

1 Development and validation of DNA Methylation scores in two European cohorts augment
2 10-year risk prediction of type 2 diabetes

3 Yipeng Cheng^{1,2}, Danni A Gadd¹, Christian Gieger^{3,4,5}, Karla Monterrubio-Gómez⁶, Yufei
4 Zhang¹, Imrich Berta¹, Michael J Stam⁷, Natalia Szlachetka⁷, Evgenii Lobzaev⁷, Nicola
5 Wrobel⁸, Lee Murphy⁸, Archie Campbell¹, Cliff Nangle¹, Rosie M Walker^{1,9}, Chloe Fawns-
6 Ritchie^{10,1}, Annette Peters^{4,5,11}, Wolfgang Rathmann^{12,5}, David J Porteous¹, Kathryn L Evans¹,
7 Andrew M McIntosh¹³, Timothy I Cannings¹⁴, Melanie Waldenberger^{3,4}, Andrea Ganna²,
8 Daniel L McCartney¹, Catalina A Vallejos^{6,15,*}, Riccardo E Marioni^{1,*}

9 ¹ Centre for Genomic and Experimental Medicine, Institute of Genetics and Cancer, University
10 of Edinburgh, Edinburgh, EH4 2XU, UK.

11 ² Institute for Molecular Medicine Finland, University of Helsinki, Helsinki, Finland

12 ³ Research Unit Molecular Epidemiology, Helmholtz Zentrum München, German Research
13 Center for Environmental Health (GmbH), Neuherberg, Germany

14 ⁴ Institute of Epidemiology, Helmholtz Zentrum München, German Research Center for
15 Environmental Health (GmbH), Neuherberg, Germany

16 ⁵ German Center for Diabetes Research (DZD), München-Neuherberg, Germany

17 ⁶ MRC Human Genetics Unit, Institute of Genetics and Cancer, University of Edinburgh,
18 Edinburgh, EH4 2XU, UK.

19 ⁷ School of Informatics, University of Edinburgh, Edinburgh, EH8 9AB, UK

20 ⁸ Edinburgh Clinical Research Facility, University of Edinburgh, Western General Hospital,
21 Edinburgh, UK

22 ⁹ Centre for Clinical Brain Sciences, University of Edinburgh, Chancellor's Building, 49 Little
23 France Crescent, Edinburgh BioQuarter, Edinburgh, EH16 4SB, UK

24 ¹⁰ Department of Psychology, University of Edinburgh, Edinburgh, EH8 9JZ, UK

25 ¹¹ DZHK (German Centre for Cardiovascular Research), partner site Munich Heart Alliance,
26 Munich, Germany

27 ¹² Institute for Biometrics and Epidemiology, German Diabetes Center, Leibniz Institute for
28 Diabetes Research at Heinrich Heine University, Düsseldorf, Germany

29 ¹³ Division of Psychiatry, Centre for Clinical Brain Sciences, University of Edinburgh,
30 Edinburgh, UK

31 ¹⁴ School of Mathematics, University of Edinburgh, Edinburgh, EH9 3FD, UK

32 ¹⁵ The Alan Turing Institute, London, UK

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

33

34 *Corresponding authors:

35 Names: Riccardo Marioni and Catalina Vallejos

36 Contact Details: riccardo.marioni@ed.ac.uk and catalina.vallejos@ed.ac.uk

37 **Abstract**

38 Type 2 diabetes mellitus (T2D) presents a major health and economic burden that could be
39 alleviated with improved early prediction and intervention. While standard risk factors have
40 shown good predictive performance, we show that the use of blood-based DNA methylation
41 information leads to a significant improvement in the prediction of 10-year T2D incidence risk.
42 Previous studies have been largely constrained by linear assumptions, the use of CpGs one-at-
43 a-time, and binary outcomes. We present a flexible approach (via an R package, *MethylPipeR*)
44 based on a range of linear and tree-ensemble models that incorporate time-to-event data for
45 prediction. Using the Generation Scotland cohort (training set $n_{\text{cases}}=374$, $n_{\text{controls}}=9,461$; test
46 set $n_{\text{cases}}=252$, $n_{\text{controls}}=4,526$) our best-performing model (Area Under the Curve (AUC)=0.872,
47 Precision Recall AUC (PRAUC)=0.302) showed notable improvement in 10-year onset
48 prediction beyond standard risk factors (AUC=0.839, PRAUC=0.227). Replication was
49 observed in the German-based KORA study ($n=1,451$, $n_{\text{cases}} = 142$, $p=1.6 \times 10^{-5}$).

50 **Introduction**

51 Diabetes mellitus is one of the most prevalent diseases in the world and a leading cause of
52 mortality. Around half a billion people live with diabetes worldwide, with type 2 diabetes
53 (T2D) making up about 90% of these cases [1]. Individuals with diabetes can suffer from
54 debilitating complications including nerve damage, kidney disease and blindness [2]. The
55 disease also increases the future risk of dementia and cardiovascular disease [3], with recent
56 studies highlighting obesity and T2D as risk factors for COVID-19 disease severity and ICU
57 admission [4]. Furthermore, risk of complications increases over time and is exacerbated if
58 blood-glucose levels are poorly managed. Despite developments in the way T2D can be
59 managed for patients, these treatments are reactive, focusing on patients that have already been
60 diagnosed. Early intervention with metformin or lifestyle changes have been shown to delay
61 onset of T2D, although they did not reduce the risk of all-cause mortality [5].

62 Beyond public health costs, T2D also presents a substantial financial burden to the NHS, with
63 estimated annual spending of £10 billion on diabetes in the UK. Around 80% of these costs are
64 for treatment of complications, many of which are preventable with early intervention [6].

65 While the mechanisms of insulin resistance in T2D are well-known, the interactions between
66 genetic and environmental factors that increase T2D susceptibility are less understood.
67 Previous T2D risk prediction models have used a range of health risk factors [7]. However,
68 these have not utilised epigenetic information. Epigenetics is the study of heritable changes to
69 DNA that do not modify its nucleotide sequence. A commonly studied form of this is DNA
70 methylation (DNAm), whereby methyl groups are attached to the DNA molecule - most
71 commonly to the 5-carbon on a cytosine in a cytosine-guanine pair (CpG). Due to its
72 involvement with gene expression and gene-environment interactions, DNAm can provide
73 dynamic predictive information for disease risk for an individual. For example, Epigenetic

74 Scores (EpiScores) built via penalised regression models have been used to show that weighted
75 linear CpG predictors can explain a substantial proportion of phenotypic variance (R^2) of
76 modifiable health factors including body mass index (BMI) (12.5%), HDL cholesterol (15.6%)
77 and smoking status (60.9%) [8]. Blood-based DNAm is of particular interest in predictive
78 modelling and biomarker development due to its comparatively non-invasive sampling
79 procedure. EpiScores have also shown the ability to explain up to 58% of variance in plasma
80 protein levels are associated with a number of incident diseases including T2D and several
81 comorbidities [9]. Epigenome-wide association studies (EWAS) have identified a number of
82 CpG sites significantly associated with T2D [10-14] as well as related risk factors such as
83 cardiovascular disease [15] and obesity [16, 17]. While these provide some predictive
84 performance for T2D prevalence, incident T2D has been less well studied. One such EWAS
85 with 563 cases and 701 controls identified 18 CpGs associated with incident T2D but did not
86 consider any prediction models [10]. Given that preventative lifestyle changes have been
87 shown to effectively reduce T2D onset [18], prediction of T2D incidence years ahead of time
88 would be greatly beneficial in stratifying populations so those at high risk can be monitored
89 and treated with early interventions.

90 Currently, most studies generating DNAm predictors consider marginal CpG effects or assume
91 only linear additive effects between CpGs. The use of predictive models that can incorporate
92 both interaction and non-linear effects could capture more complex relationships between
93 variables, resulting in greater prediction accuracy. Therefore our study aims to evaluate both
94 the additional predictive benefit that DNAm can provide for 10-year T2D risk and the
95 applicability of linear and tree-ensemble survival models.

96 Here, we use one of the world's largest studies with paired genome-wide DNAm and data
97 linkage to electronic health records (EHR), Generation Scotland (n=14,613, n=626 incident
98 T2D cases over 15 years of follow-up), to develop and validate epigenetic EpiScores for T2D.

99 We show the added contribution of these EpiScores to prediction over and above standard risk
100 factors e.g. age, sex and BMI and externally validate these results in the KORA S4 cohort.

101

102 **Methods**

103 *Reporting of Analysis Pipeline and Results*

104 To enhance reproducibility, the analysis pipeline and results presented in this study have been
105 reported via the TRIPOD checklist [19] (**Supplementary File 1**).

106 *Generation Scotland*

107 Blood-based DNAm and linked health data were obtained from Generation Scotland [20], a
108 family-structured population-based cohort. The cohort consists of 23,960 volunteers across
109 Scotland aged 18-99 years at recruitment (between 2006 and 2011), of whom over 18,000
110 currently have genome-wide DNAm data available (Illumina EPIC array). In DNAm quality
111 control, CpG sites were filtered by removing those with low bead count in $\geq 5\%$ of samples or
112 a high detection p-value (>0.05) in more than 5% of samples. Probes on X and Y chromosomes
113 were also removed. Samples were filtered by removing those with a mismatch between
114 predicted and recorded sex or $\geq 1\%$ of CpGs with detection p-value > 0.05 . Missing CpG values
115 were mean-imputed. To enable the predictors to be applied to existing cohort studies with older
116 Illumina array data, CpGs were filtered to the intersection of the 450k and EPIC array sites
117 (n=453,093 CpGs).

