

1 **Predicting cardiovascular disease in patients with mental illness**
2 **using machine learning**

3
4 Martin Bernstorff^{1,2,3}, Lasse Hansen^{1,2,3}, Kevin Kris Warnakula Olesen⁴,
5 Andreas Aalkjær Danielsen^{1,2}, Søren Dinesen Østergaard^{1,2}

6
7 ¹ Department of Affective Disorders, Aarhus University Hospital – Psychiatry, Aarhus, Denmark

8 ² Department of Clinical Medicine, Aarhus University, Aarhus, Denmark

9 ³ Center for Humanities Computing, Aarhus University, Denmark

10 ⁴ Department of Cardiology, Aarhus University Hospital, Aarhus, Denmark.
11
12
13
14

15 **Corresponding author**

16 Martin Bernstorff, MD
17 Department of Affective Disorders
18 Aarhus University Hospital - Psychiatry
19 Palle Juul-Jensens Boulevard 175
20 8200 Aarhus N
21 Denmark
22 E-mail: manber@rm.dk
23 Telephone: +45 4142 6636
24
25
26
27

This peer-reviewed article has been accepted for publication but not yet copyedited or typeset, and so may be subject to change during the production process. The article is considered published and may be cited using its DOI.

This is an Open Access article, distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is unaltered and is properly cited. The written permission of Cambridge University Press must be obtained for commercial re-use or in order to create a derivative work.

28 **Abstract**

29

30 **Background:** Cardiovascular disease (CVD) is twice as prevalent among individuals with
31 mental illness compared to the general population. Prevention strategies exist but require
32 accurate risk prediction. This study aimed to develop and validate a machine learning model
33 for predicting incident CVD among patients with mental illness using routine clinical data
34 from electronic health records.

35

36 **Methods:** A cohort study was conducted using data from 74,880 patients with 1.6 million
37 psychiatric service contacts in the Central Denmark Region from 2013 to 2021. Two machine
38 learning models (XGBoost and regularized logistic regression) were trained on 85% of the
39 data from 6 hospitals using 234 potential predictors. The best performing model was
40 externally validated on the remaining 15% of patients from another 3 hospitals. CVD was
41 defined as myocardial infarction, stroke, or peripheral arterial disease.

42

43 **Results:** The best-performing model (hyperparameter-tuned XGBoost) demonstrated
44 acceptable discrimination, with an area under the receiver operating characteristic curve of
45 0.84 on the training set and 0.74 on the validation set. It identified high-risk individuals 2.5
46 years before CVD events. For the psychiatric service contacts in the top 5% of predicted risk,
47 the positive predictive value was 5%, and the negative predictive value was 99%. The model
48 issued at least one positive prediction for 39% of patients who developed CVD.

49

50 **Conclusions:** A machine learning model can accurately predict CVD risk among patients
51 with mental illness using routinely collected electronic health record data. A decision support

52 system building on this approach may aid primary CVD prevention in this high-risk

53 population.

54

55 **Keywords:** Precision Medicine, Artificial Intelligence, Psychiatry, Cardiovascular Diseases

56

57 **Introduction**

58 CVD not only diminishes quality of life, but also contributes substantially to premature
59 mortality [1,2]. Individuals with mental illness are twice as likely to develop CVD compared
60 to the background population [3,4], and are at elevated risk of premature death due to CVD
61 [2]. This elevated risk can likely be attributed to higher prevalence of unhealthy lifestyle such
62 as poor diet, sedentary behaviour, and excessive alcohol consumption [5]. Additionally,
63 psychopharmacological treatment, antipsychotics in particular, acts as a double-edged sword
64 in the context of CVD, increasing risk due to weight gain and dysmetabolism [6], while being
65 associated with lower risk of cardiovascular disease in observational studies [7], likely via
66 beneficial effect on the underlying mental disorder.

67

68 Unfortunately, the elevated risk of CVD among those with mental illness is not reflected in
69 the administration of preventive measures, with screening for CVD occurring at 25% lower
70 rates among individuals with mental illness [3,8], and up to 88% of individuals with
71 schizophrenia with dyslipidaemia not receiving adequate treatment for the latter [9].

72 Consequently, identifying individuals with mental illness at elevated risk of CVD is a crucial
73 initial step towards implementing effective preventive strategies. However, to the best of our
74 knowledge, there is a paucity of tools designed for predicting CVD risk among patients
75 receiving treatment in psychiatric service systems.

76

77 Accurately assessing CVD risk is a multifaceted challenge. Machine learning models are
78 particularly well-suited for this task, given the presence of numerous interacting factors
79 increasing CVD risk [10], and the models' ability to capture complex relationships while
80 mitigating the impact of data idiosyncrasies [11]. Previous research has demonstrated the
81 efficacy of machine learning models in accurately predicting clinical outcomes for patients

82 with mental disorders when trained on electronic health record data. Specifically, it has been
83 possible to predict, e.g., mechanical restraint [12], progression from prediabetes to type 2
84 diabetes [13], and incidence of type 2 diabetes [14]. In line with these achievements, to aid
85 identification of patients with mental illness who may benefit from targeted intervention to
86 prevent CVD, we aimed to develop and validate a machine learning model trained on
87 electronic health record data to predict development of CVD among patients with mental
88 illness.

89

90 **Methods**

91 The methods are illustrated by panels A-I in Figure 1.

92 **Data and cohort extraction**

93 This study is based on electronic health record data from the PSYchiatric Clinical Outcome
94 Prediction (PSYCOP) cohort, which encompasses all individuals with at least one contact
95 with the Psychiatric Services of the Central Denmark Region in the period from January 1,
96 2011, and November 22, 2021. The dataset includes information from routine clinical
97 practice (i.e., there was no specific data collection for the purpose of this study) on service
98 contacts, diagnoses, medications, procedures and laboratory results from all public hospitals
99 (psychiatric as well as general hospitals) in the Central Denmark Region (Figure 1A).

100 Denmark has a tax-financed universal public healthcare system.

101

102 A flowchart illustrating the definition of the patient cohort is available as eFigure 1. For this
103 study, we restricted the cohort to patients with contacts to the Psychiatric Services of the
104 Central Denmark Region after January 1, 2013, due to data instability prior to this date
105 caused by the implementation of a new electronic health record system [15,16]. Only patients
106 aged 18 years or older were included, as the probability of developing CVD is very low in
107 those below the age of 18. Patients with known CVD, defined by meeting one of the outcome
108 criteria (see below) between January 1, 2011, and December 31, 2013, were excluded to
109 minimize issuing of predictions for prevalent cases.

110

111 **Outcome definition (cardiovascular disease)**

112 The outcome definition had three elements. First, to align with prior research, we took
113 inspiration from the outcome definition from the Systematic Coronary Risk Evaluation 2

114 (SCORE2) [17]. Specifically, we defined incident CVD as the first occurrence of a diagnosis
115 of myocardial infarction (MI) (International Classification of Diseases, 10th revision (ICD-
116 10): I21-I23 or a diagnosis of stroke (ICD-10: I6, (Figure 1B). Second, we included
117 interventions/procedures which are highly indicative of vascular disease (procedure codes are
118 available in eTable 1) to the outcome definition, namely percutaneous coronary intervention
119 (PCI), coronary artery bypass grafting (CABG), intracranial endovascular thrombolysis and
120 other intracranial endovascular surgery. Third, given the large morbidity and disability
121 burden due to peripheral arterial disease, its increasing incidence, and the potential for
122 prevention [18], we included diagnoses (ICD-10: I70.2, I73.9) and procedures (procedure
123 codes are available in eTable 1) for iliac, femoral, popliteal and distal arterial disease to the
124 outcome definition.

