. There is no mention of measurement error within the article, and so it assumed that the  
234 authors have not made use of the repeated measures in a conscious effort to reduce measurement  
235 error, but to take advantage of all the data available (which may consequently potentially minimise  
236 measurement error). The second study using repeated measures [27] used joint modelling of  
237 longitudinal measures of CA125 with the stated purpose of estimating the time trend of CA125  
238 rather than explicitly accounting for measurement error (although again, this may consequently  
239 account for measurement error).

240

241 Despite only two of the reviewed articles included repeated measurements in the modelling process,  
242 repeated measurements of at least one of the candidate predictors were actually reported to be  
243 available in 6 (18.2%) of the articles [24, 27, 30, 31, 37, 45]. Of these 6 articles that repeated  
244 measures were available, 4 of these had repeated measures that we categorised as being high risk  
245 [27, 30, 31, 45], one of which used the repeated measures within the modelling [27].

246 Predictors considered at high risk of measurement error have been grouped by key reasons for being  
247 susceptible to error in Table 2, and several examples of these predictors with more detail and  
248 related references are given in Table 3.

ACCEPTED MANUSCRIPT

249 >> insert Table 3: Key reasons for measurement error in examples of predictors at high risk of being susceptible to error<<

250

251

252 One example of a predictor at high risk of error and used in several of the final prediction models is  
253 prostate specific antigen (PSA). Roehrborn et al. [69] conclude that there is significant variability  
254 between two serum PSA measurements obtained within a short time interval, which is due to  
255 chance alone. Biomarkers such as CA125, creatinine, C-reactive protein, serum albumin and other  
256 serological markers are also likely to change if a second sample was assessed, meaning they are  
257 measured with error due to biological variability causing discrepancies away from an underlying  
258 (mean) value [71]. There is also the possibility of laboratory error being present in these biomarkers,  
259 as the equipment or methods used to take the measurements within the laboratory may not be  
260 accurate.

261 In another example of a predictor likely measured with error, Ali et al. [68] found that the depth of  
262 myometrial invasion (DMI) was different in 29% of cases when the DMI was reassessed. The area  
263 under a patients pain curve could also be measured with error as it is a subjective measurement that  
264 may be affected by various things including how the question is asked, the setting in which the  
265 question is asked or when the question is asked. It could also be subject to recall error if the patient  
266 is asked about previous days pain levels. Another example is pulse rate, where Kobayashi [63] found  
267 that error occurred when pulse rates were objectively scored for various durations (e.g. 10, 15 or 30  
268 seconds) rather than for a whole minute, so the error in a pulse rate could depend on how long the  
269 pulse was taken for.

270 A patient's primary tumour diameter is another example of a predictor which may be susceptible to  
271 being measured with error. If a histologist determined the diameter under a microscope there would  
272 be little deviation from the true value, whereas if a surgeon recorded the diameter using an  
273 endoscopy then this could be recorded with error and could have an effect on the therapy chosen to  
274 be used [72]. Another example is BMI, which again the amount of error would depend on how it was

275 measured. If measured by a clinician then there is unlikely to be much measurement error, but if  
276 measured by the patient and recalled this may be subject to error [73]. Other examples include  
277 duration of convulsions, duration of neck pain, duration on nervousness and duration of tingling  
278 which are all self-reported predictors which could be subject to imperfect recall by the patient.

279 There were also many examples of predictors that were considered to be at low risk of important  
280 error. For example, one model that aimed to identify trauma patients at high risk of pulmonary  
281 embolism included a predictor indicating if the patient arrived at the hospital by helicopter [26], and  
282 it would be unlikely this would be incorrectly classified. Other models included the patient's disease  
283 location as a predictor and again, it is unlikely that this would not be recorded correctly.

284

#### 285 Timing of measurements and relation to intended moment of model use

286 Only eight of the articles explicitly stated exactly when the intended moment of using the model  
287 would be, or exactly when the predictors used in the final model were measured. However, for the  
288 majority of the 33 included articles it was possible to make a reasonable assumption about these  
289 details. If these assumptions were indeed correct, then in 30 (90.9%) of the 33 articles, the predictor  
290 measurements were all either taken at the intended moment of model use or were available prior to  
291 this. For example, one study [42] developed a model to predict survival prognosis after surgery in  
292 patients with symptomatic metastatic spinal cord compression from non-small cell lung cancer, with  
293 the aim of being able to provide optimal treatment. Although the specific timing of the predictor  
294 measurements was not stated, the predictors were specified as preoperative characteristics. The  
295 assumption was made that the model would be intended to be used at the point when a treatment  
296 decision was being made, as it was reported that those with the most favourable survival prognosis  
297 may instead be treated with more radical surgery. Therefore, it was assumed that the preoperative  
298 characteristics considered as predictors were either measured prior to or at the point that the model  
299 would be intended to be used.



300 In another example [40], a diagnostic model was developed to predict colorectal cancer in patients  
301 selected for colonoscopy in a primary health care setting, with the aim of identifying high risk  
302 patients to reduce the time till diagnosis and hence provide more efficient treatment strategies and  
303 success. As the model is to be used to help identify high risk patients when being considered for  
304 colonoscopy, which would happen during a GP consultation, it was assumed that the model would  
305 be intended to be used during a GP consultation when considering referral for colonoscopy. The  
306 model used predictors recorded in routine care data, which would all be available at the point of  
307 care, and although the article did not state at which time the predictors were recorded, it was  
308 assumed that only measurements recorded prior to colonoscopy referral were considered in the  
309 model development.

310 In all 6 of the articles in which repeated measures were available, each of the repeated measures  
311 were recorded either at or prior to the intended moment of using the prediction model.

312 In two (6.1%) of the articles [32, 48] it was not possible to make an assumption with regards to when  
313 the predictors were measured in relation to when the model was intended to be used. In the first  
314 article [32], a prognostic model was developed to predict the specific risk of non-sentinel node  
315 metastases in women with breast cancer with the aim of preventing unnecessary axillary lymph  
316 node dissections. The model was intended to be used after diagnosis of breast cancer, and as it is to  
317 be used to prevent unnecessary axillary lymph node dissections it could be assumed that the model  
318 would be intended to be used when deciding whether to perform an intraoperative axillary lymph  
319 node dissection. Little information was given on the predictors used in the model meaning the  
320 timing of the measurements of the predictors could not be deciphered, hence it was not possible to  
321 determine whether the predictors were measured at the intended moment of using the model or  
322 not. In the second article [48], a model was developed to predict unfavourable disease in patients  
323 with prostate cancer. The aim of the model was to avoid or postpone interventions in subjects with  
324 prostate cancer of low biological potential. The article states that the model is intended to be used

325 in patients after radical prostatectomy, but who were eligible for active surveillance. The predictors  
326 included were recorded from clinical evaluation, prostatic biopsy and radical prostatectomy  
327 specimens, but the timing of the clinical evaluation and prostatic biopsy was unclear and hence it  
328 was unknown whether these were before, at, or after the intended moment of using the model.  
329 For one of the included articles [30], a classification algorithm was developed for the diagnosis of  
330 non-alcoholic fatty liver disease (NAFLD). The model was not developed to be intended to be used at  
331 a specific time but to be used to identify large scale longitudinal cohorts from electronic medical  
332 records for use in research studies.