118 This study considered DNAm data from three large subsets of the GS cohort, with 5,087 (Set
119 1), 4,450 (Set 2) and 8,877 (Set 3) individuals. Processing took place in 2017, 2019 and 2021
120 respectively. Set 1 and Set 3 included related individuals within and between each set while all
121 individuals in Set 2 were unrelated to each other and to individuals in Set 1 (genetic relationship

122 matrix (GRM) threshold <0.05). In our experiments, the training set consisted of Sets 2 and 3
123 combined and Set 1 was used as the test set. To avoid the presence of families with individuals
124 across both the training and test set, any individuals in the training set from the same family as
125 an individual in the test set were excluded from the analysis ($n_{\text{excluded}}=3,138$).

126 Participant health measures including age, BMI, sex, self-reported hypertension and family
127 (parent or sibling) history of T2D were taken at baseline (DNAm sampling) via questionnaire.
128 BMI was calculated as the individual's weight in kilograms divided by the square of their height
129 in metres. Missing values in the Set 1 health measures were treated as missing-completely-at-
130 random and the corresponding individuals were excluded ($n_{\text{Set 1}}=99$). This was not performed
131 in Sets 2 and 3 as the health measures were used for incremental modelling (Set 1 only).

132 Disease cases were ascertained through data linkage to NHS Scotland health records consisting
133 of hospital (ICD-10 codes) and GP records (Read2 codes). Prevalent cases were identified from
134 a baseline questionnaire (self-reported) or from ICD-10/Read2 codes dated prior to baseline
135 and removed from the dataset. Type 1 and juvenile cases were treated as control observations.
136 A total of 757 incident cases were observed over the follow up period (from recruitment date
137 to 01/2022) and after preprocessing, 626 cases remained. Mean time-to-T2D-onset was 5.9, 5.4
138 and 6.0 years for Sets 1, 2 and 3 respectively, with ranges of 0.2 – 14.8 (Set 1), 0.2 – 14.8 (Set
139 2) and 0.1-14.8 (Set 3) years. In GP record-derived cases, 81% of cases had a C10F. “type 2
140 diabetes mellitus” code; 12% had a C10.. “diabetes mellitus” code and 4% had a C109. “Non-
141 insulin dependent diabetes mellitus” code. The full list of included and excluded terms are
142 given in **Supplementary Table 1**.

143

144 *Composite Protein EpiScore*

145 A composite protein EpiScore model for incident T2D was trained using a set of 109 protein
146 EpiScores previously shown to associate with a range of incident diseases [9]. For each protein,
147 the EpiScore was calculated for each subject in the training and test sets. A Cox proportional-
148 hazards (Cox PH) lasso model was fit to the training set with the 109 protein EpiScores (scaled
149 within set to mean of 0 and variance 1) as features. The linear predictor from the Cox PH lasso
150 model was then used as the composite protein EpiScore in the incremental test set model.

151

152 *Direct EpiScore*

153 The direct EpiScore Cox PH lasso model for incident T2D was fit to the DNAm data in the
154 training set. Due to memory limitations in the model fitting R package (glmnet [21]), the CpGs
155 were filtered to the 200,000 sites with highest variance. The linear predictor from the Cox PH
156 lasso model was then used as the direct EpiScore in the incremental test set modelling. For tree
157 ensemble models, the Cox PH lasso-selected CpGs were used as input and the 10-year onset
158 risk was subsequently used as the direct EpiScore.

159

160 *Outcome Definition for 10-Year Onset Prediction*

161 Linkage to NHS Scotland health records provided dates for disease diagnoses from which age-
162 at-onset was calculated. Along with age at baseline (DNAm sampling), these were used to
163 calculate the time-to-event, measured in years, for each individual. For incident T2D cases and
164 controls, time-to-event was defined as the time from baseline to disease onset and censoring,
165 respectively. Controls were censored at the latest date of available GP records (09/2020).

166 While all models were trained as survival models, our primary prediction outcome was incident
167 T2D diagnosis within 10 years. Therefore, predictions on the test set were calculated using the

168 10 year onset probability (one minus survival probability). When calculating binary outcome
169 metrics, cases with time-to-event (TTE) > 10 were treated as controls. These metrics included
170 confusion matrices, area under the receiver operating characteristics curve (AUC) and area
171 under the precision-recall curve (PRAUC). After preprocessing and case thresholding (TTE >
172 10), there were 218 cases and 4,560 controls in the test set.

173 The numbers of individuals/cases and controls after each preprocessing step are shown in
174 **Supplementary Figure 1**.

175

176 *Incremental Modelling*

177 Composite protein EpiScores and direct EpiScores were generated in the training dataset using
178 different machine learning methods, via the *MethylPipeR* package (**Figure 1**), before being
179 applied to the test set in an incremental modelling approach (further detail in **Supplementary**
180 **Methods**). In the test set, a (null) risk factors-only model was defined via a Cox PH model for
181 T2D with age, sex, BMI, self-reported hypertension, and self-reported family (sibling or
182 parental) history of diabetes as predictors. A multitude of risk factors have been used in
183 previous T2D risk prediction tools, a majority of which include the set that we have used in
184 this study. While additional risk factors, such as waist-hip ratio, may also be relevant [7], we
185 selected the null model covariates based on those used in the Diabetes UK type 2 diabetes
186 ‘Know Your Risk’ tool to compare our results to an existing widely utilised tool. This was with
187 the exception of ethnicity, due to the relative homogeneity of the GS cohort. These also closely
188 match the top risk factors identified in a systematic review of previous T2D risk predictors (see
189 Figure 2 in ref. [7]).

190

191 *Penalised Regression Predictors*

192 Since the number of CpGs ($n_{\text{CpG}}=200,000$) was much greater than the number of rows in the
193 training set ($n=9,835$ after preprocessing), a regularisation method was required to reduce
194 overfitting of the Cox PH regression models.

195 Penalised regression models reduce overfitting by applying a regularisation penalty in the
196 model fitting process. This forces regression parameters to remain small, or possibly to shrink
197 them to zero. The latter allows the method to be used for variable selection, by keeping only
198 the variables with resulting non-zero coefficients.

199 Lasso penalisation was fit to the training set DNAm and protein EpiScores using glmnet [21,
200 22] via *MethylPipeR* with the best shrinkage parameter (λ) chosen by 9-fold cross-validation.

201

202 *Tree Ensemble Models*

203 Tree ensembles are non-parametric models capable of estimating complex functions using a
204 set of decision trees. Two tree ensemble approaches were used: random survival forest (RSF)
205 [23] and survival Bayesian Additive Regression Trees (sBART) [24]. Random forest [25] is an
206 ensemble machine learning model that estimates a function by averaging the output from a set
207 of independently trained decision trees. During model fitting, each tree is built using a different
208 subset of the variables from the training set to prevent individual trees from overfitting to the
209 whole dataset. sBART is a non-parametric method that estimates a function as a sum over a set
210 of regression trees. sBART incorporates the ability to model both additive and interaction
211 effects and has shown high predictive performance in comparison with similar methods [24,
212 26].

213 RSF and sBART were fit to the training set using R packages randomForestSRC (version
214 2.11.0) [27] and BART (version 2.9) [28] respectively, via *MethylPipeR*. Details on
215 hyperparameter selection are given in **Supplementary Methods**.

216 Due to computational limitations and probable overfitting in using the tree ensemble models
217 on all CpGs in the dataset, variable pre-selection was based on the coefficients in the penalised
218 Cox PH models. Each tree-ensemble model was evaluated with the features corresponding to
219 non-zero coefficients from the Cox PH lasso model.

220

221 *Evaluating Predictive Performance*

222 Survival models can be used to predict the risk of incident T2D for an arbitrary prediction
223 period. Here we focus on classification performance for the binary outcome defined by 10-year
224 T2D incidence. EpiScores were calculated as one minus 10-year survival probabilities and the
225 binary outcomes were calculated by truncating observed TTE at 10-years (see *Outcome*
226 *Definition for 10-Year Onset Prediction*).

227 AUC and PRAUC were calculated as measures of predictive performance as the discrimination
228 threshold was varied. PRAUC is more informative in situations where there is a class imbalance
229 in the test set [29]. Additionally, binary classification metrics consisting of sensitivity (recall),
230 specificity, positive predictive value (PPV/precision) and negative predictive value (NPV)
231 were calculated. These metrics require selection of a discrimination threshold to assign
232 positive/negative class predictions and have varying behaviour as this threshold is altered.
233 Therefore, each of the metrics were calculated at a range of thresholds between 0-1 in
234 increments of 0.1.

235 Model calibration was examined by comparing predicted probabilities with actual case/control
236 proportions. The test data was sorted by predicted probability and divided into bins; the mean
237 predicted probability and the proportion of cases was calculated for each bin.

