

Feature pre-selection for the development of epigenetic biomarkers

Yipeng Cheng^{1, 2}, Christian Gieger^{3, 4, 5}, Archie Campbell¹, Andrew M McIntosh⁶, Melanie Waldenberger^{3, 4}, Daniel L McCartney¹, Riccardo E Marioni^{1, *}, and Catalina A Vallejos^{7, 8, *}

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

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

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

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

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

⁶Division of Psychiatry, Centre for Clinical Brain Sciences, University of Edinburgh, Edinburgh, UK

⁷MRC Human Genetics Unit, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh, EH4 2XU, UK

⁸The Alan Turing Institute, London, UK

*Corresponding authors:

Names: Catalina Vallejos and Riccardo Marioni

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

1 Abstract

2 Over the last decade, a plethora of blood-based DNA methylation biomarkers have been
3 developed to track differences in ageing, lifestyle, health, and biological outcomes. Typ-
4 ically, penalised regression models are used to generate these predictors, with hundreds
5 or thousands of CpGs included as potential features. However, in such ultra high-
6 dimensional settings, the effectiveness of these methods may be reduced.

7 Here, we introduce Related Trait-based Feature Screening (RTFS), a method for per-
8 forming CpG pre-selection for incident disease prediction models by utilising associations
9 between CpGs and health-related continuous traits. In a comparison with commonly used
10 CpG pre-selection methods, we evaluate resulting downstream Cox proportional-hazards
11 prediction models for 10-year type 2 diabetes (T2D) onset risk in Generation Scotland
12 (n=18,414). The top performing models utilised incident T2D EWAS (AUC=0.881,
13 PRAUC=0.279) and RTFS (AUC=0.877, PRAUC=0.277). The resulting models also im-
14 prove prediction over a model using standard risk factors only (AUC=0.841, PRAUC=0.194)
15 and replication was observed in the German-based KORA study (n=4,261)

16 RTFS is a flexible and generalisable framework that can help to refine biomarker
17 development for incident disease outcomes.

18 Introduction

19 Numerous studies have shown that levels of DNA methylation (DNAm) at various CpG
20 sites can correlate with health-related traits, such as body mass index (BMI), smoking
21 status [1], and incident diseases [2, 3, 4]. DNAm is an epigenetic modification whereby
22 methyl groups are dynamically attached and removed at various genomic positions (often
23 on the cytosine of a C-G dinucleotide; CpG) throughout an individual's lifetime. Blood-
24 based DNAm is of particular interest within cohort studies as its relatively non-intrusive
25 sample procedure makes it potentially suitable for clinical biomarker development, en-
26 abling the development of risk prediction models (e.g. to predict incident disease).

27 A major challenge in developing these prediction models is the selection of relevant

28 CpG sites for use as inputs. DNAm is commonly (and affordably) ascertained through the
29 use of arrays including the Illumina Infinium HumanMethylation450 and EPIC arrays,
30 which capture methylation information for $\sim 450,000$ and $\sim 800,000$ CpG sites, respec-
31 tively [5, 6]. In contrast, cohort sizes tend to be limited to a few thousand individuals.
32 This leads to an ultra high-dimensional setting in which the number of features or pre-
33 dictors (p) is much larger than the number of observations (n).

34 A typical approach to utilising high-dimensional data involves the application of pe-
35 nalisised regression models, both for feature selection and prediction (see e.g. [1, 7, 8]).
36 However, in ultra high-dimensional settings, the effectiveness of penalised regression may
37 be reduced [9, 10, 11]. A two-stage process has previously been suggested to address this,
38 where a pre-selection (screening) step is first applied to the data, before fitting penalised
39 regression models [12, 10]. The purpose of the pre-selection is to broadly filter out irrele-
40 vant features to reduce the number of potential predictors to a size suitable for penalised
41 regression (typically of polynomial order with respect to the sample size [11]).

42 One commonly used method for CpG pre-selection is variance-based filtering, whereby
43 the top k CpGs are retained after ranking them by decreasing variance, where k is arbi-
44 trarily chosen. This method helps to remove invariant CpG sites, but its performance may
45 be problematic, particularly with small effect and sample sizes [13]. Other approaches,
46 based on the correlation of each feature with the outcome, have been proposed for con-
47 tinuous (e.g. [12]) and time-to-event data (e.g. [14]), but some of these may introduce
48 problems related to post-selection inference [15] if the same data is used for screening and
49 model fitting. An alternative is to use domain knowledge (e.g. from external data) to
50 inform the screening. One such method involves pre-selecting CpGs that have previously
51 shown associations with the outcome in Epigenome-Wide Association Studies (EWAS).
52 If the associations have been found in an independent dataset, the chance of noise from
53 spurious correlations with the outcome is reduced. However, the pre-selection is lim-
54 ited to marginal associations between each CpG site and the outcome and availability of
55 EWAS results varies depending on the outcome. Another strategy, that can bypass the
56 need for feature pre-selection, is the application of principal components analysis (PCA;

57 or other dimensionality reduction techniques) to obtain a low-dimensional set of features
58 (e.g. the first 100 principal components) to be used as inputs. The latter has been shown
59 to potentially improve out-of-sample prediction in CpG-based models [16].

60 Here, we propose a Related Trait-based Feature Screening (RTFS) pipeline, using
61 information about continuous traits that are related to the outcome of interest to perform
62 feature pre-selection. For example, to predict time-to disease incidence, we selected a
63 range of measurements (e.g. BMI, smoking, alcohol consumption) typically related to a
64 broad set of health outcomes. Feature pre-selection can be then performed by applying
65 e.g. penalised regression on the continuous traits, with lower sample size requirements
66 than time-to-event data [17, 18, 19]. Power calculations for time-to-event data typically
67 depend on the number of case events per feature, which is often small compared to the
68 overall sample size. This is in contrast to the corresponding calculations for continuous
69 traits which are based on the total number of data points per feature.

70 We apply RTFS and other popular CpG pre-selection methods to the Generation
71 Scotland (GS) [20] cohort ($n = 18,414$), one of the world's largest studies including
72 genome-wide DNAm data paired with linkage to electronic health records (EHR). We
73 compare the performance of the different pre-selection methods as well as dimensionality
74 reduction using PCA in the development of epigenetic scores (EpiScores) - weighted sums
75 of CpG methylation values - used to predict time to incident type 2 diabetes (T2D). We
76 show that RTFS is competitive with the top existing EWAS-based filtering approach,
77 leading to an increase in predictive performance above standard T2D risk factors. We
78 also show the predictive performance increases of the EpiScores compared to genetic risk
79 factors using a T2D polygenic risk score (PRS). Finally, we validated the performance of
80 resulting EpiScores derived from RTFS and incident T2D EWAS-based filtering in the
81 KORA S4 cohort [21]. All analyses and results are reported in line with the TRIPOD
82 checklist [22] for reproducibility purposes and can be found in **Supplementary File 1**.
83 Analysis scripts are provided on GitHub at <https://github.com/marioni-group/rdfs>.

84 Results

85 RTFS

86 The proposed RTFS pipeline aims to aid feature pre-selection in ultra high-dimensional
87 settings when developing prediction models for incident disease risk. We focus on time-
88 to-event outcomes, but a similar pipeline could be applied to other types of outcomes
89 (e.g. binary or counts). RTFS borrows information from a set of traits that are related
90 to the outcome — or a broad set of outcomes — of interest. The process is illustrated in
91 **Figure 1**. First, (linear) lasso regression models are trained with each of the (continuous)
92 traits as the outcome. The resulting RTFS pre-selected CpG set consists of the union of
93 CpG sites retained from any of the continuous trait models. Here, nineteen continuous
94 traits were included in the RTFS pipeline: age, glucose, total cholesterol, high-density
95 lipoprotein (HDL) cholesterol, sodium, potassium, urea, creatinine, BMI, waist-hip ratio,
96 body fat percentage, systolic blood pressure, diastolic blood pressure, heart rate, forced
97 expiratory volume (FEV), forced vital capacity (FVC), alcohol consumption, smoking
98 and general cognitive ability. All of these were recorded at baseline (see **Methods**). For
99 all the continuous traits, the in-sample predictive performance for the corresponding lasso
100 model in the test set is given in **Supplementary Table 1**.

101 Cohort summary

102 After exclusions, our data consisted of 14,531 individuals from the GS cohort (see **Meth-**
103 **ods** and **Supplementary Figure 1**). This was divided into three non-overlapping sets:
104 to train the trait-specific models (feature pre-selection set for RTFS only; $n = 5,739$) as
105 well as to train ($n = 4,158$) and test ($n = 4,634$) the incident T2D prediction model.
106 After removal of missing values in the continuous traits, the pre-selection set consisted
107 of between $n = 4,872$ and $n = 5,739$ individuals depending on the trait (see **Sup-**
108 **plementary Table 2**). Summary information for the T2D training and test sets is
109 shown in **Table 1** and **Supplementary Figure 2**. Both sets had a highly imbalanced
110 case/control distribution with 3.2% (130/4,028) and 4.6% (213/4,634) having an incident

111 T2D diagnosis in the training and test sets, respectively.

112 **Prediction performance assessment**

113 When assessing predictive performance in the test set, two types of outcomes were con-
114 sidered: prediction of time to incident T2D diagnosis and a binary outcome given by
115 whether incident T2D occurred prior to 10 years after baseline (**Methods**). Predictive
116 performance using the time to incident T2D diagnosis was assessed using C-index and
117 Brier scores. C-index measures discrimination (agreement in the ranking between pre-
118 dicted risks and observed time-to-event values across pairs of individuals) while Brier
119 scores give a measurement of both model calibration and discrimination at a given time
120 point. In our experiments, Brier scores were evaluated at all integer time points from
121 $t = 1$ to $t = 10$ years (inclusive). Binary outcome prediction performance was assessed
122 using measures of discrimination - area under the receiver operating characteristics curve
123 (AUC) and area under the precision recall curve (PRAUC). Calibration of the predic-
124 tions generated by each model was also evaluated. Other measures, such as specificity
125 and sensitivity, across a range of probability classification thresholds are also provided.

126 **Prediction of incident T2D using risk factors only**

127 A Cox proportional-hazards (Cox PH) model in the test set using established risk fac-
128 tors (age, sex, BMI, hypertension and parent/sibling history of diabetes) as covariates
129 (referred to as the risk factors-only model) had a C-index of 0.828 for time-to-event out-
130 comes (Brier scores are shown in **Supplementary Table 3**). AUC and PRAUC were
131 0.841 and 0.194, respectively, when predicting if an incident T2D diagnosis occurred prior
132 to 10 years after baseline.