333

## 334 **DISCUSSION**

335 Our review suggests that many published clinical prediction models include predictors that are  
336 susceptible to potentially important measurement error and yet this was seldom acknowledged. Of  
337 33 articles in our review only two used methods that could potentially account for measurement  
338 error by using repeated measurements in the modelling. Though the impact of ignoring  
339 measurement error in the articles reviewed is difficult to establish, it raises an important  
340 methodological consideration for future prediction model research to address, particularly as a third  
341 of the predictors used in the prediction models were categorised as being at high risk of being  
342 susceptible to measurement error. The review also found that over three-quarters of the articles  
343 included did not explicitly state the exact timing that the model is intended to be used in clinical  
344 practice, or exactly when the predictors used in the modelling development were measured.  
345 However, a reasonable assumption could be made for the majority of the articles included and,  
346 based on this, there were no articles that obviously recorded a predictor after the time it was  
347 intended to be used.

348

349 Related research

350 Measurement error has been found to generally have three main effects if not accounted for in  
351 medical research: biased or inaccurate estimates of the parameters, loss of power and masking the  
352 features of the data (making it harder to spot relationships via graphical methods) [8]. The direction  
353 and magnitude of bias from measurement error depends heavily on whether the distribution of  
354 errors for one variable depends on the actual value of the variable, the actual values of other  
355 variables, or the errors in measuring other variables [7], as well as on the true strength of  
356 association, the prevalence of the predictors [74] and whether the errors are random or systematic.  
357 Hence, the direction of bias from predictor measurement error is likely to be difficult to predict.  
358 However, failing to adjust for random measurement error could potentially lead to estimates being  
359 biased towards the null [6], which could subsequently lead to an underestimate of a patients'  
360 probability of outcome if measurement error is present in the prediction model used. Conversely,  
361 failing to account for systematic errors may change the results in different directions, which could  
362 again lead to incorrect predictions of a patients' probability of future outcome.

363 There are currently two conflicting views about whether measurement error in prediction models is  
364 an issue or not. Firstly, Carroll and colleagues [8] state that if a predictor ( $X$ ) is measured with error,  
365 and this measure ( $W$ ) is used to predict a patients outcome, then if it is this same surrogate measure  
366 of  $X$  that will be used when applying the prediction model in practice, there is little issue with using  
367  $W$  to develop the prediction model. On the other hand, a prediction model should provide the most  
368 accurate estimate possible, and if a predictor used in the development of a model is measured with  
369 error then the estimates of the predictor-outcome associations will be biased, meaning the  
370 predictions made may be untrue. Measurement error in the candidate predictors could also lead to  
371 certain predictors not being included in the final model due to the measurement error.

372 In etiologic research we are most interested in the (adjusted) estimate of a single predictor-outcome  
373 association and hence would want to minimise bias of this particular estimate. Whereas when

374 developing a prediction model, we are not predominantly interested in the individual estimates of  
375 one (or more) of the predictor-outcome associations, but in the actual absolute risk predictions  
376 calculated from the model (and the predictive performance of these risk predictions from the  
377 model). Hence, even if one (or more) of the estimates of a predictor-outcome association in a  
378 prediction model is biased due to measurement error, this may not be an issue if the model as a  
379 whole performs well in terms of the absolute risk predictions. However, measurement error in  
380 prediction models has been shown to reduce the c-statistic and increase the Brier score dramatically  
381 [13], but in that article the authors focussed on the gain in prediction performance from using error-  
382 free predictors instead of error-prone predictors, rather than the gain in prediction performance  
383 from accounting for the measurement error in the model when the true error-free values are not  
384 known. The article also only evaluated the scenario where only one error-prone predictor was  
385 included in the prediction model.

386 Another article assessed the impact of random and systematic error in self-reported height and  
387 weight on the performance of a model used to predict diabetes [14]. The authors found that random  
388 error reduced the calibration and discrimination, and biased the predicted risk upwards, whereas  
389 systematic error reduced the calibration and biased the predicted risk in the direction of the bias,  
390 but had no effect on the discrimination.

391

### 392 Strengths and Limitations

393 A strength of this review was that a clearly defined search strategy which was based on a previously  
394 published search filter [21] was used. Although this review did not include a search of every  
395 prediction model published within a certain time period due to the sheer volume of prediction  
396 models published each year [2], a search of a few of the most recently published studies was  
397 deemed appropriate to enable a general overview of the current literature and provide qualitative  
398 saturation of whether there was a susceptibility for measurement error within the predictors and

399 whether this was considered and also the timing of the predictor measurements in relation to the  
400 intended moment of using the model.

401 The reviewer's judgement had to be used and assumptions were made about the timing of  
402 measurements and when the model is intended to be used. This was due to the reviewed articles  
403 not explicitly stating these details. Based on this, all of the papers here did actually measure the  
404 predictors at the intended moment of using the model (or before), in those that it was possible to  
405 decipher this information. However, it is possible that some of these assumptions made were  
406 incorrect.

407 Another concern within prediction models in relation to predictor timing is the relevant time  
408 window, or the length of the induction period, in which the predictor of interest is causally related to  
409 the outcome. For some prediction models, certain causal factors may need to be considered from  
410 much longer ago than others, i.e. with a longer induction period. For example, if considering  
411 asbestos exposure in relation to future lung disease, the association could span back many years,  
412 whereas recent asbestos exposure may not be related to the outcome if the induction period is only  
413 relatively short, e.g. 1-2 years. On the other hand, when predicting infectious diseases, the current  
414 and recent exposure of the patient is likely to be most important, and so a relatively short induction  
415 period would be needed. Hence, the duration of follow-up of predictors prior to the intended  
416 moment of model use should be clearly specified when developing a prediction model, however we  
417 did not assess this within this review.

418 When developing a prediction model, the calendar year of time in which the measurements were  
419 made is important (relative to the calendar time of the intended moment of model use), because the  
420 precision of measurements often improves when using newer measurement methods. Using a more  
421 recent, up-to-date data set that used more improved measurement techniques to develop a  
422 prediction model would potentially provide a more relevant and better performing model than if

423 using an older dataset. While study recruitment dates are generally reported, we did not consider  
424 this in relation to when the article was published or would be intended to be used.

425 Due to many of the included studies not actually stating a complete list of all of the candidate  
426 predictors considered in the model development, only the predictors included in the final models  
427 were assessed for their susceptibility to measurement error. However, measurement error in the  
428 candidate predictors could lead to the exclusion of these predictors in the model development stage  
429 and so measurement error in these predictors could be as equally as important as measurement  
430 error in the predictors in the final models.

