

1 Can machine learning improve risk prediction of incident hypertension?

2 An internal method comparison and external validation of the

3 Framingham risk model using HUNT Study data

4

5 Filip Emil Schjerven<sup>1\*</sup>, Emma Ingeström<sup>2</sup>, Frank Lindseth<sup>1</sup>, Ingelin Steinsland<sup>3</sup>

6

7 <sup>1</sup> Department of Computer Science, Norwegian University of Science and Technology, Trondheim,

8 Norway.

9 <sup>2</sup> Department of Circulation and Medical Imaging, Norwegian University of Science and Technology,

10 Trondheim, Norway.

11 <sup>3</sup>Department of Mathematical Sciences, Norwegian University of Science and Technology, Trondheim,

12 Norway.

13

14 \* Corresponding author

15 Email: [filip.e.schjerven@ntnu.no](mailto:filip.e.schjerven@ntnu.no)

16

17

18

## 19 Abstract

20 A recent meta-review on hypertension risk models detailed that the differences in data and study-setup  
21 have a large influence on performance, meaning model comparisons should be performed using the same  
22 study data. We compared five different machine learning algorithms and the externally developed  
23 Framingham risk model in predicting risk of incident hypertension using data from the Trøndelag Health  
24 Study. The dataset yielded  $n = 23722$  individuals with  $p = 17$  features recorded at baseline before follow-  
25 up 11 years later. Individuals were without hypertension, diabetes, or history of CVD at baseline. Features  
26 included clinical measurements, serum markers, and questionnaire-based information on health and  
27 lifestyle. The included modelling algorithms varied in complexity from simpler linear predictors like  
28 logistic regression to the eXtreme Gradient Boosting algorithm. The other algorithms were Random  
29 Forest, Support Vector Machines, K-Nearest Neighbor. After selecting hyperparameters using cross-  
30 validation on a training set, we evaluated the models' performance on discrimination, calibration, and  
31 clinical usefulness on a separate testing set using bootstrapping. Although the machine learning models  
32 displayed the best performance measures on average, the improvement from a logistic regression model  
33 fitted with elastic regularization was small. The externally developed Framingham risk model performed  
34 well on discrimination, but severely overestimated risk of incident hypertension on our data. After a  
35 simple recalibration, the Framingham risk model performed as well or even better than some of the newly  
36 developed models on all measures. Using the available data, this indicates that low-complexity models  
37 may suffice for long-term risk modelling. However, more studies are needed to assess potential benefits of  
38 a more diverse feature-set. This study marks the first attempt at applying machine learning methods and  
39 evaluating their performance on discrimination, calibration, and clinical usefulness within the same study  
40 on hypertension risk modelling.

## 41 Author summary

42 Hypertension, the state of persistent high blood pressure, is a largely symptom-free medical condition  
43 affecting millions of individuals worldwide, a number that is expected to rise in the coming years. While  
44 consequences of unchecked hypertension are severe, life-style modifications have been proven to be  
45 effective in prevention and treatment of hypertension. A possible tool for identifying individuals at risk of  
46 developing hypertension has been the creation of hypertension risk scores, which calculate a probability of  
47 incident hypertension sometime in the future. We compared applying machine learning as opposed to  
48 more traditional tools for constructing risk models on a large Norwegian cohort, measuring performance  
49 by model validity and clinical usefulness. Using easily obtainable clinical information and blood  
50 biomarkers as inputs, we found no clear advantage in performance using the machine learning models.  
51 Only a few of our included inputs, namely systolic and diastolic blood pressure, age, and BMI were found  
52 to be important for accurate prediction. This suggest more diverse information on individuals, like genetic,  
53 socio-economic, or dietary information, may be necessary for machine learning to excel over more  
54 established methods. A risk model developed using an American cohort, the Framingham risk model,  
55 performed well on our data after recalibration. Our study provides new insights into machine learning may  
56 be used to enhance hypertension risk prediction.

## 57 Introduction

58 Individuals with persistently high levels of blood pressure are said to have hypertension. It is a mostly  
59 symptom-free medical condition, but it increases the risk of more severe diseases and premature death if  
60 left untreated (1). The number of hypertensive individuals is estimated to be over 1 billion worldwide, and  
61 suboptimal blood pressure accounts for around 10 % of the world's overall health expenditures (2,3). It is  
62 well-established that the risks related to hypertension can be effectively reduced through lifestyle  
63 modifications and medications (1). Hence, early identification of otherwise healthy individuals at risk of

64 hypertension has become a multidisciplinary research effort including the analysis and subsequent  
65 development of hypertension risk models.

66 Since the publication of the Framingham risk model for incident hypertension in 2009, the number of risk  
67 models has increased substantially and the topic has been reviewed multiple times (4–6). In a recent  
68 review, 52 studies and 117 risk models were identified, of which most models were developed using  
69 established statistical methods (7). Simultaneously, machine learning has emerged as a common  
70 alternative for constructing risk models. By leveraging the ability to learn more complex patterns from the  
71 data with less human intervention, machine learning has been emphasized as having the potential to  
72 construct models excelling those using less-complex methods. Machine learning has been shown to  
73 improve risk prediction for cardiovascular disease (8). Many speculate whether machine learning and  
74 other artificial intelligence (AI) methods can contribute to transforming the health-sciences (9) (10).  
75 However, this potential has been challenged in studies reviewing risk prediction models for other diseases.  
76 In these cases, applying machine learning to a problem was not synonymous with improved performance  
77 compared to more established statistical methods (11–13). When not restricting to a specific disease, a  
78 systematic review by Christodoulou et al. found that logistic regression performed just as well as machine  
79 learning for clinical prediction models when limiting themselves to studies with low risk of bias (14).

80 Considering the topic of hypertension risk models specifically, a recent systematic review identified and  
81 summarized the previous work (7). In that review, it was found that a large variation reported in  
82 discrimination performance was unrelated to the type of modelling algorithms being used. In other words,  
83 using different datasets and study setups had considerable influence on the reported performance. This  
84 implies that to confidently determine a performance advantage of one method over the other, it would be  
85 necessary to compare methods within the same study and data, keeping everything fixed except the  
86 methods under investigation. Thus, it is difficult to determine the cause of any apparent performance  
87 benefit that has been reported for one method over others in the existing literature.

88 For a prediction model to be useful to health practitioners in the real-life clinical setting, it must  
89 demonstrate acceptable performance during evaluation. A model's performance is often captured by its  
90 discriminative performance, as well as the calibration of its predictions. Furthermore, clinical usefulness  
91 has been emphasized to demonstrate that the risk models provide a clear benefit over alternatives in  
92 decision-making (15–17). In the literature on hypertension risk models, discrimination is the most  
93 frequently reported performance indicator. Although calibration is often reported for risk models  
94 developed using more established methods, few studies using machine learning have reported model  
95 calibration performance. Lastly, few studies have explored clinical usefulness for hypertension risk  
96 models using decision curve analysis like Net Benefit (7,15,17).

97 Machine learning have been used to develop hypertension risk models before, but few studies have  
98 compared their performance with low-complexity models like logistic or Cox regression (18–20). Among  
99 the studies who have, machine learning has been found to have better discrimination in some studies and  
100 worse in others (11,21–24). Calibration was assessed in only two studies applying machine learning  
101 models, and then solely by reporting a summary measure (22,25). Niu et al. evaluated machine learning  
102 models by their net benefit as well as discrimination, and displayed a large advantage of using machine  
103 learning models on a rural Chinese cohort (24).

104 Considering the many existing risk models, a new model should be a valuable contribution to the  
105 literature. External validation, as opposed to creating new risk models, has been emphasized as an equally  
106 valuable contribution to the literature and a necessity for transitioning a risk model to clinical practice.  
107 (26–28). Not only may it provide information on how well the existing model generalizes to new, unseen  
108 data, but it also serves as a benchmark for the development of any new proposed model.

109 In this work, we aim to investigate the potential performance benefit of using machine learning for  
110 predicting hypertension by conducting a thorough method comparison. Our primary aim was to assess  
111 whether more complex modelling methods provide a clear performance benefit, in terms of  
112 discrimination, calibration, and clinical usefulness, over more established statistical methods when

113 developed under equal settings using the same data. To assess the need for a new risk model and provide  
114 benchmarks for our comparison, we externally validated the Framingham risk model (4).

## 115 Results

### 116 Developed models

117 Summary statistics for our study data are provided in Table in S2 Table. There were significant  
118 differences in all features when stratified on outcome status. Between the training and testing set, BMI and  
119 Physical Activity had small, but significant differences between means or proportions between the training  
120 and testing sets. All other variables had no significant difference, see Table in S6 Table. The outcome rate  
121 of the full, training and testing set was 24.65%, 24.36% and 25.3%, respectively. Missing entries in the  
122 features were relatively low for almost all features. Missingness exceeded 10 % for ‘*Family history of*  
123 *hypertension*’ and 1% for ‘*Family history of CVD*’, ‘*Socio-economic status*’ and ‘*Physical Activity*’. In  
124 total, 4572 individuals missed one entry, 570 missed two, 73 missed three, 9 missed four, and 1 missed  
125 five. In sum, 5225 (22%) individuals had at least one missing entry.

126 Hyperparameters selected from the cross-validation procedure and the results from evaluation on the  
127 training set are listed in Table in S7 Table.

128 Applying the final models upon the testing set, we obtained the bootstrapped results given in Table 1. On  
129 average, the eXtreme Gradient Boosting (XGBoost) model performed better on the Area Under the Curve  
130 measure (AUC) and Brier score, but with largely overlapping confidence intervals compared to the elastic  
131 regression and Support Vector Machine (SVM) models. The Random Forest (RF) model excelled on the  
132 Integrated Calibration Index (ICI) and outperformed all other developed models on average.

133 **Table 1: Model results achieved on test set.**

Models	AUC (↑)	Brier (↓)	ICI (↓)
Elastic regression	0.801 [0.789, 0.813]	0.148 [0.143, 0.153]	0.023 [0.015, 0.032]
XGBoost	<b>0.803 [0.791, 0.815]</b>	<b>0.147 [0.142, 0.152]</b>	0.018 [0.010, 0.026]
RF	0.789 [0.776, 0.802]	0.152 [0.147, 0.157]	0.013 [0.006, 0.020]

SVM	0.798 [0.785, 0.810]	0.150 [0.145, 0.155]	0.029 [0.021, 0.037]
KNN	0.786 [0.773, 0.799]	0.153 [0.148, 0.158]	0.019 [0.011, 0.027]
Framingham risk model (Ext.)	0.794 [0.782, 0.805]	0.175 [0.170, 0.180]	0.127 [0.118, 0.136]
Framingham risk model, recalibrated (Ext.)	0.794 [0.782, 0.805]	0.150 [0.145, 0.154]	<b>0.011 [0.003, 0.018]</b>

Performance obtained applying the fitted models on the test set. Reported as mean and 95% confidence interval after bootstrapping. The symbols (↑) and (↓) signify increasing or decreasing values as improved performance, respectively. Models marked (Ext.) are externally developed models with or without recalibration.

134  
135 The calibration plots are given in Fig 1 with individual curves and prediction distributions given in S11  
136 Fig. The shape of the curves was similar for all developed models, with most models overestimating risk  
137 for predictions above 50%. The RF model obtained a lower ICI as it was overall well-calibrated for  
138 predictions below 50%, and far closer to the calibration reference line for most of its predictions compared  
139 to the other models.

