

Journal Pre-proof

Systematic review identifies the design and methodological conduct of studies on machine learning-based prediction models

Constanza L. Andaur Navarro, doctoral student, Johanna AA. Damen, assistant professor, Maarten van Smeden, associate professor, Toshihiko Takada, assistant professor, Steven WJ. Nijman, doctoral student, Paula Dhiman, research fellow, Jie Ma, medical statistician, Gary S. Collins, professor, Ram Bajpai, research fellow, Richard D. Riley, professor, Karel GM. Moons, professor, Lotty Hooft, professor



PII: S0895-4356(22)00300-6

DOI: <https://doi.org/10.1016/j.jclinepi.2022.11.015>

Reference: JCE 10961

To appear in: *Journal of Clinical Epidemiology*

Received Date: 25 July 2022

Revised Date: 9 October 2022

Accepted Date: 22 November 2022

Please cite this article as: Andaur Navarro CL, Damen JA, van Smeden M, Takada T, Nijman SW, Dhiman P, Ma J, Collins GS, Bajpai R, Riley RD, Moons KG, Hooft L, Systematic review identifies the design and methodological conduct of studies on machine learning-based prediction models, *Journal of Clinical Epidemiology* (2022), doi: <https://doi.org/10.1016/j.jclinepi.2022.11.015>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. 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.

© 2022 The Author(s). Published by Elsevier Inc.

REVIEW

Systematic review identifies the design and methodological conduct of studies on machine learning-based prediction models

Constanza L Andaur Navarro^{1,2} *doctoral student* (c.l.andaurnavarro@umcutrecht.nl, 0000-0002-7745-2887),

Johanna AA Damen^{1,2} *assistant professor* (j.a.a.damen@umcutrecht.nl, 0000-0001-7401-4593),

Maarten van Smeden¹, *associate professor* (m.vansmeden@umcutrecht.nl, 0000-0002-5529-1541),

Toshihiko Takada¹ *assistant professor* (t.takada@umcutrecht.nl, 0000-0002-8032-6224),

Steven WJ Nijman¹ *doctoral student* (S.W.J.Nijman@umcutrecht.nl, 0000-0001-6798-2078),

Paula Dhiman^{3,4} *research fellow* (paula.dhiman@ndorms.ox.ac.uk, 0000-0002-0989-0623),

Jie Ma³ *medical statistician* (jie.ma@csm.ox.ac.uk, 0000-0002-3900-1903),

Gary S Collins^{3,4} *professor* (gary.collins@csm.ox.ac.uk, 0000-0002-2772-2316),

Ram Bajpai⁵ *research fellow* (r.bajpai@keele.ac.uk, 0000-0002-1227-2703),

Richard D Riley⁵ *professor* (r.riley@keele.ac.uk, 0000-0001-8699-0735),

Karel GM Moons^{1,2} *professor* (k.g.m.moons@umcutrecht.nl, 0000-0003-2118-004X),

Lotty Hooft^{1,2} *professor* (l.hooft@umcutrecht.nl, 0000-0002-7950-2980)

¹Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands.

²Cochrane Netherlands, University Medical Center Utrecht, Utrecht University, Utrecht, The Netherlands.

³Center for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology & Musculoskeletal Sciences, University of Oxford, Oxford, United Kingdom.

⁴NIHR Oxford Biomedical Research Centre, Oxford University Hospitals NHS Foundation Trust, Oxford, United Kingdom.

⁵Centre for Prognosis Research, School of Medicine, Keele University, Keele, United Kingdom.

To be submitted to JCE

Word Count: 4036

Correspondence to: Constanza L Andaur Navarro, c.l.andaurnavarro@umcutrecht.nl

Julius Centre for Health Sciences and Primary Care, Universiteitsweg 100, P.O. Box 85500, 3508 GA, Utrecht, The Netherlands.

36 **ABSTRACT**

37 **Objective.** We sought to summarize the study design, modelling strategies, and performance
38 measures reported in studies on clinical prediction models developed using machine learning
39 techniques.

40 **Study Design and Setting.** We search PubMed for articles published between 01/01/2018 and
41 31/12/2019, describing the development or the development with external validation of a
42 multivariable prediction model using any supervised machine learning technique. No restrictions
43 were made based on study design, data source, or predicted patient-related health outcomes.

44 **Results.** We included 152 studies, 58 (38.2% [95%CI 30.8-46.1]) were diagnostic and 94 (61.8%
45 [95%CI 53.9-69.2]) prognostic studies. Most studies reported only the development of prediction
46 models (n=133, 87.5% [95%CI 81.3-91.8]), focused on binary outcomes (n=131, 86.2% [95%CI
47 79.8-90.8]), and did not report a sample size calculation (n=125, 82.2% [95%CI 75.4-87.5]). The
48 most common algorithms used were support vector machine (n=86/522, 16.5% [95%CI 13.5-
49 19.9]) and random forest (n=73/522, 14% [95%CI 11.3-17.2]). Values for area under the Receiver
50 Operating Characteristic curve ranged from 0.45 to 1.00. Calibration metrics were often missed
51 (n=494/522, 94.6% [95%CI 92.4-96.3]).

52 **Conclusions.** Our review revealed that focus is required on handling of missing values, methods for
53 internal validation, and reporting of calibration to improve the methodological conduct of studies
54 on machine learning-based prediction models.

55 **Systematic review registration** PROSPERO, CRD42019161764.

56 **Word count:** 201/200

57

58 **Keywords:** predictive algorithm, risk prediction, diagnosis, prognosis, development, validation.

59 **Running title:** Review on studies on ML-based prediction models

60 INTRODUCTION

61 Clinical prediction models aim to improve healthcare by providing timely information for shared
62 decision-making between clinician and their patients, risk stratification, changes in behaviour, and
63 to counsel patients and their relatives.¹ A prediction model can be defined as the (weighted)
64 combination of several predictors to estimate the likelihood or probability of the presence or
65 absence of a certain disease (diagnostic model), or the occurrence of an outcome over a time
66 period (prognostic model).² Traditionally, prediction models were developed using regression
67 techniques, such as logistic or time-to-event regression. However, in the past decade, the attention
68 and use of machine learning approaches to developing clinical prediction models has rapidly grown.

69 Machine learning can be broadly defined as the use of computer systems that fit mathematical
70 models that assume non-linear associations and complex interactions. Machine learning has a
71 wide range of potential applications in different pathways of healthcare. For example, machine
72 learning is applied in stratified medicine, triage tools, image-driven diagnosis, online consultations,
73 medication management, and to mine electronic medical records.³ Most of these applications
74 make use of supervised machine learning whereby a model is fitted to learn the conditional
75 distribution of the outcome given a set of predictors with little assumption on data distributions,
76 non-linear associations, and interactions. This model can be later applied in other but related
77 individuals to predict their (yet unknown) outcome. Support vector machines (SVM), random forests
78 (RF), and neural networks (NN) are some examples of these techniques.⁴

79 The number of studies on prediction models published in the biomedical literature increases every
80 year.^{5,6} With more healthcare data being collected and increasing computational power, we expect
81 studies on clinical prediction models based on (supervised) machine learning techniques to
82 become even more popular. Although numerous models are being developed and validated for
83 various outcomes, patients' populations, and healthcare settings, only a minority of these published
84 models are successfully implemented in clinical practice.^{7,8}

85 The use of appropriate study designs and prediction model strategies to develop or validate a
86 prediction model could improve their transportability into clinical settings.⁹ However, currently there
87 is a dearth of information about which study designs, what modelling strategies, and which
88 performance measures do studies on clinical prediction models report when choosing machine
89 learning as modelling approach.¹⁰⁻¹² Therefore, our aim was to systematically review and
90 summarise the characteristics on study design, modelling steps, and performance measures
91 reported in studies of prediction models using supervised machine learning.

92 **METHODS**

93 We followed the PRISMA 2020 statement to report this systematic review.¹³

94 **Eligibility criteria**

95 We searched via PubMed (search date 19 December 2019) for articles published between 1
96 January 2018 and 31 December 2019 (Supplemental file 1). We focused on primary studies that
97 described the development and/or validation of one or more multivariable diagnostic or prognostic
98 prediction model(s) using any supervised machine learning technique. A multivariable prediction
99 model was defined as a model aiming to predict a health outcome by using two or more predictors
100 (features). We considered a study to be an instance of supervised machine learning when reporting
101 a non-regression approach to model development. If a study reported machine learning models
102 alongside regression-based models, this was included. We excluded studies reporting only
103 regression-based approaches such as unpenalized regression (e.g., ordinary least squares or
104 maximum likelihood logistic regression), or penalized regression (e.g., lasso, ridge, elastic net, or
105 Firth's regression), regardless of whether they referred to them as machine learning. Any study
106 design, data source, study population, predictor type or patient-related health outcome was
107 considered.

108 We excluded studies investigating a single predictor, test, or biomarker. Similarly, studies using
109 machine learning or AI to enhance the reading of images or signals, rather than predicting health
110 outcomes in individuals, or studies that used only genetic traits or molecular ('omics') markers as
111 predictors, were excluded. Furthermore, we also excluded reviews, meta-analyses, conference
112 abstracts, and articles for which no full text was available via our institution. The selection was
113 restricted to humans and English-language studies. Further details about eligibility criteria can be
114 found in our protocol.¹⁴

115 **Screening and selection process**

116 Titles and abstracts were screened to identify potentially eligible studies by two independent
117 reviewers from a group of seven (CLAN, TT, SWJN, PD, JM, RB, JAAD). After selection of potentially
118 eligible studies, full text articles were retrieved and two independent researchers reviewed them
119 for eligibility; one researcher (CLAN) screened all articles and six researchers (TT, SWJN, PD, JM,
120 RB, JAAD) collectively screened the same articles for agreement. In case of any disagreement
121 during screening and selection, a third reviewer was asked to read the article in question and
122 resolve.