431 Little information was given within the included articles about any measurement error that may be  
432 present in the predictors. Without the availability of previous research on the amount of error in  
433 certain predictors, a subjective decision on whether measurement error was likely had to be made  
434 by the reviewer, although an academic GP also reviewed the list of predictors and gave their opinion  
435 on whether they would judge the predictor to be susceptible to measurement error when using in  
436 practice. One difficulty with making a decision on whether the predictor is likely to be susceptible to  
437 measurement error was that for many of the predictors it would depend on exactly how the  
438 predictor was measured, but often this level of detail is missing from the article. Despite this  
439 subjective approach to categorising measurement error, there were several predictors included in  
440 the final models that had corresponding published research suggesting they are likely to be  
441 measured with error, and this was not considered within the development of the models.

442

## 443 **CONCLUSIONS**

444 It is possible that many published prediction models include predictors that are measured with error,  
445 and this is often not accounted for or even considered. Additionally, even if the authors considered  
446 the predictors to be measured without error, either because of the way they were measured, or for

447 some other reason, this was still not stated within the articles. This suggests a need to assess under  
448 what circumstances ignoring measurement error in prediction models is a concern and whether  
449 accounting for the error will improve the predictions made and the model performance. However,  
450 researchers should be considering how susceptible to measurement error their predictors may be  
451 when developing a model.

452 Although there were no clear examples within this review of a prediction model being developed  
453 using a predictor that was measured after the intended moment of using the model, it is common in  
454 prognosis studies of recurrent and long-term conditions presenting to primary care for information  
455 on predictors (e.g. pain intensity) to be ascertained by mailed self-complete questionnaires, or  
456 personal interview and examination in research clinics several days after their index consultation  
457 [75-81]. It was found in this review that the timing of the measurements and the intended moment  
458 of using the model is often not explicitly stated, which could mean that future users of the model  
459 unknowingly estimate misleading probabilities of a patients' outcome if they are using predictors  
460 measured at a different time than those used in the model development in relation to the timing of  
461 the model use. We have previously found that displacing the collection of time-varying predictors  
462 from the intended moment of use of a prediction model can result in differences in the magnitude of  
463 predictor-outcome associations and the subsequent accuracy of the model performance [20].  
464 Hence, future prediction model research studies must clearly report the intended moment of use of  
465 the prediction model, and be explicit about whether the predictors were collected before the  
466 intended moment of use or not, and if not, justify why.

467

468

469

470

471 **LIST OF ABBREVIATIONS**

472 AUC=area under the curve; BMI=body mass index; DMI=depth of myometrial invasion; GP=general  
473 practitioner; NAFLD=non-alcoholic fatty liver disease; PSA=prostate specific antigen;  
474 TRIPOD=Transparent Reporting of a multivariable prediction model for Individual Prognosis or  
475 Diagnosis.

476

477 **DECLARATIONS OF INTEREST:** None.

478

479 **FUNDING:** RW is funded through an NIHR Research Professorship in General Practice for Christian  
480 Mallen (NIHR-RP-2014-04-026). GSC was supported by the NIHR Biomedical Research Centre,  
481 Oxford. The views expressed in this paper are those of the author(s) and not necessarily those of the  
482 NHS, the NIHR, or the Department of Health.

483

484 **AUTHORS' CONTRIBUTIONS:** RW, RR, GP and JB developed the research question. RW reviewed and  
485 extracted information from the articles, with checks from RR, GP and JB. RW drafted the first version  
486 of the manuscript, with support from RR. All authors reviewed the draft manuscript, helped to revise  
487 the manuscript, and approved the final version.

488

489 **ACKNOWLEDGEMENTS:** The authors would like to thank Dr Lorna Clarkson for assessing the list of  
490 predictors for susceptibility to measurement error.



## 491 REFERENCES

- 492 1. Hemingway H, Croft P, Perel P, Hayden JA, Abrams K, Timmis A, Briggs A, Udumyan R, Moons  
 493 KGM, Steyerberg EW *et al*: **Prognosis research strategy (PROGRESS) 1: A framework for**  
 494 **researching clinical outcomes.** *BMJ* 2013, **346**.
- 495 2. Wessler BS, Lana Lai YH, Kramer W, Cangelosi M, Raman G, Lutz JS, Kent DM: **Clinical**  
 496 **Prediction Models for Cardiovascular Disease: The Tufts PACE CPM Database.** *Circulation*  
 497 *Cardiovascular quality and outcomes* 2015, **8(4)**:368-375.
- 498 3. Steyerberg E: **Clinical Prediction Models: A Practical Approach to Development, Validation,**  
 499 **and Updating:** Springer New York; 2008.
- 500 4. Steyerberg EW, Moons KG, van der Windt DA, Hayden JA, Perel P, Schroter S, Riley RD,  
 501 Hemingway H, Altman DG: **Prognosis Research Strategy (PROGRESS) 3: prognostic model**  
 502 **research.** *PLoS medicine* 2013, **10(2)**:e1001381.
- 503 5. Moons KG, de Groot JA, Bouwmeester W, Vergouwe Y, Mallett S, Altman DG, Reitsma JB,  
 504 Collins GS: **Critical appraisal and data extraction for systematic reviews of prediction**  
 505 **modelling studies: the CHARMS checklist.** *PLoS medicine* 2014, **11(10)**:e1001744.
- 506 6. Prentice RL: **Covariate Measurement Errors and Parameter Estimation in a Failure Time**  
 507 **Regression Model.** *Biometrika* 1982, **69(2)**:331-342.
- 508 7. Rothman KJ, Greenland S, Lash TL: **Modern Epidemiology:** Wolters Kluwer Health/Lippincott  
 509 Williams & Wilkins; 2008.
- 510 8. Carroll RJ, Ruppert D, Stefanski LA, Crainiceanu CM: **Measurement Error in Nonlinear**  
 511 **Models: A Modern Perspective, Second Edition:** CRC Press; 2006.
- 512 9. Gustafson P: **Measurement Error and Misclassification in Statistics and Epidemiology:**  
 513 **Impacts and Bayesian Adjustments:** CRC Press; 2003.
- 514 10. Guolo A: **Robust techniques for measurement error correction: a review.** *Statistical*  
 515 *methods in medical research* 2008, **17(6)**:555-580.
- 516 11. Brakenhoff TB, Mitroiu M, Keogh RH, Moons KGM, Groenwold RHH, van Smeden M:  
 517 **Measurement error is often neglected in medical literature: a systematic review.** *J Clin*  
 518 *Epidemiol* 2018.
- 519 12. Grassi G, Bombelli M, Brambilla G, Trevano FQ, Dell’Oro R, Mancia G: **Total Cardiovascular**  
 520 **Risk, Blood Pressure Variability and Adrenergic Overdrive in Hypertension: Evidence,**  
 521 **Mechanisms and Clinical Implications.** *Current Hypertension Reports* 2012, **14(4)**:333-338.
- 522 13. Khudyakov P, Gorfine M, Zucker D, Spiegelman D: **The impact of covariate measurement**  
 523 **error on risk prediction.** *Statistics in medicine* 2015, **34(15)**:2353-2367.
- 524 14. Rosella LC, Corey P, Stukel TA, Mustard C, Hux J, Manuel DG: **The influence of measurement**  
 525 **error on calibration, discrimination, and overall estimation of a risk prediction model.**  
 526 *Population Health Metrics* 2012, **10(1)**:20.
- 527 15. Fox J-P, Glas CAW: **Bayesian modeling of measurement error in predictor variables using**  
 528 **item response theory.** *Psychometrika* 2003, **68(2)**:169-191.
- 529 16. Li W, Mazumdar S, Arena VC, Sussman N: **A resampling approach for adjustment in**  
 530 **prediction models for covariate measurement error.** *Computer Methods and Programs in*  
 531 *Biomedicine* 2005, **77(3)**:199-207.
- 532 17. van Walraven C, Davis D, Forster AJ, Wells GA: **Time-dependent bias was common in**  
 533 **survival analyses published in leading clinical journals.** *J Clin Epidemiol* 2004, **57(7)**:672-682.
- 534 18. Collins GS, Reitsma JB, Altman DG, Moons KGM: **Transparent Reporting of a Multivariable**  
 535 **Prediction Model for Individual Prognosis Or Diagnosis (TRIPOD): the TRIPOD statement.**  
 536 *Journal of Clinical Epidemiology*, **68(2)**:112-121.
- 537 19. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, Vickers AJ,  
 538 Ransohoff DF, Collins GS: **Transparent Reporting of a multivariable prediction model for**  
 539 **Individual Prognosis Or Diagnosis (TRIPOD): Explanation and ElaborationThe TRIPOD**  
 540 **Statement: Explanation and Elaboration.** *Annals of Internal Medicine* 2015, **162(1)**:W1-W73.