140 The Net Benefit plot is given in Fig 2 with individual curves given in S12 Fig. Assuming a threshold  
141 probability < 50%, all models displayed favorable net benefit compared to the reference options of ‘treat  
142 all’, ‘treat none’, or predicting individuals to be ‘Hypertensive’ at follow-up if they were prehypertensive  
143 at baseline. Assuming a threshold above 50%, some of the models exhibited negative benefit, i.e., that the  
144 cost of erroneous predictions is higher than the benefit. Reviewing net benefit across all thresholds, the  
145 elastic regression and XGBoost models yielded the highest, but similar, net benefit.

146 Other performance indicators from evaluations on the testing set are reported in Table in S8 Table. In  
147 short, XGBoost and elastic regression yielded the best average scores of the developed models.

## 148 External model

149 Judging by the mean AUC, the Framingham risk model performed comparably to our newly developed  
150 models, even outperforming the K-Nearest Neighbor (KNN) and RF models. However, the Framingham  
151 risk model had the worst Brier score and calibration compared to the other models. Judging by the

152 calibration plot in Fig 2, the Framingham risk model overestimated risk and increasingly so as the  
153 predicted risk became larger.

154 The Framingham risk model provided favorable Net Benefit compared to the references but was  
155 dominated by all other models. Other performance indicators for the Framingham risk model evaluations  
156 are listed in Table in S8Table. The average performance was slightly worse when evaluated on the full  
157 dataset, but within the 95% confidence interval reported for the testing set, see Table 2. After recalibration  
158 of the Framingham risk model using the training set, the recalibrated model achieved a better ICI than all  
159 other models on the testing set, see Table 1. The Brier score was comparable, but not superior, to  
160 developed models as its discriminatory performance was slightly worse than the other models. See Table  
161 in S4 Table for details on the recalibrated model.

162 **Table 2: Framingham risk model results achieved on the whole dataset.**

<b>Models</b>	<b>AUC (↑)</b>	<b>Brier Score (↓)</b>	<b>ICI (↓)</b>
Framingham risk model	0.787 [0.780, 0.793]	0.178 [0.175, 0.181]	0.133 [0.127, 0.138]

Performance obtained applying the Framingham risk model with adaptations on the whole dataset. Reported as mean and 95% confidence interval after bootstrapping. The symbols (↑) and (↓) signify increasing or decreasing values as improved performance, respectively.

163

164 **Sensitivity analysis**

165 The results from our sensitivity analysis using the LASSO regression with increasing penalty are  
166 displayed in Fig 3. Most of the coefficients were zeroed out at a low penalty. The order and penalty in  
167 which coefficients were zeroed out are given in Table in S9 Table. Diastolic blood pressure, systolic blood  
168 pressure, age, and BMI stand out as important features. The AUC performance indicator decreased slowly  
169 until these variables were eliminated. However, the Brier and ICI score became notably worse at a lower  
170 regularization penalty compared to the AUC score.

171 To compare the importance of features, we report the normalized importance gain calculated in the  
172 XGBoost and RF model from model-fitting on the training set. These are displayed in Fig 4 and reported



173 in detail in Table in S9 Table. In the XGBoost model, the feature importance was concentrated around  
174 age, diastolic blood pressure, systolic blood pressure, and BMI, with all others less important. The same  
175 top four was reported for the RF model, however, the impact of all numerical  
176 measurements were rated higher. BMI was estimated to be less important than age, diastolic blood  
177 pressure, and systolic blood pressure, but more important than the others in the LASSO regressions and by  
178 the normalized importance gain in XGBoost and RF models.

## 179 Discussion

180 In this study we have performed a thorough comparison of multiple machine learning algorithms for  
181 creating risk models for hypertension incidence. With a selection of machine learning models with  
182 distinctive characteristics, we produced well-performing models in terms of discrimination, calibration,  
183 and net benefit. This study marks the first attempt at applying machine learning methods and evaluating  
184 their performance on discrimination, calibration, and clinical usefulness within the same study. Further,  
185 this study is the first using a Norwegian cohort for constructing any risk model for hypertension incidence.  
186 In reviewing the literature, this is only the third time a Scandinavian population has been used (25,29). As  
187 these included genetic data in their models which was not present in our study data, we were not able to  
188 externally validate the model produced using other Scandinavian cohorts.

189 More complex machine learning methods performed better on average on each of AUC, Brier Score, and  
190 ICI, however, the apparent benefit was small. While the XGBoost model achieved the best mean AUC and  
191 Brier Score, the difference versus the much simpler elastic regression model was small, and far smaller  
192 than the variation induced by bootstrapping on the testing set. Comparing the AUC scores we achieved in  
193 this study versus a meta-analysis of the literature, we find that it is close to the expected mean AUC for a  
194 hypertension risk model (5–7). Relative to other studies, the small differences in discrimination between  
195 the elastic regression and the best performing machine learning model, is similar to that found in other  
196 studies (11,21,23).

197 The HUNT Study offers a dataset with large sample size. Evaluations in cross-validation and testing data  
198 yielded similar model performance indicators. This suggests that the sample size was sufficient for our  
199 study and chosen validation scheme. However, we note that we have a longer time frame than most other  
200 risk prediction models in the literature. In addition, there was a slight majority of women versus men in  
201 our data due to men being more likely to decline participation and being lost to follow-up (30). As men  
202 were more likely to develop hypertension, this might have biased the outcome rate. At the same time, the  
203 overall outcome rate of 25% in the HUNT Study is close to the age-standardized hypertension prevalence  
204 in high-income countries at 28.5%, calculated using data from 2005-2014 (31). In summary, we are  
205 confident that these factors did not affect the *relative* performance of the models, which was the focus of  
206 our study.

207 As a reference and external validation, we included the performance indicators produced by the  
208 Framingham risk model (4). However, we should be careful when comparing its results to our newly  
209 developed models, even though the same data is being applied. It is expected that an external validation on  
210 an unseen population will produce worse performance than the internal validation of a newly developed  
211 model, regardless of how the internal validation is performed (26,27,32). Further, the Framingham risk  
212 model is less complex than some of the machine learning methods applied in this study. Hence, the  
213 Framingham risk model results are more appropriately interpreted as a lower benchmark for acceptable  
214 performance for new risk models. For a fair comparison of the Framingham risk models and the newly  
215 developed model, they need to be applied to the same data in an external validation (26).

216 Performance-wise, it is notable that the Framingham model's average AUC was comparable to some of the  
217 models developed in this study. The worse indicator for calibration was more as expected (27,32). Volzke  
218 et al also found the Framingham risk model to have an acceptable discrimination, but poor calibration on  
219 their Danish cohort (25). As for clinical usefulness, the net benefit was positive, but lower than the best  
220 performing models when the risk threshold was lower than 50%. There is a negative net benefit for risk  
221 thresholds over 50% using the Framingham risk model, which is probably related to the increasing degree

222 of miscalibration seen in its calibration plot in S11 Fig. However, the Net Benefit of the newly developed  
223 models were small, null, or slightly negative in the same regions, meaning no model was particularly  
224 useful when a threshold above 50% was applied. Lastly, we note that the Framingham risk model was  
225 developed for a shorter follow-up and that there are some differences in the study cohort, see Table in S5  
226 Table. Most notable, despite the lower mean blood pressure at baseline and shorter time to follow-up, the  
227 outcome ratio in the Framingham risk model cohort was far higher, 45% vs. 25% in our study data. This  
228 might explain why the Framingham risk model overestimates risk in our external validation. The rate of  
229 hypertension in close family was also different, which we suspect is related to the difference in how the  
230 feature was recorded. For the Framingham risk model, parental hypertension was recorded from a separate  
231 cohort study, where the parents themselves were assessed for hypertension multiple times over a longer  
232 period (4). In the HUNT Study data, parental hypertension feature was recorded based on a questionnaire  
233 where individuals reported for their families, which might provide less accurate data. Lastly, we did make  
234 some amendments to the model itself to match our available data, see Table in S4 Table. This may have  
235 had an impact on our external validation of the model.

236 Model updating or recalibration may be a method for obtaining a well-performing risk model with little  
237 effort using an external model. In applying a simple recalibration using the linear predictor of the  
238 Framingham risk model we obtained a risk model with excellent calibration, see Fig G in S12 Fig. While  
239 discrimination is unaffected by our recalibration method, the recalibrated model was on average better  
240 calibrated than the new developed models (27,33). Although it discriminated slightly worse than other  
241 models, the improved calibration was the likely cause of the Net Benefit becoming on par with the best  
242 performing models in our study. Recalibrating the Framingham risk model seems to be a worthwhile  
243 alternative to develop a new risk model from scratch using the HUNT Study data.

244 In the sensitivity analysis, the LASSO regression model performed well on discrimination and calibration  
245 on the testing set, even as many coefficients were zeroed out. Despite the statistically significant  
246 differences between individuals whose blood pressure remained within the range of normotension and

247 individuals who developed hypertension at baseline, the predictive power may be low for some features.  
248 The LASSO regression model with all features except linear effects of blood pressure, age, and BMI  
249 eliminated performed well when applied on the testing set, suggesting these as important features. A  
250 possible explanation to the high performance with low features might be that the zeroed-out features have  
251 non-linear effects not captured by the LASSO regression, as opposed to a small or no linear effect.  
252 However, we saw the same features emphasized as the most important in the XGBoost and RF models.  
253 Both methods are capable of learning non-linear effects from their inputs. The RF importance did  
254 emphasize more features than those mentioned, e.g., all serum biomarkers, and we note that these  
255 constitute all the numerical features used in this study. However, the testing set performance of the RF  
256 model was not consistently better than the other models. On average, the RF model had better calibration  
257 performance but with poorer performance on discrimination. Hence, the serum biomarkers may be  
258 regarded as important by the RF model without having much predictive power.

259 There are several limitations of our study. We had available 17 diverse features, including biomarkers and  
260 other easily obtainable clinical features from the HUNT Study data. The inclusion of other types of  
261 features such as genetic information, more comprehensive socioeconomic information, or information on  
262 diet could have made a difference on model performance. This became clear in the sensitivity analysis, in  
263 which only a small subset of features was determined as important. For example, Niu et al showed that the  
264 addition of a genetic risk score improved the performance of their machine learning models, but not the  
265 less complex Cox regression model (24). However, in a recent meta-analysis, the addition of genetic  
266 information was found to not improve discrimination of the included risk models (7). The latter models  
267 were largely developed using established statistical methods which included linear effects of the genetic  
268 information. Hence, which features should be included to improve predictions of incident hypertension  
269 remain unclear and a topic of further study. The second limitation was our imputation procedure: We did  
270 not perform *multiple* imputations to capture uncertainty in the imputation method. This is suggested by  
271 guidelines for developing risk models, but we refrained from doing so due to high computational costs

272 (32). However, all models were subjected to a rigorous cross-validation scheme during development and  
273 bootstrapping of their performance on the testing set. Further, the missingness was low, with only ‘Family  
274 history of hypertension’ having a missing rate above 4%. While some variation in our imputation scheme  
275 may not have been captured, we do not think it would impact the performance of the developed models  
276 relative to each other. Third, our study is restricted to comparing the relative performance of various  
277 modelling methods for the HUNT Study data. Other relevant factors for the practical use of machine  
278 learning models are the applicability and transparency of the model (34). While many risk models are  
279 developed using relatively simple methods like logistic regression, user-friendly risk score sheets are often  
280 provided for simple application of the model for health practitioners (32). As a contrast, a computer client  
281 would be needed to calculate predictions using the XGBoost, RF, SVM, and the KNN models (35). In  
282 addition, exactly how XGBoost, RF, and SVM work for individual predictions is often complex and  
283 difficult to discern. While some argue that auxiliary methods may be helpful, other suggest avoiding use  
284 of “black-box” model predictions for high-stakes decisions (36,37). However, for these considerations to  
285 be of interest, the validation performance using more complex models should be superior to other models  
286 to justify their use. In this study, the machine learning models did not outperform models developed using  
287 less-complex methods. Lastly, we also note some limitations in the interpretation of decision curves, and  
288 how they relate to clinical usefulness. We refer to Kerr et al. and Vickers et al. for details [39], [40].