123 **Extraction of data items**

124 We selected several items from existing methodological guidelines for reporting and critical
125 appraisal of prediction model studies to build our data extraction form (CHARMS, TRIPOD,
126 PROBAST).¹⁵⁻¹⁸ Per study, we extracted the following items: characteristics of study design (e.g.

127 cohort, case-control, randomized trial) and data source (e.g. routinely collected data, registries,
128 administrative databases), study population, outcome, setting, prediction horizon, country, patient
129 characteristics, sample size (before and after exclusion of participants), number of events, number
130 of candidate and final predictors, handling of missing data, hyperparameter optimization, dataset
131 splitting (e.g. train-validation-test), method for internal validation (e.g. bootstrapping, cross
132 validation), number of models developed and/or validated, and availability of code, data and
133 model. We defined country as the location of the first author's affiliation. Per model, we extracted
134 information regarding the following items: type of algorithm used, selection of predictors, reporting
135 of variable importance, penalization techniques, reporting of hyperparameters, and metrics of
136 performance (e.g., discrimination and calibration).

137 Items were recorded by two independent reviewers. One reviewer (CLAN) recorded all items, whilst
138 the other reviewers collectively assessed all articles (CLAN, TT, SWJN, PD, JM, RB, JAAD). Articles
139 were assigned to reviewers in a random manner. To accomplish consistent data extraction, the
140 standardized data extraction form was piloted by all reviewers on five articles. Discrepancies in
141 data extraction were discussed and solved between the pair of reviewers. The full list of extracted
142 items is available in our published protocol.¹⁴

143 We extracted information on a maximum number of 10 models per article. We selected the first 10
144 models reported in the methods section of articles and extracted items accordingly in the results
145 section. For articles describing external validation or updating, we carried out a separate data
146 extraction with similar items. If studies referred to the supplemental file for detailed descriptions,
147 the items were checked in those files. Reviewers could also score an item as not applicable, not
148 reported, or unclear.

149 **Summary measures and synthesis of results**

150 Results were summarized as percentages (with confidence intervals calculated using the Wilson
151 score interval and the Wilson score continuity-corrected interval, when appropriated), medians, and
152 interquartile range (IQR), alongside a narrative synthesis. The reported number of events was
153 combined with the reported number of candidate predictors to calculate the number of events per
154 variable (EPV). Data on a model's predictive performance was summarized for the apparent
155 performance, corrected performance, and externally validated performance. We defined "apparent
156 performance" when studies reported model performance assessed in the same dataset or sample
157 in which the model was developed and in case no re-sampling methods were used; "corrected
158 performance" when studies reported model performance assessed in test dataset and/or using re-
159 sampling methods; and "externally validated performance" when studies reported model
160 performance assessed in another sample than the one use for model development. As we wanted
161 to identify the methodological conduct of studies on prediction models developed using machine
162 learning, we did not evaluate the nuances of each modelling approach or its performance, instead

163 we kept our evaluations at study level. We did not perform a quantitative synthesis of the model'
164 performance (i.e., meta-analysis), as this was beyond the scope of our review. Analysis and
165 synthesis of data was presented overall. Analyses were performed using R (version 4.1.0, R Core
166 Team, Vienna, Austria).

Journal Pre-proof

167 RESULTS

168 Among 24,814 articles retrieved, we drew a random sample of 2482 articles. After title and
169 abstract screening, 312 references potentially met the eligibility criteria. After full-text screening,
170 152 articles were included in this review: 94 (61.8% [95%CI 53.9-69.2]) prognostic and 58 (38.2%
171 [95%CI 30.8-46.1]) diagnostic prediction model studies (Figure 1). Detailed description of the
172 included articles is provided in Supplemental file 2.

173 In 152 articles, 132 (86.8% [95%CI 80.5-91.3]) studies developed prediction models and
174 evaluated their performance using an internal validation technique, 19 (12.5% [95%CI 8.2-18.7])
175 studies developed and externally validated the same ML based prediction model, and 1 (0.6%)
176 study included model development with external validation of another comparative model
177 (eventually included as development with internal validation). Eighty-seven studies (57% [95% CI
178 49.3-64.8]) were published in 2019 and 65/152 studies (42.8% [95% CI 35.2-50.7]) in 2018. The
179 three clinical fields with the most articles were oncology (n=21/152, 13.8% [95%CI 9.2–20.2]),
180 surgery (n=20/152, 13.5% [95%CI 8.7-19.5]), and neurology (n=20/152, 13.5% [95%CI 8.7-19.5]).
181 Most articles originated from North America (n=59/152, 38.8% [95%CI 31.4-46.7]), followed by
182 Asia (n=46/152, 30.3% [95%CI 23.5-38]) and Europe (n=37/152, 24.3% [95%CI 18.2-31.7]). Half
183 of the studies had a first author with a clinical affiliation (n=85/152, 56% [95%CI 48-63.6]). Other
184 characteristics are shown in Table 1.

185 Overall, 1,429 prediction models were developed (Median: 9.4 models per study, IQR: 2-8, Range:
186 1-156). As we set a limit on data extraction to 10 models per article, we evaluated 522 models.
187 The most common applied modeling techniques were support vector machine (n=86/522, 16.5%
188 [95%CI 13.5-20]), logistic regression (n=74/522, 14.2% [95%CI 11.4-17.5]), and random forest
189 (n=73/522, 14% [95%CI 11.2-17.3]). Further modelling algorithms are described in Table 2. In
190 120/152 (78.9% [95%CI 71.8-84.7]) articles, authors recommended at least one model usually
191 based on model performance (i.e., AUC).

192 *Participants*

193 Participants included in the reviewed studies were mostly recruited from secondary (n=32/152,
194 21.1% [95%CI 15.3-28.2]) and tertiary care (n=78/152, 51.3% [95%CI 43.4-59.1]) settings (Table
195 1). Approximately half of the studies involved data from one center (n=73/152, 48% [95%CI 40.2-
196 55.9]) (Table 3).

197 *Data sources*

198 The prediction models were most frequently developed using cohort data, either prospective
199 (n=50/152, 32.9% [95%CI 25.9-40.7]) or retrospective (n=48/152, 31.6% [95%CI 24.7-39.3]).
200 Electronic medical records were used in 30/152 studies (19.7% [95%CI 14.2-26.8]). Data
201 collection was conducted on average for 41.9 months (IQR 3 to 60 months) when used to develop
202 models, while for externally validation this was 44.4 months (IQR 1.75 to 42 months). In 101 out

203 of 152 studies (66.4% [95%CI 58.6-73.5]), the time horizon for the predictions was mostly
204 unspecified. However, when reported (n=51/152, 33.6% [95%CI 26.5-41.4]), the time horizon of
205 prediction ranged from 24 hours to 8 years (Table 3).

206 **Outcome**

207 Most models were developed to predict a binary outcome (n=131/152, 86.2% [95%CI 79.8-90.8]).
208 The most frequent predicted outcome was complications after a certain treatment (n=66/152,
209 43.4% [95%CI 35.8-51.4]). Mortality was also a common endpoint (n=21/152, 13.8% [95%CI 9.2-
210 20.2]) (Table 1).

211 **Candidate predictors**

212 Candidate predictors frequently involved demographics, such as age and sex (n=120/152, 78.9%
213 [95%CI 71.8-84.7]), clinical history (n=111/152, 73% [95%CI 65.5-79.4]), and blood and urine
214 parameters (n=63/152, 41.4% [95%CI 33.9-49.4]). When applicable, treatment modalities were
215 also considered as predictors (n=36/116, 31.0% [95%CI 17.6-31]). Studies included a median of
216 24 candidate predictors (IQR 13 – 112). Most studies included continuous variables as candidate
217 predictors (n=131/152, 86.2% [95%CI 79.8-90.8]). Whether continuous predictors were
218 categorized during data preparation was often unclear (n=104/152, 68.4% [95%CI 60.7-75.3])
219 (Table 4).

220 **Sample size**

221 Studies had a median sample size of 587 participants (IQR 172 – 6328). The number of events
222 across the studies had a median of 106 (IQR 50 – 364). Based on studies with available
223 information (n=28/152, 18.4% [95%CI 13.1-25.3]), a median of 12.5 events per candidate
224 predictors were used for model development (IQR 5.7 – 27.7) (Table 5). Most studies did not report
225 a sample size calculation or justification for sample size (n=125/152, 82.2% [95%CI 75.4-87.5]).
226 When sample size justification was provided, the most frequent rationale given was based on the
227 size of existing/available data used (n=16/27, 59.3% [95%CI 40.7-75.5]) (Table 3).

228 **Missing values**

229 Missing values were an explicit exclusion criterion of participants in 56 studies (n=56/152, 36.8%
230 [95%CI 29.6-44.7]). To handle missing values, complete-case analysis was the most common
231 method (n=30/152, 19.7% [95%CI 14.2-26.8]). Other methods were median imputation
232 (n=10/152, 6.6% [95%CI 3.6-11.7]), multiple imputation (n=6/152, 3.9% [95%CI 1.9-8.3]) and k-
233 nearest neighbor imputation (n=5/152, 3.3% [95%CI 1.4-7.5]). Further methods to handle missing
234 values are presented in Table 6.