- 541 20. Whittle R, Royle K-L, Jordan KP, Riley RD, Mallen CD, Peat G: **Prognosis research ideally**  
542 **should measure time-varying predictors at their intended moment of use.** *Diagnostic and*  
543 *Prognostic Research* 2017, **1**(1):1.
- 544 21. Geersing GJ, Bouwmeester W, Zuithoff P, Spijker R, Leeflang M, Moons KG: **Search filters for**  
545 **finding prognostic and diagnostic prediction studies in Medline to enhance systematic**  
546 **reviews.** *PLoS one* 2012, **7**(2):e32844.
- 547 22. Handler J: **The Importance of Accurate Blood Pressure Measurement.** *The Permanente*  
548 *Journal* 2009, **13**(3):51-54.
- 549 23. Angioli R, Capriglione S, Aloisi A, Ricciardi R, Scaletta G, Lopez S, Miranda A, Di Pinto A,  
550 Terranova C, Plotti F: **A Predictive Score for Secondary Cytoreductive Surgery in Recurrent**  
551 **Ovarian Cancer (SeC-Score): A Single-Centre, Controlled Study for Preoperative Patient**  
552 **Selection.** *Annals Of Surgical Oncology* 2015, **22**(13):4217-4223.
- 553 24. Ankerst DP, Xia J, Thompson IM, Jr., Hoefler J, Newcomb LF, Brooks JD, Carroll PR, Ellis WJ,  
554 Gleave ME, Lance RS *et al*: **Precision Medicine in Active Surveillance for Prostate Cancer:**  
555 **Development of the Canary-Early Detection Research Network Active Surveillance Biopsy**  
556 **Risk Calculator.** *European Urology* 2015, **68**(6):1083-1088.
- 557 25. Bendifallah S, Canlorbe G, Laas E, Huguet F, Coutant C, Hudry D, Graesslin O, Raimond E,  
558 Touboul C, Collinet P *et al*: **A Predictive Model Using Histopathologic Characteristics of**  
559 **Early-Stage Type 1 Endometrial Cancer to Identify Patients at High Risk for Lymph Node**  
560 **Metastasis.** *Annals Of Surgical Oncology* 2015, **22**(13):4224-4232.
- 561 26. Black SR, Howard JT, Chin PC, Starr AJ: **Toward a More Robust Prediction of Pulmonary**  
562 **Embolism in Trauma Patients. A Risk Assessment Model Based Upon 38,000 Patients.**  
563 *Journal Of Orthopaedic Trauma* 2015.
- 564 27. Chang C, Chiang AJ, Chen W-A, Chang H-W, Chen J: **A joint model based on longitudinal**  
565 **CA125 in ovarian cancer to predict recurrence.** *Biomarkers In Medicine* 2015.
- 566 28. Chen J-Y, Chen J-J, Xue J-Y, Chen Y, Liu G-Y, Han Q-X, Yang W-T, Shen Z-Z, Shao Z-M, Wu J:  
567 **Predicting Non-sentinel Lymph Node Metastasis in a Chinese Breast Cancer Population**  
568 **with 1-2 Positive Sentinel Nodes: Development and Assessment of a New Predictive**  
569 **Nomogram.** *World Journal Of Surgery* 2015, **39**(12):2919-2927.
- 570 29. Cohen MH, Hotton AL, Hershov RC, Levine A, Bacchetti P, Golub ET, Anastos K, Young M,  
571 Gustafson D, Weber KM: **Gender-Related Risk Factors Improve Mortality Predictive Ability**  
572 **of VACS Index Among HIV-Infected Women.** *Journal Of Acquired Immune Deficiency*  
573 *Syndromes (1999)* 2015, **70**(5):538-544.
- 574 30. Corey KE, Kartoun U, Zheng H, Shaw SY: **Development and Validation of an Algorithm to**  
575 **Identify Nonalcoholic Fatty Liver Disease in the Electronic Medical Record.** *Digestive*  
576 *Diseases And Sciences* 2015.
- 577 31. Côté GA, Lynch S, Easler JJ, Keen A, Vassell PA, Sherman S, Hui S, Xu H: **Development and**  
578 **Validation of a Prediction Model for Admission After Endoscopic Retrograde**  
579 **Cholangiopancreatography.** *Clinical Gastroenterology And Hepatology: The Official Clinical*  
580 *Practice Journal Of The American Gastroenterological Association* 2015, **13**(13):2323-  
581 2332.e2329.
- 582 32. Di Filippo F, Giannarelli D, Bouteille C, Bernet L, Cano R, Cunnick G, Sapino A: **Elaboration of**  
583 **a nomogram to predict non sentinel node status in breast cancer patients with positive**  
584 **sentinel node, intra-operatively assessed with one step nucleic acid amplification method.**  
585 *Journal Of Experimental & Clinical Cancer Research: CR* 2015, **34**(1):136-136.
- 586 33. Du X-J, Tang L-L, Chen L, Mao Y-P, Guo R, Liu X, Sun Y, Zeng M-S, Kang T-B, Shao J-Y *et al*:  
587 **Neoadjuvant chemotherapy in locally advanced nasopharyngeal carcinoma: Defining high-**  
588 **risk patients who may benefit before concurrent chemotherapy combined with intensity-**  
589 **modulated radiotherapy.** *Scientific Reports* 2015, **5**:16664-16664.