133 **Prediction of incident T2D using risk factors and DNAm**

134 We considered four methods for feature pre-selection (**Figure 2**; details in **Methods**):
135 filtering to sites on the 450k array (henceforth referred to as the EPIC-450k intersection);

136 filtering to the top 100k and 200k most variable CpGs; filtering to epigenome-wide signif-
137 icant CpGs from the EWAS literature (72 and 55 CpGs for incident and prevalent T2D,
138 respectively); and filtering to the 5,468 RTFS CpGs identified from the lasso models on
139 the continuous traits. We also considered applying PCA to the EPIC-450k intersection
140 and to the top 200k most variable CpGs, with the PCs explaining a cumulative variance
141 $> 95\%$ taken forward as features. This led to selecting 3,734 and 3,652 PCs, respectively.

142 The greatest C-index values were achieved from using incident T2D EWAS-based
143 filtering and RTFS (both 0.866). All C-index and Brier score values are shown in **Sup-**
144 **plementary Table 3**. Incident T2D EWAS-based filtering and RTFS resulted in the
145 lowest two Brier scores for all time points, suggesting that those methods consistently
146 performed in the top two models in terms of calibration and case/control discrimination.

147 **Table 2** shows the AUC and PRAUC values obtained from incremental Cox PH
148 models corresponding to the addition of an EpiScore, derived from each pre-selection
149 method or PCA, to the risk factors-only model. Incident T2D EWAS-based filtering
150 achieved the highest AUC (0.881) and PRAUC (0.279). Corresponding ROC curves for
151 the incident T2D EWAS-based filtering, RTFS and the risk factors-only models are shown
152 in **Figure 3**. We evaluated the robustness of this ranking by considering the number
153 of times each method was ranked in the top n methods across 1,000 bootstrap runs is
154 plotted in **Figure 4**. Incident T2D EWAS-based filtering had the highest frequency
155 of first rankings across the bootstraps in both AUC and PRAUC. RTFS also performed
156 consistently well with both methods ranking in the top three in the majority of bootstraps.

157 Differences in model calibration between the incident T2D EWAS EpiScore model,
158 RTFS EpiScore model and risk factors-only model are shown in **Supplementary Fig-**
159 **ure 4**. The incident T2D EWAS EpiScore and RTFS EpiScore models show stronger
160 calibration performance when compared to the risk factors-only model. All three models
161 plotted show underestimation of risk below a predicted probability of around 0.5 and
162 overestimation of risk otherwise.

163 **Supplementary Figure 3** shows how confusion matrix values vary across the prob-
164 ability classification threshold range for the risk factors-only, RTFS and the incident T2D

165 EWAS-based filtering EpiScore model in the test set. Overall, the incident T2D EWAS-
166 based filtering and RTFS EpiScore model improve the classification of cases with respect
167 to the risk factors-only model (increase in true positives and decrease in false negatives)
168 while showing a slight decrease in the correct classification of controls. The differences
169 in correctly classified individuals in terms of sensitivity, specificity, positive predictive
170 value (PPV) and negative predictive value (NPV) between the RTFS EpiScore and risk
171 factors-only models are also given in **Supplementary Table 4**.

172 **Comparison of incident T2D EpiScore and polygenic risk score** 173 **prediction performance**

174 To assess the added value of the EpiScores against genetic risk factors on predictive
175 performance, two additional Cox PH models were fit to the GS test set that included
176 a polygenic risk score (PRS) for incident T2D [23]. These consisted of a model using
177 the standard risk factors plus the PRS, as well as a second model which also included
178 the EpiScore derived from incident T2D EWAS-based filtering (the top performing pre-
179 selection method). These two models showed AUC values of 0.857 and 0.892 respectively.
180 PRAUC values were 0.212 and 0.302 and C-index values were 0.843 and 0.876. The PRS
181 gave a smaller increase in each of these metrics above standard risk factors compared
182 to the incident T2D EWAS EpiScore (AUC=0.881, PRAUC=0.279, C-index=0.866);
183 however, without pre-selection of CpG sites, the EpiScore gives smaller increases (EPIC-
184 450k EpiScore AUC=0.855, PRAUC=0.208, C-index=0.841. The largest increase was
185 given when using both the PRS and EpiScore in the model, showing additive increases
186 from both scores over using risk factors only.

187 **Validation of RTFS and EPIC-450k intersection EpiScores in the** 188 **KORA S4 cohort**

189 Performance in KORA S4 was only evaluated for the binary T2D incidence outcome
190 (diagnosis within 10 years of baseline date) as time to T2D diagnosis data was not avail-

191 able. The logistic risk factors-only model fit to the KORA S4 cohort showed an AUC
192 and PRAUC of 0.797 and 0.294, respectively. The logistic models including risk fac-
193 tors plus either the RTFS or EPIC-450k incident T2D EpiScore resulted in AUCs of
194 0.806 and 0.798, respectively. Corresponding PRAUC values were 0.295 and 0.293 (see
195 **Supplementary Table 5**).

196 **Overlap of pre-selected CpG sites**

197 The continuous trait lasso models in RTFS selected between 49 and 864 CpG sites per
198 trait (5,468 in the union). **Figure 5** shows the number of CpG sites selected for each
199 trait and the selection overlap between traits. This shows that the majority of CpG sites
200 were selected exclusively for a single trait. Notable overlaps were present between BMI,
201 waist-to-hip ratio and body fat as well as between systolic and diastolic blood pressure.

202 **Supplementary Figure 5** shows the number of CpGs pre-selected across all methods
203 and their overlap. Over half of the RTFS pre-selected CpGs were not in the top 200,000
204 CpGs by variance. Additionally, a small proportion of the incident and prevalent EWAS
205 CpGs overlapped with the RTFS CpGs (see **Supplementary Figure 6**)

206 **Discussion**

207 In this study, we explored the use of different feature pre-selection methods in the context
208 of ultra-high dimensional DNAm data (where the number of features largely exceeds the
209 number of observations). We introduce RTFS, which borrows information from a broad
210 set of health-related traits to identify a suitable set of CpG sites that can be used as input
211 in the development of risk prediction models for incident disease. Using type 2 diabetes as
212 a case study, we compared the performance of RTFS against a range of other commonly-
213 applied CpG pre-selection (and dimensionality reduction) approaches. Consistent with
214 [24], the inclusion of an EpiScore generally improved discrimination performance with
215 respect to the standard risk factors-only model. However, the improvement was not
216 uniform across the different methods: with only marginal improvements in the absence of

217 feature pre-selection when training the EpiScore. Our analysis also shows that EpiScores
218 can improve predictive performance compared to the use of genetic information via an
219 existing incident T2D PRS.

220 Incident T2D EWAS-based filtering resulted in the highest AUC (0.881) and PRAUC
221 (0.279) for 10-year incident disease prediction, with a notable increase in the correct
222 classification of cases and a small decrease in correct classification of controls. External
223 validation in KORA supported this, although it showed smaller improvements compared
224 to an earlier study which used a larger training set for the incident T2D EpiScores [24].

225 While filtering to significant CpG sites from an incident T2D EWAS study was the
226 highest performing model, it is reliant on the existence of large-scale EWAS studies
227 for T2D, something that may not be generally available for other diseases of interest.
228 RTFS bypasses this requirement and led to similar performance metrics (AUC = 0.877;
229 PRAUC = 0.277). It was also consistently ranked amongst the top performing models
230 in our bootstrap experiments. Additionally, the continuous traits used for RTFS were
231 primarily general health-related measures and not necessarily specific to T2D. Therefore,
232 the resulting set of RTFS CpGs may be applicable to other diseases and could potentially
233 be used as a general panel of morbidity-related sites for risk prediction.

234 We used a pre-specified set of continuous trait to perform CpG pre-selection in RTFS.
235 While we evaluated the predictive performance of each trait-specific lasso model, future
236 studies could investigate the impact of including or excluding continuous traits e.g. based
237 on a range of different performance thresholds. Additional studies could also investigate
238 other variable selection methods for RTFS continuous traits, for example using elastic-net
239 [25], as well as more general methods for risk prediction (e.g. random survival forests[26]).
240 Future studies could also consider using DNAm-based predictions for each trait directly
241 as predictors in downstream models, similar to previous approaches (e.g. the protein
242 EpiScores in [27] or the approach used to develop the GrimAge epigenetic clock [28]).

243 Access to the GS cohort enabled us to demonstrate the use of RTFS in one of the
244 largest cohorts of its kind — with three independently-processed sets of DNAm data,
245 which allowed for separate training, testing and RTFS pre-selection datasets. In addi-

246 tion, comprehensive information on incident T2D diagnoses was available through ex-
247 tensive linkage to electronic health records. Availability of both genetic and epigenetic
248 data allowed for a direct performance comparison between risk scores derived from each
249 data source and showed the benefit of using DNAm data, which can better reflect health-
250 associated changes within individuals' lifetimes. While the inclusion of DNAm resulted in
251 considerable predictive performance increases compared to using risk factors only in the
252 GS test set, these differences were small when applied to the KORA validation cohort.
253 The generalisability of our results is limited by the characteristics of the GS cohort: GS
254 participants are generally healthier, wealthier and have a different age-sex distribution
255 to the general population [29]. Similarities in these socio-demographic characteristics
256 within GS may have resulted in positive bias in the performance of RTFS. Given that the
257 models including DNAm data with and without CpG preselection both showed small per-
258 formances differences when compared to a risk factors-only model in KORA, further work
259 could explore the impact of factors such as the number of incident cases and availabil-
260 ity of primary versus secondary care data for T2D disease ascertainment. Additionally,
261 both the development and validation cohorts consisted of individuals from predominantly
262 white European ancestries. Further validation is required to evaluate the generalisability
263 of RTFS to other populations and genetic ancestries.

264 In conclusion, our study reiterated the need for pre-selection as an important step
265 in DNAm-based risk prediction models. We introduced and evaluated an effective pre-
266 selection method, RTFS, utilising information from health-related traits with the poten-
267 tial for application in predictive models for other incident diseases in future studies.

268 **Methods**

269 **Generation Scotland (GS) DNAm data**

270 The data used for this study were from the Generation Scotland (GS) cohort, recruited
271 from across Scotland between 2006 and 2011. This consists of 23,960 volunteers aged
272 18-99 at baseline (recruitment date). Of these, 18,414 have genome-wide DNAm data

273 available, ascertained from blood samples taken at baseline. DNAm quality control is
274 detailed in [24]. DNAm measurements were obtained in three large sets, processed in
275 2017 (set 1, $n = 5,087$), 2019 (set 2, $n = 4,450$) and 2021 (set 3, $n = 8,877$). Set 2 was
276 used as the training set for incident T2D and set 3 was used for feature pre-selection. Set
277 1 was used as the test set for incident T2D. Sets 1 and 3 contained related individuals
278 (genetic relationship matrix (GRM) threshold > 0.05), both within and between sets.
279 There were also related individuals between sets 2 and 3. To avoid the presence of
280 families with individuals across the training and test sets, individuals in set 3 with a
281 family member present in set 1 were excluded from the analyses ($n_{excluded} = 3,138$). To
282 maintain compatibility with previous studies using the Illumina 450K array, the CpGs
283 were filtered to those present in both the 450K and EPIC arrays (EPIC-450k intersection).