125 Data splitting

126 The data were divided into two subsets: a training dataset (85% of the data) and a test dataset
127 (15% of the data). Specifically, all visits to the Psychiatric Services in either the western or
128 eastern part of the Central Denmark Region (Aarhus, Gødstrup, Herning, Holstebro, Horsens
129 and Randers) were used for the training set, and the central part (Viborg, Silkeborg and
130 Skive) for the test-set (see Figure 1C). If a patient first had visits in one of the splits (i.e. the
131 training set or the test set), any subsequent visits in the other split was removed. This
132 guaranteed that no patient appeared in both the training and test datasets. After this point, the
133 test dataset was left aside and only used for the final evaluation of the best performing model
134 obtained during the training phase. This geographical split assessed the generalizability
135 across geography, e.g., to which extent the model could be applied without modification if a
136 new hospital was added to the region.

137

138 Prediction time filtering

139 We defined prediction times as the time of any in- or outpatient contact with the Psychiatric
140 Services (service contacts). Consequently, each patient could have multiple prediction times -
141 corresponding to their number of service contacts. We excluded prevalent cases by not
142 issuing a prediction if that patient had already met the CVD outcome criteria at the time of a
143 service contact (Figure 1D). Moreover, no prediction was made if the lookbehind window
144 (the time used for extracting predictors) included time before follow-up started on January 1,
145 2013 or if the lookahead window (the time within which to detect the outcome) of 2 years
146 extended beyond the end of follow-up, the date of moving out of the Central Denmark
147 Region, or the patient's death. These "truncations" are artifacts caused by data collection. If
148 not accounted for, they could cause the model to learn patterns that do not exist during
149 implementation, leading to discrepancies between the model's test performance and actual
150 implemented performance. In the case of a patient moving into the region, we did not issue
151 predictions for two years after the move, mirroring the wash-in for existing patients.

152

153 Predictor grouping and flattening

154 Predictors were chosen based on a recent meta-analysis of prediction models for CVD in non-
155 psychiatric settings and included demographics, laboratory results, diagnoses, antipsychotics,
156 and mood stabilizers [19]. Specifically, the following predictors were included, all
157 operationalized using routine clinical electronic health record data from the Central Denmark
158 Region: age, sex, smoking status, high- and low-density lipoprotein (HDL and LDL),
159 haemoglobin A1c (HbA1c), systolic blood pressure, diagnosis of chronic lung disease (ICD-
160 10: J40-J44*), diagnoses from all psychiatric subchapters individually (F0-F9), as well as use
161 of any one of the top 10 weight gaining antipsychotics during inpatient treatment
162 (Anatomical Therapeutic Chemical classification codes in parentheses): clozapine

163 (N05AH02), zotepine (N05AX11), olanzapine (N05AH03), sertindole (N05AE03),
164 chlorpromazine (N05AA01), iloperidone (N05AX14), quetiapine (N05AH04), paliperidone
165 (N05AX13), trifluoperazine (N05AB06), and risperidone (N05AX08), resulting in 26 eligible
166 features (Figure 1E) [20,21]. These predictors were aggregated over the lookbehind windows
167 (90, 365 and 730) days, to incorporate different temporal contexts, and with different
168 aggregation methods (min, mean, max) using the timeseriesflattener python package [22],
169 resulting in a total of 234 potential predictors (Figure 1F). For further elaboration, see the
170 Supplementary Material.

171

172 The dataset includes numerous predictors lacking values within the lookbehind window.
173 However, these absent values do not constitute missing data in the conventional sense, as
174 they are not a result of omitted data entry. Instead, the absence of data reflects the reality of
175 clinical practice. Since this absence aligns with the data available for implementation,
176 patients exhibiting such an absence should be retained in the dataset. During model training,
177 these absent values are either passed on directly (XGBoost) or imputed using the population
178 median (logistic regression).

179

180 Predictor addition by early stopping

181 The predictors were rank ordered into eight layers (see eTable 2). Models were trained
182 incrementally, adding layers until discrimination stabilized ($\Delta\text{AUROC} < 0.01$) for the last
183 two layers. The best-performing layer with the fewest features was further refined by
184 incorporating additional aggregation methods (min, max, mean) and lookbehind windows
185 (90, 365, 730 days). See the Supplementary Material for further details.

186

187 Model selection and hyperparameter tuning

188 We focused on two models: XGBoost and elastic net regularised logistic regression, due to
189 the large number of possible model configurations (Figure 1G). XGBoost was selected for its,
190 fast training, and ability to handle numerical, categorical, and missing values internally, and
191 due to the fact that gradient boosting methods generally outperform other machine learning
192 approaches on tabular data [23,24]. As simpler models are more interpretable and easier to
193 implement, logistic regression with elastic net regularisation was included as a benchmark
194 model. Logistic regression requires missing value imputation as part of pre-processing, and
195 we imputed using the median. For the elastic net penalisation to not be affected by predictor
196 units, we Z-score standardised all predictors for the logistic regression. All predictors listed
197 under “Predictor grouping and flattening” were considered for the XGBoost and elastic net
198 regularised logistic regression. As a sensitivity analysis, we trained an elastic net regularised
199 logistic regression using only predictors that mimic those from SCORE2 as closely as
200 possible with the available data (see Supplementary Table 1 for the specific predictors). All
201 models were trained using 5-fold cross-validation, with hyperparameter optimisation to
202 maximise the area under the receiver operating characteristic curve (AUROC) using the tree-
203 structured Parzen estimator algorithm in Optuna v2.10.1 (Figure 1H). Additional details,
204 including which hyperparameters were explored, are provided in the Supplementary Material.

205

206 Model evaluation

207 The model that achieved the best AUROC on the training dataset was evaluated on the
208 geographically independent (external) test dataset (Figure 1I). Performance metrics, including
209 AUROC, sensitivity, specificity, positive predictive value, and negative predictive value,
210 were calculated. Since healthcare systems are limited by available resources, and can
211 accommodate different amounts of interventions, performance metrics were calculated for
212 different predicted positive rates [25]. The predicted positive rate is the proportion of all

213 prediction times which are marked as "positive". The mean time from the first positive
214 prediction until a patient met the definition of CVD was also determined. Predictor
215 importance was estimated using information gain.

216

217 Robustness analyses

218 The stability of model prediction was assessed across patient sex, age, as well as time from
219 first visit, and month of year.

220 Post-hoc analyses

221 A model using the best performing hyperparameters was re-evaluated on a random split of
222 the entire dataset. All patients were randomly allocated (85%-15%) to either the training
223 (85%) or test set (15%), ensuring no patient overlap between the splits. This analysis assessed
224 the performance in the case where all application-sites are included in the training data.

225

226 Ethics

227 The use of electronic health record data for this study was approved by the Legal Office of
228 the Central Denmark Region in accordance with the Danish Health Care Act §46, Section 2.
229 According to the Danish Committee Act, ethical review board approval is not required for
230 studies based solely on data from electronic health records (waiver for this project: 1-10-72-
231 1-22). Data were processed and stored in accordance with the European Union General Data
232 Protection Regulation and the project is registered on the internal list of research projects
233 having the Central Denmark Region as data steward.