238

239 *Selected-CpG Comparison with EWAS Catalog*

240 The CpG sites selected by the Cox PH lasso were queried in the EWAS Catalog [30] to identify
241 traits that have previously been linked to these sites. The EWAS catalog is a database allowing
242 users to search EWAS results from existing literature. We performed a tissue-agnostic query
243 using the selected CpGs and filtered results to those with an epigenome-wide significance
244 threshold of $P < 3.6 \times 10^{-8}$ in studies with a sample size $> 1,000$ [31]. Almost all (739 out of
245 742; 99.6%) of the post-filter results were from blood-based studies. The remaining results
246 were from saliva and prefrontal cortex-based studies.

247

248 *Validation in KORA S4*

249 The Cox PH lasso model using the direct EpiScore was applied to the KORA S4 cohort [32].
250 This cohort consisted of 1,451 individuals in southern Germany, aged 25-74 years. Cohort
251 summary details are shown in **Supplementary Table 2**. Individuals with missing values in
252 the health measures were removed from the dataset. Missing CpG values in the DNAm data
253 were mean-imputed. Similar to the approach in the Generation Scotland test set, an EpiScore
254 was computed for each individual in the KORA dataset. Evaluation was then performed using
255 an incremental modelling approach. Additional cohort and methods details (such as the
256 outcome definition, follow-up period and preprocessing numbers) are provided in
257 **Supplementary methods**.

258

259 *EpiScore Prediction of Ongoing Symptomatic COVID-19/Hospitalisation*

260 The subset of the Generation Scotland cohort with reported COVID-19 infection (clinically-
261 diagnosed or positive test from linked test data), who had also participated in the CovidLife

262 study [33] were used for prediction of ongoing symptomatic COVID-19 and hospitalisation
263 from COVID-19 (n=713). Ongoing symptomatic COVID-19 cases were defined here as
264 individuals with self-reported symptoms lasting ≥ 4 weeks [34]. Hospitalisation cases were
265 defined as hospital admissions with accompanying ICD10 codes U07.1 (confirmed COVID-
266 19 test) and U07.2 (clinically diagnosed), derived from the Scottish Morbidity Records
267 (SMR01). Details of the method and summary statistics are shown in **Supplementary methods**
268 and **Supplementary Table 3**.

269

270 **Results**

271 After preprocessing, the mean time-to-onset of T2D was 5.7 and 5.9 years for the training
272 (n=374 cases) and test (n=252 cases) sets, respectively. Mean age-at-onset was also similar
273 between the training and test set at 61.2 and 60.4 years and the mean BMI for cases (at baseline)
274 was 31.7 and 32.2 kg/m². The full set of cohort summary details for cases and controls in both
275 sets are shown in **Table 1**. The machine learning prediction pipeline of the *MethylPipeR*
276 package is shown in **Figure 1**.

277

278 *Null Model for the Incremental Modelling Approach*

279 A Cox PH model in the test set with age, sex, BMI, self-reported hypertension, and family
280 history of diabetes as predictors yielded good classification metrics: AUC=0.839,
281 PRAUC=0.227.

282

283 *Incremental model using Direct EpiScore derived using Cox PH lasso, RSF and sBART*

284 In the risk factors + direct EpiScore test set model, Cox PH lasso performed the best, showing
285 an AUC and PRAUC of 0.870 and 0.299, respectively ($p=3.6 \times 10^{-27}$ for the EpiScore
286 coefficient). This corresponds to an increase of 3.1% and 7.2% over the standard risk factors
287 model.

288 Overall, the tree-ensemble models used for the direct EpiScore resulted in lower performance
289 compared to Cox PH lasso. AUC values for RSF and sBART were 0.852 and 0.840 and
290 PRAUC values were 0.247 and 0.230, respectively (**Supplementary Table 4**). P-values for the
291 EpiScore coefficients and ROC curves for all models are given in (**Supplementary Table 5**).

292 *Incremental model using Composite Protein EpiScore derived using Cox PH lasso*

293 The composite protein EpiScore showed good performance with AUC and PRAUC of 0.864
294 and 0.270, respectively (EpiScore coefficient $p=1.61 \times 10^{-18}$). The increase in PRAUC was
295 smaller for the composite protein EpiScore compared to the direct EpiScore but still a notable
296 improvement over using risk factors only.

297

298 *Incremental model using Composite Protein and Direct EpiScores*

299 The full model (risk factors + composite protein EpiScore + direct EpiScore) with a Cox PH
300 lasso direct EpiScore gave an AUC and PRAUC of 0.872 and 0.302 respectively. The full
301 models with RSF- and sBART-derived direct EpiScores showed AUCs of 0.866 and 0.864,
302 respectively. The corresponding PRAUC values were 0.273 and 0.270. Therefore, the best
303 overall model used direct and composite protein EpiScores from Cox PH lasso models. The
304 ROC and PR curves for the full models and risk factor only model are shown in **Figure 2**.

305

306 *Binary Classification Metrics and Model Calibration*

307 **Supplementary Table 6** shows how confusion matrix metrics vary for the null (risk factors
308 only) model and the Cox PH lasso model across a range of probability classification thresholds.
309 As expected, as the classification probability threshold is increased, sensitivity and NPV
310 decrease while specificity increases. The effects of these differences on the number of true
311 positives and true negatives are illustrated in **Figure 3**. The two models also show differences
312 in their calibration plots (**Supplementary Figure 2**). In addition, the difference in number of
313 correctly classified individuals between the two models are given. These are calculated
314 assuming, arbitrarily, a 10-year incidence rate of 33%, for example, in a scenario where high-
315 risk individuals have been selected for screening.

316

317 *Selected CpGs*

318 The Cox PH lasso model assigned non-zero coefficients to 145 CpGs (**Supplementary Table**
319 **7**). After filtering the EWAS Catalog by p-value ($p < 3.6 \times 10^{-8}$) [31] and sample size ($n > 1,000$),
320 119 (82%) of the model-selected CpGs were present. These CpGs corresponded to 742 entries
321 and showed epigenome-wide associations with traits including: serum HDL cholesterol, serum
322 triglycerides, smoking, C-reactive protein, BMI and age (**Supplementary Table 8**).

323

324 *Selected Protein EpiScores*

325 The composite protein EpiScore Cox PH lasso model assigned non-zero coefficients to 46
326 protein EpiScores. Details on the corresponding proteins and genes are given in
327 **Supplementary Table 9**. Out of the selected protein EpiScores, 21 have previously shown
328 associations with incident T2D [9].

329

330 *Validation in KORA S4*

331 Prediction of incident diabetes in the KORA S4 cohort using the Cox PH lasso model showed
332 good replication of direct EpiScore performance ($p=1.6 \times 10^{-5}$) with increases of 1.6% and in
333 absolute terms above the null model values for both AUC and PRAUC. Further details are
334 provided in **Supplementary Table 10**.

335

336 *EpiScore Prediction of Ongoing Symptomatic COVID-19/Hospitalisation*

337 In all models, incident T2D was predictive of hospitalisation with COVID-19 infection.
338 However, neither the composite-protein nor the direct EpiScore were predictive of the same
339 outcome (**Supplementary Table 11**). Additionally, neither the EpiScores nor incident T2D
340 were predictive of ongoing symptomatic COVID after COVID-19 infection.

341

342 **Discussion**

343 Utilising a large cohort with genome-wide epigenetic data and health records linkage to
344 longitudinal primary and secondary care, we have shown that DNAm-based predictors
345 augment standard risk factors in the prediction of incident type 2 diabetes. The best model with
346 traditional risk factors yielded an AUC of 0.839 compared to 0.872 when DNAm was also
347 considered and the PRAUC increased from 0.227 to 0.305. Using a variety of linear and non-
348 linear survival models, we showed that overall, the Cox PH lasso model produced the most
349 predictive direct EpiScore. A composite protein EpiScore also notably increased predictive
350 performance. The direct EpiScore also showed good external validation performance in the
351 KORA S4 cohort. Beyond the T2D analysis presented here, we have developed the
352 *MethylPipeR* R package, along with a user-interface *MethylPipeR-UI*, to facilitate reproducible

353 machine learning time-to-event and binary prediction using DNAm or other types of high-
354 dimensional omics data.

355 Determining a 'best' model is complicated and depends on the trade-off that a user wishes to
356 make. Here, we optimised AUC and PRAUC but binary classification metrics vary by method
357 and/or classification threshold. When using classifiers in clinical settings, decisions need to be
358 made about the number of patients that can be recommended for intervention as well as the
359 acceptable proportion of false positives and false negatives. We showed an increase in correct
360 identification of positives/negatives at varying probability thresholds when adding direct and
361 composite EpiScores above standard risk factors. For instance, given an (arbitrary) incidence
362 rate of 33% (commonly used as a cut-off for high-risk of T2D) [35] over 10 years in a sample
363 of 10,000 individuals, our best model would correctly classify an additional 448 individuals
364 compared to the risk factors only model at a threshold of 0.2 (**Supplementary Table 6**). Given
365 the costs of treating T2D-related complications, our study gives evidence for possible benefits
366 of EpiScores on public health that could also alleviate the financial burden to the NHS. In
367 addition, an assessment of calibration is also critical [36]. Investigation of these related criteria
368 could assist in deciding an optimal threshold given clinical constraints and provide a more
369 comprehensive assessment of model predictions than AUCs or metrics at the commonly-
370 utilised threshold of 0.5.