284 A range of traits were also recorded at baseline via questionnaire or clinical appoint-
285 ment. These included (units listed within parenthesis): age (*years*), glucose (*millimoles*
286 *per litre; mmol/L*), total cholesterol (*mmol/L*), high-density lipoprotein (HDL) choles-
287 terol (*mmol/L*), sodium (*mmol/L*), potassium (*mmol/L*), urea (*mmol/L*), creatinine
288 (*mmol/L*), BMI (kg/m^2), waist-hip ratio, body fat percentage, systolic blood pressure
289 (*millimetres of mercury; mmHg*), diastolic blood pressure (*mmHg*), heart rate (*beats*
290 *per minute; bpm*), forced expiratory volume (FEV) (*L*), forced vital capacity (FVC) (*L*),
291 alcohol consumption (*units/week*), smoking (*pack years*) and general cognitive ability.
292 The latter was defined as the first unrotated principal component from a PCA of four
293 cognitive tests (logical memory, digit symbol, verbal fluency and vocabulary), scaled to
294 mean of 0 and standard deviation of 1 [30].

295 **Continuous trait preprocessing**

296 The 19 baseline traits listed above were used as the continuous traits for RTFS. These
297 were processed separately within each data set, to remove outliers, and to regress out
298 age and sex (after trait-specific transformation, if applied). Trait-specific transformation
299 included adding 1 to each value of alcohol consumption and smoking, prior to a natural
300 log-transform. Glucose and BMI were also log-transformed. Outliers were defined as

301 points greater than 4 standard deviations away from the mean. This is with the exception
302 of BMI for which outliers were defined as a BMI < 18 or BMI > 50. A linear regression
303 model with *age*, *age*² (to include non-linear effects with age) and *sex* as covariates was
304 then fit to each continuous trait. The resulting residuals were kept for further analyses.
305 For FEV and FVC, height was also included in the linear regression. Missing values in
306 each continuous trait were treated as missing-at-random and corresponding individuals
307 were removed from the training set when fitting the predictive model for the respective
308 trait. The number of missing values for each trait is given in **Supplementary Table 2**.

309 **Time to incident Type 2 Diabetes (T2D)**

310 History of disease diagnoses (prevalent and incident) was ascertained via data linkage
311 to NHS Scotland health records. Secondary care (hospital) records from January 1980
312 to April 2022 were available for all subjects, with disease diagnoses encoded using ICD-
313 9/10. Due to restricted consent from data controllers, only partial linkage to primary care
314 (general practice; GP) records was available (a subset of general practice centres were
315 unable to provide data): only available for 35% ($n = 3,191$, $n_{training} = 1,421$, $n_{test} =$
316 $1,770$) of individuals in the incident T2D training and test sets. Primary care records
317 cover the period from January 1980 to October 2020 and use Read2 codes to record
318 disease diagnoses.

319 Hospital record-derived prevalent and incident T2D cases were defined as individuals
320 with an E11* ICD-10 code or 250.0/250.1 ICD-9 code. GP record-derived cases were
321 defined using a set of diabetes-related Read2 codes. A full list of ICD-9/10 and Read2
322 codes is provided in **Supplementary Table 6**. Type 1 and juvenile diabetes cases were
323 treated as controls (no T2D). Additional prevalent cases were identified from self-reported
324 history in a baseline questionnaire. All prevalent cases were removed. For incident cases,
325 time-to-event (*years*) was calculated as the time from baseline to disease onset (first T2D
326 record) for cases, and to censoring for controls. Controls were censored at the latest date
327 of available hospital records (April 2022) or time-to-death, whichever happened sooner.

328 For the individuals with both primary and secondary care records, a comparison be-

329 tween time-to-event outcomes derived from hospital and GP records was used to assess
330 possible delays in hospital diagnoses. As a sensitivity analysis (RTFS only), the end
331 of GP follow-up (October 2020) was also considered as a censoring date for those indi-
332 viduals (in the absence of a hospital diagnosis). **Supplementary Table 7** shows the
333 AUCs, PRAUCs and incremental Cox PH model coefficient estimates from this analy-
334 sis. Differences between the two outcome derivations were minor in terms of both the
335 discrimination metrics and coefficient estimates.

336 **T2D risk factors**

337 Risk factors used in the incremental EpiScore models included age, sex, BMI, hyperten-
338 sion and parent/sibling history of diabetes. Hypertension and parent/sibling history of
339 diabetes were defined as self-reported in the baseline questionnaires. While many T2D
340 risk factors have been identified, we based these on the most utilised factors in existing
341 risk scores according to [31]. These five risk factors were used as variables in the risk
342 factors-only model.

343 **Related Trait-based Feature Screening (RTFS)**

344 Linear lasso [32] was applied to each continuous trait (after pre-processing) using set 3.
345 Lasso is a penalised regression method which shrinks regression coefficients to be small,
346 forcing some to be exactly equal to zero. As such, it performs feature selection by keeping
347 only the features with non-zero coefficients. The strength of the penalty is controlled by
348 a hyper-parameter λ . Five-fold cross-validation was used to select λ , to minimise the
349 mean squared-error of out-of-sample predictions. Lasso models were fit using the glmnet
350 R package version 4.1-1 [33]. For computational efficiency, the top 200,000 sites with the
351 highest marginal variance were used as inputs. The union of lasso-selected CpGs from
352 the final set of continuous trait models (hereafter referred to as the RTFS set) were used
353 as input to predict time-to T2D incidence.

354 As RTFS performs feature pre-selection via linear models applied to a set of continuous
355 traits, it implicitly assumes that each trait can be predicted by a linear combination of

356 CpGs. The predictive ability of DNAm for each trait was quantified in a test set (GS set
357 1) using the percentage of variance explained R^2 .

358 **Alternative approaches for feature pre-selection**

359 Initially, the EPIC-450k intersection set was used without pre-selection. Two commonly
360 used approaches for feature pre-selection were then considered as an alternative to RTFS:

361 **Highest variance.** The per-feature variance is calculated, and the top p features with
362 the highest variance in set 3 are pre-selected. For the T2D analysis, we used $p = 100,000$
363 and $p = 200,000$.

364 **EWAS-based filtering.** Existing EWAS analysis for incident or prevalent disease
365 were used to pre-select CpGs. For the T2D analysis, CpGs identified to be statistically sig-
366 nificant by two recent large meta analyses for incident and prevalent T2D were included.
367 The first study [34] consisted of five European cohorts ($N_{cases} = 1,250$, $N_{controls} = 1,950$)
368 and identified 76 differentially methylated CpG sites ($p < 1.1 \times 10^{-7}$) for incident T2D. Af-
369 ter filtering these to those present in the EPIC-450k intersection, 72 CpG sites remained.
370 The second [35] consisted of four European cohorts ($N_{cases} = 340$, $N_{controls} = 3,088$)
371 identifying 58 differentially methylated CpG sites ($p < 1.0 \times 10^{-5}$) for prevalent T2D (55
372 post EPIC-450k intersection filtering). The full list of CpG sites identified from the two
373 studies is shown in **Supplementary Table 8**.

374 **Dimensionality reduction**

375 As an alternative to feature pre-selection, we also explored whether dimensionality re-
376 duction techniques can be used to create a low-dimensional set of features to be used as
377 input when predicting T2D. Here, we focus on PCA (as in [16]). In this study, we applied
378 PCA (in set 2) to the 450k-EPIC intersection and the top 200k CpGs by variance. PCs
379 were ordered by the variance explained in set 2 and the top PCs required to explain 95%
380 of the variance were kept for the final T2D model.

381 **Incident T2D EpiScore**

382 Using the CpGs identified by each pre-selection method (or top PCs, where appropri-
383 ate), a Cox PH elastic-net model [36] was fit to the set 2 DNAm data (training set) using
384 time-to-T2D incidence as the outcome. Similar to lasso, elastic-net provides a regularised
385 model fit, reducing overfitting. The strength of the regularization is controlled by hyper-
386 parameters λ and α . If $\alpha > 0$, the model performs feature selection by setting a subset of
387 coefficients to 0. Hyperparameters were selected using 9-fold cross-validation. Lambda
388 was optimised using the `cv.glmnet` function. Alpha was selected by testing values between
389 0 and 1 (inclusive) in increments of 0.1 and selecting the value which maximised mean
390 partial-likelihood across the nine folds. The linear predictor from the resulting Cox PH
391 elastic-net model was defined as an incident T2D EpiScore.

392 **T2D Incremental Modelling**

393 The incident T2D EpiScores obtained after applying each feature pre-selection (or dimen-
394 sionality reduction) method were subsequently applied to set 1 (test set) in an incremental
395 modelling approach. Firstly, a risk factors-only model was fit using a Cox PH model and a
396 set of known T2D risk factors (listed above) as predictors. For each pre-selection method,
397 Cox PH models were then fit using the same variables as the risk factors-only model plus
398 the corresponding EpiScore. Full details on the incremental modelling calculations are
399 given in the **Supplementary Note**.

400 **Incident T2D Polygenic Risk Score**

401 To compare the differences in predictive performance of the EpiScores to genetic risk
402 factors, two additional Cox PH models were fit in the test. The first included the standard
403 risk factors plus a T2D PRS [23]. The second included these same variables plus the top
404 performing incident T2D EpiScore (incident T2D EWAS-based filtering).

405 Predictive performance evaluation

406 Predictive performance for each of the Cox PH models above was evaluated on the test
407 set. Two types of prediction outcomes were used: time-to-T2D diagnosis and a binary
408 outcome defined by whether a T2D diagnosis was recorded within 10 years from baseline.

409 For the time-to-T2D outcome, C-index and Brier scores were calculated using the
410 SurvMetrics R package (version 0.5.0) [37]. C-index gives a measure of discrimination
411 for a model, defined as proportion of concordant pairs of individuals predicted by the
412 model. This value is between 0 and 1 (inclusive) with higher scores representing better
413 discrimination. A pair of individuals is concordant if the individual with the smaller
414 time-to-event is given a greater risk by the model. The Brier score measures both dis-
415 crimination and calibration, calculated as the mean square difference between the true
416 classes (i.e. whether a T2D diagnosis has occurred) and the predicted probabilities at a
417 given time point. Brier scores range between 0 and 1 (inclusive) and lower scores rep-
418 resent better discrimination and calibration. Brier scores were evaluated at each integer
419 time point from $t = 1$ to $t = 10$.

420 For the binary 10-year T2D onset outcome, predictions were calculated as one minus
421 the estimated 10-year survival probability. This calculation was based on the Breslow
422 estimator [38] for the cumulative baseline hazard. This calculation is detailed in the
423 **Supplementary Note**. Censored individuals were defined as controls when assessing
424 predictive performance. Discrimination metrics including area under the receiver operat-
425 ing characteristics curve (AUC) and the area under the precision-recall curve (PRAUC)
426 were compared.