234

235 Data and code sharing

236 The code for all analyses is available on GitHub: <https://github.com/Aarhus-Psychiatry->
237 [Research/psycop-](https://github.com/Aarhus-Psychiatry-Research/psycop-common/tree/7cc7ad912e638957e983a1af2a6df0f474aa6345/psycop/projects/t2d)
238 [common/tree/7cc7ad912e638957e983a1af2a6df0f474aa6345/psycop/projects/t2d](https://github.com/Aarhus-Psychiatry-Research/psycop-common/tree/7cc7ad912e638957e983a1af2a6df0f474aa6345/psycop/projects/t2d)
239

240 **Results**

241 The eligible cohort consisted of 27,954 patients with a total of 364,791 psychiatric service
242 contacts (prediction times). Demographic and clinical information on the cohort is reported in
243 Table 1. Patients in the train- and test data were broadly similar, with median ages of 35.2
244 and 35.9 years, and proportions of females of 54.9% and 58.0%, respectively. Among the
245 27,954 patients, 524 (2.0%) experienced a CVD event. The incidence of CVD was slightly
246 higher in the test data compared to the training data (2.2% vs. 1.8%). The incidence of CVD
247 spiked around the end of the wash-out period, after which it declined (eFigure 2). For each
248 hpredictor, the proportion of prediction times using the fallback value is described in eTable
249 3.

250 Figure 2A presents the results of the model training. The XGBoost model using only
251 predictor layers 1+2 (sex, age, LDL, systolic blood pressure, smoking (pack-years) and
252 smoking (daily/occasionally/prior/never) achieved an AUROC of 0.84 (95% CI: 0.83;
253 0.84). Incorporating additional lookbehinds or aggregation methods did not enhance
254 model performance. Furthermore, the inclusion of further predictor layers did not
255 increase the AUROC materially or statistically significantly (see eTable 4). The SCORE2-
256 like elastic net regularised logistic regression model performed comparably, with an
257 AUROC of 0.83 (95% CI: 0.83; 0.83).

258 Figure 2B shows the results for the XGBoost model with a 5-year lookahead window applied
259 to the test data. It achieved an AUROC of 0.74 (95% CI: 0.73; 0.75). Figure 2C shows the
260 resulting confusion matrix at a predicted positive rate of 5% with a positive predictive value
261 of 5% and a negative predictive value of 99%, reflecting that for every twenty positive
262 predictions, one prediction was followed by CVD within 5 years. At this predicted positive
263 rate, the sensitivity at the level of prediction times (contacts to the Psychiatric Services) was

264 19%, and 39% of all patients who developed CVD were predicted positive at least once
265 (Table 2). Figure 2C shows that, for patients experiencing a CVD event, the model's
266 probability of flagging them as positive (high risk) increases as the prediction time
267 approaches the CVD event. Figure 2D shows the time from a patient's first positive
268 prediction until they experienced the CVD event. The model marked patients as being at high
269 risk an average of 1.4 years before the CVD event.

270

271 Supplementary Table 3 lists prediction by information gain for the best-performing XGBoost
272 model (layers 1+2). The most important predictor was age, followed by smoking
273 (daily/occasionally/prior/never), sex, systolic blood pressure, smoking (pack-years), and
274 LDL-cholesterol.

275 Figure 3 highlights that the model was stable across sex, age, and month of year. When
276 calculating model performance within specific age bins, it dropped markedly, which is
277 expected given the relative importance of increasing age for prediction. Model
278 performance also dropped somewhat for patients having been in the system for longer,
279 perhaps indicating a decreasing predictor-sampling-frequency over time (most
280 diagnostic workup in the initial hospital contacts).

281 Post-hoc analyses

282 When training (85% split) and evaluating (15% split) the model on a random split of the
283 entire dataset, it obtained an AUROC of 0.84 on the test data, identical to the cross-validated
284 performance in the training data.

285

286 **Discussion**

287 In this study, we explored the feasibility of developing a machine learning model trained on
288 routine clinical data from electronic health records to predict the development of CVD in
289 patients with mental illness. An XGBoost model based only on layers 1+2 (sex, age, LDL,
290 systolic blood pressure, smoking (pack-years) and smoking (daily/occasionally/prior/never)
291 achieved an AUROC of 0.74 in the test set at the level of individual service contacts, with a
292 PPV of 5% and an NPV of 99%. For patients who developed CVD and were identified by the
293 model, the median time from initial positive prediction to CVD diagnosis was 1.4 years. This
294 relatively simple model, in which the predictors overlap substantially with those from
295 SCORE2, offers easy implementation in psychiatric services with less comprehensive
296 electronic health record systems [26]. Notably, in spite of the theoretical improvements
297 stemming from the use of machine learning, logistic regression with elastic net penalisation
298 performed as well as the more complex XGBoost. This implies that, for prediction of CVD
299 with a well-established aetiology, simpler models may be sufficient.

300 A substantial decline in model performance was observed when evaluating on the test
301 set (from an AUROC of 0.84 during cross-validation on the training set to an AUROC of
302 0.74 on the test set). Of note, the training and test sets comprised data from different
303 psychiatric hospitals within the Psychiatric Services of the Central Denmark region. This
304 suggests that substantial distribution shifts can occur even within a relatively
305 homogeneous population sharing geographical proximity, healthcare infrastructure,
306 and clinical protocols, which is further supported by the relative lack of performance
307 difference between training and test when performing a random split of the data (from
308 an AUROC of 0.84 during cross-validation on the training set to an AUROC of 0.84 on the
309 test set). These shifts may be due to variations in patient demographics and/or in data

310 collection between hospitals – despite geographical proximity. More broadly, this lends
311 credence to the argument that external validation should not be considered an absolute
312 prerequisite for scientific publication or model evaluation. Instead, it is proposed that
313 models should undergo rigorous testing within the specific population which they are
314 targeting [27].

315 Adding information on psychiatric diagnoses by subchapter and antipsychotics (predictor
316 layer 4) did not improve predictive performance. We hypothesise that this is either due to the
317 relatively crude granularity with which these predictors were included, or that their effects are
318 mediated by predictors were already included in the model (e.g. LDL, systolic blood
319 pressure, HbA1c). If diagnoses and antipsychotics affect CVD risk mostly through these
320 variables, they will add no further information. Moreover, the use of antipsychotics results in
321 better treatment of the underlying disease, perhaps resulting in more health-promoting
322 behaviour. In observational studies, antipsychotic use is associated with a lower risk of
323 cardiovascular mortality [7].

324

325 To our knowledge, this is the first study to predict the onset of CVD specifically in patients
326 with mental illness based on routine clinical EHR data from psychiatric services.

327 Consequently, comparisons can only be made to studies from other settings/populations.

328 Osborn et al. trained a CVD prediction model specifically for patients with severe mental
329 illness in a primary care setting, including diagnoses and use of antipsychotics as potential
330 predictors [10]. The final model (PRIMROSE) was based on age, gender, height, weight,
331 systolic blood pressure, diabetes, smoking, body mass index, lipid profile, social deprivation,
332 severe mental illness diagnosis, prescriptions of antidepressants, antipsychotics, and reports
333 of heavy alcohol use. It achieved a C-statistic of 0.78, compared to 0.76 of the Framingham
334 risk score (including weights from age, sex, current smoking, total cholesterol, HDL

335 cholesterol, systolic blood pressure, and blood pressure medications). Quadackers et al.
336 compared multiple model's absolute risk estimates for psychiatric inpatient populations,
337 namely SCORE (blood pressure, age, sex, smoking, total cholesterol, and geographical
338 region), the Framingham risk score and PRIMROSE (described above) [28]. They found very
339 low agreement between the methods, with the Framingham risk score estimating risks 5-10
340 times higher than SCORE, arguing that it overestimates risk because the risk of CVD was
341 higher at the time of model development than it is now. This indicates the need for re-
342 calibrating models if they are used in markedly different populations than those in which they
343 were developed – one example being patients with mental illness.