289 The main strength of our study is that we have taken great care in ensuring that all aspects of model-fitting  
290 and evaluation were as similar as possible for all modelling methods. This ensures that any variation seen  
291 in results or feature importance between methods are due to the characteristics of the modelling methods  
292 themselves. Further, a large sample size ensures that our fitted models were robust, with similar  
293 performance in both cross-validation and testing. Furthermore, we evaluated our models with respect to  
294 discrimination, calibration, and clinical usefulness. The similarity in our scheme for model-fitting and  
295 evaluation allows us to feel confident about comparing the performance of the developed models relative  
296 to each other. Lastly, we evaluated the externally developed Framingham risk model as an alternative to

297 developing a new risk model from scratch. However, for application of the developed models as risk  
298 prediction models and to assess their generalization to new data, a separate validation on new unseen data  
299 would be necessary (26,27).

300 In conclusion, we have developed and compared the performance of hypertension risk models using  
301 different machine learning methods. While more complex methods displayed good discrimination and  
302 calibration, they did not consistently outperform a logistic regression model fitted with elastic  
303 regularization. We found the externally developed Framingham risk model to produce almost as good  
304 discrimination scores as the newly developed models. The original model overestimated hypertension risk  
305 in the HUNT Study, but this was amended by a simple recalibration to our data. In our sensitivity analysis,  
306 the features age, systolic blood pressure, diastolic blood pressure, and BMI was found to be particularly  
307 important compared to the other included features.

## 308 **Materials and methods**

### 309 **Data**

310 A dataset was derived from The Trøndelag Health (HUNT) Study, based in the now former county of  
311 Nord-Trøndelag in Norway. The HUNT Study constitutes a large population database for medical and  
312 health-related research based in four health surveys over four decades (38). Specifically, baseline data was  
313 collected from HUNT2 (1995-1997) with endpoint derived from the follow-up in HUNT3 (2006-2008).

314 We included individuals (>20 years of age) participating in both surveys:

- 315 – With complete information on blood pressure measurements and use of blood pressure medication  
316 at baseline and follow-up,
- 317 – without missing information on diabetes or history of cardiovascular disease (CVD) at baseline,
- 318 – with a blood pressure below the hypertension threshold and being free of blood pressure  
319 medication, CVD, and diabetes at baseline.

320 Blood pressure measurements in the HUNT Study were performed three times per survey, with the initial  
321 measurement used to calibrate the measurement device (38). The recorded pressure was the average of  
322 recording two and three. Hypertension status was determined following the ESC/ESH guidelines, i.e., a  
323 systolic pressure above 140, diastolic pressure above 90, or usage of blood pressure medication (39). The  
324 process of applying exclusion criteria and dataflow is shown in S10 Fig. In total, 23 722 individuals were  
325 found eligible for this study. The features available for our study are well-established risk factors of  
326 hypertension and CVD and commonly used in risk modelling of incident hypertension (7,39). We  
327 estimated physical activity by a novel physical activity metric, Personal Activity Intelligence (PAI). The  
328 PAI algorithm converts self-reported leisure time physical activity to an average weekly PAI score for the  
329 last year (40–43). The HUNT Study protocol have been described in detail by Åsvold et al. and more  
330 information about how features were collected can be found in Table in S1 Table and at [https://hunt-](https://hunt-db.medisin.ntnu.no/hunt-db/#/)  
331 [db.medisin.ntnu.no/hunt-db/#/](https://hunt-db.medisin.ntnu.no/hunt-db/#/) (38). All participants provided informed written informed consent. This  
332 study was approved by the Regional Committee on Medical and Health Research Ethics of Norway (REK;  
333 22902; 2018/1824). Data can be obtained upon approval from REK and HUNT Research Centre. For more  
334 information see: [www.ntnu.edu/hunt/data](http://www.ntnu.edu/hunt/data).

335 The eligible cohort was stratified on outcome status (normotension; hypertension) and described by  
336 summary statistics and missing rate in Table in S2 Table. We applied unpaired t-tests or chi-square tests as  
337 appropriate to detail significant differences between those whose blood pressure remained within range of  
338 normotension or developed hypertension.

### 339 Model development

340 To minimize the risk of providing overoptimistic results, we used a thorough development and validation  
341 scheme. First, we divided the available dataset randomly into a training and testing set by a 7:3 ratio. We  
342 applied unpaired t-test and the chi-square test for evaluating differences between the training and test set.  
343 Second, we applied a 4-fold cross-validation scheme on the training set to select hyperparameters for our  
344 modelling methods, shown in Fig 5. The combination of hyperparameters that produced the best mean test

345 fold performance during cross-validation was selected for each method separately. Lastly, to produce the  
346 final model for each method, the model was fitted using the whole training set with the selected  
347 hyperparameters, see Fig 6.

## 348 **Model validation**

349 The final models were applied on the testing set to evaluate their performance, see Fig 7. To account for  
350 variations in the data, we used a simple bootstrap scheme: The testing set was resampled with replacement  
351 100 times, measuring the performance of each model on each resampled testing set. For each final model,  
352 the performance indicators were summarized by its mean and standard deviation from the 100 evaluations.

353 To ensure equal conditions for our modelling methods, we used the same data folds in cross-validation of  
354 all methods and measured model performance on the same bootstrapped testing sets. Any variation  
355 induced from the random sampling of data is therefore equal for all modelling methods, which allows for  
356 an accurate comparison of model results.

## 357 **Preprocessing and missing data**

358 As part of the model development and evaluation, the data was preprocessed by standardizing the  
359 numerical features and one-hot-encoding the categorical or nominal features. In the case of missing  
360 feature values, we used a bagged decision tree model to impute the missing values (44). The imputation  
361 model was blinded from the outcome and learned to predict the missing feature using the other features in  
362 the dataset. Both preprocessing and imputation requires determining some parameters: For  
363 standardization, the mean and standard deviation on each feature need to be determined. For imputation,  
364 the bagged decision trees models must be fitted. To avoid 'data-leakage', these parameters were learned  
365 without using data the models were evaluated on. In short, the parameters and imputation method were  
366 determined in the cross-validation scheme using only the training data folds before being applied on the  
367 testing fold. Likewise, when fitting the final models, the standardization parameters and imputation model  
368 was learned using only the training set.



## 369 Modelling methods

370 We included several popular and frequently used machine learning algorithms as modelling methods. In  
371 addition, we included the logistic regression method, with and without elastic regularization. The learning  
372 algorithms used were the logistic regression with and without elastic regularization (referred to as elastic  
373 regression), eXtreme Gradient Boosting algorithm (XGBoost), regularized Random Forest (RF), Support  
374 Vector Machine (SVM), and K-Nearest Neighbor (KNN) (44–48). Notably, we did not include any  
375 machine learning method using neural networks algorithms, as a recent review suggest that they display  
376 less than optimal performance on tabular data, which we employ in this study (49). Features were included  
377 in the models without defining any interactions, for simplicity.

378 Depending on the algorithm, the different methods included necessitated the selection of multiple  
379 hyperparameters. For each hyperparameter, a sensible range of possible values was defined, and search  
380 strategies was employed to select a value from these ranges. For the XGBoost, RF, and SVM modelling  
381 methods we employed a sampling of these ranges. In total, 70, 30, and 30 combinations of  
382 hyperparameters were sampled for the XGBoost, RF, and SVM methods, respectively, and used as  
383 candidates in the cross-validation scheme. For elastic regression and KNN, we employed a grid search  
384 where all combinations were trialed in cross-validation. The rationale behind different sampling strategies  
385 was motivated by the higher training time and computational cost of XGBoost, RF, and SVM methods  
386 compared to elastic regression and KNN. The hyperparameters, ranges and search strategies are described  
387 in Table in S3 Table.

## 388 External models

389 We searched the literature for existing risk models allowing validation: Using similar features to those  
390 available in the HUNT Study data; reporting model performance; and suitable for the 11-year follow-up  
391 period between baseline and outcome. From this, we found four articles detailing such models: One being  
392 the Framingham risk model, and the three others being refitted versions of the Framingham risk model (4).  
393 We therefore included the Framingham risk model in our full model evaluation, i.e., calculated its

394 performance on the test set along with our newly developed models. As it is expected that an external  
395 model will be poorly calibrated, we perform a simple recalibration using the linear predictor of the  
396 Framingham risk model (27,33). We learned the recalibration parameter on the training set and reported  
397 the recalibrated models performance on the test set. We did not perform any further recalibration, e.g., re-  
398 estimating coefficients. This was due to already including a simple linear predictor in our hyperparameter  
399 search for the elastic regression, i.e., elastic regression with ‘lambda’ hyperparameter set to zero, which  
400 reduce to unregularized logistic regression. We also applied the original Framingham risk model on the  
401 whole dataset, after imputation. The model, with recalibration and adaptations made for it to suit our data  
402 can be found in Table in S4 Table. We compare key aspects of the data from the Framingham risk model  
403 development study and the HUNT Study in Table in S5 Table.

#### 404 [Sensitivity analysis](#)

405 As sensitivity analysis, we fitted Lasso logistic regression models on the training data with increasing  
406 regularization penalty and evaluated their performance on the test set (50). We employ all features as in  
407 the model development, without interactions. With increasing penalty, less-important coefficients in the  
408 model-fitting will be shrunk towards zero. By looking at which degree of penalty the different coefficients  
409 are eliminated, we may see which features are of most importance in model-fitting for the logistic  
410 regression model. To compare, we investigated the feature importance calculated during model-fitting on  
411 the training data for the XGBoost and RF methods.