427 Confusion matrix metrics were also assessed by calculating the number of true/false
428 positives/negatives using the ten-year onset probabilities and a range of discrimination
429 thresholds between 0 and 1, in increments of 0.1. Calibration was assessed by plotting
430 loess calibration curves using the `valProbggplot` function in the `CalibrationCurves` R
431 package (version 2.0.0) [39]. These show the observed event proportions plotted against
432 the predicted event probabilities.

433 To assess the robustness of the relative rankings for the pre-selection methods, the

434 incremental modelling was repeated in 1,000 bootstrap samples of the test set for each
435 EpiScore model. For each bootstrap sample, the EpiScore and risk factors-only models
436 (**Table 2**) were ranked based on their AUC and PRAUC estimates. The number of times
437 that each method was included in the top n ranks was calculated ($n = 1, \dots, 10$).

438 **Overlap of pre-selected CpG sites**

439 The sets of CpGs selected across the continuous trait lasso models were analysed using
440 an UpSet plot (UpSetR R package, version 1.4.0 [40]) showing, the number of CpG
441 sites selected across all of the traits in each combination of continuous traits. The same
442 visualisation method was used for analysing the overlap between CpGs selected using
443 each pre-selection method.

444 **Validation of RTFS and EPIC-450k intersection EpiScores in the** 445 **KORA S4 cohort**

446 The incident T2D EpiScores derived from the RTFS and EPIC-450k intersection CpGs
447 were validated in a subset of the German-based KORA S4 cohort, which consisted of
448 1,451 individuals aged 25-74 years and recruited in southern Germany. The subset was
449 defined by individuals with DNAm and incident T2D data available, after removing
450 prevalent cases at baseline. Missing CpG values in the DNAm data were mean-imputed
451 and individuals with missing health measures were removed from the dataset.

452 The prediction outcome was defined as the occurrence of a T2D diagnosis within 10
453 years after individuals' baseline date. A time-to-event outcome was not used for validation
454 as time-to-T2D diagnosis data was not available in KORA S4.

455 Validation was performed using incremental logistic models. Firstly, a risk factors-
456 only model was fit to the KORA S4 subset using age, sex, BMI, hypertension and parent
457 history of diabetes as variables. Then, two additional logistic models were fit using
458 the risk factors plus each of the RTFS and EPIC-450k intersection EpiScores. Prediction
459 performance was evaluated by calculating AUC and PRAUC for each of the three models.

460 The incident T2D EWAS EpiScore was not validated in KORA S4 as the correspond-
461 ing EWAS meta-analysis included KORA participants.

462 Additional details of participant follow-up, ascertainment of incident T2D diagnoses
463 and preprocessing numbers are provided in the **Supplementary Note**.

464 **Code and data sharing**

465 Analysis scripts for this study are available at <https://github.com/marioni-group/rdfs>.
466 According to the terms of consent for Generation Scotland participants, access to data
467 must be reviewed by the Generation Scotland Access Committee. Applications should
468 be made to access@generationscotland.org. The informed consent given by the KORA
469 S4 study participants does not cover data posting in public databases. However, data
470 are available upon request from the KORA Project Application Self-Service Tool. Data
471 requests can be submitted online (<https://epi.helmholtz-muenchen.de/>) and are subject
472 to approval by the KORA board.

473 **Ethics**

474 All components of Generation Scotland received ethical approval from the NHS Tayside
475 Committee on Medical Research Ethics (REC Reference Number: 05/S1401/89). Gener-
476 ation Scotland has also been granted Research Tissue Bank status by the East of Scotland
477 Research Ethics Service (REC Reference Number: 20-ES-0021), providing generic ethical
478 approval for a wide range of uses within medical research. Written, informed consent was
479 provided by Generation Scotland participants.

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

484 Acknowledgements

485 This research was funded in whole, or in part, by the Wellcome Trust [104036/Z/14/Z,
486 108890/Z/15/Z, 216767/Z/19/Z]. For the purpose of open access, the author has ap-
487 plied a CC BY public copyright licence to any Author Accepted Manuscript version
488 arising from this submission. Generation Scotland received core support from the Chief
489 Scientist Office of the Scottish Government Health Directorates (CZD/16/6) and the
490 Scottish Funding Council (HR03006) and is currently supported by the Wellcome Trust
491 (216767/Z/19/Z). DNA methylation profiling of the Generation Scotland samples was
492 carried out by the Genetics Core Laboratory at the Edinburgh Clinical Research Facil-
493 ity, Edinburgh, Scotland and was funded by the Medical Research Council UK and the
494 Wellcome Trust (Wellcome Trust Strategic Award "Stratifying Resilience and Depres-
495 sion Longitudinally" (STRADL; Reference 104036/Z/14/Z). The DNA methylation data
496 assayed for Generation Scotland was partially funded by a 2018 NARSAD Young Inves-
497 tigator Grant from the Brain & Behavior Research Foundation (Ref: 27404; awardee: Dr
498 David M Howard) and by a JMAS SIM fellowship from the Royal College of Physicians
499 of Edinburgh (Awardee: Dr Heather C Whalley). Y.C. is supported by the University of
500 Edinburgh and University of Helsinki joint PhD program in Human Genomics. C.A.V.
501 is a Chancellor's Fellow funded by the University of Edinburgh. R.E.M. is supported by
502 Alzheimer's Society major project grant AS-PG-19b-010.

503 The KORA study was initiated and financed by the Helmholtz Zentrum München
504 – German Research Center for Environmental Health, which is funded by the German
505 Federal Ministry of Education and Research (BMBF) and by the State of Bavaria. Fur-
506 thermore, KORA research has been supported within the Munich Center of Health Sci-
507 ences (MC-Health), Ludwig-Maximilians-Universität, as part of LMUinnovativ and is
508 supported by the DZHK (German Centre for Cardiovascular Research). The KORA
509 study is funded by the Bavarian State Ministry of Health and Care through the research
510 project DigiMed Bayern (www.digimed-bayern.de).

511 Conflicts of interest

512 R.E.M is an advisor to the Epigenetic Clock Development Foundation and Optima Part-
 513 ners. All other authors declare no competing interests.

	Training		Test	
	Cases	Controls	Cases	Controls
n	130	4,028	213	4,421
Time-to-event (Years to Onset or Censoring)	6.2 (3.1)	12.3 (2.1)	6.1 (3.4)	12.7 (1.8)
Age (Baseline)	58.3 (10.9)	51.1 (13.2)	55.6 (9.8)	48.1 (14.1)
Age (Onset or Censoring)	64.5 (11.0)	63.4 (13.1)	61.8 (10.1)	60.8 (13.8)
Sex (Male)	67 (51.5%)	1,724 (42.8%)	108 (50.7%)	1,645 (37.2%)
BMI (kg/m ²)	31.8 (6.3)	26.5 (4.7)	32.2 (6.0)	26.6 (5.1)
Self-reported Parent or Sibling Diabetes	55 (42.3%)	700 (17.4%)	95 (44.6%)	861 (19.5%)
Self-reported Hypertension	51 (39.2%)	555 (13.8%)	86 (40.4%)	597 (13.5%)
GP records available	62 (46.3%)	1,359 (32.8%)	110 (54.5%)	1,660 (36.1%)

Table 1: Summary details for the incident T2D training and test sets.

Incremental Model Variables (in addition to risk factors)	AUC	PRAUC	C-index	Alpha	Lambda
Risk Factors (RF) only	0.841	0.194	0.828	NA	NA
RF + PCA EPIC-450k EpiScore	0.849	0.206	0.837	0.2	0.325
RF + PCA Top 200k by Variance EpiScore	0.853	0.205	0.841	0	9.716
RF + EPIC-450k EpiScore	0.855	0.208	0.841	0.8	0.014
RF + PRS	0.857	0.212	0.843	NA	NA
RF + Top 100k by Variance EpiScore	0.864	0.215	0.852	0.7	0.012
RF + Prevalent T2D EWAS EpiScore	0.869	0.255	0.858	0.6	0.004
RF + Top 200k by Variance EpiScore	0.872	0.233	0.860	0.5	0.017
RF + Prevalent and Incident T2D EWAS EpiScore	0.873	0.262	0.858	1	0.002
RF + RTFS EpiScore	0.877	0.277	0.866	0.5	0.012
RF + Incident T2D EWAS EpiScore	0.881	0.279	0.866	0.5	0.006
RF + PRS + Incident T2D EWAS EpiScore	0.892	0.302	0.876	0.5	0.006

Table 2: Incremental modelling performance metrics for each pre-selection / PCA result calculated in the GS test set.

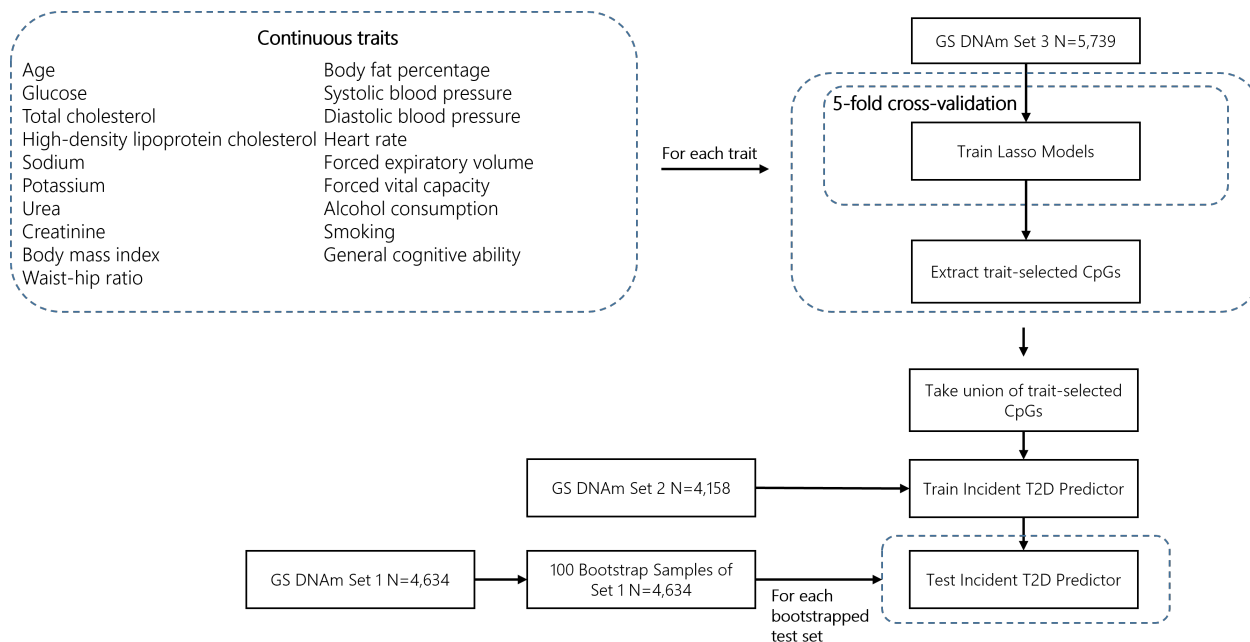


Figure 1: The RTFS pipeline applied to Generation Scotland.