344

345 Outside the context of patients with mental illness/psychiatric services, a recent meta-analysis
346 found 16 studies comparing machine-learning models to traditional statistical models for
347 prediction of CVD [19]. In aggregate, the point estimate of the machine-learning methods
348 was marginally better, with a C-statistic of 0.77 (0.74-0.81) vs. 0.76 (0.73-0.79) for
349 traditional statistical models. However, they also find that their implementation is rare and
350 uncertain, arguing that “the impact of missing or unavailable variables and different baseline
351 characteristics on model performance when applied cross-institutionally is unclear”. Indeed,
352 implementing a model based on research cohorts can be challenging, because information on
353 predictors is often not collected as part of routine clinical care, and/or the model assumes that
354 all predictors are available at the time(s) of prediction. we intentionally used only readily
355 available routine clinical data from electronic health records.

356

357 If the model developed in this study were to be implemented in the Psychiatric Services of
358 the Central Denmark Region, positive CVD predictions could be automatically presented to
359 healthcare staff via the EHR system, enabling them to initiate appropriate interventions at the

360 level of the individual patient. The specific interventions will depend on the situation. As a
361 first step, more information should typically be gathered, including blood pressure, and a full
362 cardiovascular risk profile. Based on these measurements, patients should be treated
363 according to guidelines [29]. Notably, lifestyle interventions do not appear to be cost-
364 effective in this population, with a large randomised trial of patients with schizophrenia
365 finding no effect [30,31], and a meta-analysis of trials finding only a clinically insignificant
366 change to BMI (-0.63 kg/m²) [32]. Pharmacological interventions, such as statins and
367 antihypertensive drugs, may be more successful, as they require smaller changes to daily life.
368 Another candidate, smoking cessation medication (e.g. bupropion), is as effective among
369 patients with severe mental illness as in the general population, but underutilised [29,33].

370

371 There are limitations to this study that should be considered by the reader. First, prevalent
372 cases of CVD can be misclassified as incident, leading to a false spike in incidence at the
373 beginning of the follow-up period. We mitigated this by employing a 2-year wash-in period.
374 We found that, for most CVD events, incidence was decreasing after the wash-in period.
375 There are multiple potential reasons for this finding. Specifically, it may reflect a true drop in
376 incidence as studies show decreasing incidence rates of CVD in Denmark, but these drops are
377 insufficient to fully explain the trend [34,35]. As such, it cannot be ruled out that some part of
378 the events we detect are prevalent cases. This is, however, unlikely to cause harm to patients,
379 as prevalent cases also need prevention of further events, but it may have inflated the
380 prediction estimates. Second, this study does not address potential effects of implementing
381 the developed model. When prediction models are implemented, they should affect
382 behaviour, for example by inducing further testing or treatment. Specifically, implementing a
383 CVD prediction model would likely induce more relevant LDL- and blood-pressure
384 measurements. These model-induced measurements should improve the next prediction

385 issued by the model, meaning that predictions following a positive prediction are likely less
386 accurate in the present dataset than they would be following implementation. Third, many
387 important variables for CVD, such as physical activity, dietary habits, or waist circumference,
388 are not collected with sufficient regularity as part of current clinical practice and could not be
389 included in the model. If they had been available, the model would likely perform with
390 greater accuracy. Fourth, since, most patients who experienced an event in the test set had a
391 stroke (71.6%) the model is less likely to generalise to cohorts where stroke is less prevalent.
392 However, given that the important features for the model are very general CVD features, we
393 would expect meaningful generalisation. Finally, machine learning models vary markedly in
394 their generalisability. We used routine clinical data from a system with universal healthcare
395 and observed performance differences between departments within the same regional
396 Psychiatric Services. Therefore, direct transfer of the model to other healthcare system would
397 probably yield suboptimal predictions. However, the approach is likely to be generalisable,
398 and retraining the model on data from other settings using the same architecture may allow
399 for transferability.

400

401 In conclusion, a machine learning model trained on routine clinical data from electronic
402 health records can predict development of CVD among patients with mental illness at a level
403 that may make clinical implementation as a decision support tool feasible. Specifically, the
404 model may help clinicians identifying which patients will benefit from primary preventative
405 initiatives. Moving forward, we see two main tasks arising from this work. First, we will
406 work towards testing the feasibility of implementing the model as a clinical decision support
407 tool in the Psychiatric Services of the Central Denmark Region. Second, as we believe the
408 model may hold potential for broader application, we aim to conduct external validation in
409 independent samples.

410

411 **Acknowledgement Section**

412

413 Author contributions

414 The study was conceptualized and designed by MD, KKWO, AAD and SDØ. The coding and
415 statistical analyses were carried out by MB with assistance from LH. All authors contributed
416 to the interpretation of the results. MB wrote the first draft of the manuscript, which was
417 subsequently revised for important intellectual content by the remaining authors. All authors
418 approved the final version of the manuscript prior to submission.

419

420 The authors thank Bettina Nørreremark from Aarhus University Hospital – Psychiatry for
421 assistance with extraction of data and Mathias Brønd Sørensen from the Business Intelligence
422 Office, Central Denmark Region, for assistance in building the infrastructure for model
423 training.

424

425 **Data availability**

426 According to Danish law, the personally sensitive data used in this study is only available for
427 research projects conducted by employees in the Central Denmark Region following approval
428 from the Legal Office under the Central Denmark Region (in accordance with the Danish
429 Health Care Act §46, Section 2).

430

431 **Funding**

432 The study is supported by grants from the Lundbeck Foundation (grant number: R344-2020-
433 1073), the Danish Cancer Society (grant number: R283-A16461), the Central Denmark
434 Region Fund for Strengthening of Health Science (grant number: 1-36-72-4-20), the Danish
435 Agency for Digitisation Investment Fund for New Technologies (grant number 2020-6720),
436 and by CAAIR; Collaboration on Applied AI Research and Digital Innovation in Healthcare.
437 SDØ reports further funding from the Lundbeck Foundation (grant number: R358-2020-
438 2341), the Novo Nordisk Foundation (grant number: NNF20SA0062874) and Independent
439 Research Fund Denmark (grant numbers: 7016-00048B and 2096-00055A). KKWO is
440 supported by a grant from the Danish Cardiovascular Academy (grant no. CPD5Y-2022001-
441 HF), which is funded by the Danish Heart Association and the Novo Nordisk Foundation.
442 The funders played no role in study design, collection, analysis or interpretation of data, the
443 writing of the report or the decision to submit the paper for publication.

444

445 **Conflicts of interest**

446 Danielsen has received a speaker honorarium from Otsuka Pharmaceutical. SDØ received the
447 2020 Lundbeck Foundation Young Investigator Prize. SDØ owns/has owned units of mutual
448 funds with stock tickers DKIGI, IAIMWC, SPIC25KL and WEKAFKI, and owns/has owned
449 units of exchange traded funds with stock tickers BATE, TRET, QDV5, QDVH, QDVE,
450 SADM, IQQH, USPY, EXH2, 2B76, IS4S, OM3X and EUNL. The remaining authors report
451 no conflicts of interest.

452

453 **References**

454

455 1. Roth Gregory A., Mensah George A., Johnson Catherine O., Addolorato Giovanni,
456 Ammirati Enrico, Baddour Larry M., et al. Global Burden of Cardiovascular Diseases and
457 Risk Factors, 1990–2019. *J Am Coll Cardiol.* 2020 Dec 22;76(25):2982–3021.