235 **Class imbalance**

236 In our sample, 27/152 (17.8% [95%CI 12.5-24.6]) studies applied at least one method to
237 purportedly address class imbalance, that is – when one class of the outcome outnumbers the

238 other class (Table 7). The most applied technique was Synthetic Minority Over-sampling Technique
239 (SMOTE), a method that combines oversampling the minority class with undersampling the majority
240 class.^{19,20}

241 **Modelling algorithms**

242 Tree-based methods were applied in 166/522 (31.8% [95%CI 27.9-36]) models with random forest
243 being the most popular (n=73/522, 14% [95%CI 11.2-17.3]). Alongside machine learning
244 algorithm, unpenalized regression methods (n=101/522, 19.3% [95%CI 16.1-23.1]), and
245 particularly logistic regression (n=74/522, 14.2 [95%CI 11.4-17.5]) were often applied. Few
246 studies reported models built with penalized regression (n=29/522, 5.6% [95%CI 3.8-8]). NNs
247 (n=74/522, 14.2% [95%CI 11.4-17.5]) and Naïve Bayes (n=22/522, 4.2% [95%CI 2.7-6.4]) were
248 also applied in our sample of articles.

249 **Selection of predictors**

250 The strategy to build models was unclear in 168 out of 522 models (32.2% [95%CI 28.2-36.4]).
251 Most models reported a data-driven approach for model building (n=192/522, 36.8% [95%CI 32.7-
252 41.1]). One study reported the use of recursive feature elimination for model building (n=3/522,
253 0.6% [95%CI 0.1-1.8]). Selection of candidate predictors based on univariable predictor-outcome
254 associations was used in 27/522 (5.2% [95%CI 3.5-7.5]) of the models. Further details on
255 modelling strategies are presented in Table 8. Of the three studies that reported time-to-event
256 outcomes none reported how they dealt with censoring.

257 **Variable importance and hyperparameters**

258 Variable importance scores show insight into how much each variable contributed to the prediction
259 model.²¹ For 316/522 (60.5% [95%CI 56.2-64.7]) models, authors did not provide these scores,
260 while in 115/522 (22% [95%CI 18.6-25.9]) models these scores were reported without specifying
261 the methods applied to obtain such calculations (Table 8). When reported, the mean decrease in
262 node impurity was the most popular method (n=31/522, 5.9% [95%CI 4.1-8.4]). Hyperparameters
263 (including default settings) were reported in 160/552 (30.7% [95%CI 26.8-34.8]) models.
264 Strategies for hyperparameter optimization were described in 44/152 studies (28.9% [95%CI 22.3-
265 36.3]). The most common method reported was cross-validation (n=15/152) [9.9% [95%CI 6.1-
266 15.6]. Nine studies (n=9/152, 5.9% [95%CI 3.1-10.9]) split their dataset into a validation set for
267 hyperparameter tuning (Table 7).

268 **Performance metrics**

269 Most models used measures of the area under the Receiver Operating Characteristic curve
270 (AUC/ROC or the concordance (c)-statistic) (n=358/522, 68.6% [95%CI 64.4-72.5]) to describe the
271 discriminative ability of the model (Table 9). A variety of methods were used to describe the
272 agreement between predictions and observations (i.e., calibration), the most frequent being a
273 calibration plot (n=23/522, 4.4% [95%CI 2.9-6.6]), calibration slope (n=17/522, 3.3% [95%CI 2-

274 5.3]), and calibration intercept (n=16/522, 3.1% [95%CI 1.8-5]). However, for the large majority no
275 calibration metrics were reported (n=494/522, 94.6% [95%CI 92.2-96.3]). Decision curve analysis
276 was reported for two models (n=2/522, 0.4% [95%CI 0.1-1.5]).²² We also found overall metrics
277 such as classification accuracy (n=324/522, 62.1% [95%CI 57.8-66.2]) and F1-score (n=79/522,
278 15.1% [95%CI 12.2-18.6]).

279 ***Uncertainty quantification***

280 In 53/152 (34.9% [95% CI 22.8-42.7]) studies, discrimination was reported without precision
281 estimates (i.e., confidence intervals or standard errors). Likewise, 7/152 (4.6% [95%CI 2.2-9.2])
282 studies reported model calibration without precision estimates.

283 ***Predictive performance***

284 Most models achieved discriminative ability better than chance (i.e., AUC 0.5) with a median
285 apparent AUC of 0.82 (IQR 0.75-0.90; range 0.45 to 1.00), while internally validated AUC was also
286 0.82 (IQR: 0.74-0.89; range 0.46 to 0.99). For external validation, the median AUC was 0.73 (IQR:
287 0.70-0.78, range: 0.51-0.88). For calibration and overall performance metrics, see Table 10.

288 ***Internal validation***

289 In total, 86/152 studies (56.6% [95%CI 48.6-64.2]) internally validated their models, most often
290 splitting the dataset into a training and test set. The train-test sets were often split randomly
291 (n=49/86, 57% [95%CI 46.4-66.9]) and in a few studies a temporal (non-random) split was applied
292 (n=9/86, 10.5% [95%CI 5.6-18.7]). The proportion of the data used for test sets ranged from 10%
293 to 50% of the total dataset. Seventy studies also performed cross-validation (46.1% [95%CI 38.3-
294 54]) with ten studies reporting nested cross-validation (6.6% [95%CI 3.6-11.7]). Out of five studies
295 performing bootstrapping (n=5/152, 3.3% [95%CI 1.4-7.5]), one reported 250 iterations, three
296 reported 1000 iterations and one did not report the number of iterations. For further details see
297 Table 3.

298 ***External validation***

299 Few studies (n=19/152, 12.5% [95%CI 8.2-18.7]) performed an external validation. Eleven studies
300 (n=11/19, 57.9% [95%CI 36.3-76.9]) used data from independent cohorts and eight (n=8/19,
301 42.1% [95%CI 23.1-63.7]) used subcohorts within the main cohort to validate their developed
302 models. From the independent cohorts, three studies (n=3/19, 15.8% [95%CI 5.5-37.6]) used data
303 from a different country. Five studies (n=5/19, 26.3% [95%CI 11.8-48.8]) described an external
304 validation based on temporal differences on the inclusion of participants. Seven studies (36.8%
305 [95%CI 19.1-59]) reported differences and similarities in definitions between the development and
306 validation data.

307 ***Model availability***

308 Some studies shared their prediction model either as a web-calculator or worked example
309 (n=31/152, 20.4% [95%CI 14.8-27.5]). Furthermore, in a minority of studies datasets and code
310 were accessible through repositories, which were shared as supplemental material (n=18/152,
311 11.8% [95%CI 7.6-17.9]; n=13/152, 8.6% [95%CI 5.1-14.1]). Details in Table 1.

Journal Pre-proof

312 DISCUSSION

313 Principal findings

314 In this study, we evaluated the study design, data sources, modelling steps, and performance
315 measures in studies on clinical prediction models using machine learning across. The methodology
316 varied substantially between studies, including modelling algorithms, sample size, and
317 performance measures reported. Unfortunately, longstanding deficiencies in reporting and
318 methodological conduct previously seen in studies with a regression-based approach, were also
319 extensively found in our sample of studies on machine learning models.^{9,23}

320 The spectrum of supervised machine learning techniques is quite broad.^{24,25} In this study, the most
321 popular modelling algorithms were tree-based methods (RF in particular) and SVM. RF is an
322 ensemble of random trees trained on bootstrapped sub-sets of the dataset.²⁶ On the other hand,
323 SVM first map each data point into a feature space to then identify the hyperplane that separates
324 the data items into two classes while maximizing the marginal distance for both classes and
325 minimizing the classification errors.²⁷ Several studies also applied regression-based methods (LR
326 in particular) as benchmark to compare against the predictive performance of machine learning-
327 based models.

328 Various other well-known methodological issues in prediction model research need to be further
329 discussed. Our reported estimate on EPV is likely to be overestimated given than we were unable
330 to calculate it based on number of parameters, and instead we used only the number of candidate
331 predictors. A simulation study concluded that modern modelling techniques such as SVM and RF
332 might even require 10 times more events.²⁸ Hence, the sample size in most studies on prediction
333 models using machine learning remains relatively low. Furthermore, splitting datasets persists as
334 a method for internal validation (i.e., testing), reducing even more the actual sample size for model
335 development and increasing the risk of overfitting.^{29,30} Whilst AUC was a frequently reported metric
336 to assess predictive performance, prediction calibration or prediction error was often overlooked.³¹
337 Moreover, a quarter of studies in our sample corrected for class imbalance without reporting
338 recalibration, although recent research has shown that correcting for class imbalance may lead to
339 poor calibration and thus, prediction errors.³² Finally, therapeutic interventions were rarely
340 considered as predictors in the prognostic models, although these can affect the accuracy and
341 transportability of models.³³

342 Variable importance scores, tuning of hyperparameters, and data preparation (i.e., data pre-
343 processing) are items closely related to machine learning prediction models. We found that most
344 studies reporting variable importance scores did not specify the calculation method. Data
345 preparation steps (i.e., data quality assessment, cleaning, transformation, reduction) were often
346 not described in enough transparent detail. Complete-case analysis remains a popular method to
347 handle missing values in machine learning based models. Detailed description and evaluation on

348 how missing values were handled in our included studies has been provided elsewhere.³⁴ Last, only
349 one third of models reported their hyperparameters settings, which is needed for reproducibility
350 purposes.