#### 412 [Performance indicators](#)

413 A risk models performance is primarily quantified by indicators for discrimination and calibration. In this  
414 study, we evaluated the models by the *Area Under the receiver operator Curve* (AUC). The AUC is a  
415 frequently used in literature on risk models for hypertension (7,15). We also captured the models overall  
416 performance by the Brier score, which is a proper scoring function (15,51). Further, calibration was  
417 assessed both graphically using smoothed calibration plots, and summarized using the *Integrated*  
418 *Calibration Index* (ICI). The ICI measures the deviation of the smoothed calibration curve of a model

419 versus a perfect, straight, diagonal calibration line, weighted by the distribution of the model's predictions.  
420 A low ICI score indicates that the model is well calibrated for the risk percentiles the model frequently  
421 provided predictions for in the data (52,53). Lastly, we evaluate the clinical usefulness of our models  
422 compared to sensible benchmarks graphically by presenting the *Net Benefit* plot derived from the testing  
423 set. The net benefit complements the smoothed calibration curve in assessing clinical usefulness (15,16).  
424 The benchmarks compared against in the Net Benefit plot was 1) treat all individuals as 'Hypertensive', 2)  
425 treat all as 'Normotensive', and 3) predict and treat an individual as 'Hypertensive' if they were  
426 prehypertensive, i.e., systolic BP above 130 or diastolic BP above 80 mmHg, at baseline, and  
427 'Normotensive' otherwise.

428 In addition, we also reported some performance indicators frequently used in machine learning literature:  
429 The area under the Precision-Recall curve, the F1 measure, sensitivity, specificity, positive predictive  
430 value, negative predictive value, and the Matthews correlation coefficient (8). For the performance  
431 indicators where predictions needed to be either "Normotensive" or "Hypertensive", we assigned all  
432 individual predictions below the outcome rate (24.36%) of the training data as 'Normotensive', and above  
433 as 'Hypertensive'.

434 As a common criteria for choosing the optimal models during cross-validation, we used the Brier Score, as  
435 it is a proper scoring rule regardless of modelling method (54). Hence, the optimal set of hyperparameters  
436 during cross-validation was the one producing models with the lowest Brier Score.

## 437 [Software and reporting](#)

438 We have used the RStudio IDE and R for implementing the modelling algorithms and data-processing.  
439 The following R-packages were used: *skimr* for data exploration, *randomForest*, *RRF*, *glmnet*, *Matrix*,  
440 *xgboost*, *plyr*, *kermlab*, *class* and *caret* for modelling algorithms, *recipes* for preprocessing, *glmnet* for  
441 sensitivity analysis, and *ggplot2*, *ggextra*, *patchwork*, *ggh4x* and *dcurves* for graphics (55–70).

442 To ensure a high degree of transparency and quality in reporting, we strived to follow the *Transparent*  
443 *Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis* (TRIPOD) (32). We  
444 note that a guideline is being developed for prognostic prediction model studies based on artificial  
445 intelligence, which may have been more relevant for our study (71). However, the guideline was not  
446 available at the time of writing. A form detailing adherence to the TRIPOD guidelines for the  
447 development of the new models and the validation of the Framingham risk model has been attached as  
448 Supporting Information, see Files S13 and S14.

## 449 Acknowledgments

450 The HUNT Study is a collaboration between HUNT Research Centre (Faculty of Medicine and Health  
451 Sciences, Norwegian University of Science and Technology, Trøndelag County Council, Central Norway  
452 Regional Health Authority, and the Norwegian Institute of Public Health. We thank the participants and  
453 management team of the HUNT Study. We also thank the staff at HUNT Cloud who aided us with tools  
454 for data storage and analysis.

## 455 Online resources

456 To encourage dissemination of the risk models developed in this study, we created an online resource,  
457 [https://github.com/filsch/hypertension\\_prediction\\_models\\_hunt\\_study](https://github.com/filsch/hypertension_prediction_models_hunt_study), where multiple risk models and  
458 auxiliary functions are provided for easy utilization by external researchers. Some example data is  
459 provided to ensure the right formatting of data to be used in the models. The risk models included in the  
460 resource are the XGBoost, RF and the elastic regression model. In addition, a logistic regression model  
461 using fewer features derived from the sensitivity analysis, the adapted Framingham risk model, and the  
462 adapted Framingham risk model recalibrated on the HUNT Study data are also included.

## 463 Funding

464 The author(s) received no specific funding for this work.

## 465 References

- 466 1. Carretero OA, Oparil S. Essential Hypertension: Part I: Definition and Etiology. *Circulation*. 2000 Jan  
467 25;101(3):329–35.
- 468 2. Zhou B, Bentham J, Di Cesare M, Bixby H, Danaei G, Cowan MJ, et al. Worldwide trends in blood  
469 pressure from 1975 to 2015: a pooled analysis of 1479 population-based measurement studies with  
470 19·1 million participants. *The Lancet*. 2017 Jan;389(10064):37–55.
- 471 3. Gaziano TA, Bitton A, Anand S, Weinstein MC. The global cost of nonoptimal blood pressure. *Journal*  
472 *of Hypertension*. 2009 Jul;27(7):1472–7.
- 473 4. Parikh NI, Pencina MJ, Wang TJ, Benjamin EJ, Lanier KJ, Levy D, et al. A risk score for predicting near-  
474 term incidence of hypertension: The Framingham Heart Study. Vol. 148, *Annals of Internal Medicine*.  
475 2008. p. 102–10.
- 476 5. Sun D, Liu J, Xiao L, Liu Y, Wang Z, Li C, et al. Recent development of risk-prediction models for  
477 incident hypertension: An updated systematic review. Vol. 12, *PLoS One*. 2017. p. e0187240.
- 478 6. Echouffo-Tcheugui JB, Batty GD, Kivimaki M, Kengne AP. Risk Models to Predict Hypertension: A  
479 Systematic Review. Vol. 8, *Plos One*. 2013.
- 480 7. Chowdhury MZI, Naeem I, Quan H, Leung AA, Sikdar KC, O’Beirne M, et al. Prediction of  
481 hypertension using traditional regression and machine learning models: A systematic review and  
482 meta-analysis. Palazón-Bru A, editor. *PLoS ONE*. 2022 Apr 7;17(4):e0266334.
- 483 8. Weng SF, Reps J, Kai J, Garibaldi JM, Qureshi N. Can machine-learning improve cardiovascular risk  
484 prediction using routine clinical data? Liu B, editor. *PLoS ONE*. 2017 Apr 4;12(4):e0174944.
- 485 9. Chen JH, Asch SM. Machine Learning and Prediction in Medicine — Beyond the Peak of Inflated  
486 Expectations. *New England Journal of Medicine*. 2017 Jun;376(26):2507–9.
- 487 10. Beam AL, Kohane IS. Big Data and Machine Learning in Health Care. *JAMA*. 2018 Apr 3;319(13):1317.
- 488 11. Nusinovici S, Tham YC, Chak Yan MY, Wei Ting DS, Li J, Sabanayagam C, et al. Logistic regression was  
489 as good as machine learning for predicting major chronic diseases. *Journal of Clinical Epidemiology*.  
490 2020 Jun;122:56–69.
- 491 12. Lynam AL, Dennis JM, Owen KR, Oram RA, Jones AG, Shields BM, et al. Logistic regression has similar  
492 performance to optimised machine learning algorithms in a clinical setting: application to the  
493 discrimination between type 1 and type 2 diabetes in young adults. *Diagn Progn Res*. 2020  
494 Dec;4(1):6.
- 495 13. Gravesteijn BY, Nieboer D, Ercole A, Lingsma HF, Nelson D, van Calster B, et al. Machine learning  
496 algorithms performed no better than regression models for prognostication in traumatic brain  
497 injury. *Journal of Clinical Epidemiology*. 2020 Jun;122:95–107.

- 498 14. Christodoulou E, Ma J, Collins GS, Steyerberg EW, Verbakel JY, Calster BV. A systematic review shows  
499 no performance benefit of machine learning over logistic regression for clinical prediction models.  
500 *Journal of Clinical Epidemiology*. 2019 Jun;110:12–22.
- 501 15. Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, et al. Assessing the  
502 performance of prediction models: a framework for traditional and novel measures. *Epidemiology*.  
503 2010 Jan;21(1):128–38.
- 504 16. Vickers AJ, Elkin EB. Decision curve analysis: a novel method for evaluating prediction models. *Med*  
505 *Decis Making*. 2006 Dec;26(6):565–74.
- 506 17. Kappen TH, van Klei WA, van Wolfswinkel L, Kalkman CJ, Vergouwe Y, Moons KGM. Evaluating the  
507 impact of prediction models: lessons learned, challenges, and recommendations. *Diagn Progn Res*.  
508 2018 Dec;2(1):11.
- 509 18. Ramezankhani A, Kabir A, Pournik O, Azizi F, Hadaegh F. Classification-based data mining for  
510 identification of risk patterns associated with hypertension in Middle Eastern population: A 12-year  
511 longitudinal study. Vol. 95, *Medicine*. 2016.
- 512 19. Sakr S, Elshawi R, Ahmed A, Qureshi WT, Brawner C, Keteyian S, et al. Using machine learning on  
513 cardiorespiratory fitness data for predicting hypertension: The Henry Ford Exercise Testing (FIT)  
514 Project. Vol. 13, *Plos One*. 2018.
- 515 20. Silva GFS, Fagundes TP, Teixeira BC, Chiavegatto Filho ADP. Machine Learning for Hypertension  
516 Prediction: a Systematic Review. *Curr Hypertens Rep [Internet]*. 2022 Jun 22 [cited 2022 Sep 8];  
517 Available from: <https://link.springer.com/10.1007/s11906-022-01212-6>
- 518 21. Dritsas E, Fazakis N, Kocsis O, Fakotakis N, Moustakas K. Long-Term Hypertension Risk Prediction  
519 with ML Techniques in ELSA Database. In: Simos DE, Pardalos PM, Kotsireas IS, editors. *Learning and*  
520 *Intelligent Optimization [Internet]*. Cham: Springer International Publishing; 2021 [cited 2022 Jun  
521 29]. p. 113–20. (Lecture Notes in Computer Science; vol. 12931). Available from:  
522 [https://link.springer.com/10.1007/978-3-030-92121-7\\_9](https://link.springer.com/10.1007/978-3-030-92121-7_9)
- 523 22. Xu F, Zhu JC, Sun N, Wang L, Xie C, Tang QX, et al. Development and validation of prediction models  
524 for hypertension risks in rural Chinese populations. Vol. 9, *Journal of Global Health*. 2019.
- 525 23. Kanegae H, Suzuki K, Fukatani K, Ito T, Harada N, Kario K. Highly precise risk prediction model for  
526 new-onset hypertension using artificial intelligence techniques. Vol. 22, *Journal of Clinical*  
527 *Hypertension*. 2020. p. 445–50.
- 528 24. Niu M, Wang Y, Zhang L, Tu R, Liu X, Hou J, et al. Identifying the predictive effectiveness of a genetic  
529 risk score for incident hypertension using machine learning methods among populations in rural  
530 China. *Hypertension Research*. 2021 Nov 1;44(11):1483–91.
- 531 25. Volzke H, Fung G, Ittermann T, Yu SP, Baumeister SE, Dorr M, et al. A new, accurate predictivemodel  
532 for incident hypertension. Vol. 31, *Journal of Hypertension*. 2013. p. 2142–50.
- 533 26. Ramspek CL, Jager KJ, Dekker FW, Zoccali C, van Diepen M. External validation of prognostic models:  
534 what, why, how, when and where? *Clinical Kidney Journal*. 2021 Feb 3;14(1):49–58.