- 590 34. Dua A, Algodí MM, Furlough C, Ray H, Desai SS: **Development of a scoring system to**  
591 **estimate mortality in abdominal aortic aneurysms management.** *Vascular* 2015, **23**(6):586-  
592 591.
- 593 35. Englum BR, Pura J, Reed SD, Roman SA, Sosa JA, Scheri RP: **A Bedside Risk Calculator to**  
594 **Preoperatively Distinguish Follicular Thyroid Carcinoma from Follicular Variant of Papillary**  
595 **Thyroid Carcinoma.** *World Journal Of Surgery* 2015, **39**(12):2928-2934.
- 596 36. Faget C, Taourel P, Charbit J, Ruyer A, Alili C, Molinari N, Millet I: **Value of CT to predict**  
597 **surgically important bowel and/or mesenteric injury in blunt trauma: performance of a**  
598 **preliminary scoring system.** *European Radiology* 2015, **25**(12):3620-3628.
- 599 37. Horn SD, Barrett RS, Fife CE, Thomson B: **A Predictive Model for Pressure Ulcer Outcome:**  
600 **The Wound Healing Index.** *Advances In Skin & Wound Care* 2015, **28**(12):560-572.
- 601 38. Kaymakçalan MD, Xie W, Albiges L, North SA, Kollmannsberger CK, Smoragiewicz M, Kroeger  
602 N, Wells JC, Rha S-Y, Lee JL *et al*: **Risk factors and model for predicting toxicity-related**  
603 **treatment discontinuation in patients with metastatic renal cell carcinoma treated with**  
604 **vascular endothelial growth factor-targeted therapy: Results from the International**  
605 **Metastatic Renal Cell Carcinoma Database Consortium.** *Cancer* 2015.
- 606 39. Koller L, Rothgerber D-J, Sulzgruber P, El-Hamid F, Förster S, Wojta J, Goliash G, Maurer G,  
607 Niessner A: **History of previous bleeding and C-reactive protein improve assessment of**  
608 **bleeding risk in elderly patients (≥ 80 years) with myocardial infarction.** *Thrombosis And*  
609 *Haemostasis* 2015, **114**(5):1085-1091.
- 610 40. Koning NR, Moons LMG, Büchner FL, Helsper CW, Ten Teije A, Numans ME: **Identification of**  
611 **patients at risk for colorectal cancer in primary care: an explorative study with routine**  
612 **healthcare data.** *European Journal Of Gastroenterology & Hepatology* 2015, **27**(12):1443-  
613 1448.
- 614 41. Kusamura S, Torres Mesa PA, Cabras A, Baratti D, Deraco M: **The Role of Ki-67 and Pre-**  
615 **cytoreduction Parameters in Selecting Diffuse Malignant Peritoneal Mesothelioma**  
616 **(DMPM) Patients for Cytoreductive Surgery (CRS) and Hyperthermic Intraperitoneal**  
617 **Chemotherapy (HIPEC).** *Annals Of Surgical Oncology* 2015.
- 618 42. Lei M, Liu Y, Tang C, Yang S, Liu S, Zhou S: **Prediction of survival prognosis after surgery in**  
619 **patients with symptomatic metastatic spinal cord compression from non-small cell lung**  
620 **cancer.** *BMC Cancer* 2015, **15**(1):853-853.
- 621 43. Matsuo K, Jung CE, Hom MS, Gualtieri MR, Randazzo SC, Kanao H, Yessaian AA, Roman LD:  
622 **Predictive Factor of Conversion to Laparotomy in Minimally Invasive Surgical Staging for**  
623 **Endometrial Cancer.** *International Journal Of Gynecological Cancer: Official Journal Of The*  
624 *International Gynecological Cancer Society* 2015.
- 625 44. Nykanen DG, Forbes TJ, Du W, Divekar AA, Reeves JH, Hagler DJ, Fagan TE, Pedra CA, Fleming  
626 GA, Khan DM *et al*: **CRISP: Catheterization RiSk score for pediatrics: A Report from the**  
627 **Congenital Cardiac Interventional Study Consortium (CCISC).** *Catheterization And*  
628 *Cardiovascular Interventions: Official Journal Of The Society For Cardiac Angiography &*  
629 *Interventions* 2015.
- 630 45. Olmedilla L, Lisbona CJ, Pérez-Peña JM, López-Baena JA, Garutti I, Salcedo M, Sanz J, Tisner  
631 M, Asencio JM, Fernández-Quero L *et al*: **Early Measurement of Indocyanine Green**  
632 **Clearance Accurately Predicts Short-Term Outcomes After Liver Transplantation.**  
633 *Transplantation* 2015.
- 634 46. Resch JE, Brown CN, Macciocchi SN, Cullum CM, Blueitt D, Ferrara MS: **A Preliminary**  
635 **Formula to Predict Timing of Symptom Resolution for Collegiate Athletes Diagnosed With**  
636 **Sport Concussion.** *Journal Of Athletic Training* 2015.
- 637 47. Rosenkrantz AB, Ream JM, Nolan P, Rusinek H, Deng F-M, Taneja SS: **Prostate Cancer: Utility**  
638 **of Whole-Lesion Apparent Diffusion Coefficient Metrics for Prediction of Biochemical**  
639 **Recurrence After Radical Prostatectomy.** *AJR American Journal Of Roentgenology* 2015,  
640 **205**(6):1208-1214.