351 **Comparison to previous studies**

352 Regression methods were not our focus (as we did not define them to be machine learning
353 methods), but other reviews including both approaches show similar issues with methodological
354 conduct and reporting.^{12,35-37} Missing data, sample size, calibration, and model availability remain
355 largely neglected aspects.^{7,12,37-40} A review looking at the trends of prediction models using
356 electronic health records (EHR) observed an increase in the use of ensemble models from 6% to
357 19%.⁴¹ Another detailed review on prediction models for hospital readmission shows that the use
358 of algorithms such as SVM, RF, and NN increased from none to 38% over the last five years.¹⁰
359 Methods to correct for class imbalance in datasets concerning EHR increased from 7% to 13%.⁴¹

360 **Strengths and limitations of this study**

361 In this comprehensive review, we summarized the study design, data sources, modelling strategies,
362 and reported predictive performance in a large and diverse sample of studies on clinical prediction
363 model studies. We focused on all types of studies on clinical prediction models rather than on a
364 specific type of outcome, population, clinical specialty, or methodological aspect. We appraised
365 studies published almost three years ago and thus, it is possible that further improvements might
366 have raised. However, improvements in methodology and reporting are usually small and slow even
367 when longer periods are considered.⁴² Hence, we believe that the results presented in this
368 comprehensive review still largely apply to the current situation of studies on machine learning-
369 based prediction models. Given the limited sample, our findings can be considered a representative
370 rather than exhaustive description of studies on machine learning models.

371 Our data extraction was restricted to what was reported in articles. Unfortunately, few articles
372 reported the minimum information required by reporting guidelines, thereby hampering data-
373 extraction.²³ Furthermore, terminology differed between papers. For example, the term 'validation'
374 was often used to describe tuning, as well as testing (i.e., internal validation). An issue already
375 observed by a previous review of studies on deep learning models.⁴³ This shows the need to
376 harmonize the terminology for critical appraisal of machine learning models.⁴⁴ Our data extraction
377 form was based mainly on the items and signaling questions from TRIPOD and PROBAST. Although
378 both tools were primarily developed for studies on regression-based prediction models, most items
379 and signaling questions were largely applicable for studies on machine learning-based models, as
380 well.

381 **Implication for researchers, editorial offices, and future research**

382 In our sample, it is questionable whether studies ultimately aimed to improve clinical care.⁴⁵ Aim,
383 clinical workflow, outcome format, prediction horizon, and clinically relevant performance metrics

384 received very little attention. The importance of applying optimal methodology and transparent
385 reporting in studies on prediction models has been intensively and extensively stressed by
386 guidelines and meta-epidemiological studies.⁴⁶⁻⁴⁸ Researchers can benefit from TRIPOD and
387 PROBAST, as these provide guidance on best practices for prediction model study design, conduct
388 and reporting regardless of their modelling technique.^{16,17,46,47} However, special attention is
389 required on extending the recommendations to include areas such as data preparation, tunability,
390 fairness, and data leakage. In this review, we have provided evidence on the use and reporting of
391 methods to correct for class imbalance, data preparation, data splitting, and hyperparameter
392 optimization. PROBAST-AI and TRIPOD-AI, both extensions to artificial intelligence (AI) or machine
393 learning based prediction models are underway.^{44,49} As machine learning continues to emerge as
394 a relevant player in healthcare, we recommend researchers and editors to reinforce a minimum
395 standard on methodological conduct and reporting to ensure further transportability.^{16,17,46,47}

396 We identified that studies covering the general population (e.g., for personalized screening),
397 primary care settings, and time-to-event outcomes are underrepresented in current research.
398 Similarly, only a relatively small proportion of the studies evaluated (validated) their prediction
399 model on a different dataset (i.e., external validation).⁵⁰ In addition, the poor availability of the
400 developed models hampers further independent validation, an important step before their
401 implementation in clinical practice. Sharing the code and ultimately the clinical prediction model is
402 a fundamental step to create trustworthiness on AI and machine learning for clinical application.⁵¹

403 **CONCLUSIONS**

404 Our study provides a comprehensive overview of the applied study designs, data sources, modelling
405 steps, and performance measures used. Special focus is required in areas such as handling of
406 missing values, methods for internal validation, and reporting of calibration to improve the
407 methodological conduct of studies on prediction models developed using machine learning
408 techniques.

409 Registration and protocol

410 This review was registered in PROSPERO (CRD42019161764). The study protocol can be accessed in doi:10.1136/bmjopen-2020-
411 038832.

412 Support

413 GSC is funded by the National Institute for Health Research (NIHR) Oxford Biomedical Research Centre (BRC) and by Cancer Research
414 UK program grant (C49297/A27294). PD is funded by the NIHR Oxford BRC. RB is affiliated to the National Institute for Health and Care
415 Research (NIHR) Applied Research Collaboration (ARC) West Midlands. The views expressed are those of the authors and not necessarily
416 those of the NHS, NIHR, or Department of Health and Social Care. None of the funding sources had a role in the design, conduct,
417 analyses, or reporting of the study or in the decision to submit the manuscript for publication.

418 Competing interests

419 None

420 Availability of data, code, and other materials

421 Articles that support our findings are publicly available. Template data collection forms, detailed data extraction on all included studies,
422 and analytical code are available upon reasonable request.

423 Supplemental Material

424 Supplemental file 1. Search strategy

425 Supplemental file 2. Characteristics and citation of included articles

426 Acknowledgements

427 The authors would like to thank and acknowledge the support of René Spijker, information specialist. The peer-reviewers are thanked
428 for critically reading the manuscript and suggesting substantial improvements.

429 Authors' contributions

430 **Constanza L. Andaur Navarro:** Conceptualization, Methodology, Investigation, Data Curation, Formal analysis, Writing - original draft,
431 Writing - review & editing; **Johanna A.A. Damen:** Conceptualization, Methodology, Investigation, Writing - review & editing, Supervision;

432 **Maarten van Smeden:** Conceptualization, Writing - review & editing; **Toshihiko Takada:** Investigation, Writing - review & editing.

433 **Steven WJ Nijman:** Investigation, Writing - review & editing; **Paula Dhiman:** Conceptualization, Methodology, Investigation, Writing - review
434 & editing; Jie Ma: Investigation, Writing - review & editing; **Gary S Collins:** Conceptualization, Methodology, Writing - review & editing; **Ram**

435 **Bajpai:** Investigation, Writing - review & editing; **Richard D Riley:** Conceptualization, Methodology, Writing - review & editing; **Karel GM**

436 **Moons:** Conceptualization, Methodology, Writing - review & editing, Supervision; **Lotty Hoofst:** Conceptualization, Methodology, Writing -
437 review & editing, Supervision

438 Ethical approval

439 Not required for this work.

References

1. Moons KGM, Royston P, Vergouwe Y, Grobbee DE, Altman DG. Prognosis and prognostic research: What, why, and how? *BMJ*. 2009;338(7706):1317-1320. doi:10.1136/bmj.b375
2. van Smeden M, Reitsma JB, Riley RD, Collins GS, Moons KG. Clinical prediction models: diagnosis versus prognosis. *J Clin Epidemiol*. 2021;132:142-145. doi:10.1016/j.jclinepi.2021.01.009
3. Meskó B, Görög M. A short guide for medical professionals in the era of artificial intelligence. *npj Digit Med*. 2020;3(1). doi:10.1038/s41746-020-00333-z
4. Uddin S, Khan A, Hossain ME, Moni MA. Comparing different supervised machine learning algorithms for disease prediction. *BMC Med Inform Decis Mak*. 2019;19(1). doi:10.1186/S12911-019-1004-8
5. Macleod MR, Michie S, Roberts I, et al. Biomedical research: Increasing value, reducing waste. *Lancet*. 2014;383(9912):101-104. doi:10.1016/S0140-6736(13)62329-6
6. Jong Y de, Ramspek CL, Zoccali C, Jager KJ, Dekker FW, Diepen M van. Appraising prediction research: a guide and meta-review on bias and applicability assessment using the Prediction model Risk Of Bias ASsessment Tool (PROBAST). *Nephrology*. 2021:1-9. doi:10.1111/NEP.13913
7. Damen JAAG, Hooft L, Schuit E, et al. Prediction models for cardiovascular disease risk in the general population: Systematic review. *BMJ*. 2016;353. doi:10.1136/BMJ.l2416
8. Collins GS, Mallett S, Omar O, Yu LM. Developing risk prediction models for type 2 diabetes: A systematic review of methodology and reporting. *BMC Med*. 2011;9(1):103. doi:10.1186/1741-7015-9-103
9. Andaur Navarro CL, Damen JAA, Takada T, et al. Risk of bias in studies on prediction models developed using supervised machine learning techniques: systematic review. *BMJ*. 2021;375:n2281. doi:10.1136/bmj.n2281
10. Artetxe A, Beristain A, Graña M. Predictive models for hospital readmission risk: A systematic review of methods. *Comput Methods Programs Biomed*. 2018;164:49-64. doi:10.1016/j.cmpb.2018.06.006
11. Stafford IS, Kellermann M, Mossotto E, Beattie RM, MacArthur BD, Ennis S. A systematic review of the applications of artificial intelligence and machine learning in autoimmune diseases. *npj Digit Med*. 2020;3(1). doi:10.1038/s41746-020-0229-3
12. Dhiman P, Ma J, Andaur Navarro CL, et al. Methodological conduct of prognostic prediction models developed using machine learning in oncology: a systematic review. *BMC Med Res Methodol*. 2022;22(1):1-16. doi:10.1186/s12874-021-01469-6
13. Page MJ, McKenzie JE, Bossuyt PM, et al. The PRISMA 2020 statement: An updated guideline for reporting systematic reviews. *BMJ*. 2021;372. doi:10.1136/bmj.n71
14. Andaur Navarro CL, Damen JAAG, Takada T, et al. Protocol for a systematic review on the methodological and reporting quality of prediction model studies using machine learning techniques. *BMJ Open*. 2020;10(11):1-6. doi:10.1136/bmjopen-2020-038832
15. Moons KGM, de Groot JAH, Bouwmeester W, et al. Critical Appraisal and Data Extraction for Systematic Reviews of Prediction Modelling Studies: The CHARMS Checklist. *PLoS Med*. 2014;11(10). doi:10.1371/journal.pmed.1001744
16. Moons KGM, 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. doi:10.7326/M14-0698
17. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): The TRIPOD Statement. *Ann Intern Med*. 2015;162(1):55. doi:10.7326/M14-0697
18. Luo W, Phung D, Tran T, et al. Guidelines for developing and reporting machine learning predictive models in biomedical research: A multidisciplinary view. *J Med Internet Res*. 2016;18(12). doi:10.2196/jmir.5870
19. Chawla N V, Bowyer KW, Hall LO, Kegelmeyer WP. SMOTE: Synthetic Minority Over-sampling Technique. *J Artif*