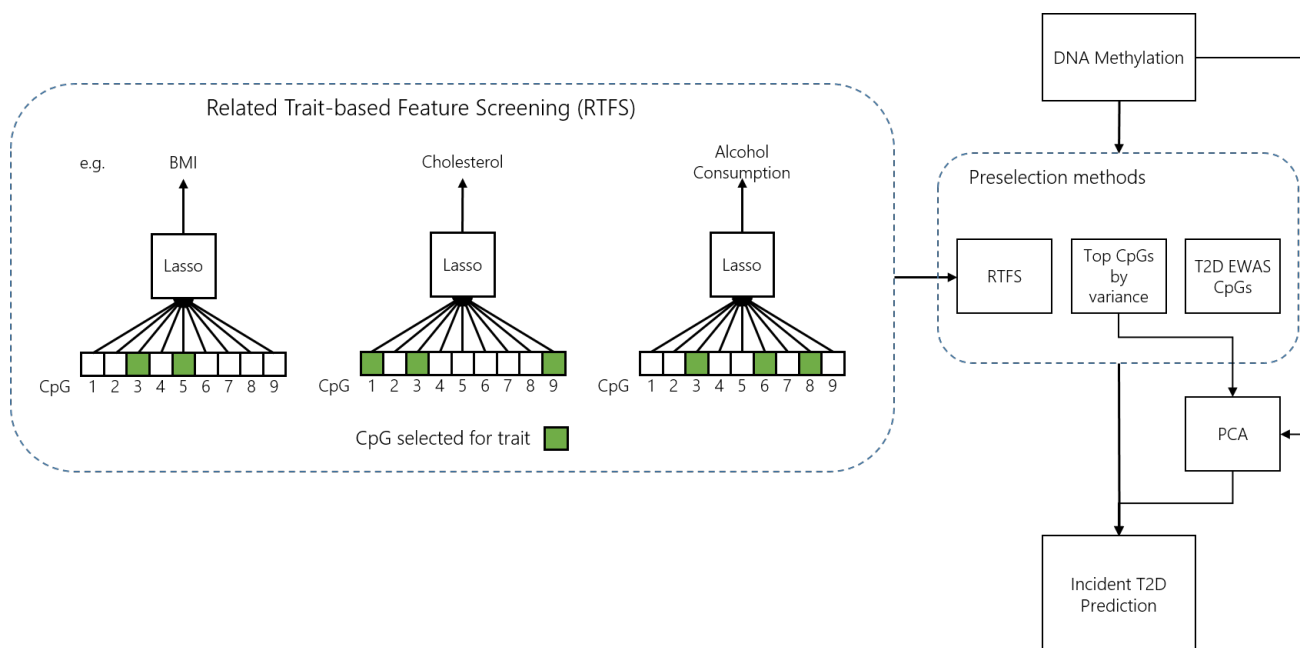


Figure 2: An overview of the pre-selection comparison pipeline.

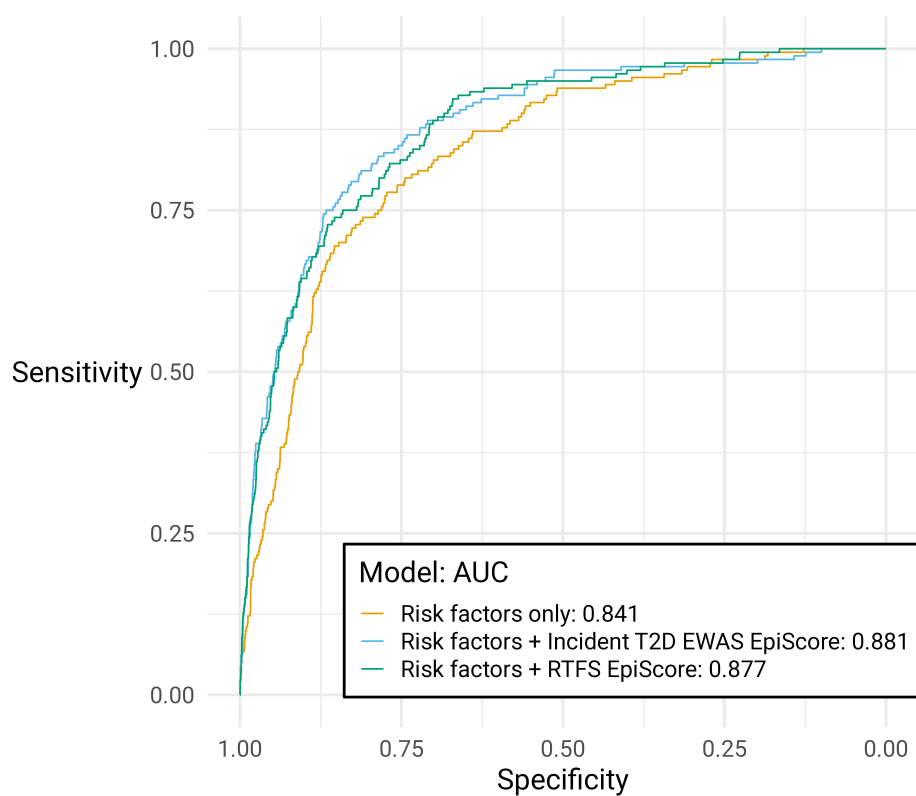


Figure 3: ROC curves for incremental incident T2D models. Results are shown for the model including risk factors only in addition to the models using RTFS and incident T2D EWAS-based filtering.

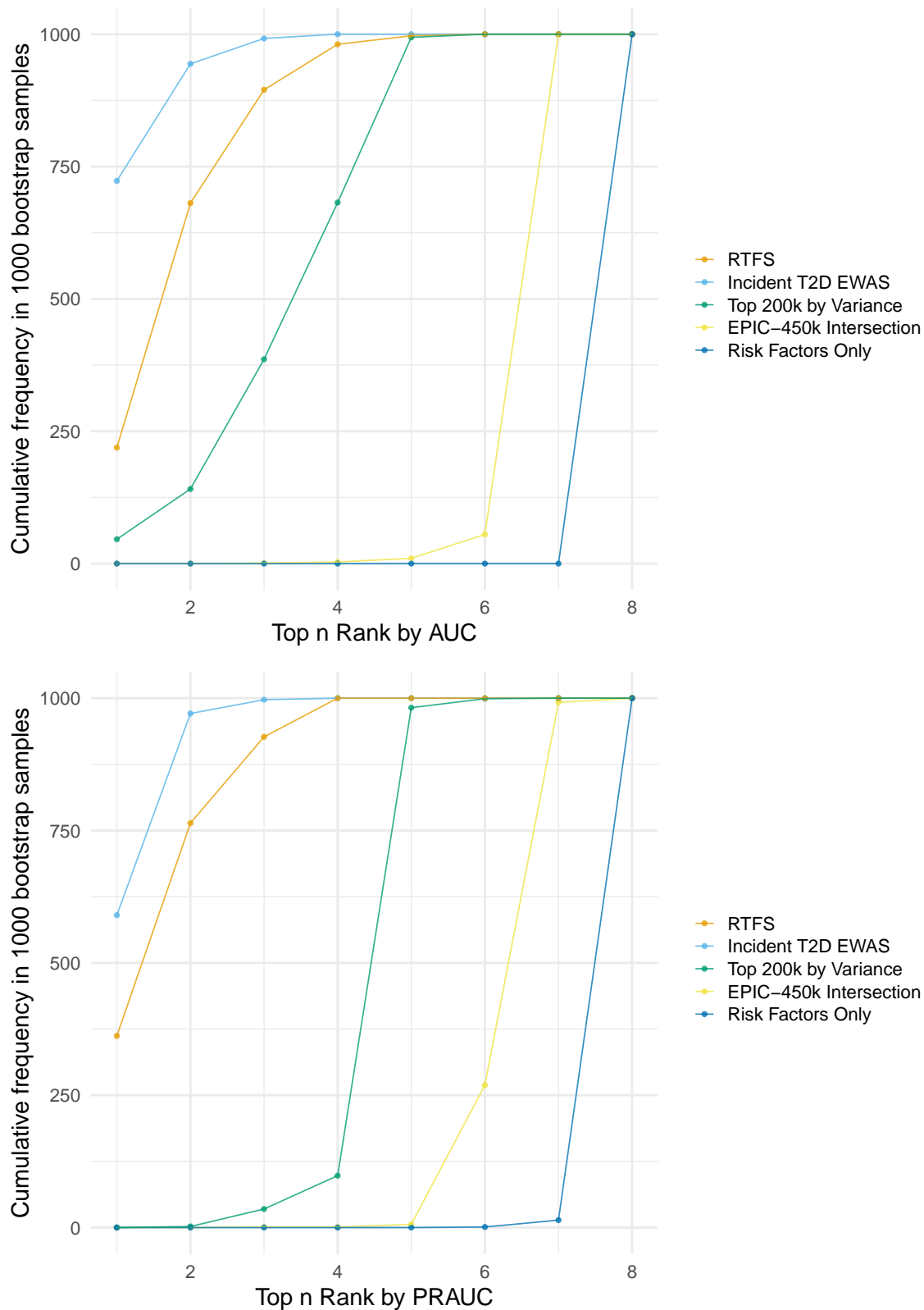


Figure 4: Rank order cumulative frequencies for AUC (top) and PRAUC (bottom) values for pre-selection methods across 1,000 bootstrap samples of the test set. The plots show, for each method, the number of bootstraps in which the method ranked in the top n in terms of their respective AUC/PRAUC. Models shown include: Related Trait-based Feature Selection (RTFS), Incident T2D EWAS, Top 200k by Variance, EPIC-450k Intersection, and Risk Factors Only.