458 2. Erlangsen A, Andersen PK, Toender A, Laursen TM, Nordentoft M, Canudas-Romo V.
459 Cause-specific life-years lost in people with mental disorders: A nationwide, register-
460 based cohort study. *Lancet Psychiatry.* 2017 Dec;4(12):937–45.

461 3. Solmi M, Fiedorowicz J, Poddighe L, Delogu M, Miola A, Høye A, et al. Disparities in
462 Screening and Treatment of Cardiovascular Diseases in Patients With Mental Disorders
463 Across the World: Systematic Review and Meta-Analysis of 47 Observational Studies.
464 *Am J Psychiatry.* 2021 Sep;178(9):793–803.

465 4. Rødevand L, Steen NE, Elvsåshagen T, Quintana DS, Reponen EJ, Mørch RH, et al.
466 Cardiovascular risk remains high in schizophrenia with modest improvements in bipolar
467 disorder during past decade. *Acta Psychiatr Scand.* 2019 Apr;139(4):348–60.

468 5. Scott D, Happell B. The High Prevalence of Poor Physical Health and Unhealthy Lifestyle
469 Behaviours in Individuals with Severe Mental Illness. *Issues Ment Health Nurs.* 2011
470 Aug 19;32(9):589–97.

471 6. Rohde C, Köhler-Forsberg O, Nierenberg AA, Østergaard SD. Pharmacological treatment
472 of bipolar disorder and risk of diabetes mellitus: A nationwide study of 30,451 patients.
473 *Bipolar Disord.* 2023 Feb 8;

474 7. Taipale H, Tanskanen A, Mehtälä J, Vattulainen P, Correll CU, Tiihonen J. 20-year follow-
475 up study of physical morbidity and mortality in relationship to antipsychotic treatment
476 in a nationwide cohort of 62,250 patients with schizophrenia (FIN20). *World*
477 *Psychiatry.* 2020;19(1):61–8.

478 8. Mitchell AJ, Delaffon V, Vancampfort D, Correll CU, Hert MD. Guideline concordant
479 monitoring of metabolic risk in people treated with antipsychotic medication:
480 systematic review and meta-analysis of screening practices. *Psychol Med.* 2012
481 Jan;42(1):125–47.

482 9. Nasrallah HA, Meyer JM, Goff DC, McEvoy JP, Davis SM, Stroup TS, et al. Low rates of
483 treatment for hypertension, dyslipidemia and diabetes in schizophrenia: Data from the
484 CATIE schizophrenia trial sample at baseline. *Schizophr Res.* 2006 Sep;86(1–3):15–22.

485 10. Osborn DPJ, Hardoon S, Omar RZ, Holt RIG, King M, Larsen J, et al. Cardiovascular Risk
486 Prediction Models for People With Severe Mental Illness: Results From the Prediction
487 and Management of Cardiovascular Risk in People With Severe Mental Illnesses
488 (PRIMROSE) Research Program. *JAMA Psychiatry.* 2015 Feb 1;72(2):143–51.

- 489 11. Song X, Mitnitski A, Cox J, Rockwood K. Comparison of Machine Learning Techniques
490 with Classical Statistical Models in Predicting Health Outcomes. *MEDINFO 2004*.
491 2004;736–40.
- 492 12. Danielsen AA, Fenger MHJ, Østergaard SD, Nielbo KL, Mors O. Predicting mechanical
493 restraint of psychiatric inpatients by applying machine learning on electronic health
494 data. *Acta Psychiatr Scand*. 2019 Aug;140(2):147–57.
- 495 13. Cahn A, Shoshan A, Sagiv T, Yesharim R, Goshen R, Shalev V, et al. Prediction of
496 progression from pre-diabetes to diabetes: Development and validation of a machine
497 learning model. *Diabetes Metab Res Rev*. 2020 Feb;36(2):e3252.
- 498 14. Bernstorff M, Hansen L, Enevoldsen K, Damgaard J, Hæstrup F, Perfalk E, et al.
499 Development and validation of a machine learning model for prediction of type 2
500 diabetes in patients with mental illness. *Acta Psychiatr Scand*. 2024 Apr 4;acps.13687.
- 501 15. Hansen L, Enevoldsen K, Bernstorff M, Perfalk E, Danielsen AA, Nielbo KL, et al. Lexical
502 stability of psychiatric clinical notes from electronic health records over a decade. *Acta*
503 *Neuropsychiatr*. 2023 Aug 25;1–11.
- 504 16. Bernstorff M, Hansen L, Perfalk E, Danielsen AA, Østergaard SD. Stability of diagnostic
505 coding of psychiatric outpatient visits across the transition from the second to the third
506 version of the Danish National Patient Registry. *Acta Psychiatr Scand*. 2022;146(3):272–
507 83.
- 508 17. SCORE2 working group and ESC Cardiovascular risk collaboration, Hageman S, Pennells
509 L, Ojeda F, Kaptoge S, Kuulasmaa K, et al. SCORE2 risk prediction algorithms: New
510 models to estimate 10-Year risk of cardiovascular disease in Europe. *Eur Heart J*. 2021
511 Jul 1;42(25):2439–54.
- 512 18. Kim MS, Hwang J, Yon DK, Lee SW, Jung SY, Park S, et al. Global burden of peripheral
513 artery disease and its risk factors, 1990–2019: a systematic analysis for the Global
514 Burden of Disease Study 2019. *Lancet Glob Health*. 2023 Oct 1;11(10):e1553–65.
- 515 19. Liu W, Laranjo L, Klimis H, Chiang J, Yue J, Marschner S, et al. Machine-learning versus
516 traditional approaches for atherosclerotic cardiovascular risk prognostication in
517 primary prevention cohorts: a systematic review and meta-analysis. *Eur Heart J - Qual*
518 *Care Clin Outcomes*. 2023 Jun 1;9(4):310–22.
- 519 20. Rotella F, Cassioli E, Calderani E, Lazzeretti L, Ragghianti B, Ricca V, et al. Long-term
520 metabolic and cardiovascular effects of antipsychotic drugs. A meta-analysis of
521 randomized controlled trials. *Eur Neuropsychopharmacol*. 2020 Mar 1;32:56–65.
- 522 21. WHOCC - ATC/DDD Index [Internet]. [cited 2023 Apr 12]. Available from:
523 https://www.whooc.no/atc_ddd_index/
- 524 22. Bernstorff M, Enevoldsen K, Damgaard J, Danielsen A, Hansen L. timeseriesflattener: A
525 Python package for summarizing features from (medical) time series. *J Open Source*
526 *Softw*. 2023 Mar 29;8(83):5197.