- Intell Res*. Published online 2002:321-357.
20. Seiffert C, Khoshgoftaar TM, Van Hulse J, Napolitano A. RUSBoost: A hybrid approach to alleviating class imbalance. *IEEE Trans Syst Man, Cybern Part A Systems Humans*. 2010;40(1):185-197. doi:10.1109/TSMCA.2009.2029559
 21. Fisher A, Rudin C, Dominici F. All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously. *J Mach Learn Res*. 2019;20:1-81.
 22. Van Calster B, Wynants L, Verbeek JFM, et al. Reporting and Interpreting Decision Curve Analysis: A Guide for Investigators. *Eur Urol*. 2018;74(6):796-804. doi:10.1016/j.eururo.2018.08.038
 23. Andaur Navarro CL, Damen JAA, Takada T, et al. Completeness of reporting of clinical prediction models developed using supervised machine learning: a systematic review. *BMC Med Res Methodol*. 2022;22(1):12. doi:10.1186/s12874-021-01469-6
 24. Beam AL, Kohane IS. Big data and machine learning in health care. *JAMA - J Am Med Assoc*. 2018;319(13):1317-1318. doi:10.1001/jama.2017.18391
 25. Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. 2nd editio. Springer-Verlag; 2009.
 26. Breiman L. Random forests. *Mach Learn*. 2001;45(1):5-32. doi:10.1023/A:1010933404324
 27. Scholkopf B, Smola AJ. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press; 2001.
 28. Ploeg T Van Der, Austin PC, Steyerberg EW. Modern modelling techniques are data hungry : a simulation study for predicting dichotomous endpoints. *BMC Med Res Methodol*. 2014;14:137. doi:10.1186/1471-2288-14-137
 29. Steyerberg EW, Harrell Jr FE, Borsboom GJ, Eijkemans RM, Vergouwe Y, Habbema J. Internal validation of predictive models: Efficiency of some procedures for logistic regression analysis. *Proc ICED 2007, 16th Int Conf Eng Des*. 2007;DS 42:774-781.
 30. Steyerberg E, Jr FH. Prediction models need appropriate internal, internal-external, and external validation. *J Clin Epidemiol*. 2016;69:245-247. doi:10.1016/j.jclinepi.2015.04.005
 31. Van Calster B, McLernon DJ, Van Smeden M, et al. Calibration: The Achilles heel of predictive analytics. *BMC Med*. 2019;17(1):1-7. doi:10.1186/s12916-019-1466-7
 32. Goorbergh R van den, van Smeden M, Timmerman D, Van Calster B. The harm of class imbalance corrections for risk prediction models: illustration and simulation using logistic regression. *J Am Med Informatics Assoc*. 2022;00(0):1-10.
 33. Pajouheshnia R, Damen JAAG, Groenwold RHH, Moons KGM, Peelen LM. Treatment use in prognostic model research: a systematic review of cardiovascular prognostic studies. *Diagnostic Progn Res*. 2017;1(1):1-10. doi:10.1186/s41512-017-0015-0
 34. Nijman SWJ, Leeuwenberg AM, Beekers I, et al. Missing data is poorly handled and reported in prediction model studies using machine learning: a literature review. *J Clin Epidemiol*. 2022;142:218-229. doi:10.1016/j.jclinepi.2021.11.023
 35. Dhiman P, Ma J, Andaur Navarro C, et al. Reporting of prognostic clinical prediction models based on machine learning methods in oncology needs to be improved. *J Clin Epidemiol*. 2021;138:60-72. doi:10.1016/j.jclinepi.2021.06.024
 36. Heus P, Reitsma JB, Collins GS, et al. Transparent Reporting of Multivariable Prediction Models in Journal and Conference Abstracts: TRIPOD for Abstracts. *Ann Intern Med*. 2020;173(1):43. doi:10.7326/M20-0193
 37. Collins GS, De Groot JA, Dutton S, et al. External validation of multivariable prediction models: A systematic review of methodological conduct and reporting. *BMC Med Res Methodol*. 2014;14(1):40. doi:10.1186/1471-2288-14-40
 38. Wynants L, Van Calster B, Collins GS, et al. Prediction models for diagnosis and prognosis of covid-19:

- Systematic review and critical appraisal. *BMJ*. 2020;369. doi:10.1136/bmj.m1328
39. Heus P, Damen JAAG, Pajouheshnia R, et al. Poor reporting of multivariable prediction model studies: Towards a targeted implementation strategy of the TRIPOD statement. *BMC Med*. 2018;16(1):1-12. doi:10.1186/s12916-018-1099-2
 40. Bouwmeester W, Zuithoff NPA, Mallett S, et al. Reporting and methods in clinical prediction research: A systematic review. *PLoS Med*. 2012;9(5). doi:10.1371/journal.pmed.1001221
 41. Yang C, Kors JA, Ioannou S, et al. Trends in the conduct and reporting of clinical prediction model development and validation: a systematic review. *J Am Med Informatics Assoc*. 2022;00(0):1-7. doi:10.1093/jamia/ocac002
 42. Zamanipoor Najafabadi AH, Ramspek CL, Dekker FW, et al. TRIPOD statement: a preliminary pre-post analysis of reporting and methods of prediction models. *BMJ Open*. 2020;10(9):e041537. doi:10.1136/bmjopen-2020-041537
 43. Kim DW, Jang HY, Ko Y, et al. Inconsistency in the use of the term “validation” in studies reporting the performance of deep learning algorithms in providing diagnosis from medical imaging. *PLoS One*. 2020;15(9 September):1-10. doi:10.1371/journal.pone.0238908
 44. Collins GS, Dhiman P, Andaur Navarro CL, et al. Protocol for development of a reporting guideline (TRIPOD-AI) and risk of bias tool (PROBAST-AI) for diagnostic and prognostic prediction model studies based on artificial intelligence. *BMJ Open*. 2021;11(e048008):1-7. doi:10.1136/bmjopen-2020-048008
 45. Vollmer S, Mateen BA, Bohner G, et al. Machine learning and artificial intelligence research for patient benefit: 20 critical questions on transparency, replicability, ethics, and effectiveness. *BMJ*. 2020;368:1-12. doi:10.1136/bmj.l6927
 46. Wolff RF, Moons KGM, Riley RD, et al. PROBAST: A tool to assess the risk of bias and applicability of prediction model studies. *Ann Intern Med*. 2019;170(1):51-58. doi:10.7326/M18-1376
 47. Moons KGM, Wolff RF, Riley RD, et al. PROBAST: A tool to assess risk of bias and applicability of prediction model studies: Explanation and elaboration. *Ann Intern Med*. 2019;170(1):W1-W33. doi:10.7326/M18-1377
 48. Damen JAAG, Debray TPA, Pajouheshnia R, et al. Empirical evidence of the impact of study characteristics on the performance of prediction models: A meta-epidemiological study. *BMJ Open*. 2019;9(4):1-12. doi:10.1136/bmjopen-2018-026160
 49. Collins GS, Moons KG. Reporting of artificial intelligence prediction models. *Lancet*. 2019;393. doi:10.1016/S0140-6736(19)30235-1
 50. Altman DG, Royston P. What do we mean by validating a prognostic model? *Stat Med*. 2000;19(4):453-473. doi:10.1002/(SICI)1097-0258(20000229)19:4<453::AID-SIM350>3.0.CO;2-5
 51. Van Calster B, Wynants L, Timmerman D, Steyerberg EW, Collins GS. Predictive analytics in health care: how can we know it works? *J Am Med Informatics Assoc*. 2019;26(12):1651-1654. doi:10.1093/JAMIA/OCZ130