371 Several CpGs from the direct EpiScore were previously identified as epigenome-wide
372 significant correlates of traits commonly linked to T2D [14, 17, 37-41]. Future work could
373 investigate overlap between these and time-to-event EWAS studies. Further studies could also
374 include DNAm predictors for traditional risk factors that are included in the null model, such
375 as BMI [8].

376 Limitations include the relatively small number of disease cases in the dataset, the limited
377 hyperparameter optimisation performed for sBART and the relatively simple variable pre-
378 selection method for tree-ensemble methods. Given the lower performance of these methods
379 compared to the best models in this study, there is potential for additional improvement in
380 predictive performance with further investigation of more advanced feature (pre-)selection.
381 This is particularly important when we consider that the pre-selection step utilised linear
382 models prior to the non-linear model fitting. The model fitting and pre-selection were also
383 performed using the same training set which may have introduced issues associated to post-
384 selection inference [42, 43]. In addition, factors such as overfitting, related individuals in the
385 test set and batch effects between the three rounds of DNAm data processing may all have an
386 effect on test-set AUC. Finally, a small proportion of the linkage codes used to define diabetes
387 included broad terms that were non-specific to T2D; however, late age of onset in these
388 individuals meant there was a high likelihood that they had developed T2D. EpiScores for
389 T2D-associated proteins have also been shown to replicate incident T2D-protein associations
390 within this sample [9] suggesting that the case definitions we use capture biological signals
391 relevant to T2D.

392 There are numerous strengths to our study. Firstly, the models used capture relationships
393 between CpGs as well as time-to-event information, which is not possible using traditional
394 EWAS methods. Secondly, data linkage to health care measures provided comprehensive T2D
395 incidence data in a very large cohort study, Generation Scotland. Validation performance in
396 the KORA cohort also strengthened evidence for the applicability of the models to other
397 populations. Finally, the R package, *MethylPipeR*, encourages reproducibility and allows
398 others to develop similar predictors on new data with minimal setup.

399 In conclusion, we have demonstrated the potential for DNA methylation data to provide notable
400 improvement in predictive performance for incident T2D, as compared to traditional risk

401 factors (age, sex, BMI, hypertension, and family history). We evaluated different models with
402 a systematic approach and presented a framework with the ability to generalise to other traits
403 and datasets for training and testing predictors in future studies.

404

405 **Declarations**

406 *Ethics approval and consent to participate*

407 All components of Generation Scotland received ethical approval from the NHS Tayside
408 Committee on Medical Research Ethics (REC Reference Number: 05/S1401/89). Generation
409 Scotland has also been granted Research Tissue Bank status by the East of Scotland Research
410 Ethics Service (REC Reference Number: 20-ES-0021), providing generic ethical approval for
411 a wide range of uses within medical research.

412 The KORA studies were approved by the Ethics Committee of the Bavarian Medical
413 Association (Bayerische Landesärztekammer; S4: #99186) and were conducted according to
414 the principles expressed in the Declaration of Helsinki. All study participants gave their written
415 informed consent.

416

417 *Availability of Data and Material*

418 According to the terms of consent for Generation Scotland participants, access to data must be
419 reviewed by the Generation Scotland Access Committee. Applications should be made to
420 access@generationscotland.org.

421 All code is available with open access at the following Gitlab repository:

422 <https://github.com/marioni-group>

423 *MethylPipeR* (version 1.0.0) is available at: <https://github.com/marioni-group/MethylPipeR>

424 *MethylPipeR-UI* is available at: <https://github.com/marioni-group/MethylPipeR-UI>

425 The informed consents given by KORA study participants do not cover data posting in public
426 databases. However, data are available upon request from KORA Project Application Self-
427 Service Tool (<https://epi.helmholtz-muenchen.de/>). Data requests can be submitted online and
428 are subject to approval by the KORA Board.

429

430 *Competing Interests*

431 R.E.M has received a speaker fee from Illumina and is an advisor to the Epigenetic Clock
432 Development Foundation. A.M.M has previously received speaker fees from Janssen and
433 Illumina and research funding from The Sackler Trust. L.M has received payment from
434 Illumina for presentations and consultancy. All other authors declare no competing interests.

435

436 *Acknowledgements*

437 **This research was funded in whole, or in part, by the Wellcome Trust [104036/Z/14/Z,**
438 **108890/Z/15/Z, 216767/Z/19/Z]. For the purpose of open access, the author has applied a**
439 **CC BY public copyright licence to any Author Accepted Manuscript version arising from**
440 **this submission.**

441 Generation Scotland received core support from the Chief Scientist Office of the Scottish
442 Government Health Directorates (CZD/16/6) and the Scottish Funding Council (HR03006) and
443 is currently supported by the Wellcome Trust (216767/Z/19/Z). DNA methylation profiling of
444 the Generation Scotland samples was carried out by the Genetics Core Laboratory at the
445 Edinburgh Clinical Research Facility, Edinburgh, Scotland and was funded by the Medical

446 Research Council UK and the Wellcome Trust (Wellcome Trust Strategic Award “STratifying
447 Resilience and Depression Longitudinally” (STRADL; Reference 104036/Z/14/Z). The DNA
448 methylation data assayed for Generation Scotland was partially funded by a 2018 NARSAD
449 Young Investigator Grant from the Brain & Behavior Research Foundation (Ref: 27404;
450 awardee: Dr David M Howard) and by a JMAS SIM fellowship from the Royal College of
451 Physicians of Edinburgh (Awardee: Dr Heather C Whalley). Y.C. is supported by the
452 University of Edinburgh and University of Helsinki joint PhD program in Human Genomics.
453 D.A.G. is supported by funding from the Wellcome Trust 4-year PhD in Translational
454 Neuroscience—training the next generation of basic neuroscientists to embrace clinical research
455 [108890/Z/15/Z]. C.A.V. is a Chancellor’s Fellow funded by the University of Edinburgh.
456 D.L.Mc.C. and R.E.M. are supported by Alzheimer’s Research UK major project grant ARUK-
457 PG2017B–10. R.E.M. is supported by Alzheimer’s Society major project grant AS-PG-19b-
458 010.

459 Recruitment to the CovidLife study was facilitated by SHARE- the Scottish Health Research
460 Register and Biobank.

461 SHARE is supported by NHS Research Scotland, the Universities of Scotland and the Chief
462 Scientist Office of the Scottish Government.

463 The KORA study was initiated and financed by the Helmholtz Zentrum München – German
464 Research Center for Environmental Health, which is funded by the German Federal Ministry
465 of Education and Research (BMBF) and by the State of Bavaria. Furthermore, KORA research
466 has been supported within the Munich Center of Health Sciences (MC-Health), Ludwig-
467 Maximilians-Universität, as part of LMUinnovativ and is supported by the DZHK (German
468 Centre for Cardiovascular Research). The KORA study is funded by the Bavarian State

469 Ministry of Health and Care through the research project DigiMed Bayern (www.digimed-
470 bayern.de).

471 **Supplementary Methods**

472 *MethylPipeR and MethylPipeR-UI*

473 The analysis pipeline was implemented via a new R package, *MethylPipeR*, along with
474 accompanying user interface, for systematic and reproducible development of complex trait
475 and incident disease predictors. *MethylPipeR* provides functionality for tasks such as model
476 fitting, prediction and performance evaluation as well as automatic logging of experiments and
477 trained models. This is complemented by *MethylPipeR-UI* which provides an interface to the
478 R package functionality while removing the need to write scripts. While *MethylPipeR* was
479 applied to incident T2D prediction with DNA methylation in our experiments, the package is
480 designed for generalised development of predictive models and is applicable to a wide range
481 of omics data and target traits. *MethylPipeR* and *MethylPipeR-UI* are publicly available at
482 <https://github.com/marioni-group/MethylPipeR> and [https://github.com/marioni-](https://github.com/marioni-group/MethylPipeR-UI)
483 [group/MethylPipeR-UI](https://github.com/marioni-group/MethylPipeR-UI) respectively. **Supplementary Figure 3** shows an example from the
484 *MethylPipeR-UI* interface including functionality such as data upload, specification of model
485 options and visualisation of model diagnostics.

486

487 *Methylation Risk Scores and Incremental Modelling*

488 Composite protein EpiScore models were fit to the training set using 109 protein EpiScores as
489 features [9]. This trained model was used to create a composite protein EpiScore for each
490 individual in the test set. Similarly, direct predictors were fit to the training set using 200,000
491 CpGs as features and used to create a direct EpiScore for each test set individual.

492 The following Cox PH models were then fit to the test set to assess the difference in metrics
493 such as the area under the curve (AUC) and area under the precision-recall curve (PRAUC): a
494 risk factors only model; risk factors + protein EpiScore; risk factors + direct EpiScore and a

495 full model (risk factors + protein EpiScore + direct EpiScore). To obtain 10-year survival
496 probabilities in the test set, the cumulative baseline hazard at $t=10$ was calculated for each
497 model using the Breslow estimator [44].

498 *Calculating Predictions*

499 To obtain predictions from the incremental models, the following should be applied.

500 Direct EpiScore

501 For the Cox lasso model, the direct EpiScore can be obtained by summing over the CpG values
502 weighted by their corresponding coefficients (for all non-zero coefficients), as given in

503 **Supplementary Table 12.**

$$504 \quad \text{Direct EpiScore} = \beta_1 CpG_1 + \beta_2 CpG_2 + \dots + \beta_p CpG_p$$

505 For sBART and RSF, the direct EpiScore requires use of the model object to obtain predictions
506 on new data. These can be found in <https://github.com/marioni-group/> and utilised with
507 *MethylPipeR*'s predict function.