- 641 48. Russo GI, Castelli T, Favilla V, Reale G, Urzì D, Privitera S, Fragalà E, Cimino S, Morgia G:  
642 **Performance of biopsy factors in predicting unfavorable disease in patients eligible for**  
643 **active surveillance according to the PRIAS criteria.** *Prostate Cancer And Prostatic Diseases*  
644 2015, **18**(4):338-342.
- 645 49. Shaikh AY, Esa N, Martin-Doyle W, Kinno M, Nieto I, Floyd KC, Browning C, Ennis C, Donahue  
646 JK, Rosenthal LS *et al*: **Addition of B-Type Natriuretic Peptide to Existing Clinical Risk Scores**  
647 **Enhances Identification of Patients at Risk for Atrial Fibrillation Recurrence After**  
648 **Pulmonary Vein Isolation.** *Critical Pathways In Cardiology* 2015, **14**(4):157-165.
- 649 50. Siegel CA, Horton H, Siegel LS, Thompson KD, Mackenzie T, Stewart SK, Rice PW, Stempak  
650 JM, Dezfoli S, Haritunians T *et al*: **A validated web-based tool to display individualised**  
651 **Crohn's disease predicted outcomes based on clinical, serologic and genetic variables.**  
652 *Alimentary Pharmacology & Therapeutics* 2015.
- 653 51. Spolverato G, Vitale A, Cucchetti A, Popescu I, Marques HP, Aldrighetti L, Gamblin TC,  
654 Maithel SK, Sandroussi C, Bauer TW *et al*: **Can hepatic resection provide a long-term cure**  
655 **for patients with intrahepatic cholangiocarcinoma?** *Cancer* 2015, **121**(22):3998-4006.
- 656 52. Suh YJ, Hong YJ, Lee H-J, Hur J, Kim YJ, Lee HS, Hong SR, Im DJ, Kim YJ, Park CH *et al*:  
657 **Prognostic value of SYNTAX score based on coronary computed tomography angiography.**  
658 *International Journal Of Cardiology* 2015, **199**:460-466.
- 659 53. Tada H, Takanashi J-I, Okuno H, Kubota M, Yamagata T, Kawano G, Shiihara T, Hamano S-I,  
660 Hirose S, Hayashi T *et al*: **Predictive score for early diagnosis of acute encephalopathy with**  
661 **biphasic seizures and late reduced diffusion (AESD).** *Journal Of The Neurological Sciences*  
662 2015, **358**(1-2):62-65.
- 663 54. Takahashi N, Leng S, Kitajima K, Gomez-Cardona D, Thapa P, Carter RE, Leibovich BC,  
664 Sasiwimonphan K, Sasaguri K, Kawashima A: **Small (< 4 cm) Renal Masses: Differentiation of**  
665 **Angiomyolipoma Without Visible Fat From Renal Cell Carcinoma Using Unenhanced and**  
666 **Contrast-Enhanced CT.** *AJR American Journal Of Roentgenology* 2015, **205**(6):1194-1202.
- 667 55. Zhou J, Zhou Y, Cao S, Li S, Wang H, Niu Z, Chen D, Wang D, Lv L, Zhang J *et al*: **Multivariate**  
668 **logistic regression analysis of postoperative complications and risk model establishment of**  
669 **gastrectomy for gastric cancer: A single-center cohort report.** *Scandinavian Journal Of*  
670 *Gastroenterology* 2016, **51**(1):8-15.
- 671 56. Scott J, Huskisson EC: **Accuracy of subjective measurements made with or without previous**  
672 **scores: an important source of error in serial measurement of subjective states.** *Annals of*  
673 *the Rheumatic Diseases* 1979, **38**(6):558-559.
- 674 57. Daoust R, Sirois MJ, Lee JS, Perry JJ, Griffith LE, Worster A, Lang E, Paquet J, Chauny JM,  
675 Emond M: **Painful Memories: Reliability of Pain Intensity Recall at 3 Months in Senior**  
676 **Patients.** *Pain research & management* 2017, **2017**:5983721.
- 677 58. Tso E, Elson P, Vanlente F, Markman M: **The "real-life" variability of CA-125 in ovarian**  
678 **cancer patients.** *Gynecologic oncology* 2006, **103**(1):141-144.
- 679 59. Tuxen MK, Soletormos G, Petersen PH, Schioler V, Dombernowsky P: **Assessment of**  
680 **biological variation and analytical imprecision of CA 125, CEA, and TPA in relation to**  
681 **monitoring of ovarian cancer.** *Gynecologic oncology* 1999, **74**(1):12-22.
- 682 60. Peake M, Whiting M: **Measurement of Serum Creatinine – Current Status and Future Goals.**  
683 *Clinical Biochemist Reviews* 2006, **27**(4):173-184.
- 684 61. Reinhard M, Erlandsen EJ, Randers E: **Biological variation of cystatin C and creatinine.**  
685 *Scandinavian journal of clinical and laboratory investigation* 2009, **69**(8):831-836.
- 686 62. Macy EM, Hayes TE, Tracy RP: **Variability in the measurement of C-reactive protein in**  
687 **healthy subjects: implications for reference intervals and epidemiological applications.**  
688 *Clinical Chemistry* 1997, **43**(1):52-58.
- 689 63. Kobayashi H: **Effect of measurement duration on accuracy of pulse-counting.** *Ergonomics*  
690 2013, **56**(12):1940-1944.

- 691 64. Sawers L: **Measuring and modelling concurrency**. *Journal of the International AIDS Society*  
692 2013, **16**(1):17431.
- 693 65. Delanaye P, Schaeffner E, Ebert N, Cavalier E, Mariat C, Krzesinski J-M, Moranne O: **Normal**  
694 **reference values for glomerular filtration rate: what do we really know?** *Nephrology*  
695 *Dialysis Transplantation* 2012, **27**(7):2664-2672.
- 696 66. Braga F, Ferraro S, Mozzi R, Panteghini M: **The importance of individual biology in the**  
697 **clinical use of serum biomarkers for ovarian cancer**. In: *Clinical Chemistry and Laboratory*  
698 *Medicine (CCLM)*. vol. 52; 2014: 1625.
- 699 67. Polley M-YC, Leung SCY, McShane LM, Gao D, Hugh JC, Mastropasqua MG, Viale G, Zabaglo  
700 LA, Penault-Llorca F, Bartlett JMS *et al*: **An International Ki67 Reproducibility Study**. *JNCI*  
701 *Journal of the National Cancer Institute* 2013, **105**(24):1897-1906.
- 702 68. Ali A, Black D, Soslow RA: **Difficulties in assessing the depth of myometrial invasion in**  
703 **endometrial carcinoma**. *International journal of gynecological pathology : official journal of*  
704 *the International Society of Gynecological Pathologists* 2007, **26**(2):115-123.
- 705 69. Roehrborn CG, Pickens GJ, Carmody T, 3rd: **Variability of repeated serum prostate-specific**  
706 **antigen (PSA) measurements within less than 90 days in a well-defined patient population**.  
707 *Urology* 1996, **47**(1):59-66.
- 708 70. Winkel P, Statland BE, Bokelund H: **Factors contributing to intra-individual variation of**  
709 **serum constituents: 5. Short-term day-to-day and within-hour variation of serum**  
710 **constituents in healthy subjects**. *Clin Chem* 1974, **20**(12):1520-1527.
- 711 71. Braga F, Panteghini M: **Generation of data on within-subject biological variation in**  
712 **laboratory medicine: An update**. *Critical Reviews in Clinical Laboratory Sciences* 2016,  
713 **53**(5):313-325.
- 714 72. Mori H, Kobara H, Tsushimi T, Nishiyama N, Fujihara S, Masaki T: **Unavoidable Human Errors**  
715 **of Tumor Size Measurement during Specimen Attachment after Endoscopic Resection: A**  
716 **Clinical Prospective Study**. *PloS one* 2015, **10**(4):e0121798.
- 717 73. Hill A, Roberts J: **Body mass index: a comparison between self-reported and measured**  
718 **height and weight**. *Journal of public health medicine* 1998, **20**(2):206-210.
- 719 74. Jurek AM, Greenland S, Maldonado G, Church TR: **Proper interpretation of non-differential**  
720 **misclassification effects: expectations vs observations**. *International journal of*  
721 *epidemiology* 2005, **34**(3):680-687.
- 722 75. Wardenaar KJ, Conradi HJ, de Jonge P: **Data-driven course trajectories in primary care**  
723 **patients with major depressive disorder**. *Depression and anxiety* 2014, **31**(9):778-786.
- 724 76. Von Korff M, Deyo RA, Cherkin D, Barlow W: **Back pain in primary care. Outcomes at 1 year**.  
725 *Spine* 1993, **18**(7):855-862.
- 726 77. Scheele J, Luijsterburg PA, Ferreira ML, Maher CG, Pereira L, Peul WC, van Tulder MW,  
727 Bohnen AM, Berger MY, Bierma-Zeinstra SM *et al*: **Back complaints in the elders (BACE);**  
728 **design of cohort studies in primary care: an international consortium**. *BMC musculoskeletal*  
729 *disorders* 2011, **12**:193.
- 730 78. Radanov BP, di Stefano G, Schnidrig A, Ballinari P: **Role of psychosocial stress in recovery**  
731 **from common whiplash [see comment]**. *Lancet (London, England)* 1991, **338**(8769):712-  
732 715.
- 733 79. Licht-Strunk E, Beekman AT, de Haan M, van Marwijk HW: **The prognosis of undetected**  
734 **depression in older general practice patients. A one year follow-up study**. *Journal of*  
735 *affective disorders* 2009, **114**(1-3):310-315.
- 736 80. Hermsen LA, Leone SS, van der Windt DA, Smalbrugge M, Dekker J, van der Horst HE:  
737 **Functional outcome in older adults with joint pain and comorbidity: design of a**  
738 **prospective cohort study**. *BMC musculoskeletal disorders* 2011, **12**:241.
- 739 81. Diehm C, Darius H, Pittrow D, Schwertfeger M, Tepohl G, Haberl RL, Allenberg JR, Burghaus I,  
740 Trampisch HJ: **Prognostic value of a low post-exercise ankle brachial index as assessed by**  
741 **primary care physicians**. *Atherosclerosis* 2011, **214**(2):364-372.