- 535 27. Moons KGM, Kengne AP, Grobbee DE, Royston P, Vergouwe Y, Altman DG, et al. Risk prediction  
536 models: II. External validation, model updating, and impact assessment. *Heart*. 2012 May  
537 1;98(9):691.
- 538 28. Altman DG, Vergouwe Y, Royston P, Moons KGM. Prognosis and prognostic research: validating a  
539 prognostic model. *BMJ*. 2009 May 28;338(may28 1):b605–b605.
- 540 29. Fava C, Sjogren M, Montagnana M, Danese E, Almgren P, Engstrom G, et al. Prediction of Blood  
541 Pressure Changes Over Time and Incidence of Hypertension by a Genetic Risk Score in Swedes. Vol.  
542 61, *Hypertension*. 2013. p. 319-+.
- 543 30. Hofman AC, Espeland L, Steinsland I, Ingeström EML. A Shared Parameter Model for Systolic Blood  
544 Pressure Accounting for Data Missing Not at Random in the HUNT Study [Internet]. arXiv; 2022 [cited  
545 2022 Sep 27]. Available from: <http://arxiv.org/abs/2203.16602>
- 546 31. Mills KT, Bundy JD, Kelly TN, Reed JE, Kearney PM, Reynolds K, et al. Global Disparities of  
547 Hypertension Prevalence and Control: A Systematic Analysis of Population-Based Studies From 90  
548 Countries. *Circulation*. 2016 Aug 9;134(6):441–50.
- 549 32. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et al. Transparent  
550 Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD):  
551 Explanation and Elaboration. *Ann Intern Med*. 2015 Jan 6;162(1):W1–73.
- 552 33. Steyerberg EW, Borsboom GJJM, van Houwelingen HC, Eijkemans MJC, Habbema JDF. Validation and  
553 updating of predictive logistic regression models: a study on sample size and shrinkage. *Stat Med*.  
554 2004 Aug 30;23(16):2567–86.
- 555 34. Watson DS, Krutzinna J, Bruce IN, Griffiths CE, McInnes IB, Barnes MR, et al. Clinical applications of  
556 machine learning algorithms: beyond the black box. *BMJ*. 2019 Mar 12;l886.
- 557 35. Van Calster B, Wynants L, Timmerman D, Steyerberg EW, Collins GS. Predictive analytics in health  
558 care: how can we know it works? *Journal of the American Medical Informatics Association*. 2019 Dec  
559 1;26(12):1651–4.
- 560 36. Molnar C, Casalicchio G, Bischl B. Interpretable Machine Learning -- A Brief History, State-of-the-Art  
561 and Challenges. In 2020 [cited 2022 Sep 12]. p. 417–31. Available from:  
562 <http://arxiv.org/abs/2010.09337>
- 563 37. Rudin C. Stop explaining black box machine learning models for high stakes decisions and use  
564 interpretable models instead. *Nat Mach Intell*. 2019 May;1(5):206–15.
- 565 38. Åsvold BO, Langhammer A, Rehn TA, Kjellvik G, Grøntvedt TV, Sjørgjerd EP, et al. Cohort Profile  
566 Update: The HUNT Study, Norway. *International Journal of Epidemiology* [Internet]. 2022 May 17  
567 [cited 2022 Jul 7]; Available from: [https://academic.oup.com/ije/advance-](https://academic.oup.com/ije/advance-article/doi/10.1093/ije/dyac095/6586600)  
568 [article/doi/10.1093/ije/dyac095/6586600](https://academic.oup.com/ije/advance-article/doi/10.1093/ije/dyac095/6586600)
- 569 39. Williams B, Mancia G, Spiering W, Agabiti Rosei E, Azizi M, Burnier M, et al. 2018 ESC/ESH Guidelines  
570 for the management of arterial hypertension. *European Heart Journal*. 2018 Sep 1;39(33):3021–104.

- 571 40. Kurtze N, Rangul V, Hustvedt BE, Flanders WD. Reliability and validity of self-reported physical  
572 activity in the Nord-Trøndelag Health Study (HUNT 2). *Eur J Epidemiol*. 2007;22(6):379–87.
- 573 41. Kieffer SK, Nauman J, Syverud K, Selboskar H, Lydersen S, Ekelund U, et al. Association between  
574 Personal Activity Intelligence (PAI) and body weight in a population free from cardiovascular disease  
575 - The HUNT study. *Lancet Reg Health Eur*. 2021 Jun;5:100091.
- 576 42. Nauman J, Nes BM, Zisko N, Revdal A, Myers J, Kaminsky LA, et al. Personal Activity Intelligence  
577 (PAI): A new standard in activity tracking for obtaining a healthy cardiorespiratory fitness level and  
578 low cardiovascular risk. *Progress in Cardiovascular Diseases*. 2019 Mar;62(2):179–85.
- 579 43. Nes BM, Gutvik CR, Lavie CJ, Nauman J, Wisløff U. Personalized Activity Intelligence (PAI) for  
580 Prevention of Cardiovascular Disease and Promotion of Physical Activity. *The American Journal of*  
581 *Medicine*. 2017 Mar;130(3):328–36.
- 582 44. James G, Witten D, Hastie T, Tibshirani R. An Introduction to Statistical Learning [Internet]. New  
583 York, NY: Springer New York; 2013 [cited 2022 Jun 29]. (Springer Texts in Statistics; vol. 103).  
584 Available from: <http://link.springer.com/10.1007/978-1-4614-7138-7>
- 585 45. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the 22nd ACM  
586 SIGKDD International Conference on Knowledge Discovery and Data Mining [Internet]. 2016 [cited  
587 2022 Jun 29]. p. 785–94. Available from: <http://arxiv.org/abs/1603.02754>
- 588 46. Breiman L. Random Forests. *Machine Learning*. 2001;45(1):5–32.
- 589 47. Zou H, Hastie T. Regularization and variable selection via the elastic net. *J Royal Statistical Soc B*.  
590 2005 Apr;67(2):301–20.
- 591 48. Cortes C, Vapnik V. Support-vector networks. *Mach Learn*. 1995 Sep;20(3):273–97.
- 592 49. Borisov V, Leemann T, Seßler K, Haug J, Pawelczyk M, Kasneci G. Deep Neural Networks and Tabular  
593 Data: A Survey [Internet]. arXiv; 2022 [cited 2022 Jun 29]. Available from:  
594 <http://arxiv.org/abs/2110.01889>
- 595 50. Tibshirani R. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society*  
596 *Series B (Methodological)*. 1996;58(1):267–88.
- 597 51. Brier GW. Verification of forecasts expressed in terms of probability. *Monthly weather review*.  
598 1950;78(1):1–3.
- 599 52. Austin PC, Steyerberg EW. The Integrated Calibration Index (ICI) and related metrics for quantifying  
600 the calibration of logistic regression models. *Statistics in Medicine*. 2019 Sep 20;38(21):4051–65.
- 601 53. Calster BV, McLernon DJ, Smeden M van, Wynants L, Steyerberg EW. Calibration: the Achilles  
602 heel of predictive analytics. *BMC Medicine*. 2019 Dec;17(1).
- 603 54. Gneiting T, Raftery AE. Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the*  
604 *American Statistical Association*. 2007 Mar;102(477):359–78.



- 605 55. Waring E, Quinn M, McNamara A, Rubia EA de la, Zhu H, Ellis S. skimr: Compact and Flexible  
606 Summaries of Data [Internet]. 2022. Available from: <https://CRAN.R-project.org/package=skimr>
- 607 56. Liaw A, Wiener M. Classification and Regression by randomForest. R News. 2002;2(3):18–22.
- 608 57. Deng H. Guided Random Forest in the RRF Package. arXiv:13060237. 2013;
- 609 58. Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models via Coordinate  
610 Descent. Journal of Statistical Software. 2010;33(1):1–22.
- 611 59. Bates D, Maechler M, Jagan M. Matrix: Sparse and Dense Matrix Classes and Methods [Internet].  
612 2022. Available from: <https://CRAN.R-project.org/package=Matrix>
- 613 60. Chen T, He T, Benesty M, Khotilovich V, Tang Y, Cho H, et al. xgboost: Extreme Gradient Boosting  
614 [Internet]. 2021. Available from: <https://github.com/dmlc/xgboost>
- 615 61. Wickham H. The Split-Apply-Combine Strategy for Data Analysis. Journal of Statistical Software.  
616 2011;40(1):1–29.
- 617 62. Karatzoglou A, Smola A, Hornik K. kernlab: Kernel-Based Machine Learning Lab [Internet]. 2022.  
618 Available from: <https://CRAN.R-project.org/package=kernlab>
- 619 63. Venables WN, Ripley BD. Modern Applied Statistics with S. 4. ed., [Nachdr.]. New York: Springer;  
620 2010. 495 p. (Statistics and computing).
- 621 64. Kuhn M. caret: Classification and Regression Training [Internet]. 2022. Available from:  
622 <https://CRAN.R-project.org/package=caret>
- 623 65. Kuhn M, Wickham H. recipes: Preprocessing and Feature Engineering Steps for Modeling [Internet].  
624 2022. Available from: <https://CRAN.R-project.org/package=recipes>
- 625 66. Wickham H. ggplot2: Elegant Graphics for Data Analysis [Internet]. Springer-Verlag New York; 2016.  
626 Available from: <https://ggplot2.tidyverse.org>
- 627 67. Attali D, Baker C. ggExtra: Add Marginal Histograms to “ggplot2”, and More “ggplot2” Enhancements  
628 [Internet]. 2022. Available from: <https://CRAN.R-project.org/package=ggExtra>
- 629 68. Pedersen TL. patchwork: The Composer of Plots [Internet]. 2022. Available from: <https://CRAN.R-project.org/package=patchwork>
- 631 69. Brand T van den. ggh4x: Hacks for “ggplot2” [Internet]. 2022. Available from: <https://CRAN.R-project.org/package=ggh4x>
- 633 70. Sjoberg DD. dcurves: Decision Curve Analysis for Model Evaluation [Internet]. 2022. Available from:  
634 <https://CRAN.R-project.org/package=dcurves>
- 635 71. Collins GS, Dhiman P, Andaur Navarro CL, Ma J, Hooft L, Reitsma JB, et al. Protocol for development  
636 of a reporting guideline (TRIPOD-AI) and risk of bias tool (PROBAST-AI) for diagnostic and prognostic  
637 prediction model studies based on artificial intelligence. BMJ Open. 2021 Jul;11(7):e048008.

638 72. Vickers AJ, van Calster B, Steyerberg EW. A simple, step-by-step guide to interpreting decision curve  
639 analysis. *Diagn Progn Res.* 2019 Dec;3(1):18.

640

## 641 **Figure captions**

642 **Fig 1. Calibration curves calculated on the test set.**

643 Calculated as mean calibration curve after pointwise bootstrapping.

644 **Fig 2. Standardized Net Benefit calculated on the test set.**

645 Standardized Net Benefit displays the benefit of applying a model on a population-level and has an upper  
646 bound of 1 with no lower bound. A model should be preferred over another if it dominates over all the  
647 relevant risk thresholds. See Vickers et al. for further details on interpretation of Net Benefit (72).