- 527 23. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: Proceedings of the
528 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining
529 [Internet]. 2016 [cited 2023 Feb 17]. p. 785–94. Available from:
530 <http://arxiv.org/abs/1603.02754>
- 531 24. Grinsztajn L, Oyallon E, Varoquaux G. Why do tree-based models still outperform deep
532 learning on tabular data? [Internet]. arXiv; 2022 [cited 2023 Feb 17]. Available from:
533 <http://arxiv.org/abs/2207.08815>
- 534 25. Wornow M, Xu Y, Thapa R, Patel B, Steinberg E, Fleming S, et al. The Shaky Foundations
535 of Clinical Foundation Models:
- 536 26. Sculley D, Holt G, Golovin D, Davydov E, Phillips T, Ebner D, et al. Machine Learning: The
537 High Interest Credit Card of Technical Debt. In: SE4ML: Software Engineering for
538 Machine Learning (NIPS 2014 Workshop). 2014.
- 539 27. Sperrin M, Riley RD, Collins GS, Martin GP. Targeted validation: validating clinical
540 prediction models in their intended population and setting. *Diagn Progn Res*. 2022 Dec
541 22;6(1):24.
- 542 28. Quadackers D, Liemburg E, Bos F, Doornbos B, Risselada A, PHAMOUS investigators, et
543 al. Cardiovascular risk assessment methods yield unequal risk predictions: a large
544 cross-sectional study in psychiatric secondary care outpatients. *BMC Psychiatry*. 2023
545 Jul 24;23(1):536.
- 546 29. Polcwiartek C, O’Gallagher K, Friedman DJ, Correll CU, Solmi M, Jensen SE, et al. Severe
547 mental illness: cardiovascular risk assessment and management. *Eur Heart J*. 2024 Mar
548 27;45(12):987–97.
- 549 30. Speyer H, Christian Brix Nørgaard H, Birk M, Karlsen M, Storch Jakobsen A, Pedersen K,
550 et al. The CHANGE trial: no superiority of lifestyle coaching plus care coordination plus
551 treatment as usual compared to treatment as usual alone in reducing risk of
552 cardiovascular disease in adults with schizophrenia spectrum disorders and abdominal
553 obesity. *World Psychiatry Off J World Psychiatr Assoc WPA*. 2016 Jun;15(2):155–65.
- 554 31. Jakobsen AS, Speyer H, Nørgaard HCB, Karlsen M, Birk M, Hjorthøj C, et al. Effect of
555 lifestyle coaching versus care coordination versus treatment as usual in people with
556 severe mental illness and overweight: Two-years follow-up of the randomized CHANGE
557 trial. *PLOS ONE*. 2017 Oct 6;12(10):e0185881.
- 558 32. Speyer H, Jakobsen AS, Westergaard C, Nørgaard HCB, Pisinger C, Krogh J, et al. Lifestyle
559 Interventions for Weight Management in People with Serious Mental Illness: A
560 Systematic Review with Meta-Analysis, Trial Sequential Analysis, and Meta-Regression
561 Analysis Exploring the Mediators and Moderators of Treatment Effects. *Psychother
562 Psychosom*. 2019 Sep 13;88(6):350–62.

- 563 33. Tsoi DT yin, Porwal M, Webster AC. Efficacy and safety of bupropion for smoking
564 cessation and reduction in schizophrenia: systematic review and meta-analysis. *Br J*
565 *Psychiatry*. 2010 May;196(5):346–53.
- 566 34. Skajaa N, Adelborg K, Horváth-Puhó E, Rothman KJ, Henderson VW, Casper Thygesen L,
567 et al. Nationwide Trends in Incidence and Mortality of Stroke Among Younger and
568 Older Adults in Denmark. *Neurology*. 2021 Mar 30;96(13):e1711–23.
- 569 35. Schmidt M, Andersen LV, Friis S, Juel K, Gislason G. Data Resource Profile: Danish Heart
570 Statistics. *Int J Epidemiol*. 2017 Oct 1;46(5):1368–1369g.