ACCEPTED MANUSCRIPT

1 Table 1: Data extracted from each article

Design and aim	<ul style="list-style-type: none"> <li>• Prognostic versus diagnostic prediction model</li> <li>• Intended scope of the review <ul style="list-style-type: none"> <li>– Clinical area</li> <li>– Aim of prediction model (e.g. inform therapeutic decision making, inform referral or withholding from invasive diagnostic testing, inform patients of probability of event)</li> </ul> </li> <li>• Source of data (e.g. cohort, case-control, randomised trial or registry data)</li> </ul>
Outcomes to be predicted	<ul style="list-style-type: none"> <li>• Definition and method for measurement of outcome</li> <li>• Type of outcome (e.g. single or combined endpoints; binary or time to event)</li> </ul>
Candidate predictors	<ul style="list-style-type: none"> <li>• Number and type of predictors (e.g. demographics, patient history, physical examination, additional testing, disease characteristics)</li> <li>• Definition and method for measurement of candidate predictors</li> <li>• Timing of predictor measurement</li> <li>• Handling of predictors in the modelling (e.g. continuous, linear, non-linear transformations or categorised)</li> </ul>
Sample size	<ul style="list-style-type: none"> <li>• Number of participants</li> <li>• Number of outcomes/events</li> <li>• Number of outcomes/events in relation to the number of candidate predictors (events per variable)</li> </ul>
Missing data	<ul style="list-style-type: none"> <li>• How much missing data</li> <li>• Handling of missing data (e.g. complete-case analysis, imputation, or other methods)</li> </ul>
Model development	<ul style="list-style-type: none"> <li>• Modelling method (e.g. logistic or survival)</li> <li>• Method for selection of predictors for inclusion in multivariable modelling</li> <li>• Method for selection of predictors during multivariable modelling</li> </ul>
Intended moment of using the model & timing of predictor measurements	<ul style="list-style-type: none"> <li>• Intended moment of use</li> <li>• Timing of the measurement of predictors included in the final model, and whether it matched the intended moment of using the model</li> </ul>
Measurement error of predictors	<ul style="list-style-type: none"> <li>• Susceptibility to measurement error for the predictors included in the final model</li> <li>• Whether measurement error was accounted for and, if so, how</li> </ul>
Model performance	<ul style="list-style-type: none"> <li>• Calibration (e.g. calibration slope, calibration plot, Hosmer-Lemeshow test)</li> <li>• Discrimination (e.g. c-statistic, D-statistic, log-rank)</li> <li>• Classification measures (e.g. sensitivity, specificity, predictive values, net reclassification improvement)</li> </ul>
Model evaluation	<ul style="list-style-type: none"> <li>• Method used for testing model performance: internal (e.g. random split of data, resampling methods, none) or external (e.g. temporal, geographical, different setting, different investigators)</li> <li>• In case of poor validation, whether the model was adjusted or updated (e.g. intercept recalibrated, predictor effects adjusted, new predictors added)</li> </ul>

2

3 *Table 2: Predictors in final models at high risk of measurement error*

<b>Key reasons for being at high risk of error</b>	<b>Predictors included in final models</b>
Fluctuations in human samples/ biological variability	Serum albumin, Serologic markers, Prostate Specific Antigen (PSA) density, Prostate Specific Antigen (PSA), Ki-67, Human epididymis protein 4 (HE4), Glomerular filtration rate, Emergency room pulse rate, CRUSADE score, C-reactive protein, Creatinine on admission, CA125, Ascites
Inaccuracy of measurement instruments	Body Mass Index (BMI), Myometrial invasion depth, Emergency room pulse rate, Creatinine on admission, Weight, Ascites, International normalised ratio (INR1) , Infection/bioburden
Imperfect recall	Body Mass Index (BMI), Duration of convulsions, Duration of drowsiness, Duration of neck pain, Duration of nervousness, Duration of tingling, History of transactional sex, Area under pain curve, Congestive heart failure, Weight, Previous bleeding, Endoscopic retrograde cholangiopancreatography (ERCP) time, Time developing motor deficits, ImpACT total symptom score, Eastern Cooperative Oncology Group (ECOG) performance status, Depression, Number of non-major comorbidities, Systemic illness/organ failure
Subjective nature of measures	Abdominal pain, Tumour stage, Suboptimal pelvic examination or enlarged uterus during preoperative evaluation, Area under pain curve, Hypertension, Clinical stage, Malnutrition, Obesity, Procedure risk category, Pressure ulcer stage, ImpACT total symptom score, Eastern Cooperative Oncology Group (ECOG) performance status, Depression, Pre-catheterisation diagnosis
Laboratory or measurer error	Tumour stage, Suboptimal pelvic examination or enlarged uterus during preoperative evaluation, Myometrial invasion depth, CRUSADE score, CA125, Histologic grade, Primary tumour diameter, Clinical stage, Residual tumour, Endoscopic Retrograde Cholangiopancreatography (ERCP) Time, Tumour size, Pressure ulcer stage, Ascites, International normalised ratio (INR1) , Peritoneal Cancer Index, Infection/bioburden, Operating time and age, Wound (ulcer) age at first encounter