508 Composite Protein EpiScore

509 The composite protein EpiScore can be obtained by summing over the protein EpiScores
510 weighted by their corresponding coefficients (for all non-zero coefficients), as given in

511 **Supplementary Table 12.**

$$512 \quad \text{Composite Protein EpiScore} = \beta_1 PES_1 + \beta_2 PES_2 + \dots + \beta_q PES_q$$

513 where q is the number of protein EpiScores with a non-zero coefficient.

514 Each protein EpiScore is calculated as the sum over CpG values weighted by their
515 corresponding coefficients (for all non-zero coefficients); the coefficients for each of these can

516 be found in the supplementary materials in Gadd et al. [9]

517 Incremental Cox PH

518 The linear predictor in an incremental Cox PH model can be calculated by multiplying the
519 relevant EpiScores and risk factors by their corresponding coefficient, given in
520 **Supplementary Table 12**. The cumulative baseline hazard is calculated using the `basehaz.gbm`
521 function in the `gbm` R package version 2.1.8 [45]. The 10-year survival probabilities can then
522 be calculated as $S(t) = \exp[-(\Lambda_0(t))]^{\exp(l)}$ at $t = 10$, where $\Lambda_0(t)$ is the cumulative baseline
523 hazard and l is the linear predictor. The 10-year onset probability is therefore $1 - S(10)$.

524 *Penalised Cox Proportional-hazards Regression*

525 Since the number of features (200,000) was much greater than the number of individuals in the
526 training set (=9,835 after data preprocessing), a regularisation method was required to reduce
527 overfitting of the Cox PH models.

528 Cox PH models with lasso penalisation was fit to the training set using `glmnet` (version 4.1-1)
529 [21, 22] via *MethylPipeR*.

530 Hyperparameter optimisation and cross validation (CV) were used to select the λ that
531 minimised the partial-likelihood for each Cox PH model corresponding to the minimum partial-
532 likelihood. The training set was divided equally into nine partitions. For each pre-selected value
533 of λ , nine models were fit, each using eight of the partitions as the training set and the ninth for
534 prediction. The mean partial likelihood over the nine models was then calculated. The model
535 using the λ that minimised this was chosen to evaluate on the test set. Nine-fold CV was used
536 to balance the advantages provided by using a greater number of folds with the limitations of
537 the number of cases in each fold as well as required computation time. In addition, individuals
538 belonging to the same family were assigned to the same fold and each contained a similar
539 number of cases to avoid folds with too few cases.

540

541 *Random Forest*

542 Random forest [25] is an ensemble machine learning model that estimates a function by
543 averaging the output from a set of independently trained decision trees. During model fitting,
544 each tree is built using a different subset of the variables and the training set to prevent
545 individual trees from overfitting to the whole dataset. In addition, random survival forests
546 adapts the original method to incorporate right-censored time-to-event data [23].

547 The hyperparameters corresponding to the number of trees (ntrees), the number of variables
548 considered at each tree split (mtry) and the minimum terminal node size (nodesize) were
549 selected using a grid-search CV method. **Supplementary Table 13** shows the set of values that
550 were considered for each hyperparameter. The R package randomForestSRC (version 2.11.0)
551 [27] was used for fitting random survival forests via *MethylPipeR*.

552

553 *Bayesian Additive Regression Trees*

554 Bayesian Additive Regression Trees (BART) [26] is a nonparametric method that estimates a
555 function as a sum over a set of regression trees. BART incorporates the ability to model both
556 additive and interaction effects and has shown high predictive performance in comparison with
557 similar methods. To reduce overfitting and model uncertainty in parameters and predictions,
558 BART uses prior distributions over tree-related parameters. Posterior estimates are obtained in
559 a Bayesian framework through Markov Chain Monte Carlo (MCMC).

560 A variant of BART for survival analysis [24] was used for 10-year onset prediction. 1,000
561 posterior samples for model parameters were kept after 500 burn-in samples. These were used
562 to generate 10-year survival probabilities on the test set. This resulted in 1,000 survival

563 probabilities for each individual in the test set, the mean of which was used as their survival
564 prediction. MCMC convergence was assessed using Geweke's diagnostic (gewekediag
565 function in the BART R package). This was calculated on the BART output (yhat.test with
566 $t \leq 10$ [24]) for the 1,000 posterior samples using 30 randomly chosen individuals in the test
567 set. Most Z-values (77%) fell within the interval [-1.96, 1.96] suggesting convergence.

568 Due to the computation time requirements of MCMC sampling and the apparent robustness of
569 BART to hyperparameter misspecification [26], the sBART model was run with
570 hyperparameters set to default. This was performed using the R package BART (version 2.9)
571 [28] via *MethylPipeR*.

572

573 *Validation of Best Performing Model in KORA S4*

574 The present analyses are based on a subsample of the participants of the KORA S4 study.
575 KORA (Cooperative Health Research in the Region of Augsburg) is a research platform
576 performing population-based surveys and subsequent follow-ups in the region of Augsburg in
577 Southern Germany [46]. Participants were of German nationality, aged between 25-74 years,
578 50% female, and sampled from the population registers in the study area where main place of
579 residence was: Ausburg city town, county Ausburg or county Aichach-Friedberg. Each
580 participant completed a health questionnaire, providing details on health status and medication.
581 Blood samples were also taken for assaying of omics data. KORA S4 participants were
582 recruited between 25/10/1999-28/04/2001. This study used a subsample of the 1,451
583 participants of the KORA S4 study with DNAm and incident T2D data available and no
584 prevalent diabetes at baseline.

585 The best performing model selected for the Generation Scotland cohort (Cox PH lasso) was
586 used for prediction of incident T2D in the KORA S4 cohort.

587 For diabetes morbidity, the data are limited to a follow-up of 10 years - starting from KORA
588 S4 recruitment. For incident T2D all prevalent diabetics as well all other diabetes types except
589 T2D cases are excluded. Age, BMI, hypertension, sex as well as self-reported family (mother
590 or father) history of T2D were taken at the baseline of KORA S4. BMI was calculated as the
591 individual's weight in kg divided by the square of their height in metres.