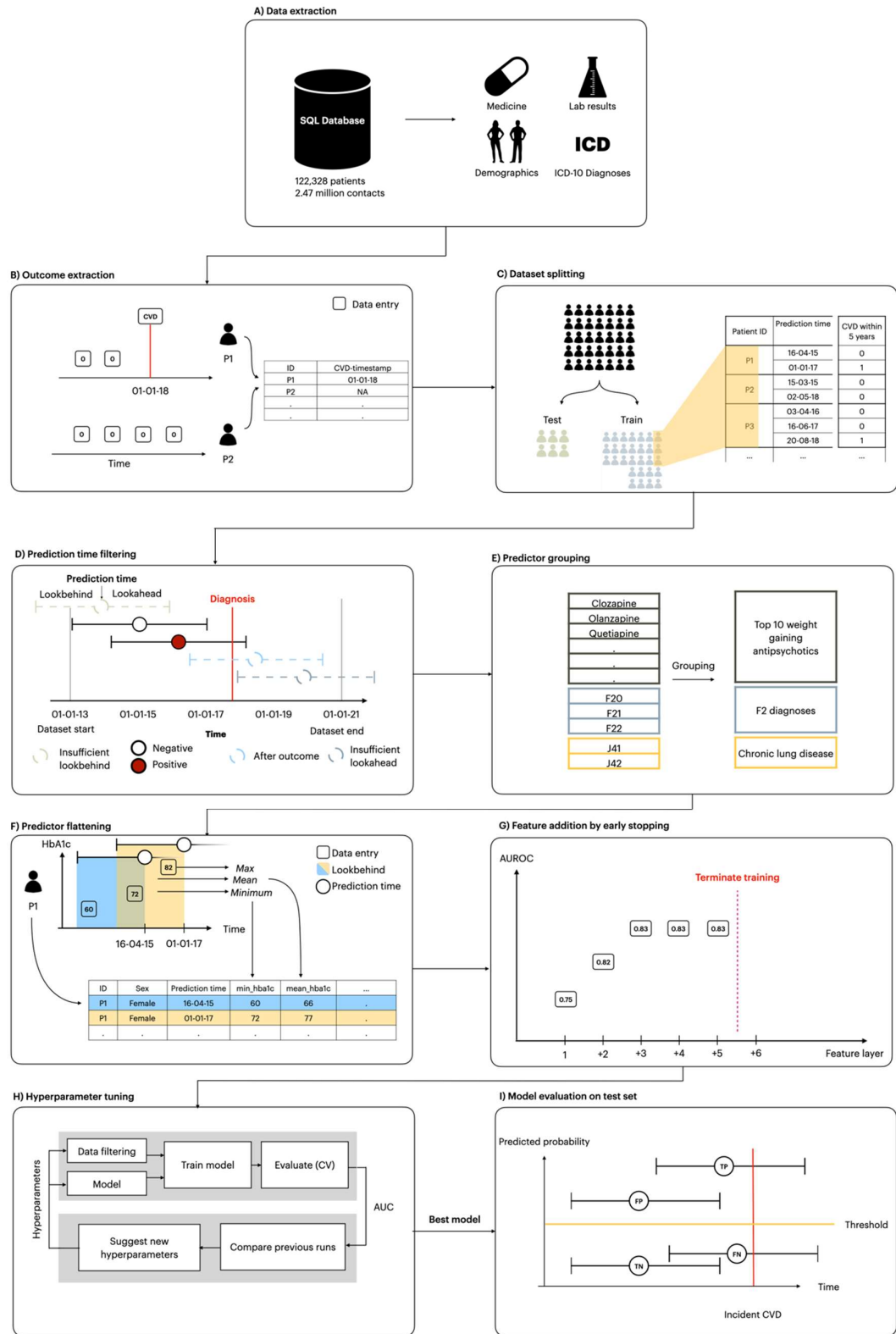
571

572

573

574
575
576

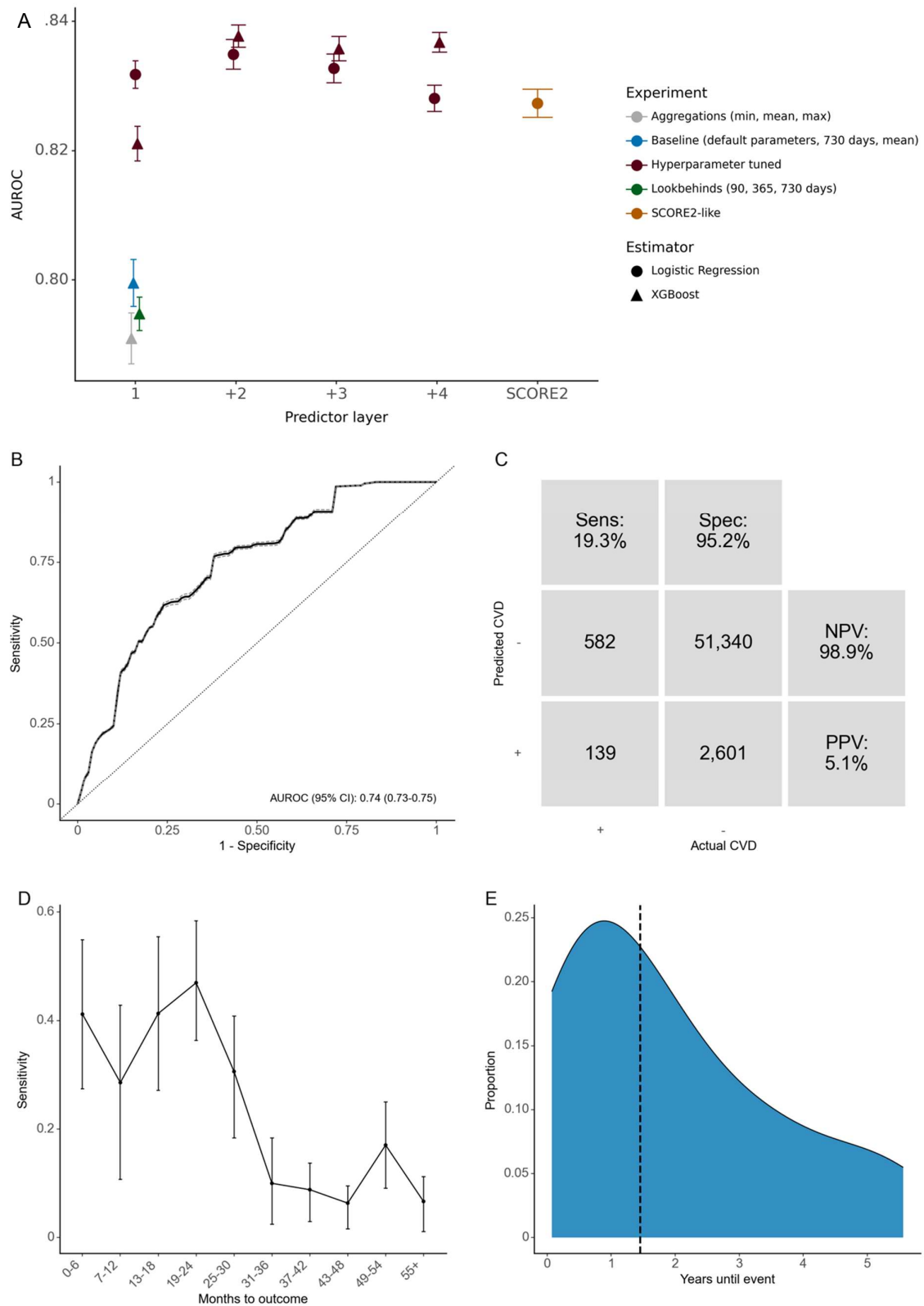
Figure 1: Extraction of data and outcome, dataset splitting, prediction time filtering, specification of predictors and flattening, model training, testing and evaluation



577
578

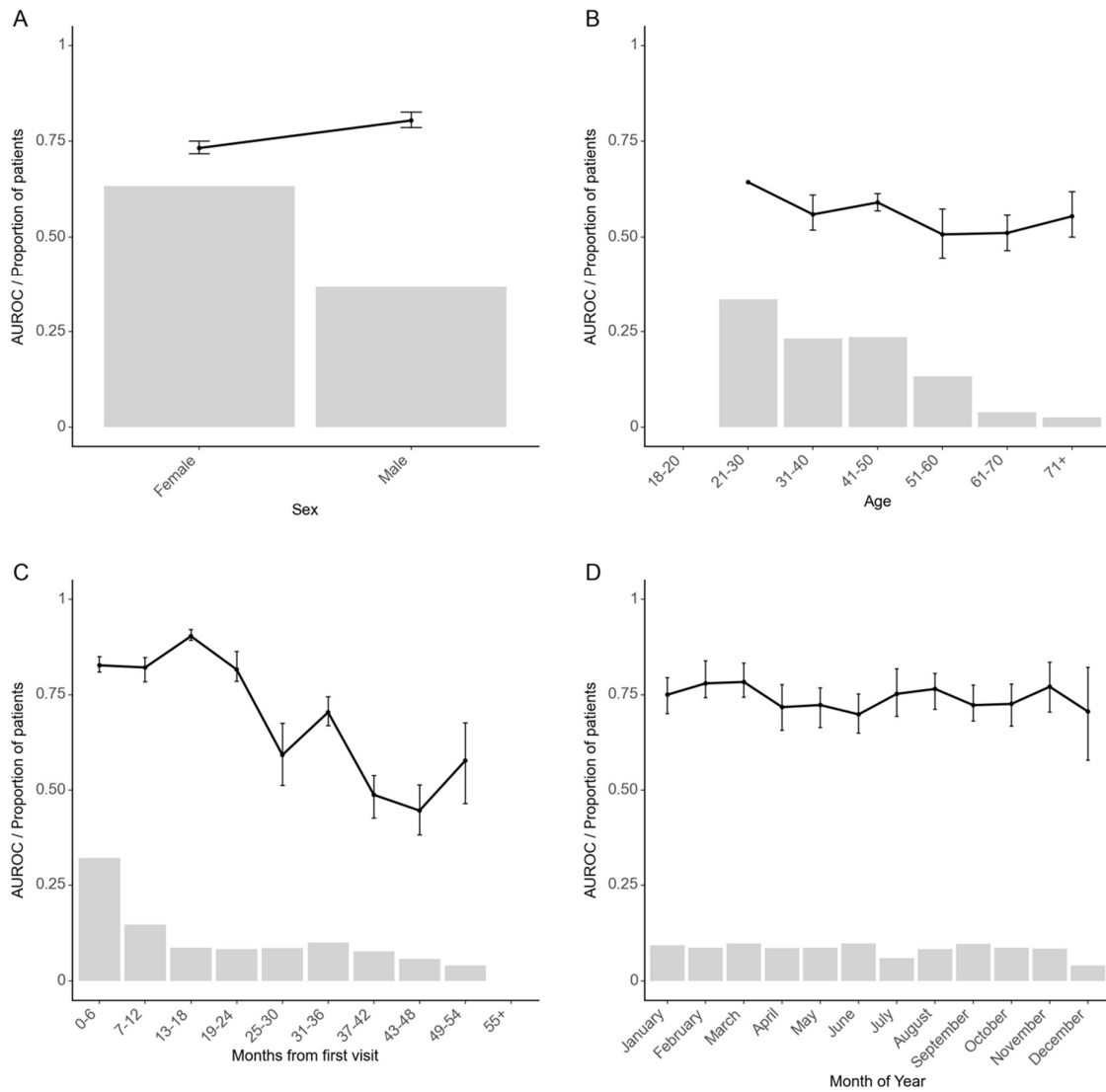
579 A: Data was extracted from the electronic health records
580 B: Potential CVD was identified
581 C: The dataset obtained is split geographically into an independent training dataset (85%) and test dataset (15%) with no
582 patient being present in both groups.
583 D: Prediction times were removed if their lookbehind window extended beyond the start of the dataset or their lookahead
584 extended beyond the end of the dataset. Prediction times were also removed after a patient developed CVD.
585 E: Predictors were grouped.
586 F: Predictors for each prediction time were extracted by aggregating the variables within the lookbehind with multiple
587 aggregation functions. As a result, each row in the dataset represents a specific prediction time with a column for each
588 predictor.
589 G: Predictor layers were added until model performance no longer improved.
590 H: Models were trained and optimized on the training set using 5-fold cross-validation. Hyperparameters were tuned to
591 optimize AUROC.
592 I: The best candidate model was evaluated on the independent test set. True positive predictions were those with predicted
593 probabilities above the decision threshold and the patient having a CVD event within the lookahead window. False positive
594 predictions were those where the model's predicted probability was above the decision threshold, but the patient did not have
595 a CVD event within the lookahead window. False negatives had predicted probabilities below the threshold, but the patient
596 had a CVD event within the lookahead window. True negatives had predicted probabilities below the threshold, and the
597 patient did not have a CVD event within the lookahead window.
598