4

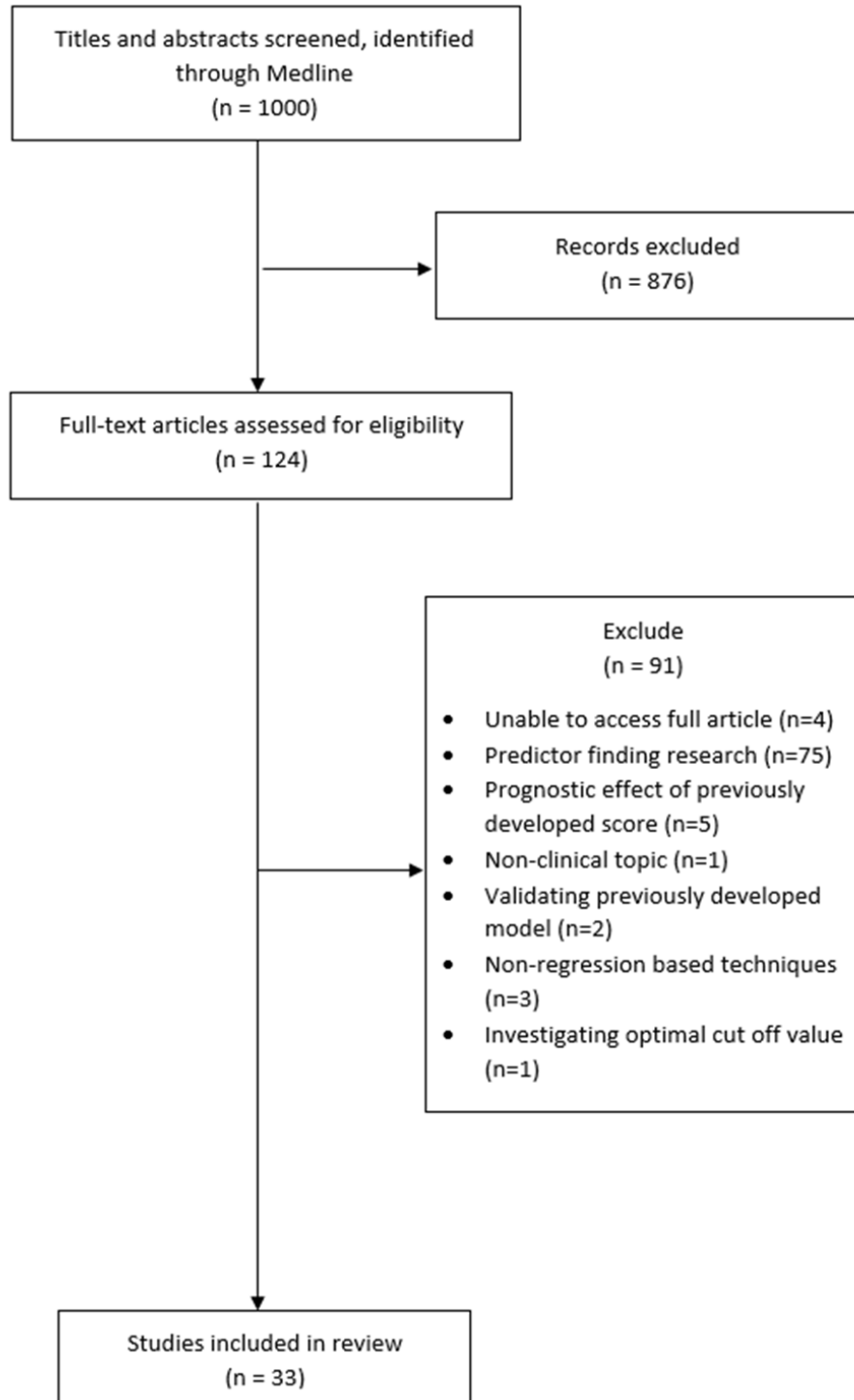


Table 3: Key reasons for measurement error in examples of predictors at high risk of being susceptible to error

Predictor	Key Reasons	Explanation
Area under pain curve	Subjective/subject to recall	Requires patient to report pain, which is a subjective measure and could report the same pain differently at a different time/by a different method or if previous scores were not provided [56], and recall incorrectly [57]
CA125	Biological variability/ laboratory error	Assay imprecision can contribute considerably to result variations in a conventional laboratory setting [58] and changes can occur due to normal biological variation [59]
Creatinine on admission	Biological variability/ inaccuracy of measurement method	Bias and imprecision may occur by use of different measurement methods [60] and changes can occur due to normal biological variation [61]
C-reactive protein	Biological variability	Within-individual variability exists, so a second sample may produce different results [62]
CRUSADE score	Biological variability/ measurer error	May get a different value if calculated again shortly afterwards as includes measures that vary and may be affected by measurer error such as of blood pressure [22]
Emergency room pulse rate	Biological variability/ inaccuracy of measurement method	May change if measured a couple of minutes later and there may be error depending on how long the measurer counted for [63]
History of transactional sex	Imperfect recall	Patient may not be truthful about history [64]
Glomerular filtration rate	Biological variability	A second sample may produce different results due to biological variation [65]
Human epididymis protein 4 (HE4)	Biological variability	A second sample may produce different results [66]
Ki-67	Biological variability	A second sample may produce different results and differences may be present from different laboratories [67]
Myometrial invasion depth	Measurer error/ inaccuracy of measurement method	Results may be different when reassessed [68]

Prostate specific antigen (PSA)	Biological variability	A second sample may produce different results [69]
Serum albumin	Biological variability	A second sample may produce different results [70]
Tumour stage	Subjective/measurer error	May get a different result from different assessors dependent on experience level or areas of speciality

Figure 1: Flow chart of included studies



ACCEPTED MANUSCRIPT

**WHAT IS NEW?****Key findings**

- Many published prediction models include predictors that are susceptible to measurement error and this measurement error is not being acknowledged or accounted for in the development of the models.
- Most prediction model articles do not explicitly state the intended moment of model use, or exactly when the predictors used in the model development were measured.

**What this adds to what is known**

- Reporting of measurement error and intended moment of model use is poor in prediction model studies.

**What is the implication, what should change now?**

- There is a need to identify circumstances where ignoring measurement error in prediction models is consequential and whether accounting for the error will improve the predictions.
- Future prediction model research studies must clearly report the intended moment of use of the prediction model, and be explicit about when the predictors were measured.