Figure 1. Flowchart of included studies

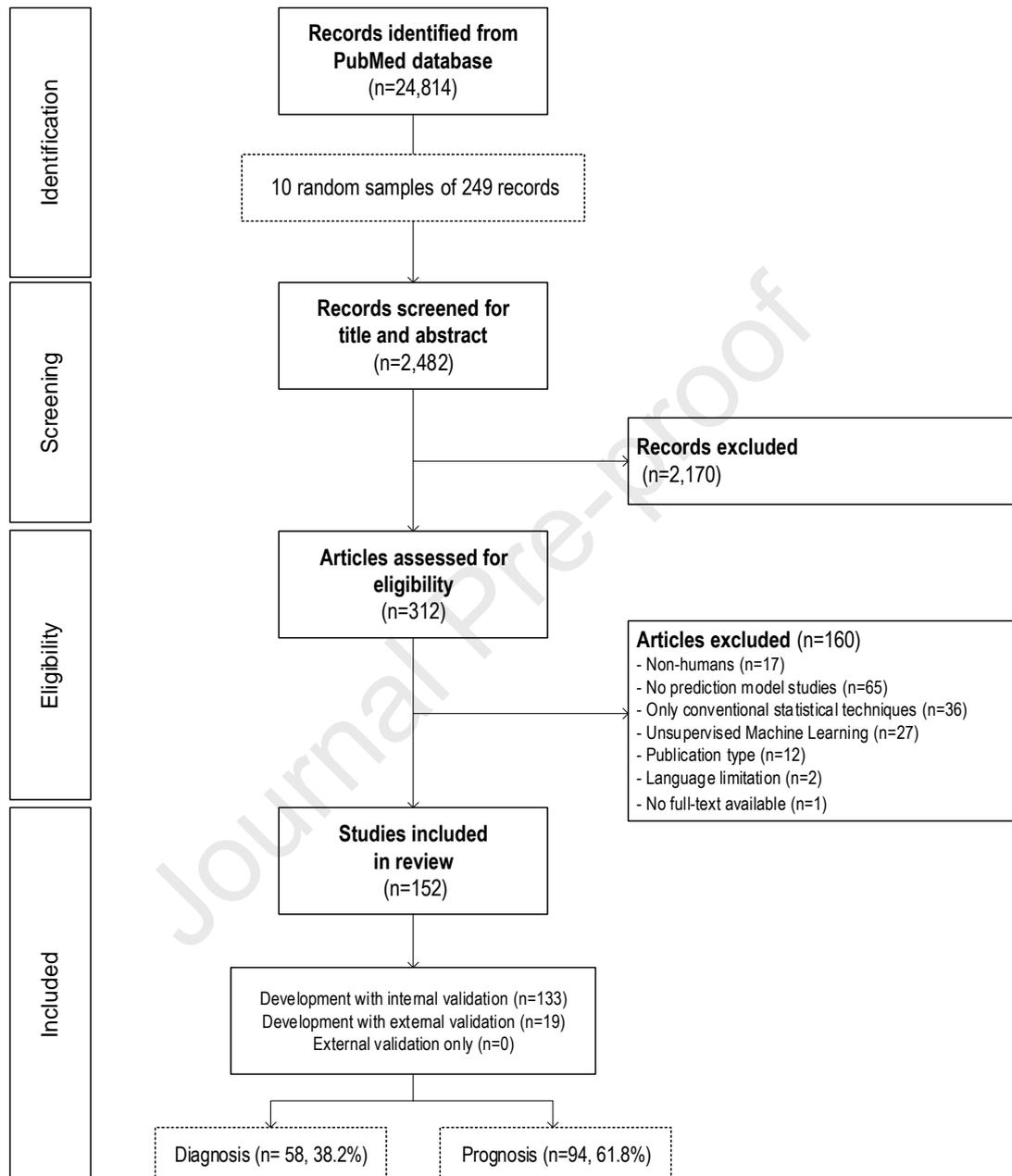


Table 1. General characteristics of included studies

		Total (n=152)
		n (%) [95% CI]
Study aim		
	Diagnosis	58 (38.2) [30.8-46.1]
	Prognosis	94 (61.8) [53.9-69.2]
Study type		
	Model development only	133 (87.5) [81.3-91.8]
	Model development with external validation	19 (12.5) [8.2-18.7]
Outcome aim		
	Classification	120 (78.9) [71.8-84.7]
	Risk probabilities	32 (21.0) [80.5-91.3]
Setting[†]		
	General population	17 (11.2) [7.1-17.2]
	Primary care	15 (9.9) [6.1-15.6]
	Secondary care	32 (21.1) [15.3-28.2]
	Tertiary care	78 (51.3) [43.4-59.1]
	Unclear	13 (8.6) [5.1-14.1]
Outcome format		
	Continuous	7 (4.6) [2.2-9.2]
	Binary	131 (86.2) [79.8-90.8]
	Multinomial	7 (4.6) [2.2-9.2]
	Ordinal	2 (1.3) [0.4-4.7]
	Time-to-event	3 (2.0) [0.7-5.6]
	Count	2 (1.3) [0.4-4.7]
Type of outcome		
	Death	21 (13.8) [9.2-20.2]
	Complications	65 (42.8) [35.2-50.7]
	Disease detection	30 (19.7) [14.2-26.8]
	Disease recurrence	9 (5.9) [3.1-10.9]
	Survival	3 (2.0) [0.7-5.6]
	Readmission	4 (2.6) [1-6.6]
	Other ^a	20 (13.2) [8.7-19.5]
Mentioning of reporting guidelines[†]		
	TRIPOD	8 (5.3) [2.7-10]
	STROBE	3 (2.0) [0.7-5.6]
	Other ^b	5 (3.3) [1.4-7.5]
	None	139 (91.4) [85.9-94.9]
Model availability[†]		
	Repository for data	18 (11.8) [7.6-17.9]
	Repository for code	13 (8.6) [5.1-14.1]
	Model presentation ^c	31 (20.4) [14.8-27.5]
	None	121 (79.6) [72.5-85.2]

[†] Counts are absolute numbers with column percentages in parentheses. The percentages sometimes do not add up to 100% because studies reported more than one option.

^a This includes length of stay, medication dose, patient's disposition, order type, lesion extension, laboratory results, cancer stage, treatment option, attendance, equipment usage, operative time.

^b Guidelines for developing and reporting machine learning models in biomedical research (n=2), STARD (n=2), BRISQ (n=1).

^oThis includes simplified scoring rule, chart, nomogram, online calculator, or worked examples.

Journal Pre-proof

Table 2. Modelling algorithms for all extracted models

Modelling algorithm	All extracted models (n=522)	
	n (%) [95%CI]	
Unpenalized regression models	101 (19.3) [16.1-23.1]	
Ordinary least squares regression ^a	27 (5.2)	[3.5-7.5]
Maximum likelihood logistic regression	74 (14.2)	[11.4-17.5]
Penalized regression models	29 (5.6) [3.8-8]	
Elastic Net	9 (1.7)	[0.8-3.4]
LASSO	13 (2.5)	[1.4-4.3]
Ridge	7 (1.3)	[0.6-2.9]
Tree-based models	166 (31.8) [28-36]	
Decision trees (e.g., CART) ^c	46 (8.8)	[6.6-11.7]
Random forest ^d	73 (14)	[11.2-17.3]
Extremely randomized trees	1 (0.2)	[0.01-1.2]
Regularized Greedy Forest	1 (0.2)	[0.01-1.2]
Gradient boosting machine ^e	34 (6.5)	[4.6-9.1]
XGBoost	11 (2.1)	[1.1-3.9]
Neural Network (incl. deep learning) ^b	75 (14.4) [11.5-17.7]	
Support Vector Machine	86 (16.5) [13.5-20]	
Naïve Bayes	22 (4.2) [2.7-6.4]	
K-nearest neighbor	15 (2.9) [1.7-4.8]	
Superlearner ensembles	14 (2.7) [1.5-4.6]	
Other ^f	10 (1.9) [1-3.6]	
Unclear	4 (0.8) [0.2-2.1]	

CART, Classification And Regression Tree; LASSO, Least Absolute Shrinkage and Selection Operator; XGBoost, extreme gradient boosting; CI, confidence interval.

^a Discriminant analysis, generalized additive models (GAM), partial least squares were extracted as OLS regression.

^b Multilayer perceptron, denseNet, convolutional, recurrent, and Bayesian neural networks were extracted as neural networks.

^c This includes conditional inference tree (n=3), optimal tree (n=1).

^d This includes Random Survival Forest (n=2).

^e This includes lightGBM (n=1), adaBoost (n=8), catBoost (n=1), logitboost (n=1), RUSBoost (n=1), and stochastic (n=1).

^f This includes bayesian network (n=3), rule-based classifier (n=1), highly predictive signatures (n=1), Kalman filtering (n=1), fuzzy soft set (n=1), adaptive neuro-fuzzy inference system (n=1), stochastic gradient descent (n=1), fully corrective binning (n=1).

Table 3. Study design of included studies, stratified by type of prediction model study

	Total (n=152)	Development only (n=133)	Development with external validation (n=19)
	n (%) [95%CI]	n (%) [95%CI]	n (%) [95%CI]
Data sources ^{†a}			
Prospective cohort	50 (32.9) [25.9-40.7]	43 (32.3) [25-40.7]	7 (36.8) [19.1-59]
Retrospective cohort	48 (31.6) [24.7-39.3]	45 (33.8) [26.3-42.2]	4 (21.1) [8.5-43.3]
Randomized Controlled Trial	3 (2.0) [0.7-5.6]	2 (1.5) [0.4-5.3]	1 (5.3) [0.3-24.6]
EMR	30 (19.7) [14.2-26.8]	28 (21.1) [15-28.7]	0
Registry	18 (11.8) [7.6-17.9]	15 (11.3) [7-17.8]	4 (21.1) [8.5-43.3]
Administrative claims	4 (2.6) [1-6.6]	4 (3.0) [1.2-7.5]	0
Case-control	18 (11.8) [7.6-17.9]	15 (11.3) [7-17.8]	3 (15.8) [5.5-37.6]
Number of centers			
	110 (72.4)	98 (73.7)	12 (63.2)
Median [IQR] (range)	1 [1-3], 1 to 51920	1 [1-3], 1 to 712	1 [1-10], 1 to 51920
Follow-up (months) [‡]			
	47 (30.9)	39 (29.2)	8 (42.1)
Median [IQR] (range)	41.9 [3-60], 0.3 to 307	43.6 [4.5-60], 0.3 to 307	33.5 [1.75-42], 1 to 144
Predictor horizon (months) [‡]			
	49 (32.2) [25.3-40]	61 (45.9) [37.6-54.3]	7 (36.8)
Median [IQR] (range)	8.5[1-36], 0.03 to 120	6 [1-33.5], 0.03 to 120	36 [6.5-60], 1 to 60
Sample size justification			
	27 (17.8) [12.5-24.6]	24 (18.0) [12.4-25.4]	3 (15.8)
Power	5 (18.5) [8.2-36.7]	5 (20.8) [9.2-40.5]	0
Justified time interval	5 (18.5) [8.2-36.7]	3 (12.5) [4.3-31]	2 (66.7)
Size of existing/available data	16 (59.3) [40.7-75.5]	15 (62.5) [42.7-78.8]	1 (33.3)
Events per variable	1 (3.7) [0.2-18.3]	1 (4.2) [0.2-20.2]	0
Internal validation [†]			
Split sample with test set	86 (56.6) [48.6-64.2]	NA	NA
(Random) split	49 (57) [46.4-66.9]		
(Non-random) split	9 (10.5) [5.6-18.7]		
Split ^b	28 (32.6) [23.6-43]		
Bootstrapping	5 (3.3) [1.4-7.5]	NA	NA
With test set	3 (60.0) [23.1-88.2]		
With cross-validation	1 (20) [1-62.4]		
Cross-validation	70 (46.1) [38.3-54]	NA	NA
Non-nested (single)	32 (45.7) [34.6]		
Nested	10 (14.3) [7.9-24.3]		
With test set	24 (34.3) [24.2-46]		
External validation [†]			
Chronological	NA	NA	5 (26.3) [11.8-48.8]
Geographical	NA	NA	3 (15.8) [5.5-37.6]
Independent dataset	NA	NA	11 (57.9) [36.3-76.9]
Fully independent dataset	NA	NA	8 (42.1) [23.1-63.7]