648 **Fig 3. Sensitivity analysis with LASSO regression.**

649 (A): Coefficient sizes versus penalty in LASSO regression fitted on the training set. Ten last coefficients  
650 to be zeroed out are colored and named. Note, numerical features have been standardized, i.e., coefficient  
651 sizes correspond to the increase of one standard deviation of the training set, see Table in S6 Table. (B):  
652 AUC, Brier score and ICI calculated on test set using LASSO models with coefficients as in Panel A, with  
653 panels corresponding on the X-axis. Black mean line and red 95% confidence interval derived from  
654 pointwise bootstrapping. Note, the performance scores as all coefficients are zeroed out have been  
655 cropped out as the AUC becomes 0.5, Brier Score 0.1875 and ICI undefined.

656 **Fig 4. Feature-importance calculated by Random Forest and XGBoost**

657 Permutation importance calculated by XGBoost and RF models after model fitting on the training set,  
658 sorted by importance in Random Forest. Normalized relative to 'Age', which was the highest ranked  
659 feature in both models. Coloring of features follow the legend of Fig A in S3 Fig.

660 **Fig 5. Cross-validation scheme for selecting hyperparameters.**

661 The Key Performance Indicator (KPI) used to select hyperparameters was the mean out-of-fold Brier  
662 Score.

663 **Fig 6. Model-fitting and learning of imputation and preprocessing parameters.**

664 Each method was fitted using the optimal hyperparameters selected in cross-validation.

665 **Fig 7. Bootstrap scheme for evaluating performance on the test set.**

## 666 Supporting information captions

667 **S1 Table. Variable names used to construct features, as named in HUNT databank.**

668 **S2 Table. Feature distributions for all data and subdivided by outcome status.**

669 **S3 Table. Hyperparameters, candidate values, and search strategy per modelling method.**

670 **S4 Table. Adaptations and full model equation of the original Framingham risk model and after  
671 recalibration to the HUNT Study data.**

672 **S5 Table. Study and data characteristics used in developing the Framingham risk model and in this  
673 study.**

674 **S6 Table. Feature and outcome distributions for all data and subdivided by training and test set.**

675 **S7 Table. Hyperparameters selected with out-of-fold performance in cross-validation on training  
676 set.**

677 **S1 Table. Auxiliary performance measures calculated on the test set.**

678 **S9 Table. Results from sensitivity analysis in tabular form.**

679 **S10 Fig. Dataflow on applying exclusion criteria.**

680 The flow of datapoints relative to the application of exclusion criteria on the available data from the  
681 HUNT Study.

682 **S11 Fig. Individual calibration curves for all models using the test set data.**

683 Displayed as black mean curves and red shaded 95% confidence interval, calculated using pointwise  
684 bootstrapping. Dashed line indicates a perfect calibration line. Histogram of predictions on the test set  
685 displayed on top of plot, color-coded by outcome status. Color coding of curves corresponds to legend in  
686 Fig 1 and Fig 2.

687 **S12 Fig. Individual net benefit curves for all models using the test set data.**

688 Displayed as black mean curves and red shaded 95% confidence interval, calculated using pointwise  
689 bootstrapping. Color coding of curves corresponds to legend in Fig 1 and Fig 2.

690 **S13 File. TRIPOD form for the development of models using the HUNT Study data.**

691 **S14 File. TRIPOD form for the external validation of the Framingham risk model.**

692 **S15 File. R scripts for analysis and plotting, and datafile for plotting figures and tabular data.**