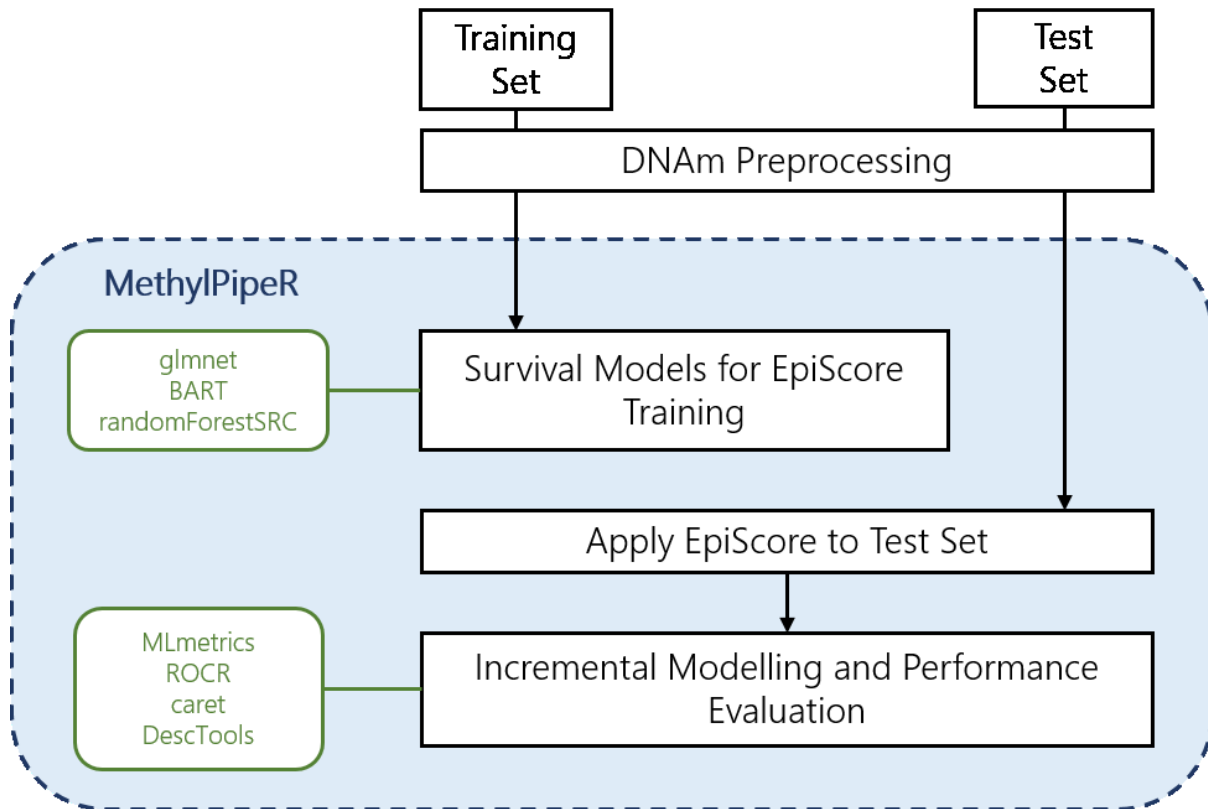
592

593 *EpiScore Prediction of Ongoing Symptomatic COVID-19/Hospitalisation*

594 Self-reported COVID-19 phenotypes were available in a subset of individuals from the
595 Generation Scotland DNA methylation sample who had also participated in the CovidLife
596 surveys (n=2,399) [33]. Ongoing symptomatic COVID-19 phenotypes were ascertained from
597 CovidLife survey 3, (n=1,802 Generation Scotland participants with DNAm profiled), where
598 participants were asked about the total overall time they experienced symptoms in their
599 first/only episode of illness, as well as the whole of their COVID-19 illness. Ongoing
600 symptomatic COVID-19 was defined here as symptoms persisting for at least 4 weeks from
601 infection and is correct as of February 2021, when the survey 3 was administered.
602 Hospitalisation information, derived from the Scottish Morbidity Records (SMR01), was used
603 to obtain COVID-19 hospital admissions using ICD10 codes U07.1 (lab-confirmed COVID-
604 19 diagnosis), and U07.2 (clinically-diagnosed COVID-19). Hospitalisation data is correct as
605 of February 2021. Logistic regression was used to assess the predictive performance of the
606 T2D EpiScores in relation to ongoing symptomatic COVID-19 and severe COVID-19 (i.e.
607 hospitalisation), adjusting for the standard risk factors and incident T2D before a positive
608 COVID-19 test.

609

610 **Figure 1.** The prediction pipeline and functionality provided in *MethylPipeR*



611

612

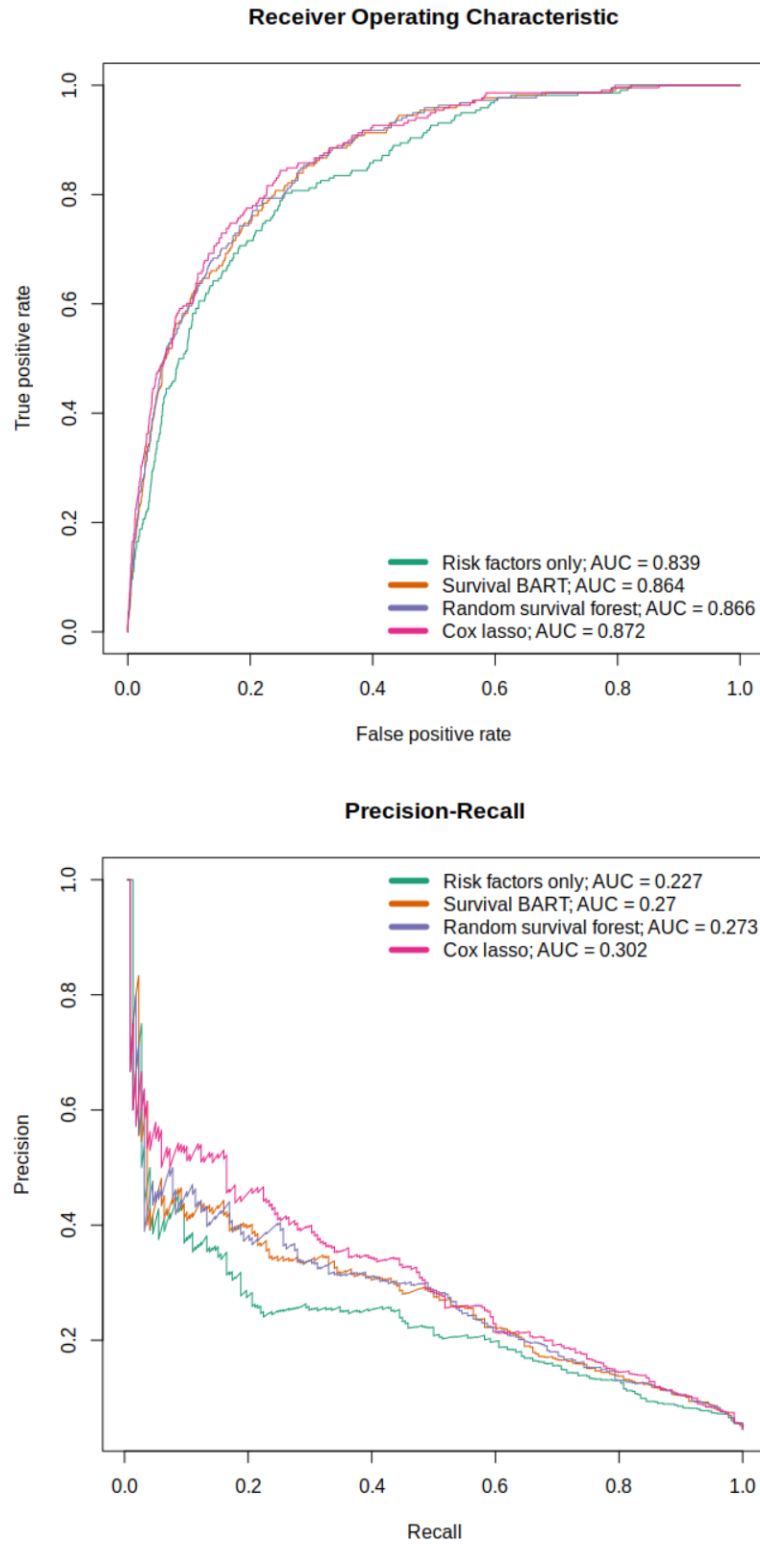
613

614

615 **Figure 2.** ROC and PR curves for the full models with Cox PH Lasso, Random Survival
616 Forest and Survival BART direct EpiScores.

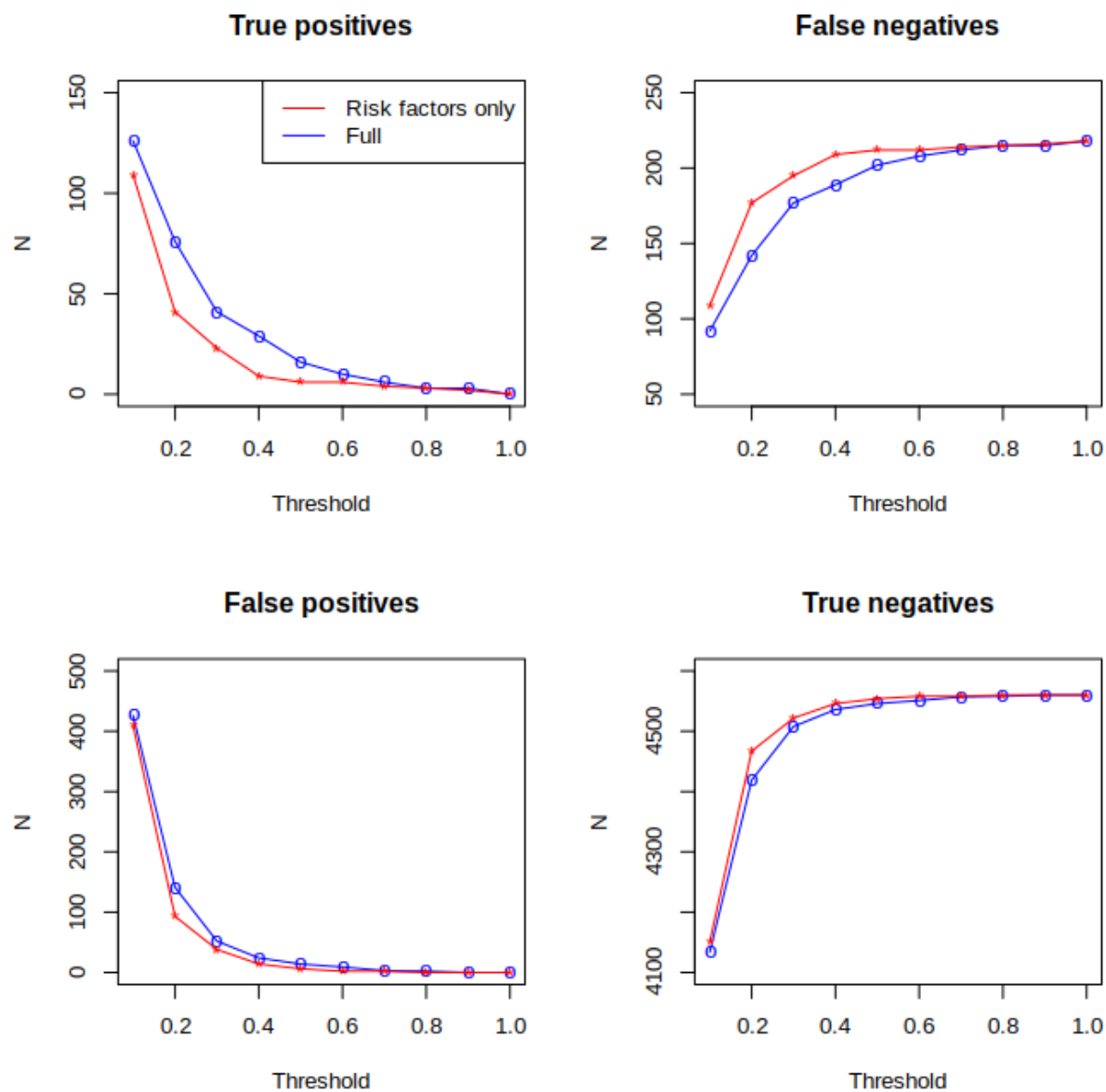
617

618



619

621 **Figure 3.** Confusion matrix plot of true/false positives/negatives for the risk factors only and
622 full model in the Generation Scotland test dataset. (Full model uses direct and composite
623 protein EpiScores from Cox PH Lasso.)