[†]Counts are absolute numbers, with column percentages in parentheses. The percentages sometimes do not add up to 100% because studies reported more than one measure. We report then the raw percentages. NA, not applicable.

[‡]We collected the longest follow-up and longest prediction horizon.

^aData sources also included surveys (n=2), cross-sectional studies (n=2).

^bUnclear whether split sample was performed random or non-random.

Table 4. Predictors in included studies

		Total (n=152)
		n (%) [95% CI]
Type of candidate predictors[†]		
	Demography	120 (78.9) [71.8-84.7]
	Clinical history	111 (73.0) [65.5-79.4]
	Physical examination	0
	Blood or Urine parameters	63 (41.4) [33.9-49.4]
	Imaging	49 (32.2) [23.3-40]
	Genetic risk score	7 (4.6) [2.2-9.2]
	Pathology	16 (10.5) [6.6-16.4]
	Scale score	31 (20.4) [14.8-27.5]
	Questionnaires	0
Treatment as candidate predictor		
	Yes	36 (23.7) [17.6-31]
	No	80 (52.6) [44.7-60.4]
	Not applicable	36 (23.7) [17.6-31]
Continuous variables as candidate predictors		
	Yes	131 (86.2) [79.8-90.8]
	Unclear	17 (11.2) [7.1-17.2]
A-priori selection of candidate predictors[‡]		
	Yes	63 (41.4) [33.9-49.4]
	No	47 (30.9) [24.1-38.7]
	Unclear	42 (27.6) [21.1-35.2]
Methods to handle continuous predictors^{†‡}		
	Linear (no change)	13 (8.6) [5.1-14.1]
	Non-linear (planned)	2 (1.3) [0.4-4.7]
	Non-linear (unplanned)	4 (2.6) [1-6.6]
	Categorised (some)	16 (10.5) [6.6-16.4]
	Categorised (all)	18 (11.8) [7.6-17.9]
	Unclear	104 (68.4) [60.7-75.3]
Categorization of continuous predictors[‡]		
	Data dependent	4 (2.6) [1-6.6]
	No rationale	17 (11.2) [7.1-17.2]
	Based on previous literature or standardization	13 (8.6) [5.1-14.1]
	Not reported	118 (77.6) [70.4-83.5]

[†]Counts are absolute numbers, with column percentages in parentheses. The percentages sometimes do not add up to 100% because studies can report more than one measure.

[‡]As data preparation

Table 5. Sample size of included studies (n=152)

	Total (n=152)	
	n (%)	Median [IQR], range
Initial sample size	93 (61.2)	999 [272-24522], 8 to 1093177
External validation ^a	13 (68.4)	318 [90-682], 19 to 1113656
Final sample size	151 (99.3)	587 [172-6328], 8 to 594751
Model development	83 (54.6)	641 [226-10512], 5 to 392536
Internal validation ^b	83 (54.6)	230 [75-2892], 2 to 202215
External validation ^a	18 (94.7)	293 [71-1688], 19 to 59738
Initial number of events	10 (6.6)	66 [15-207], 15 to 4370
External validation ^a	1 (5.3)	107
Final number of events	37 (24.3)	106 [50-364], 15 to 7543
Model development	19 (13.2)	156 [47-353], 10 to 5054
Internal validation ^b	19 (13.2)	35 [26-109], 4 to 2489
External validation ^a	4 (21.1)	250 [121-990], 107 to 2834
Number of candidate predictors	119 (78.3)	24 [13-112], 2 to 39212
Number of included predictors	90 (59.2)	12 [7-23], 2 to 570
Events per candidate predictor ^c	28 (18.4)	12.5 [5.7-27.7], 1.2 to 754.3

^a External validation was performed in 19 studies.

^b Combines all internal validation methods, e.g., split sample, cross validation, bootstrapping.

^c For model development.

Table 6. Handling of missing values, stratified by study type

	Total (n=152)	Development only (n=133)	Development with external validation (n=19)
	n (%) [95%CI]	n (%) [95%CI]	n (%) [95%CI]
Missingness as exclusion criteria for participants			
Yes	56 (36.8) [29.6-44.7]	51 (38.3) [30.5-46.8]	2 (10.5) [2.9-31.4]
Unclear	36 (23.7) [17.6-31]	33 (24.8) [18.2-32.8]	6 (31.6) [15.4-54]
Number of patients excluded	36 (23.7) [17.6-31]	34 (25.6) [18.9-33.6]	0
Median [IQR] (range)	191 [19-4209], 1 to 627180	224 [16-4699], 1 to 627180	0
Methods of handling missing data[†]			
No missing data	4 (2.6) [1-6.6]	3 (2.3) [0.8-6.4]	1 (5.3) [0.3-24.6]
No imputation	4 (2.6) [1-6.6]	4 (3) [1.2-7.5]	0
Complete case-analysis	30 (19.7) [14.2-26.8]	28 (21.1) [15-28.7]	2 (10.5) [2.9-31.4]
Mean imputation	4 (2.6) [1-6.6]	3 (2.3) [0.8-6.4]	1 (5.3) [0.3-24.6]
Median imputation	10 (6.6) [3.6-11.7]	10 (7.5) [4.1-13.3]	0
Multiple imputation	6 (3.9) [1.8-8.3]	6 (4.5) [2.1-9.5]	0
K-nearest neighbor imputation	5 (3.3) [1.4-7.5]	5 (3.8) [1.6-8.5]	0
Replacement with null value	3 (2.0) [0.7-5.6]	1 (0.8) [0-4.1]	2 (10.5) [2.9-31.4]
Last value carried forward	4 (2.6) [1-6.6]	4 (3) [1.2-7.5]	0
Surrogate variable	1 (0.7) [0-3.6]	1 (0.8) [0-4.1]	0
Random forest imputation	4 (2.6) [1-6.6]	3 (2.3) [0.8-6.4]	1 (5.3) [0.3-24.6]
Categorization	3 (2) [0.7-5.6]	2 (1.5) [0.4-5.3]	1 (5.3) [0.3-24.6]
Unclear	6 (3.9) [1.8-8.3]	5 (3.8) [1.6-8.5]	1 (5.3) [0.3-24.6]
Presentation of missing data			
Not summarized	129 (84.9) [78.3-89.7]	114 (85.7) [78.8-90.7]	16 (84.2) [62.4-94.5]
Overall	6 (3.9) [1.8-8.3]	4 (3) [1.2-7.5]	2 (10.5) [2.9-31.4]
By all final model variables	3 (2) [0.7-5.6]	3 (2.3) [0.8-6.4]	0
By all candidate predictors	13 (8.6) [5.1-14.1]	11 (8.3) [4.7-14.2]	1 (5.3) [0.3-24.6]
By number of variables	1 (0.7) [0-3.6]	1 (0.8) [0-4.1]	0

[†]Counts are absolute numbers, with column percentages in parentheses. The percentages sometimes do not add up to 100% because studies can report more than one technique.

Table 7. Machine learning aspects in the included studies

		Total (n=152)
		n (%) [95% CI]
Data preparation[†]		58 (38.2) [30.8-46.1]
	Cleaning	21 (36.2) [25.1-49.1]
	Aggregation	6 (10.3) [4.8-20.8]
	Transformation	6 (10.3) [4.8-20.8]
	Sampling	2 (3.4) [1-11.7]
	Standardization/Scaling	11 (19) [10.9-30.9]
	Normalization	22 (37.9) [26.6-50.8]
	Integration	0
	Reduction	12 (20.7) [12.3-32.8]
	Other ^a	9 (15.5) [8.4-26.9]
Data splitting		86 (56.6) [48.6-64.2]
	Train-test set	77 (50.7) [42.8-58.5]
	Train-validation-test set	9 (5.9) [3.1-10.9]
Dimensionality reduction techniques		9 (5.9) [3.1-10.9]
	CART	1 (11.1) [0.6-43.5]
	Principal component analysis	3 (33.3) [12.1-64.6]
	Factor analysis	1 (11.1) [0.6-43.5]
	Image decomposition	1 (11.1) [0.6-43.5]
Class imbalance[†]		27 (17.8) [12.5-24.6]
	Random undersampling	4 (14.8) [5.9-32.5]
	Random oversampling	5 (18.5) [8.2-36.7]
	SMOTE	11 (40.7) [24.5-59.3]
	RUSBoost	1 (3.7) [0.2-18.3]
	Other ^b	7 (25.9) [13.2-44.7]
Strategy for hyperparameter optimization[†]		44 (28.9) [22.3-36.6]
	Grid search (no further details)	5 (3.3) [1.4-7.5]
	Cross-validated grid search	14 (9.2) [5.6-14.9]
	Randomized grid search	1 (0.7) [0-3.6]
	Cross-validation	15 (9.9) [6.1-15.6]
	Manual search	1(0.7) [0-3.6]
	Pre-defined values/default	3 (2) [0.7-5.6]
	Bayesian optimization	2 (1.3) [0.4-4.7]
	Tree-structured parzen estimator method	1(0.7) [0-3.6]
	Unclear	4 (2.6) [1-6.6]

[†]Counts are absolute numbers, with column percentages in parentheses. The percentages sometimes do not add up to 100% because studies can report more than one measure. CART - Classification And Regression Tree.