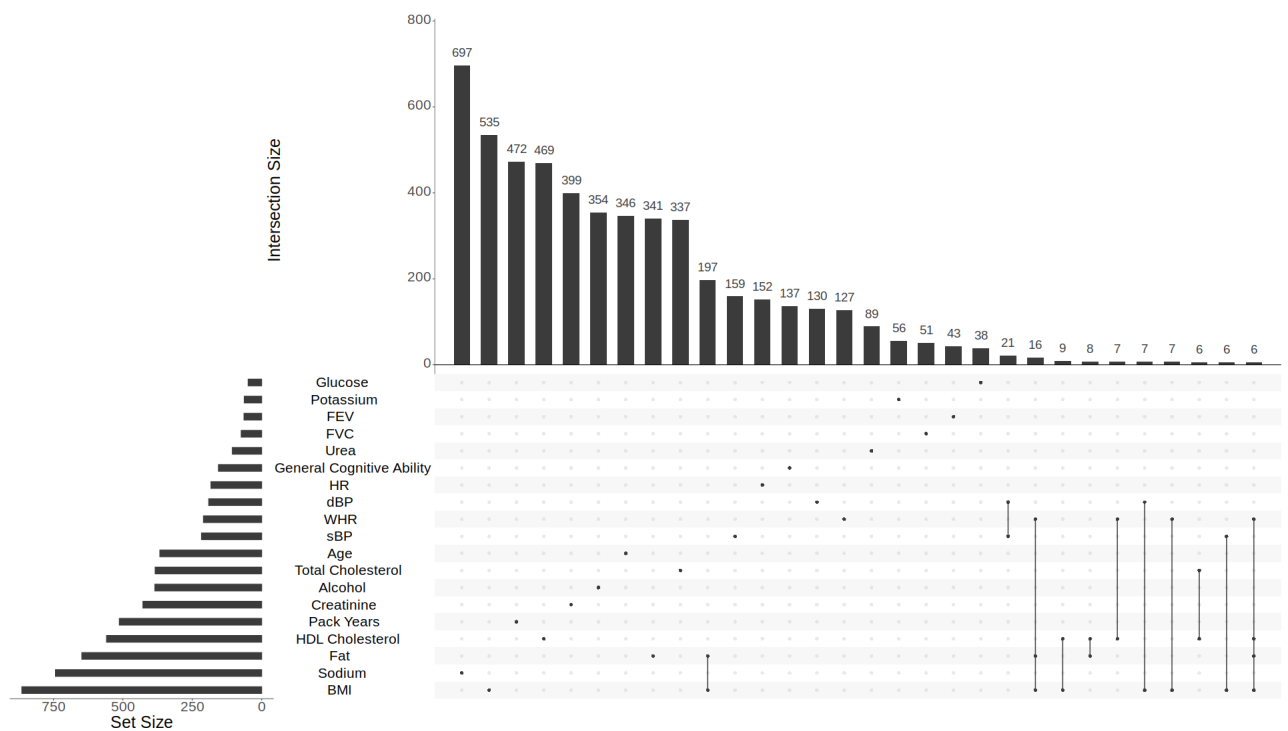


Figure 5: UpSet plot showing number of CpGs selected for each continuous trait and overlaps between traits. The frequency of the top 30 trait combinations are shown. Each column represents the number of CpGs pre-selected for the corresponding specific combination of traits. This was generated with the "distinct" option, meaning the presence or absence of a point in a column explicitly corresponds to the presence or absence of the corresponding trait in the set

514 References

- 515 [1] Daniel L McCartney et al. “Epigenetic prediction of complex traits and death”. In:
516 *Genome biology* 19 (2018), pp. 1–11.
- 517 [2] Yujing Xia, Alison Brewer, and Jordana T Bell. “DNA methylation signatures of
518 incident coronary heart disease: findings from epigenome-wide association studies”.
519 In: *Clinical Epigenetics* 13.1 (2021), pp. 1–16.
- 520 [3] Eliza Fraszczyk et al. “Epigenome-wide association study of incident type 2 dia-
521 betes: a meta-analysis of five prospective European cohorts”. In: *Diabetologia* 65.5
522 (2022), pp. 763–776.
- 523 [4] Robert F. Hillary et al. “Blood-based epigenome-wide analyses on the prevalence
524 and incidence of nineteen common disease states”. In: *medRxiv* (2023). DOI: 10 .
525 1101 / 2023 . 01 . 10 . 23284387. eprint: [https://www.
526 medrxiv.org/content/early/2023/01/11/2023.01.10.23284387.
527 medrxiv.org/content/early/2023/01/11/2023.01.10.23284387.](https://www.medrxiv.org/content/early/2023/01/11/2023.01.10.23284387.full.pdf)
- 528 [5] Marina Bibikova et al. “High density DNA methylation array with single CpG site
529 resolution”. In: *Genomics* 98.4 (2011), pp. 288–295.
- 530 [6] Ruth Pidsley et al. “Critical evaluation of the Illumina MethylationEPIC BeadChip
531 microarray for whole-genome DNA methylation profiling”. In: *Genome biology* 17.1
532 (2016), pp. 1–17.
- 533 [7] Steve Horvath and Kenneth Raj. “DNA methylation-based biomarkers and the epi-
534 genetic clock theory of ageing”. In: *Nature Reviews Genetics* 19.6 (2018), pp. 371–
535 384.
- 536 [8] Elena Bernabeu et al. “Refining epigenetic prediction of chronological and biological
537 age”. In: *Genome Medicine* 15.1 (2023), pp. 1–15.
- 538 [9] Hansheng Wang. “Forward regression for ultra-high dimensional variable screen-
539 ing”. In: *Journal of the American Statistical Association* 104.488 (2009), pp. 1512–
540 1524.

- 541 [10] Jianqing Fan, Richard Samworth, and Yichao Wu. “Ultrahigh dimensional feature
542 selection: beyond the linear model”. In: *The Journal of Machine Learning Research*
543 10 (2009), pp. 2013–2038.
- 544 [11] Jianqing Fan and Rui Song. “Sure independence screening in generalized linear
545 models with NP-dimensionality”. In: (2010).
- 546 [12] Jianqing Fan and Jinchi Lv. “Sure independence screening for ultrahigh dimensional
547 feature space”. In: *Journal of the Royal Statistical Society: Series B (Statistical*
548 *Methodology)* 70.5 (2008), pp. 849–911.
- 549 [13] Joanna Zhuang, Martin Widschwendter, and Andrew E Teschendorff. “A compar-
550 ison of feature selection and classification methods in DNA methylation studies
551 using the Illumina Infinium platform”. In: *BMC bioinformatics* 13 (2012), pp. 1–
552 14.
- 553 [14] Jianqing Fan, Yang Feng, and Yichao Wu. “High-dimensional variable selection
554 for Cox’s proportional hazards model”. In: *Borrowing strength: Theory powering*
555 *applications—a Festschrift for Lawrence D. Brown*. Vol. 6. Institute of Mathematical
556 Statistics, 2010, pp. 70–87.
- 557 [15] Jason D Lee and Jonathan E Taylor. “Exact post model selection inference for
558 marginal screening”. In: *Advances in neural information processing systems* 27
559 (2014).
- 560 [16] Albert T Higgins-Chen et al. “A computational solution for bolstering reliability
561 of epigenetic clocks: Implications for clinical trials and longitudinal tracking”. In:
562 *Nature aging* 2.7 (2022), pp. 644–661.
- 563 [17] Richard D Riley et al. “Minimum sample size for developing a multivariable predic-
564 tion model: Part I—Continuous outcomes”. In: *Statistics in medicine* 38.7 (2019),
565 pp. 1262–1275.
- 566 [18] Richard D Riley et al. “Minimum sample size for developing a multivariable predic-
567 tion model: PART II—binary and time-to-event outcomes”. In: *Statistics in medicine*
568 38.7 (2019), pp. 1276–1296.

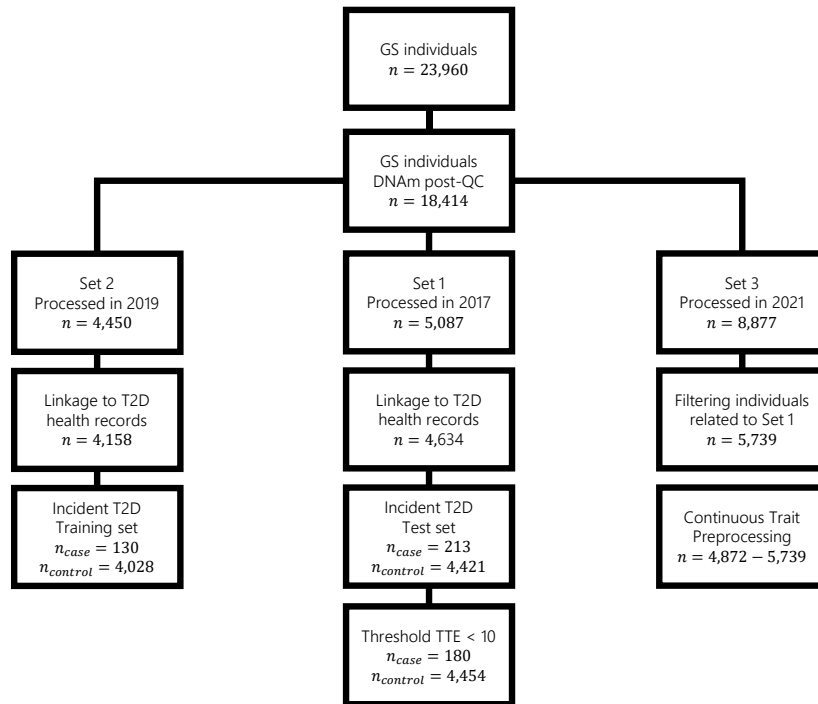
- 569 [19] Emmanuel O Ogundimu, Douglas G Altman, and Gary S Collins. “Adequate sample
570 size for developing prediction models is not simply related to events per variable”.
571 In: *Journal of clinical epidemiology* 76 (2016), pp. 175–182.
- 572 [20] Blair H Smith et al. “Generation Scotland: the Scottish Family Health Study; a
573 new resource for researching genes and heritability”. In: *BMC medical genetics* 7.1
574 (2006), pp. 1–9.
- 575 [21] H-E Wichmann et al. “KORA-gen-resource for population genetics, controls and a
576 broad spectrum of disease phenotypes”. In: *Das Gesundheitswesen* 67.S 01 (2005),
577 pp. 26–30.
- 578 [22] Gary S Collins et al. “Transparent reporting of a multivariable prediction model for
579 individual prognosis or diagnosis (TRIPOD): the TRIPOD statement”. In: *Annals
580 of internal medicine* 162.1 (2015), pp. 55–63.
- 581 [23] Anubha Mahajan et al. “Multi-ancestry genetic study of type 2 diabetes highlights
582 the power of diverse populations for discovery and translation”. In: *Nature genetics*
583 54.5 (2022), pp. 560–572.
- 584 [24] Yipeng Cheng et al. “Development and validation of DNA Methylation scores in two
585 European cohorts augment 10-year risk prediction of type 2 diabetes”. In: *Nature
586 Aging* (2023), pp. 1–9.
- 587 [25] Hui Zou and Trevor Hastie. “Regularization and variable selection via the elastic
588 net”. In: *Journal of the Royal Statistical Society Series B: Statistical Methodology*
589 67.2 (2005), pp. 301–320.
- 590 [26] Hemant Ishwaran et al. “Random survival forests”. In: (2008).
- 591 [27] Danni A Gadd et al. “Epigenetic scores for the circulating proteome as tools for
592 disease prediction”. In: *Elife* 11 (2022), e71802.
- 593 [28] Ake T Lu et al. “DNA methylation GrimAge strongly predicts lifespan and healthspan”.
594 In: *Aging (albanay NY)* 11.2 (2019), p. 303.