624

625

626

627

628 **Table 1.** Summary information for the Generation Scotland training and test sets. Summary
629 information is mean (SD) or n (%).

	Training		Test	
	Cases	Controls	Cases	Controls
n	374	9461	252	4526
Time-to-event (Years to Onset or Censoring)	5.7 (3.4)	11.1 (1.8)	5.9 (3.4)	11.3 (1.7)
Age (Onset or Censoring)	61.2 (10.7)	58.1 (14.6)	60.4 (9.4)	59.2 (13.9)
Sex (Male)	184 (49.2)	3903 (41.3)	133 (52.8)	1681 (37.1)
BMI (kg/m ²)	31.7 (5.7)	26.3 (4.8)	32.2 (6.2)	26.5 (5.0)
Self-reported Parent or Sibling Diabetes	137 (36.6)	1553 (16.4)	105 (41.7)	858 (19.0)
Self-reported Hypertension	117 (31.3)	1022 (10.8)	90 (35.7)	575 (12.7)

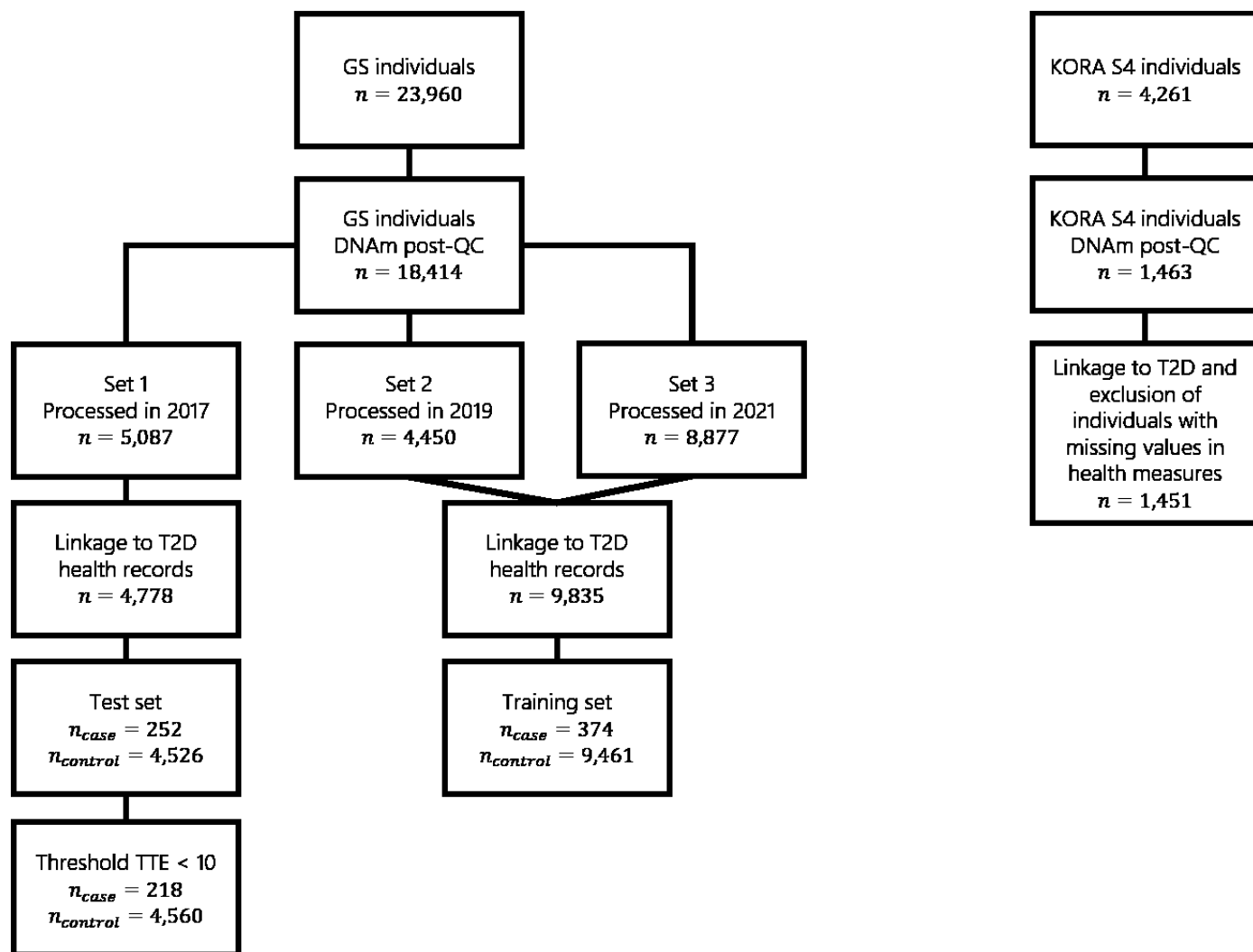
630

631

632

633

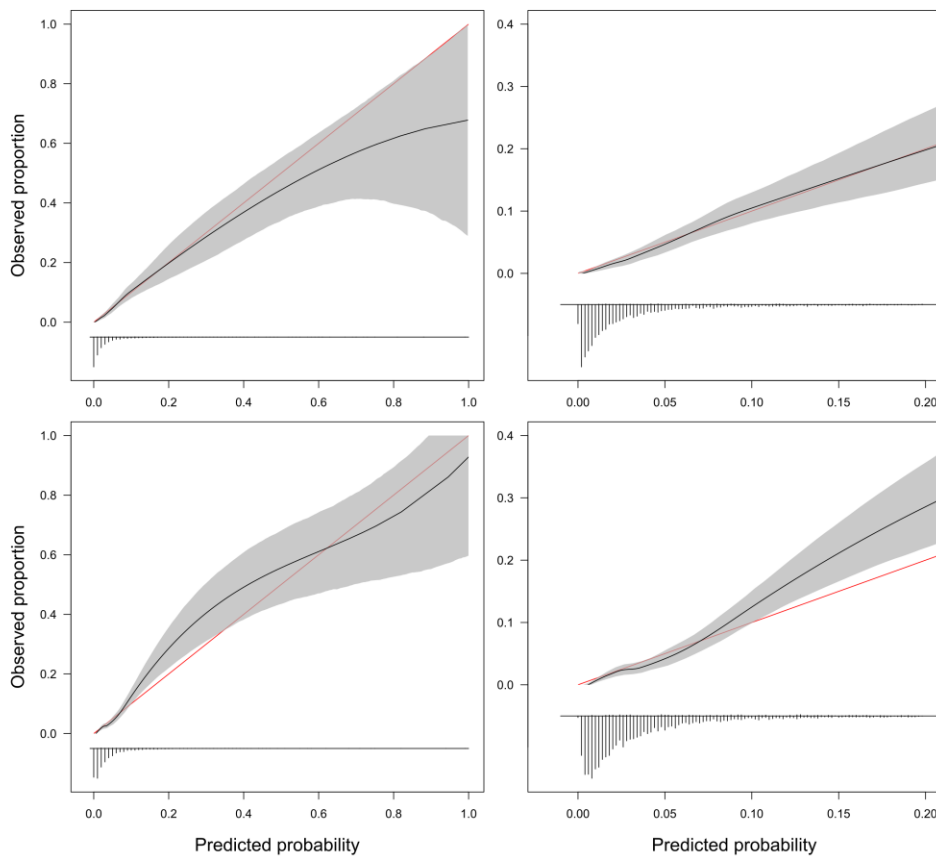
634 **Supplementary Figure 1.** Preprocessing steps for Generation Scotland with number of
635 individuals/cases and controls after each step.



636

637

638 **Supplementary Figure 2.** Calibration plots for the full model (risk factors + composite
639 protein EpiScore + Cox PH lasso direct EpiScore) (top-left) and the risk factors only model
640 (bottom-left). The grey area shows 95% confidence intervals calculated from 2000 bootstrap
641 samples. The ideal calibration line (observed = predicted) is shown in red. The histogram
642 shows the distribution of predicted probabilities. The wider confidence intervals at higher
643 predicted probabilities are due to the small number of predictions in those ranges. Most
644 predictions are low in the probability range, emphasised in the zoomed-in plots (top-right and
645 bottom-right).



646

647

648

649

650 **Supplementary Figure 3.** An example from the *MethylPipeR-UI* Shiny app.