^aThis includes matching, augmentation, noise filtering, merging, splitting, binning.

^bThis includes matching, resampling, class weighting, inverse class probability.

Table 8. Model building of all included studies

	Total (n=522)	
	n (%) [95% CI]	
Selection of predictors		
Stepwise	8 (1.5)	[0.7-3.1]
Forward selection	31 (5.9)	[4.1-8.4]
Backward selection	5 (1)	[0.4-2.4]
All predictors	72 (13.8)	[11.1-17.1]
All significant in univariable analysis	27 (5.2)	[3.5-7.5]
Embedded in learning process	192 (36.8)	[32.7-41.1]
Other	19 (3.6)	[2.3-5.7]
Unclear	168 (32.2)	[28.2-36.4]
Hyperparameter tuning reported		
Yes	160 (30.7)	[26.7-34.8]
No	283 (54.2)	[49.8-58.5]
Not applicable/Unclear	79 (15.1)	[12.2-18.6]
Variable importance reported		
Mean decrease in accuracy	26 (5)	[3.3-7.3]
Mean decrease in node impurity	31 (5.9)	[4.1-8.4]
Weights/correlation	10 (1.9)	[1-3.6]
Gain information	24 (4.6)	[3-6.9]
Unclear method	115 (22)	[18.6-25.9]
None	316 (60.5)	[56.2-64.7]
Penalization methods used		
None	481 (92.1)	[89.4-94.2]
Uniform shrinkage	3 (0.6)	[0.1-1.8]
Penalised estimation	27 (5.2)	[3.5-7.5]
Other	11 (2.1)	[1.1-3.9]

Counts are absolute numbers, with column percentages in parentheses. The percentages sometimes do not add up to 100% because studies can report more than one measure.

^a Both unpenalized and penalized regression models.

^b Only for penalized regression techniques.

Table 9. Performance measures reported, stratified by model development and validation

		All extracted models (n=522)	
		n (%) [95% CI]	
		DEV	VAL
Calibration†			
	Calibration plot	23 (4.4) [2.9-6.6]	1 (0.2) [0.01-1.2]
	Calibration slope	17 (3.3) [2.5-3]	1 (0.2) [0.01-1.2]
	Calibration intercept	16 (3.1) [1.8-5]	1 (0.2) [0.01-1.2]
	Calibration in the large	1 (0.2) [0.01-1.2]	0
	Calibration table	1 (0.2) [0.01-1.2]	0
	Kappa	10 (1.9) [1-3.6]	0
	Observed/expected ratio	1 (0.2) [0.01-1.2]	0
	Homer-Lemeshow statistic	4 (0.8) [0.3-2.1]	0
	None	494 (94.6) [92.3-96.3]	
Discrimination			
	AUC/ AUC-ROC	349 (66.9) [62.6-70.9]	46 (8.8) [6.6-11.7]
	C-statistic	9 (1.7) [0.8-3.4]	0
	None	164 (31.4) [27.5-35.6]	
Classification †			
	NRI	9 (1.7) [0.8-3.4]	0
	Sensitivity/Recall	239 (45.8) [41.5-50.2]	30 (5.7) [4-8.2]
	Specificity	193 (37) [32.8-41.3]	22 (4.2) [2.7-6.4]
Decision-analytic †			
	Decision Curve Analysis	2 (0.4) [0.01-1.5]	0
	IDI	1 (0.2) [0.01-1.2]	0
Overall †			
	R2	14 (2.7) [1.5-4.6]	0
	Brier score	19 (3.6) [2.3-5.7]	6 (1.1) [0.5-2.6]
	Predictive values*	160 (30.7) [26.8-34.8]	10 (1.9) [1-3.6]
	AUC difference	2 (0.4) [0.01-1.5]	0
	Accuracy**	234 (44.8) [40.5-49.2]	26 (5) [3.4-7.3]
	F1-score	79 (15.1) [12.2-18.6]	0
	Mean square error	21 (4) [2.6-6.2]	0
	Misclassification rate	9 (1.7) [0.8-3.4]	0
	Mathew's correlation coefficient	5 (1) [0.4-2.4]	0
	AUPR	21 (4) [2.6-6.2]	0

† Counts are absolute numbers, with column percentages in parentheses. The percentages sometimes do not add up to 100% because studies can report more than one performance measure.

*This includes models reporting positive predictive value as precision.

**This includes models reporting balance accuracy.

DEV, developed model; VAL, validation; AUC-ROC, area under the receiver operation characteristic curve; NRI, net reclassification index; IDI, integrated discrimination improvement; AUPR, area under the precision-recall curve; CI, confidence interval.

Table 10. Predictive performance of all extracted models*

		All extracted models (n=522)					
		Apparent performance		Corrected performance**		Externally validated performance	
		Reported, n (%)	Median [IQR], range	Reported, n (%)	Median [IQR], range	Reported, n (%)	Median [IQR], range
Calibration							
	Slope	11 (1.9)	1.05 [1.02 - 1.07], 0.53 to 1.46	15 (2.9)	1.3 [1-4], 0.52 to 17.6	4 (0.8)	9.9 [7.87-12.8], 5.7 to 17.6
	Intercept	10 (1.9)	0.07 [0.05 - 0.12], -0.08 to 2.32	15 (2.9)	-0.01 [-1.85 - 0.15], -8.3 to 2.74	4 (0.8)	-4.5 [-5.7 - -3.8], -8.3 to -3
	Calibration-in-the-large	1 (0.2)	-0.008	0		0	
	Observed:expected ratio	1 (0.2)	0.993	4 (0.8)	0.99 [0.98 - 1.01], 0.98 to 1.04	0	
	Homer-Lemeshow	2 (0.2)	Not significant	0		0	
	Pearson chi-square	1 (0.2)	Not significant	0		0	
	Mean Calibration Error	4 (0.8)	0.81 [0.7 - 0.88], 0.51 to 0.99	0		0	
Discrimination							
	AUC	249 (47.7)	0.82 [0.74-0.90], 0.45 to 1.00	154 (29.5)	0.82 [0.74-0.90], 0.46 to 0.99	46 (8.8)	0.82 [0.73-0.98], 0.52 to 0.97
Accuracy							
		128 (24.5)	79.8 [72.6-89.8], 44.2 to 100	117 (22.4)	81.4 [76-89.9], 17.8 to 97.5	9 (1.7)	70 [64-87], 55 to 90
Sensitivity							
		156 (29.9)	74 [58.6-87.8], 0 to 100	103 (19.7)	80 [66.3-89.7], 14.8 to 100	12 (2.3)	77.5 [63.9-83.5], 0.7 to 91
Specificity							
		122 (23.4)	82.2 [73.3-89.7], 17 to 100	80 (15.3)	83.2 [73.6-90.8], 46.6 to 100	10 (1.9)	74.4 [64.8-86.7], 42 to 90.5

* Counts are absolute numbers with column percentages in parentheses. The percentages sometimes do not add up to 100% because some studies did not report performance measure for all models pre-specified.

**We considered corrected performance only when authors stated results as such. Otherwise, performance measures were considered apparent performance by default.

What is new?

Key Findings

- Design and methodological conduct of studies on clinical prediction models based on machine learning vary substantially.

What this adds to what was known?

- Studies on clinical prediction models based on machine learning suffered from poor methodology and reporting similar to studies using regression approaches.

What is the implication and what should change now?

- Methodologies for model development and validation should be more carefully designed and reported to avoid research waste.
- More attention is needed to missing data, internal validation procedures, and calibration.
- Methodological guidance for studies on prediction models based on machine learning techniques is urgently needed.

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre-proof

Authors' contributions

Constanza L. Andaur Navarro: Conceptualization, Methodology, Investigation, Data Curation, Formal analysis, Writing - original draft, Writing - review & editing; **Johanna A.A. Damen:** Conceptualization, Methodology, Investigation, Writing - review & editing, Supervision; **Maarten van Smeden:** Conceptualization, Writing - review & editing; **Toshihiko Takada:** Investigation, Writing - review & editing. **Steven WJ Nijman:** Investigation, Writing - review & editing; **Paula Dhiman:** Conceptualization, Methodology, Investigation, Writing - review & editing; **Jie Ma:** Investigation, Writing - review & editing; **Gary S Collins:** Conceptualization, Methodology, Writing - review & editing; **Ram Bajpai:** Investigation, Writing - review & editing; **Richard D Riley:** Conceptualization, Methodology, Writing - review & editing; **Karel GM Moons:** Conceptualization, Methodology, Writing - review & editing, Supervision; **Lotty Hoof:** Conceptualization, Methodology, Writing - review & editing, Supervision

Journal Pre-proof