Figure 2. Results from model training of all models (A) and on geographically independent (external/test) data (B-E)



A) Results of experiments across aggregation methods (mean vs. min, mean and max), lookbehinds (730 days vs. 90, 365 and 730 days), predictor layers (1, +2, +3, +4) and hyperparameter tuning. Note that results for each layer also includes the features of the prior layers. **B)** Receiver operating characteristics (ROC) curve. **C)** Confusion matrix. PPV: Positive predictive value. NPV: Negative predictive value. **D)** Sensitivity by months from prediction time to event, stratified by desired predicted positive rate (PPR). Note that the numbers do not match those in Table 1, since all prediction times with insufficient lookahead distance have been dropped. **E)** Time (months) from the first positive prediction to the patient developing CVD at a 5% predicted positive rate (PPR).

Figure 3. Robustness of the best performing model on geographically independent (external/test) data



Robustness of the model across stratifications. The line is the area under the receiver operating characteristics curve. Bars represent the proportion of prediction times in each bin. Error bars are 95%-confidence intervals from 100-fold bootstrap.

Table 1. Descriptive statistics for service contacts (A) and patients (B) that were eligible for prediction.**A. Service contacts**

	Train	Test
Service contacts, n	310127	54664
Demographics		
Age, median [Q1,Q3]	35.2 [25.9,46.7]	35.9 [25.1,47.3]
Female, n (%)	185681 (59.9)	34579 (63.3)
Smoking (pack-years), mean (SD)	30.5 (75.3)	25.1 (92.8)
Smoking (daily/occasionally/prior/never), median [Q1,Q3]	2.0 [1.0,4.0]	3.0 [1.0,4.0]
BMI, median [Q1,Q3]	25.6 [22.1,30.2]	25.7 [22.0,30.2]
Height (cm), median [Q1,Q3]	171.0 [165.0,178.5]	170.8 [165.0,178.0]
Weight (kg), median [Q1,Q3]	77.0 [64.5,91.4]	76.5 [63.9,91.2]
Diagnoses		
Angina, n (%)	2355 (0.8)	355 (0.6)
Atrial fibrillation, n (%)	1822 (0.6)	453 (0.8)
Chronic kidney failure, n (%)	805 (0.3)	149 (0.3)
Chronic lung disease, n (%)	2307 (0.7)	819 (1.5)
F0 - Organic disorders, n (%)	8357 (2.7)	1245 (2.3)
F1 - Substance abuse, n (%)	32767 (10.6)	4387 (8.0)
F2 - Psychotic disorders, n (%)	49889 (16.1)	6171 (11.3)
F3 - Mood disorders, n (%)	115999 (37.4)	20048 (36.7)
F4 - Neurotic and stress-related, n (%)	94095 (30.3)	13865 (25.4)
F5 - Eating and sleeping disorders, n (%)	13689 (4.4)	2068 (3.8)
F6 - Personality disorders, n (%)	47249 (15.2)	7185 (13.1)
F7 - Mental retardation, n (%)	5778 (1.9)	320 (0.6)
F8 - Developmental disorders, n (%)	9584 (3.1)	1687 (3.1)
F9 - Child and adolescent disorders, n (%)	45151 (14.6)	11018 (20.2)
Type 1 diabetes, n (%)	1865 (0.6)	308 (0.6)
Type 2 diabetes, n (%)	6291 (2.0)	1009 (1.8)
Lab results		
HDL, mean (SD)	1.4 (0.4)	1.4 (0.4)
HbA1c, mean (SD)	35.7 (7.0)	35.2 (6.9)
LDL, mean (SD)	2.9 (0.9)	2.9 (0.9)
Systolic blood pressure, median [Q1,Q3]	126.8 [117.5,137.8]	125.2 [117.0,136.0]

Total cholesterol, mean (SD)		4.9 (1.0)	4.8 (1.0)
Medications			
Antihypertensives, n (%)		692 (0.2)	70 (0.1)
Top 10 weight gaining antipsychotics, n (%)		74900 (24.2)	10709 (19.6)
Outcomes			
Incident CVD, n (%)		2885 (0.9)	721 (1.3)
By subtype, n (group-%)	CABG	15 (0.5)	8 (1.0)
	MI	608 (18.8)	75 (9.3)
	PAD	82 (2.5)	70 (8.7)
	PCI	626 (19.3)	37 (4.6)
	Stroke	1909 (58.9)	618 (76.5)

B. Patients

		Train	Test
Patients, n		23584	4370
Female, n (%)		12946 (54.9)	2535 (58.0)
Incident CVD, n (%)		430 (1.8)	94 (2.2)
By subtype, n (group-%)	CABG	6 (1.4)	<5
	MI	70 (16.1)	14 (13.7)
	PAD	13 (3.0)	8 (7.8)
	PCI	66 (15.2)	6 (5.9)
	Stroke	280 (64.4)	73 (71.6)

Cohort demographics by split after preprocessing. For filtering steps, see eFigure 1. Definitions are available in eTable 3. CVD: Cardiovascular disease. MI: Myocardial infarction. PCI: Percutaneous coronary intervention. PAD: Peripheral artery disease. CABG: Coronary artery bypass grafting. Note that < 5 is required by Danish Data Legislation.

Table 2. Performance by predicted positive rate for the best performing model (XGBoost) with 5 years of lookahead on the test set.

Predicted positive rate	True prevalence	PPV	NPV	Sensitivity	Specificity	FPR	FNR	Accuracy	TP	TN	FP	FN	% of all patients with CVD captured	Median years from first positive to CVD
1.0%	1.3%	5.6%	98.7%	1.0%	95.7%	4.3%	99.0%	97.8%	31	53,417	524	690	7.4%	2.7
5.0%		5.1%	98.9%	4.8%	80.7%	19.3%	95.2%	94.2%	139	51,340	2,601	582	39.4%	2.5
10.0%		3.3%	98.9%	9.8%	75.2%	24.8%	90.2%	89.3%	179	48,647	5,294	542	48.9%	2.6
20.0%		3.6%	99.2%	19.6%	45.6%	54.4%	80.4%	80.1%	392	43,383	10,558	329	70.2%	2.8

Predicted positive rate: The proportion of contacts predicted positive by the model. Since the model outputs a predicted probability, this is a threshold set during evaluation.

True prevalence: The proportion of contacts that qualified for CVD within the lookahead window.

PPV: Positive predictive value.

NPV: Negative predictive value.

FPR: False positive rate.

FNR: False negative rate.

TP: True positives. Numbers are service contacts.

TN: True negatives. Numbers are service contacts.

FP: False positives. Numbers are service contacts.

FN: False negatives. Numbers are service contacts.

% of all patients with CVD captured: Percentage of all patients who developed CVD, who had at least one positive prediction.

Median years from first positive to CVD: For all patients with at least one true positive, the number of years from their first positive prediction to having developed CVD.