651

652

653 References

- 654 1. Saeedi, P., et al., *Global and regional diabetes prevalence estimates for 2019 and projections*
655 *for 2030 and 2045: Results from the International Diabetes Federation Diabetes Atlas.*
656 *Diabetes research and clinical practice*, 2019. **157**: p. 107843.
- 657 2. Gregg, E.W., N. Sattar, and M.K. Ali, *The changing face of diabetes complications.* *The lancet*
658 *Diabetes & endocrinology*, 2016. **4**(6): p. 537-547.
- 659 3. Biessels, G.J. and F. Despa, *Cognitive decline and dementia in diabetes mellitus: mechanisms*
660 *and clinical implications.* *Nature Reviews Endocrinology*, 2018. **14**(10): p. 591-604.
- 661 4. McGurnaghan, S.J., et al., *Risks of and risk factors for COVID-19 disease in people with*
662 *diabetes: a cohort study of the total population of Scotland.* *The Lancet Diabetes &*
663 *Endocrinology*, 2021. **9**(2): p. 82-93.
- 664 5. Lee, C.G., et al., *Effect of metformin and lifestyle interventions on mortality in the diabetes*
665 *prevention program and diabetes prevention program outcomes study.* *Diabetes care*, 2021.
666 **44**(12): p. 2775-2782.
- 667 6. Keng, M.J., et al., *Impact of achieving primary care targets in type 2 diabetes on health*
668 *outcomes and healthcare costs.* *Diabetes, Obesity and Metabolism*, 2019. **21**(11): p. 2405-
669 2412.
- 670 7. Collins, G.S., et al., *Developing risk prediction models for type 2 diabetes: a systematic review*
671 *of methodology and reporting.* *BMC medicine*, 2011. **9**(1): p. 1-14.
- 672 8. McCartney, D.L., et al., *Epigenetic prediction of complex traits and death.* *Genome biology*,
673 2018. **19**(1): p. 1-11.
- 674 9. Gadd, D.A., et al., *Epigenetic scores for the circulating proteome as tools for disease prediction.*
675 *Elife*, 2022. **11**: p. e71802.
- 676 10. Cardona, A., et al., *Epigenome-wide association study of incident type 2 diabetes in a British*
677 *population: EPIC-Norfolk study.* *Diabetes*, 2019. **68**(12): p. 2315-2326.
- 678 11. Meeks, K.A., et al., *Epigenome-wide association study in whole blood on type 2 diabetes*
679 *among sub-Saharan African individuals: findings from the RODAM study.* *International journal*
680 *of epidemiology*, 2019. **48**(1): p. 58-70.
- 681 12. Walaszczyk, E., et al., *DNA methylation markers associated with type 2 diabetes, fasting*
682 *glucose and HbA 1c levels: a systematic review and replication in a case-control sample of the*
683 *Lifelines study.* *Diabetologia*, 2018. **61**(2): p. 354-368.
- 684 13. Al Muftah, W.A., et al., *Epigenetic associations of type 2 diabetes and BMI in an Arab*
685 *population.* *Clinical epigenetics*, 2016. **8**(1): p. 1-10.
- 686 14. Chambers, J.C., et al., *Epigenome-wide association of DNA methylation markers in peripheral*
687 *blood from Indian Asians and Europeans with incident type 2 diabetes: a nested case-control*
688 *study.* *The lancet Diabetes & endocrinology*, 2015. **3**(7): p. 526-534.
- 689 15. Nakatochi, M., et al., *Epigenome-wide association of myocardial infarction with DNA*
690 *methylation sites at loci related to cardiovascular disease.* *Clinical epigenetics*, 2017. **9**(1): p.
691 1-9.
- 692 16. Wang, X., et al., *An epigenome-wide study of obesity in African American youth and young*
693 *adults: novel findings, replication in neutrophils, and relationship with gene expression.* *Clinical*
694 *epigenetics*, 2018. **10**(1): p. 1-9.
- 695 17. Wahl, S., et al., *Epigenome-wide association study of body mass index, and the adverse*
696 *outcomes of adiposity.* *Nature*, 2017. **541**(7635): p. 81-86.
- 697 18. Haw, J.S., et al., *Long-term sustainability of diabetes prevention approaches: a systematic*
698 *review and meta-analysis of randomized clinical trials.* *JAMA internal medicine*, 2017. **177**(12):
699 p. 1808-1817.
- 700 19. Collins, G.S., et al., *Transparent reporting of a multivariable prediction model for individual*
701 *prognosis or diagnosis (TRIPOD): the TRIPOD statement.* *Journal of British Surgery*, 2015.
702 **102**(3): p. 148-158.

- 703 20. Smith, B.H., et al., *Cohort Profile: Generation Scotland: Scottish Family Health Study (GS: SFHS). The study, its participants and their potential for genetic research on health and illness.* International journal of epidemiology, 2013. **42**(3): p. 689-700.
- 704
- 705
- 706 21. Simon, N., et al., *Regularization paths for Cox's proportional hazards model via coordinate descent.* Journal of statistical software, 2011. **39**(5): p. 1.
- 707
- 708 22. Friedman, J., T. Hastie, and R. Tibshirani, *Regularization paths for generalized linear models via coordinate descent.* Journal of statistical software, 2010. **33**(1): p. 1.
- 709
- 710 23. Ishwaran, H., et al., *Random survival forests.* Annals of Applied Statistics, 2008. **2**(3): p. 841-860.
- 711
- 712 24. Sparapani, R.A., et al., *Nonparametric survival analysis using Bayesian additive regression trees (BART).* Statistics in medicine, 2016. **35**(16): p. 2741-2753.
- 713
- 714 25. Breiman, L., *Random forests.* Machine learning, 2001. **45**(1): p. 5-32.
- 715 26. Chipman, H.A., E.I. George, and R.E. McCulloch, *BART: Bayesian additive regression trees.* The Annals of Applied Statistics, 2010. **4**(1): p. 266-298.
- 716
- 717 27. Ishwaran, H. and U. Kogalur, *Fast unified random forests for survival, regression, and classification (RF-SRC).* R package version, 2019. **2**(1).
- 718
- 719 28. Sparapani, R., C. Spanbauer, and R. McCulloch, *Nonparametric machine learning and efficient computation with bayesian additive regression trees: the BART R package.* Journal of Statistical Software, 2021. **97**(1): p. 1-66.
- 720
- 721
- 722 29. Saito, T. and M. Rehmsmeier, *The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets.* PloS one, 2015. **10**(3): p. e0118432.
- 723
- 724
- 725 30. Battram, T., et al., *The EWAS Catalog: a database of epigenome-wide association studies.* 2021.
- 726
- 727 31. Saffari, A., et al., *Estimation of a significance threshold for epigenome-wide association studies.* Genetic epidemiology, 2018. **42**(1): p. 20-33.
- 728
- 729 32. Wichmann, H.-E., et al., *KORA-gen-resource for population genetics, controls and a broad spectrum of disease phenotypes.* Das Gesundheitswesen, 2005. **67**(S 01): p. 26-30.
- 730
- 731 33. Fawns-Ritchie, C., et al., *CovidLife: a resource to understand mental health, well-being and behaviour during the COVID-19 pandemic in the UK.* Wellcome Open Research, 2021. **6**(176): p. 176.
- 732
- 733
- 734 34. Shah, W., et al., *Managing the long term effects of covid-19: summary of NICE, SIGN, and RCGP rapid guideline.* bmj, 2021. **372**.
- 735
- 736 35. Ekoe, J.-M., R. Goldenberg, and P. Katz, *Screening for diabetes in adults.* Canadian journal of diabetes, 2018. **42**: p. S16-S19.
- 737
- 738 36. Van Calster, B., et al., *Calibration: the Achilles heel of predictive analytics.* BMC medicine, 2019. **17**(1): p. 1-7.
- 739
- 740 37. Demerath, E.W., et al., *Epigenome-wide association study (EWAS) of BMI, BMI change and waist circumference in African American adults identifies multiple replicated loci.* Human molecular genetics, 2015. **24**(15): p. 4464-4479.
- 741
- 742
- 743 38. Mendelson, M.M., et al., *Association of body mass index with DNA methylation and gene expression in blood cells and relations to cardiometabolic disease: a Mendelian randomization approach.* PLoS medicine, 2017. **14**(1): p. e1002215.
- 744
- 745
- 746 39. Sayols-Baixeras, S., et al., *Identification and validation of seven new loci showing differential DNA methylation related to serum lipid profile: an epigenome-wide approach. The REGICOR study.* Human molecular genetics, 2016. **25**(20): p. 4556-4565.
- 747
- 748
- 749 40. Braun, K.V., et al., *Epigenome-wide association study (EWAS) on lipids: the Rotterdam Study.* Clinical epigenetics, 2017. **9**(1): p. 1-11.
- 750
- 751 41. Kriebel, J., et al., *Association between DNA methylation in whole blood and measures of glucose metabolism: KORA F4 study.* PloS one, 2016. **11**(3): p. e0152314.
- 752

- 753 42. Lee, J.D., et al., *Exact post-selection inference, with application to the lasso*. The Annals of
754 Statistics, 2016. **44**(3): p. 907-927.
- 755 43. Taylor, J. and R. Tibshirani, *Post-selection inference for-penalized likelihood models*. Canadian
756 Journal of Statistics, 2018. **46**(1): p. 41-61.
- 757 44. Lin, D., *On the Breslow estimator*. Lifetime data analysis, 2007. **13**(4): p. 471-480.
- 758 45. Ridgeway, G., *Generalized Boosted Models: A guide to the gbm package*. Update, 2007. **1**(1):
759 p. 2007.
- 760 46. Holle, R., et al., *KORA-a research platform for population based health research*. Das
761 Gesundheitswesen, 2005. **67**(S 01): p. 19-25.
- 762