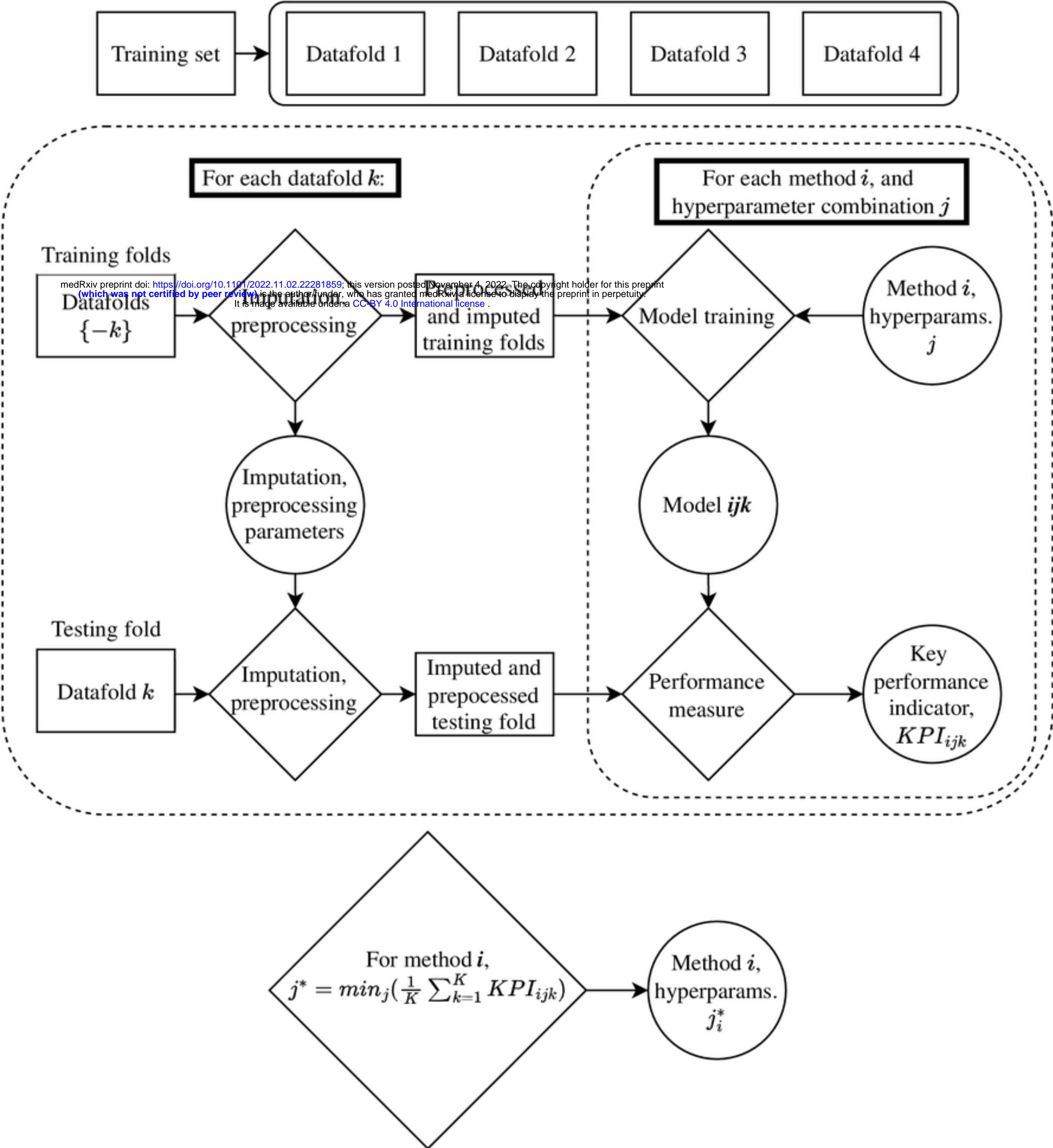


Fig 5

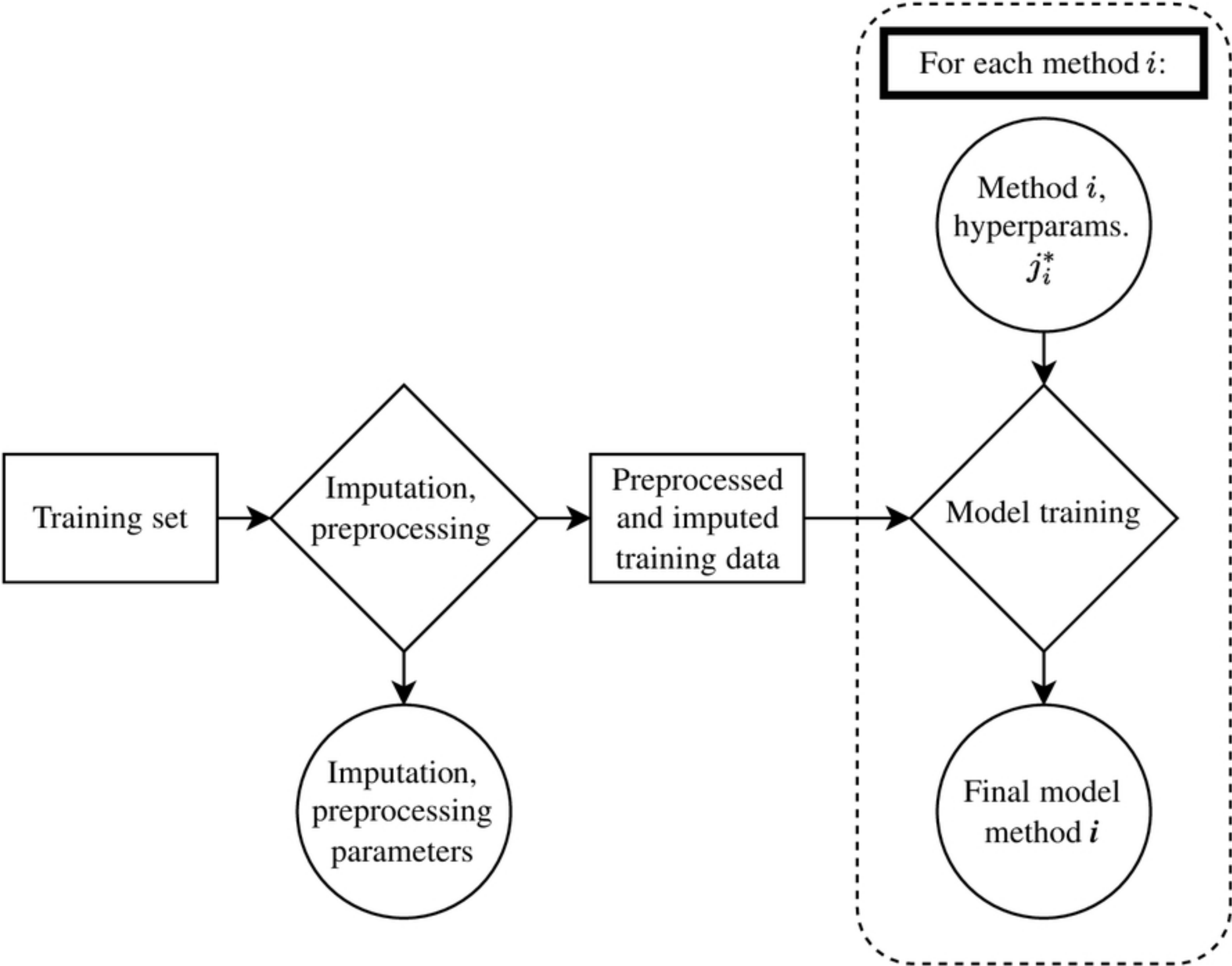


Fig 6

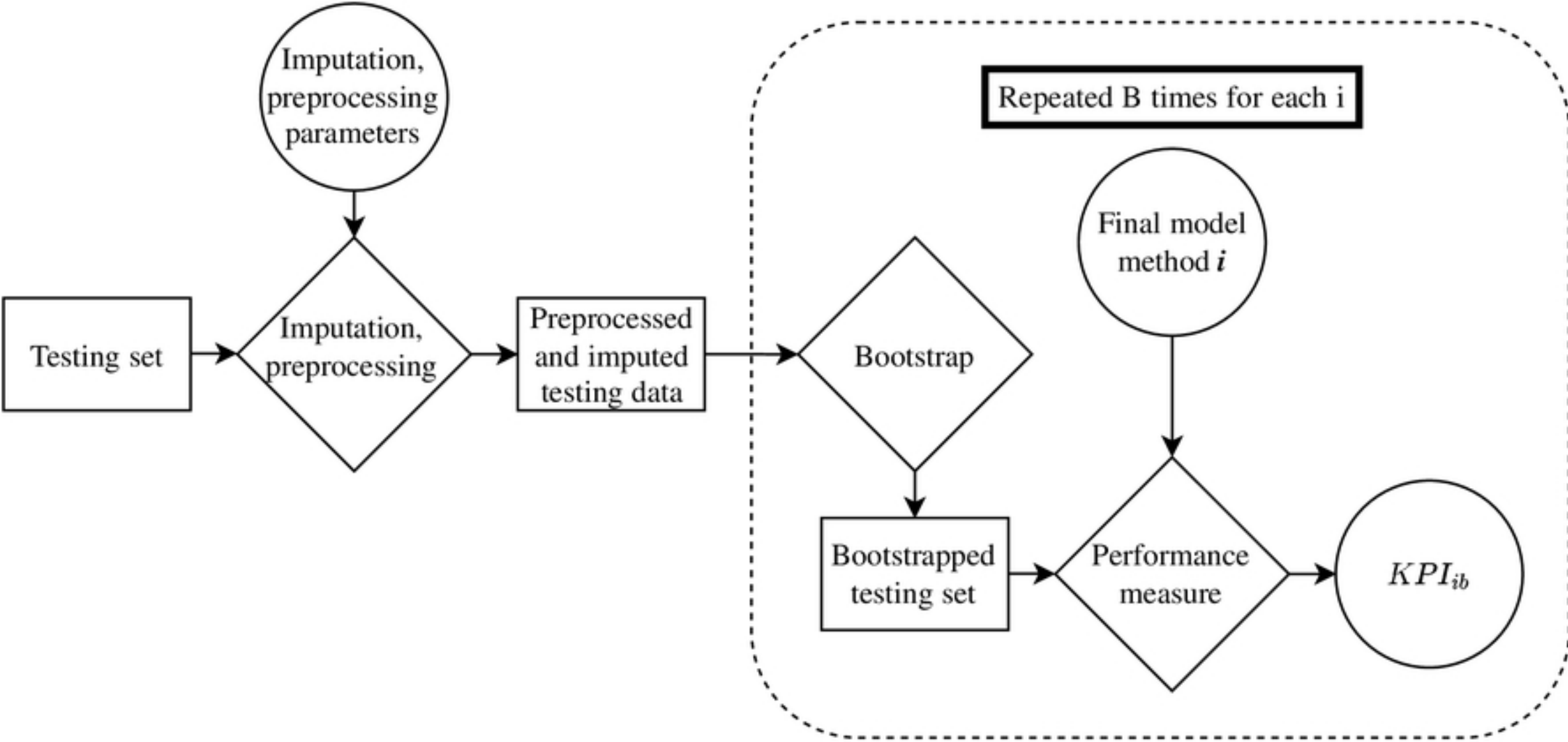


Fig 7

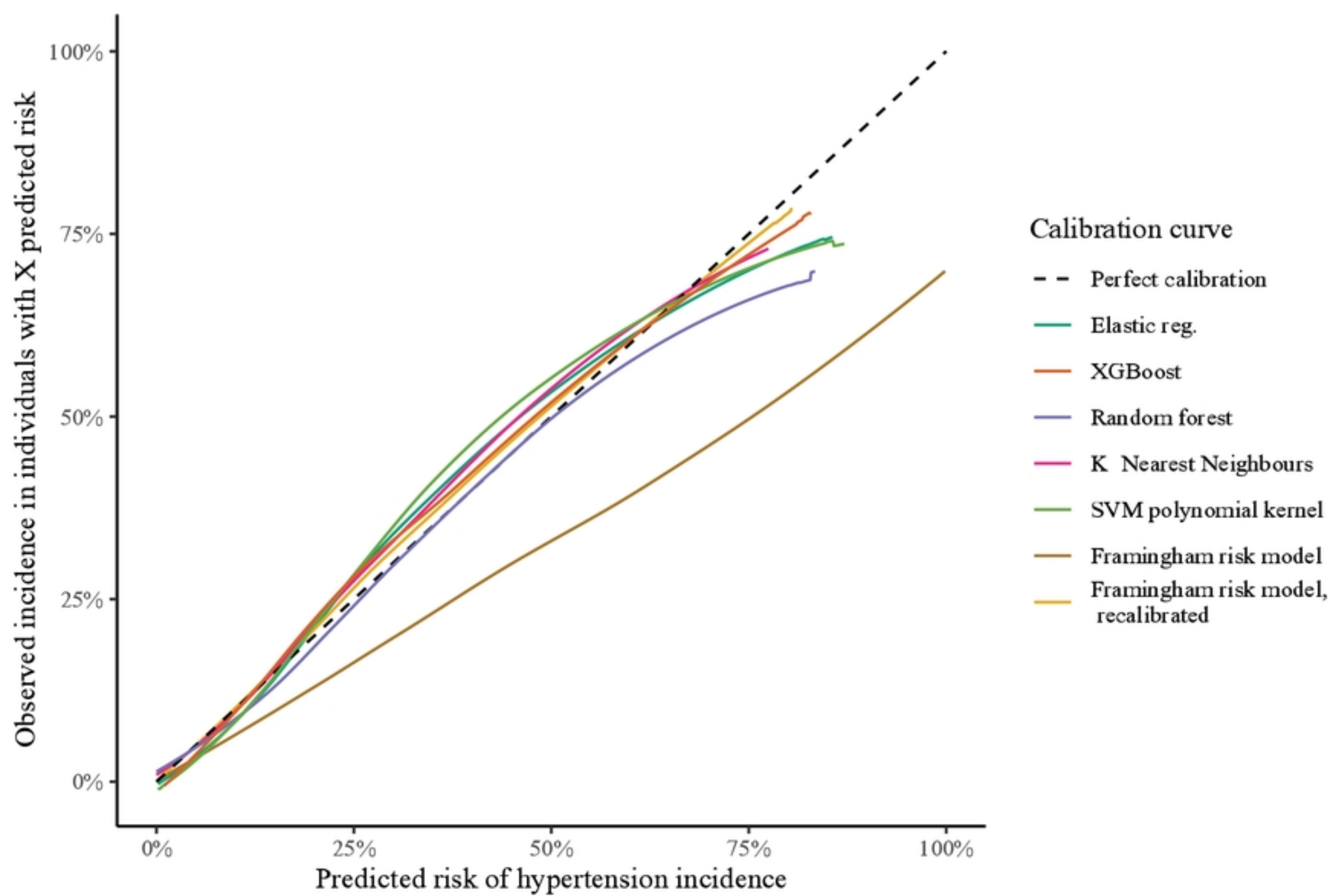