- 595 [29] Blair H Smith et al. “Cohort Profile: Generation Scotland: Scottish Family Health
596 Study (GS: SFHS). The study, its participants and their potential for genetic re-
597 search on health and illness”. In: *International journal of epidemiology* 42.3 (2013),
598 pp. 689–700.
- 599 [30] Daniel L McCartney et al. “Blood-based epigenome-wide analyses of cognitive abil-
600 ities”. In: *Genome Biology* 23.1 (2022), p. 26.
- 601 [31] Gary S Collins et al. “Developing risk prediction models for type 2 diabetes: a
602 systematic review of methodology and reporting”. In: *BMC medicine* 9.1 (2011),
603 pp. 1–14.
- 604 [32] Robert Tibshirani. “Regression Shrinkage and Selection via the Lasso”. In: *Journal*
605 *of the Royal Statistical Society. Series B (Methodological)* 58.1 (1996), pp. 267–
606 288. ISSN: 00359246. URL: <http://www.jstor.org/stable/2346178> (visited on
607 04/02/2023).
- 608 [33] Jerome Friedman, Trevor Hastie, and Rob Tibshirani. “Regularization paths for
609 generalized linear models via coordinate descent”. In: *Journal of statistical software*
610 33.1 (2010), p. 1.
- 611 [34] Eliza Fraszczyk et al. “Epigenome-wide association study of incident type 2 dia-
612 betes: a meta-analysis of five prospective European cohorts”. In: *Diabetologia* 65.5
613 (2022), pp. 763–776.
- 614 [35] Diana L Juvinao-Quintero et al. “DNA methylation of blood cells is associated with
615 prevalent type 2 diabetes in a meta-analysis of four European cohorts”. In: *Clinical*
616 *Epigenetics* 13 (2021), pp. 1–14.
- 617 [36] Noah Simon et al. “Regularization paths for Cox’s proportional hazards model via
618 coordinate descent”. In: *Journal of statistical software* 39.5 (2011), p. 1.
- 619 [37] Hanpu Zhou et al. “SurvMetrics: An R package for Predictive Evaluation Metrics
620 in Survival Analysis.” In: *R J.* 14.4 (2023), pp. 252–263.
- 621 [38] DY Lin. “On the Breslow estimator”. In: *Lifetime data analysis* 13 (2007), pp. 471–
622 480.

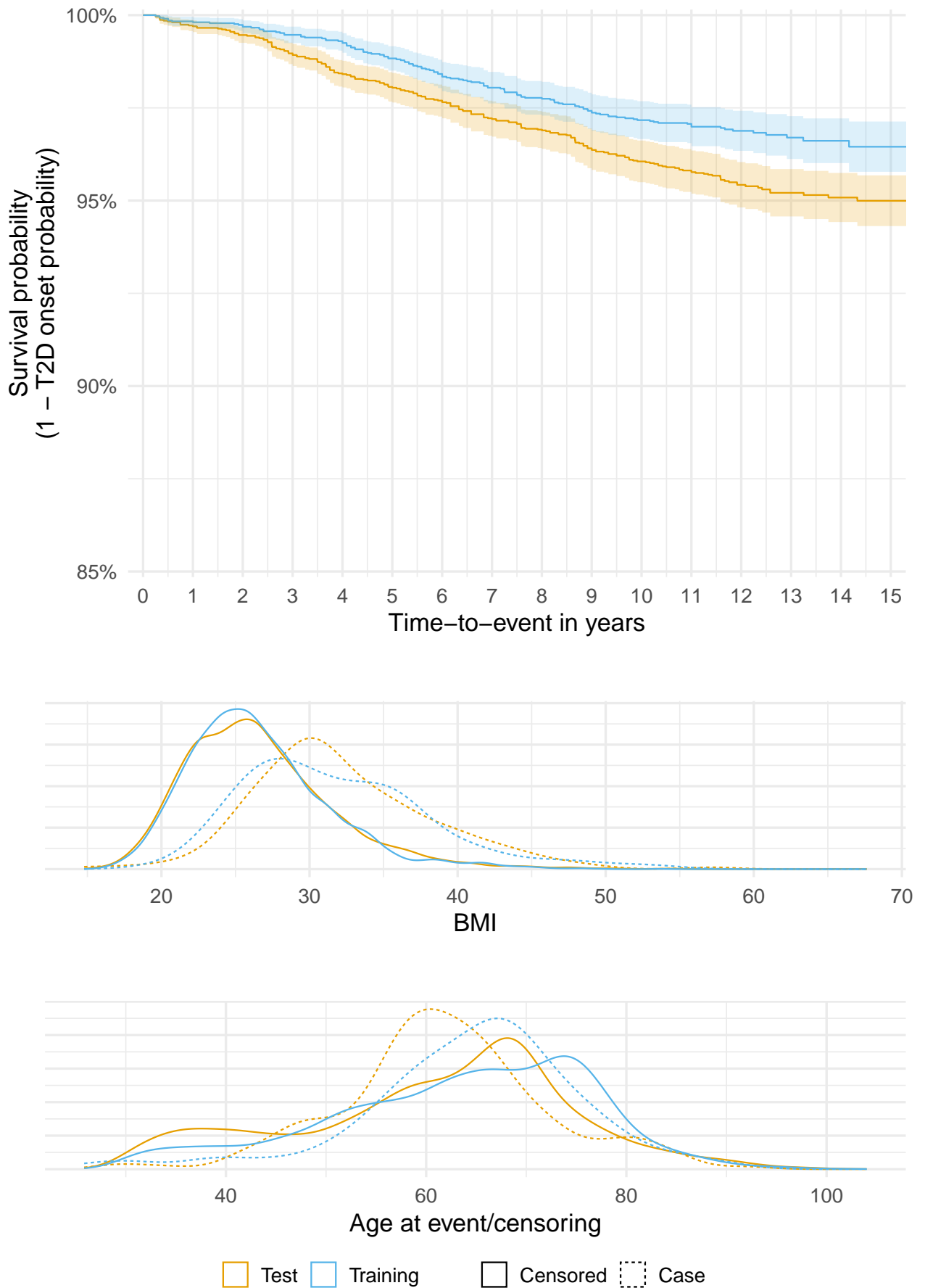
- 623 [39] Ben Van Calster et al. “A calibration hierarchy for risk models was defined: from
624 utopia to empirical data”. In: *Journal of clinical epidemiology* 74 (2016), pp. 167–
625 176.
- 626 [40] Jake R Conway, Alexander Lex, and Nils Gehlenborg. “UpSetR: an R package for
627 the visualization of intersecting sets and their properties”. In: *Bioinformatics* 33.18
628 (2017), pp. 2938–2940.

629 Supplementary Materials

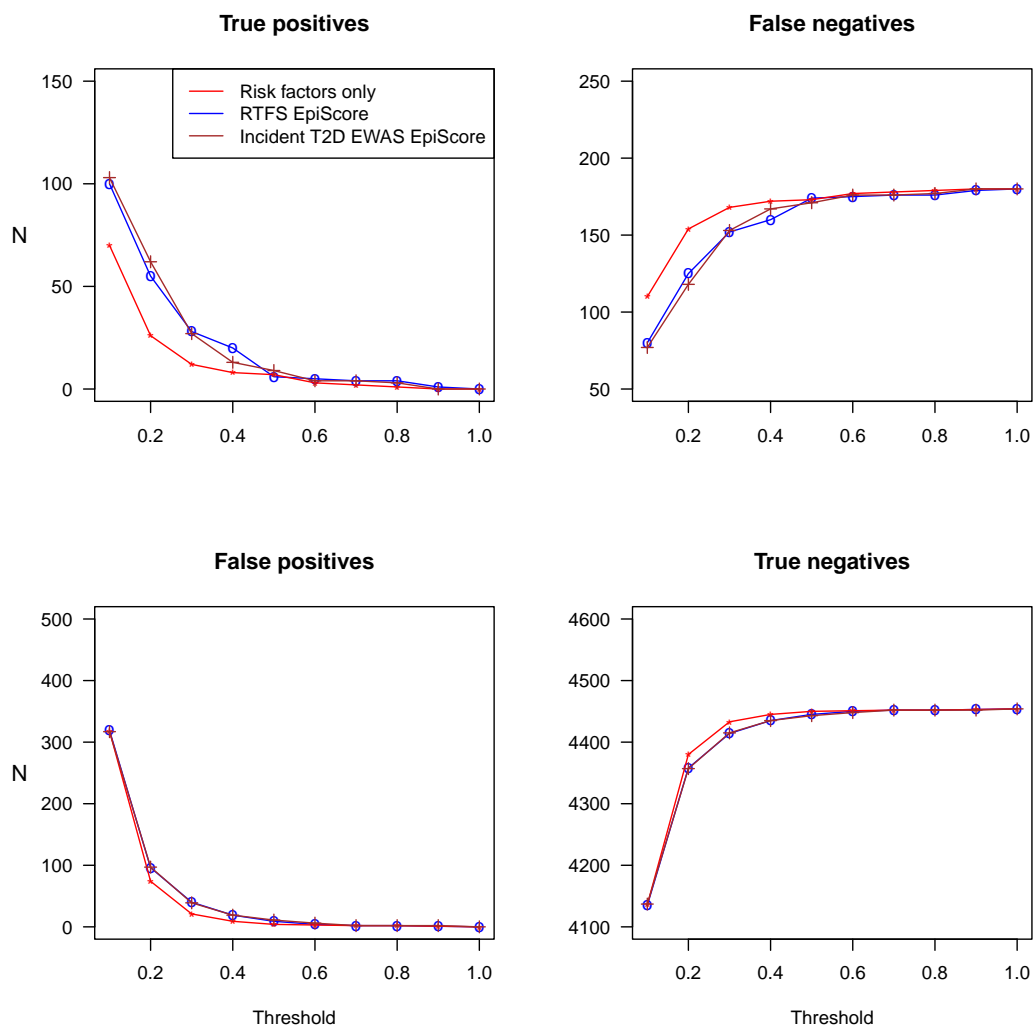
Preprocessing steps for Generation Scotland with
number of individuals/cases and controls after each step.