Fig 1



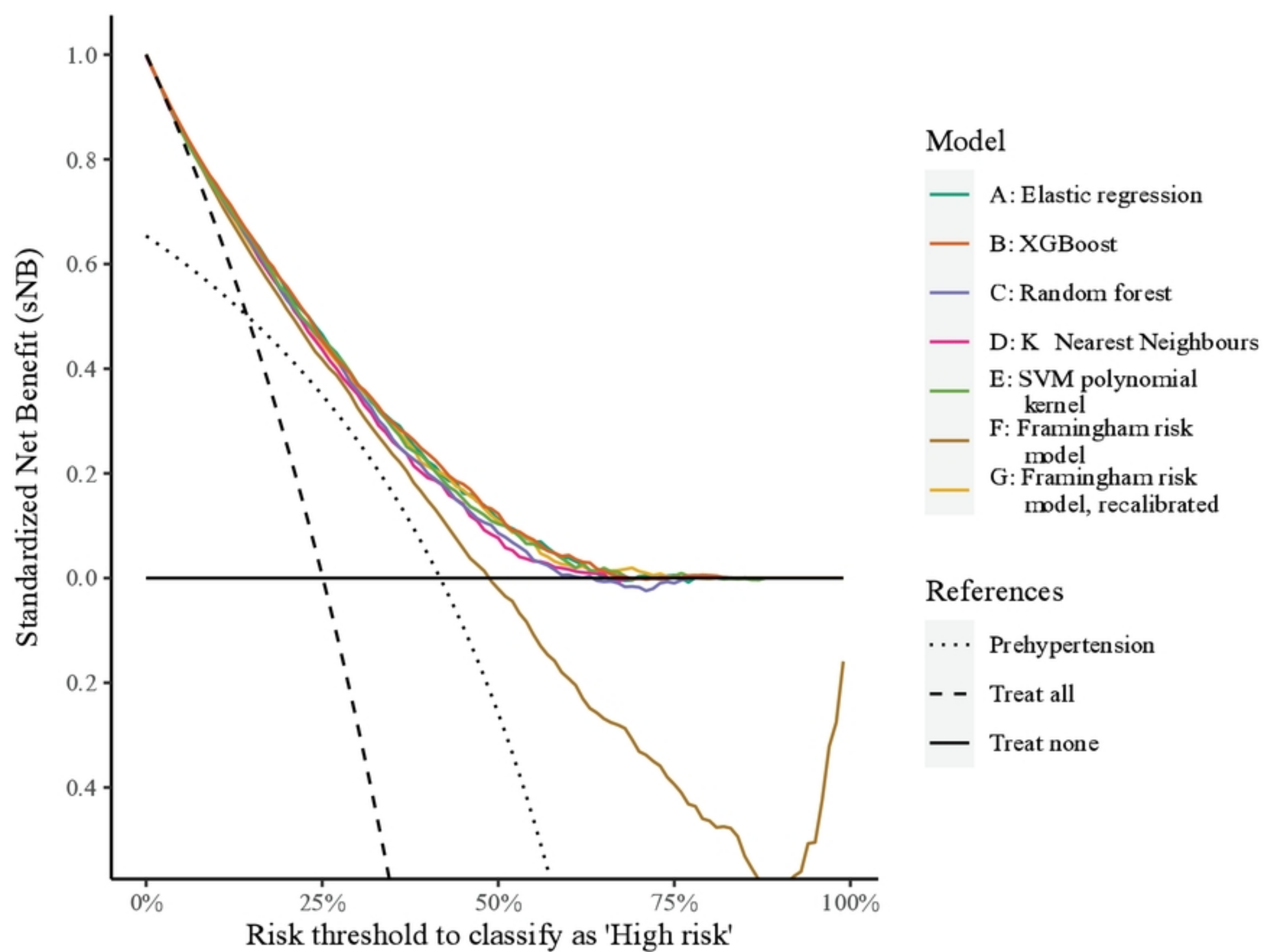


Fig 2

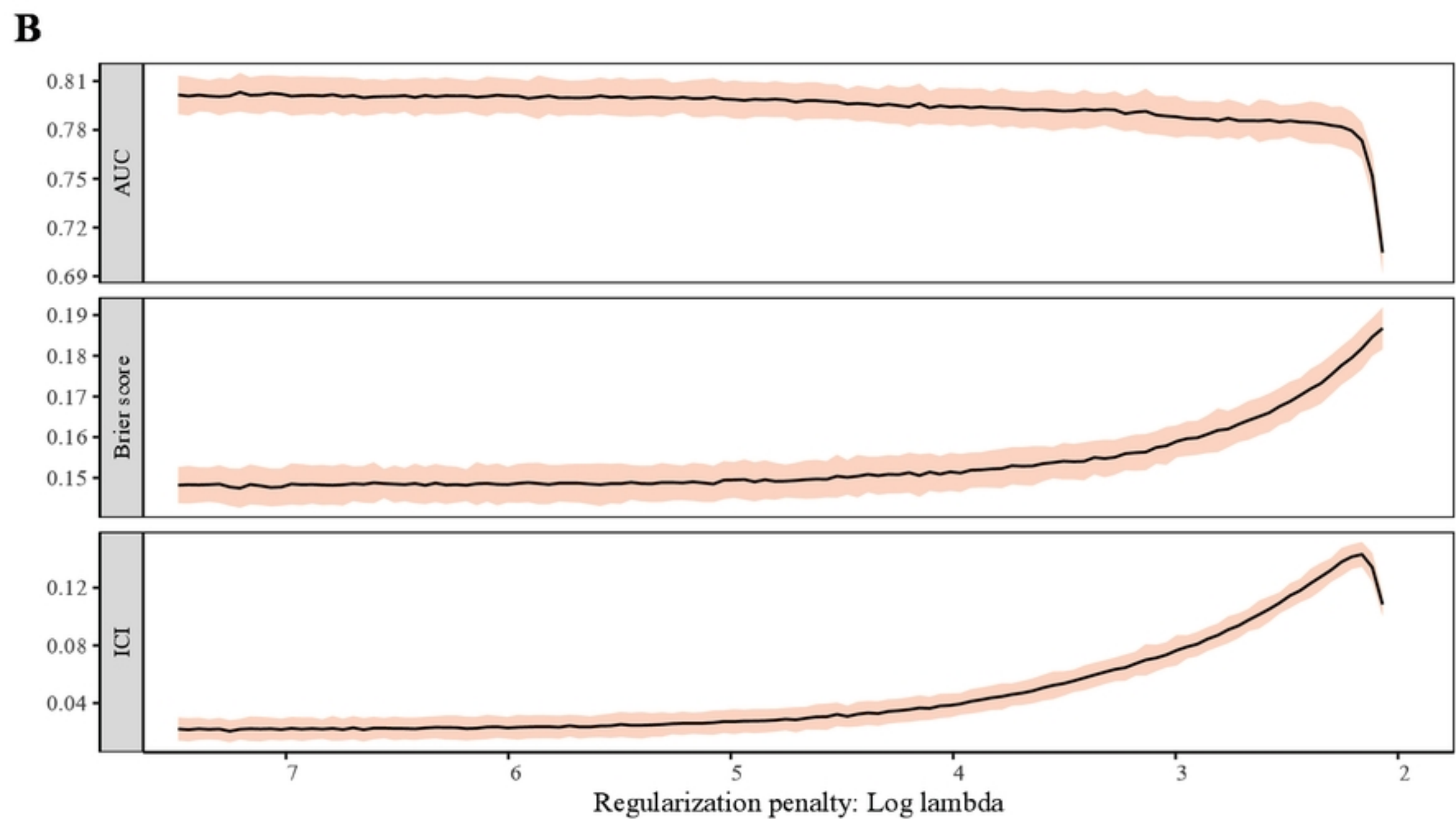
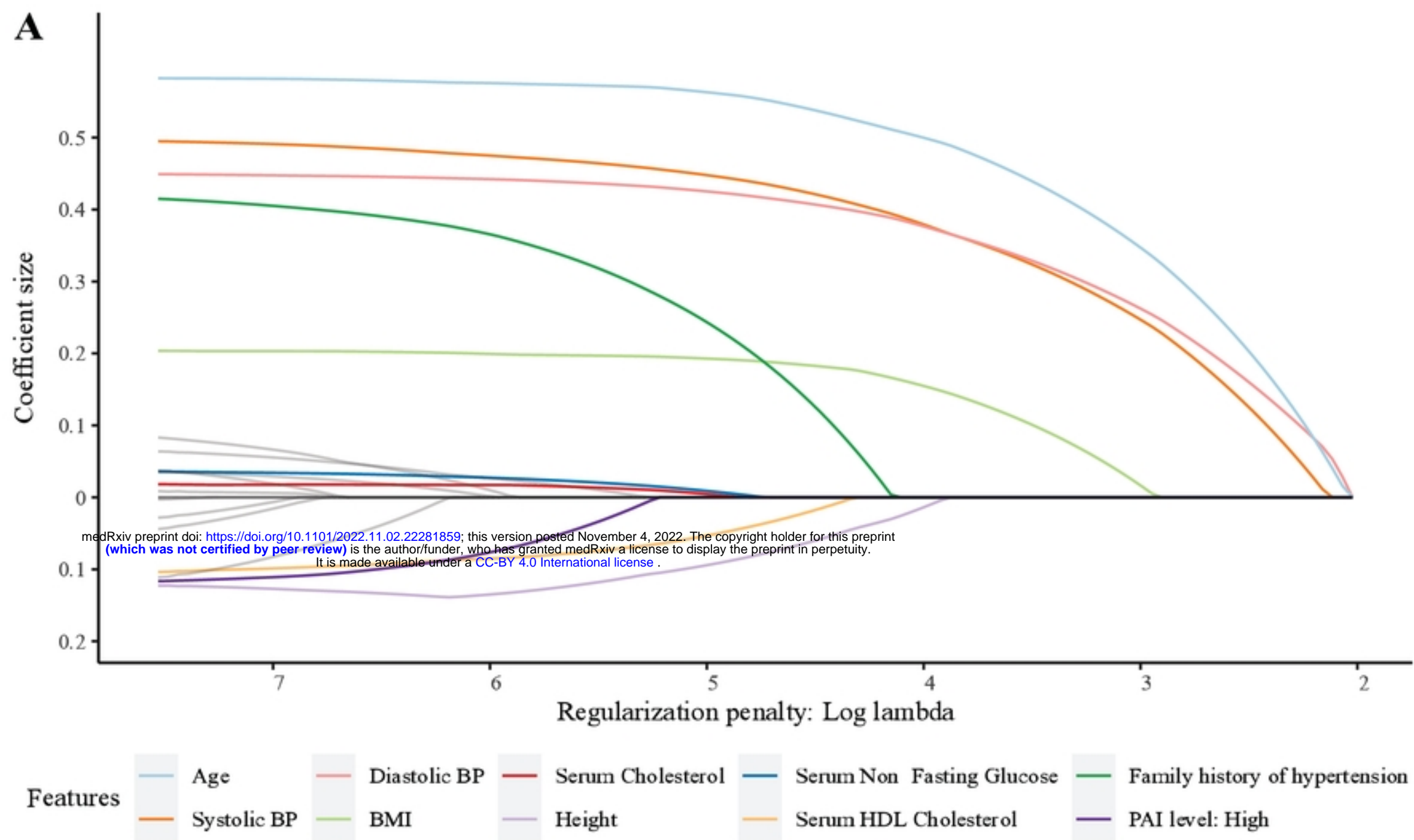


Fig 3

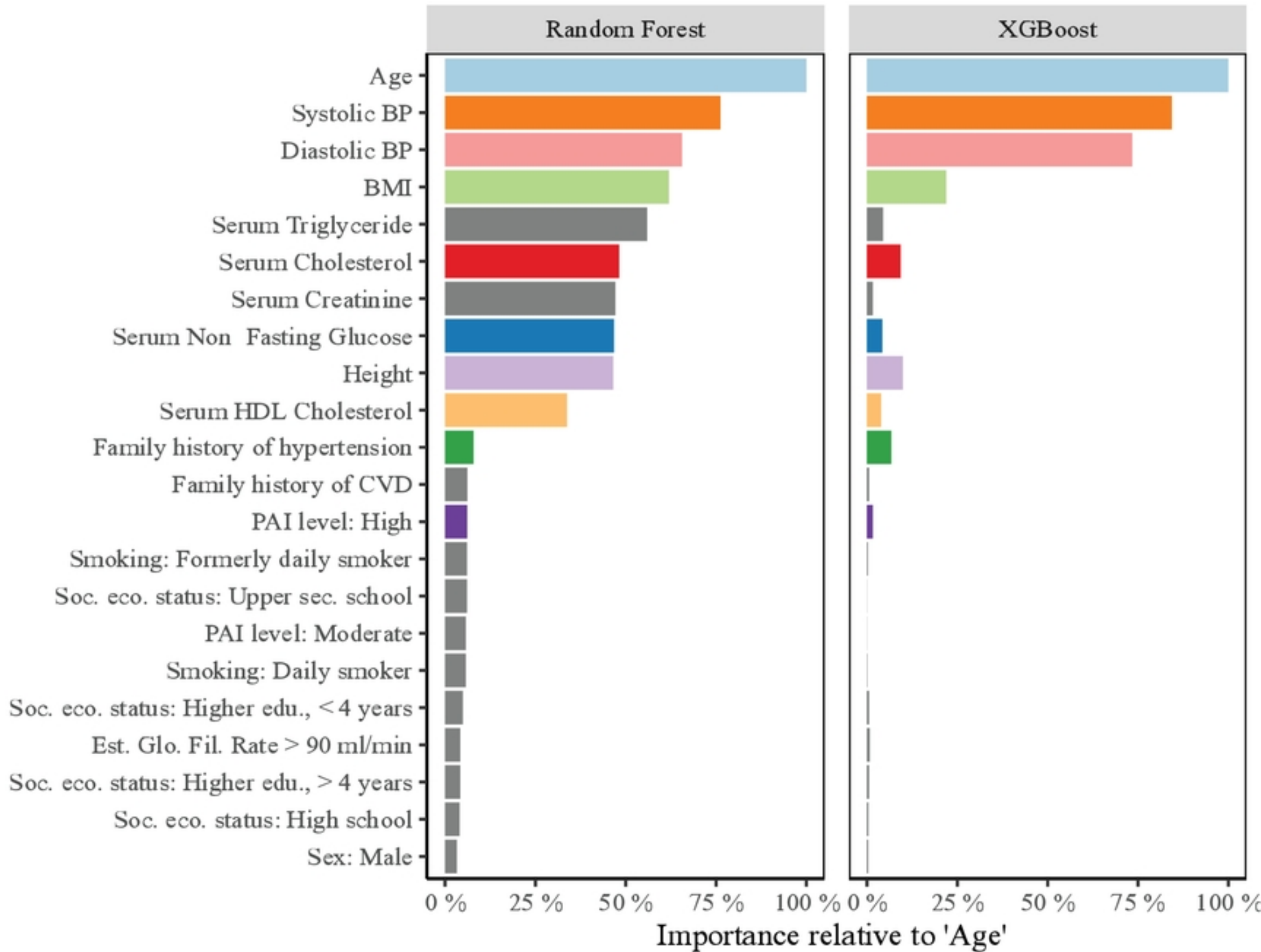


Fig 4