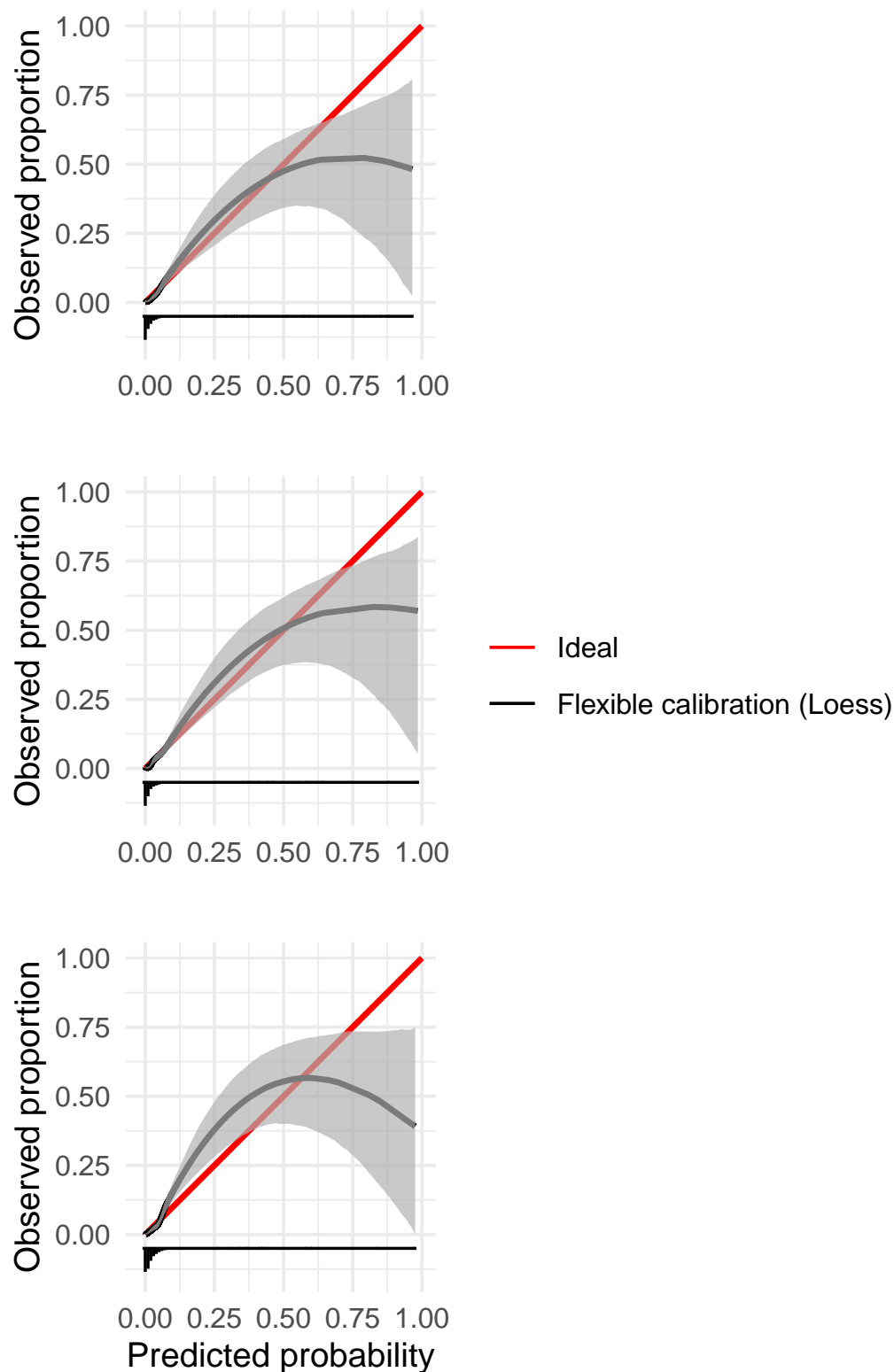
Supplementary Figure 1: Dataset numbers at each study pipeline processing step



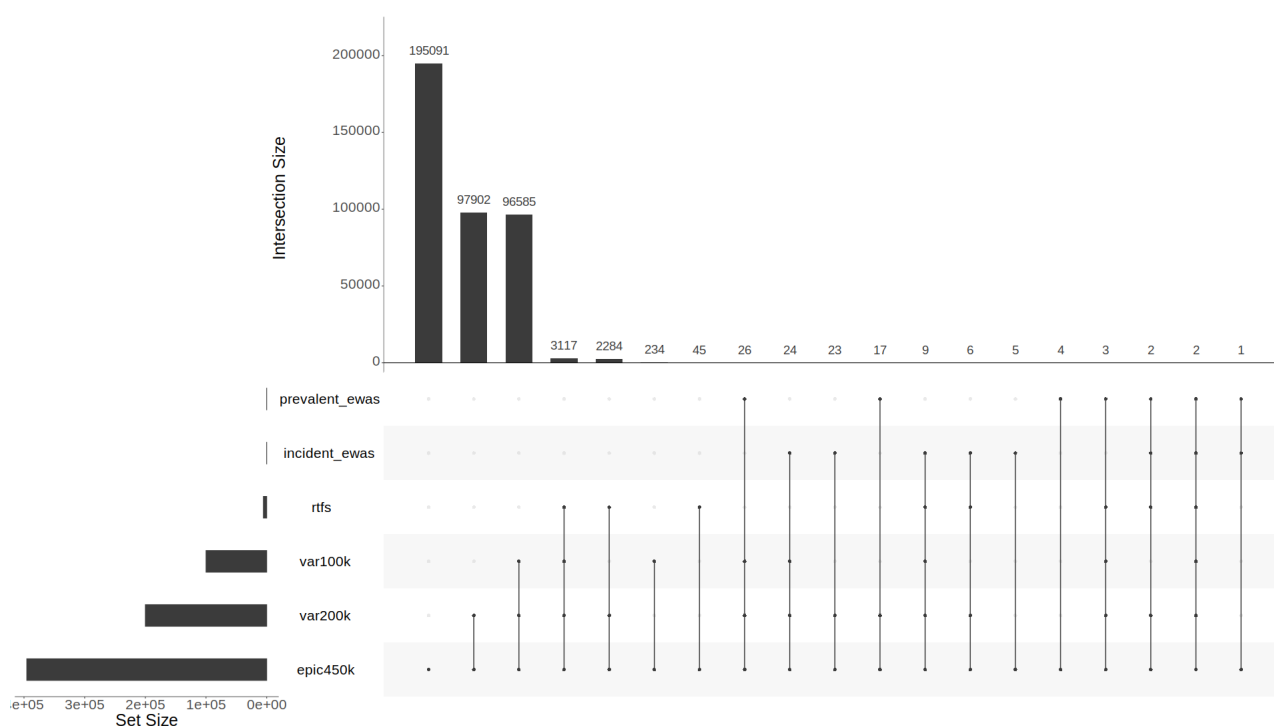
Supplementary Figure 2: Kaplan-Meier curves (top) and density plots for BMI and age at event/censoring (middle and bottom, respectively) for the incident T2D training and test sets



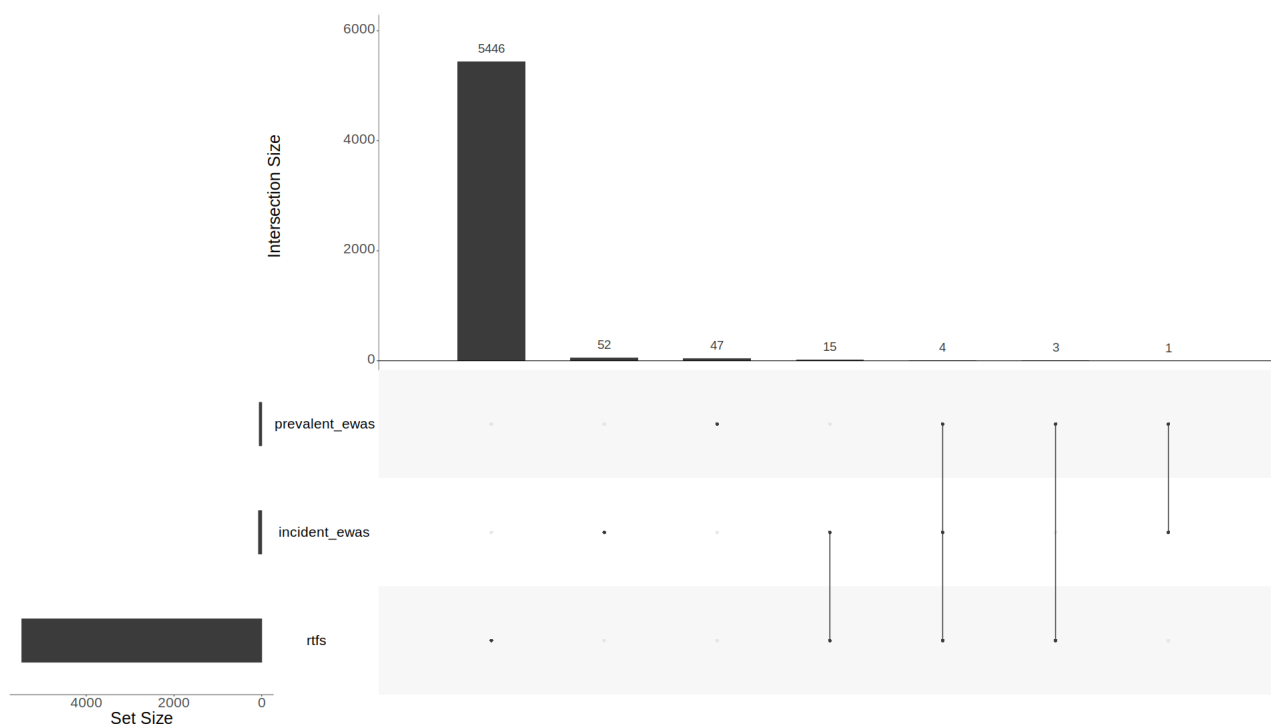
Supplementary Figure 3: Confusion matrix metrics across the probability threshold range 0-1 for the RTFS, incident T2D EWAS-based filtering and risk factors-only models.



Supplementary Figure 4: Calibration curves for the incident T2D EWAS EpiScore model (top), RTFS EpiScore model (middle) and risk factors-only model (bottom). The latter shows weaker calibration performance; the fitted calibration curve (black) is overall further from the perfect calibration line (red). All models show underestimation of risk below a predicted probability of around 0.5 and overestimation of risk otherwise.



Supplementary Figure 5: UpSet plot showing number of CpGs selected using each pre-selection method and overlaps between methods. Each column represents the number of CpGs pre-selected by the corresponding specific combination of methods. This was generated with the "distinct" option, meaning the presence or absence of a point in a column explicitly corresponds to the presence or absence of the corresponding method in the set.



Supplementary Figure 6: UpSet plot showing number of CpGs selected using the RTFS and EWAS-based pre-selection methods and overlaps between methods. Each column represents the number of CpGs pre-selected by the corresponding specific combination of methods. This was generated with the "distinct" option, meaning the presence or absence of a point in a column explicitly corresponds to the presence or absence of the corresponding method in the set.