

1 **Full title: Latent class regression improves the predictive acuity and**  
2 **clinical utility of survival prognostication amongst chronic**  
3 **heart failure patients**

4 †[John L Mbotwa](#),<sup>1,2,6</sup> ‡[Marc de Kamps](#),<sup>1,3</sup> †Paul D Baxter,<sup>1,2</sup> \*[George TH Ellison](#),<sup>1,2,4</sup> \*[Mark S](#)  
5 [Gilthorpe](#)<sup>1,2,5</sup>

6 <sup>1</sup>Leeds Institute for Data Analytics, University of Leeds, Leeds, UK; <sup>2</sup>Faculty of Medicine & Health,  
7 University of Leeds, Leeds, UK; <sup>3</sup>School of Computing, University of Leeds, Leeds, UK; <sup>4</sup>Centre for  
8 Data Innovation, University of Central Lancashire, Preston, UK; <sup>5</sup>The Alan Turing Institute, London,  
9 UK, <sup>6</sup>Department of Applied Studies, Malawi University of Science and Technology, Malawi

10 \*Joint senior author.

11

12 † This author analysed data and drafted the original manuscript.

13 \* These authors reviewed and redrafted the manuscript.

14 ‡ These authors reviewed the manuscript and provided suggestions for further improvement of the  
15 manuscript before the final manuscript was produced.

16

17 Correspondence

18 Email: [J.L.Mbotwa@leeds.ac.uk](mailto:J.L.Mbotwa@leeds.ac.uk)

19

20

21

22

23

24

## 25 **Abstract**

26 The present study aimed to compare the predictive acuity of latent class regression (LCR) modelling  
27 with: standard generalised linear modelling (GLM); and GLMs that include the membership of  
28 subgroups/classes (identified through prior latent class analysis; LCA) as alternative or additional  
29 candidate predictors. Using real world demographic and clinical data from 1,802 heart failure patients  
30 enrolled in the UK-HEART2 cohort, the study found that univariable GLMs using LCA-generated  
31 subgroup/class membership as the sole candidate predictor of survival were inferior to standard  
32 multivariable GLMs using the same four covariates as those used in the LCA. The inclusion of the LCA  
33 subgroup/class membership together with these four covariates as candidate predictors in a  
34 multivariable GLM showed no improvement in predictive acuity. In contrast, LCR modelling resulted  
35 in a 10-14% improvement in predictive acuity and provided a range of alternative models from which  
36 it would be possible to balance predictive acuity against entropy to select models that were optimally  
37 suited to improve the efficient allocation of clinical resources to address the differential risk of the  
38 outcome (in this instance, survival). These findings provide proof-of-principle that LCR modelling can  
39 improve the predictive acuity of GLMs and enhance the clinical utility of their predictions. These  
40 improvements warrant further attention and exploration, including the use of alternative techniques  
41 (including machine learning algorithms) that are also capable of generating latent class structure while  
42 determining outcome predictions, particularly for use with large and routinely collected clinical  
43 datasets, and with binary, count and continuous variables. [245/300 words]

## 44 **Key words**

45 Latent Class Regression; Prediction; Generalised Linear Modelling; Latent Class Analysis; Heart  
46 Failure; UK-HEART2 Cohort

## 47 **Short title**

48 Predicting survival of heart failure patients using latent class regression [66/100 characters]

## 49 **Introduction**

### 50 **The limited acuity and clinical utility of Generalised Linear Models (GLMs)**

51 The potential utility of predictive modelling, using routinely collected data for diagnosis,  
52 prognostication and health service planning, is one of five ‘novel capabilities’ that Wang et al. [1]  
53 identified as pertinent to the application of data analytics in medicine and health. Long before John  
54 Mashey first applied the term ‘Big Data’ in this context during the late 1990s [2], generalised linear  
55 models (GLMs) were used to develop clinical ‘risk scores’ based on much smaller scale datasets [3].  
56 Indeed, clinical prediction models (CPMs) remain popular for prognostication and are in widespread  
57 usage to this day, not least in cardiovascular medicine [4,5].

58 While CPMs and their wider utility remain contentious (beyond strict prognostication, and particularly  
59 in prevention [6-9]), many of the standard statistical modelling techniques commonly used are on  
60 clinical datasets that remain relatively small – at least when compared to contemporary notions of ‘Big  
61 Data’ [2]. A substantial statistical weakness of the commonest of these (generalised linear models;  
62 GLMs) as a predictive tool is that they often fail to make *full* use of the joint information available  
63 amongst *all* candidate predictor variables. This is because these models rarely explore nonlinear  
64 relationships and interactions. Moreover, even when analysts optimally parameterise the candidate  
65 predictors available, and carefully consider all possible interaction terms between these, the clinical  
66 utility of GLMs is typically limited to predictions made at the population level [6,10], while predictions  
67 at the individual level often lack precision (and with it, utility).

68 Although more sophisticated machine learning techniques may overcome the rigidity of GLMs and  
69 analysts’ tendency to ignore (or overlook) nonlinear relationships and interactions, population-level  
70 predictions generated using cutting edge machine learning techniques will still be more reliable than  
71 individual-level predictions. Indeed, this bald fact applies to all prediction modelling techniques,  
72 including those underpinning contemporary claims of ‘personalised’ or ‘precision medicine’ [11]. It is  
73 therefore critical to recognise that while it is possible to determine what proportion of any given

74 population will experience a specified outcome with a reasonable degree of accuracy, all such models  
75 provide less accuracy in determining outcomes for each *individual* within that population.

76 Meanwhile, a further epistemological consideration that commonly arises in CPMs (and elsewhere) is  
77 the mistaken belief that the coefficient estimates of covariates included/retained in the model indicate  
78 the extent to which each covariate contributes to the model's overall prediction (of the model's specified  
79 outcome). This belief is mistaken because each covariate's coefficient estimate is generated *conditional*  
80 on the adjustment of all other covariates in the model, such that the contribution of any one covariate is  
81 merged with that of all other covariates included in that model. For this reason, what the coefficient  
82 estimate of each covariate actually represents is the residual relationship between that covariate and the  
83 outcome *subject to* the joint contributions made by all other covariates in the model *considered*  
84 *simultaneously*.

85 This situation is further complicated where any of the covariates included in the model reflect events  
86 that occurred contemporaneously with or even after the specified outcome – in this instance, the joint  
87 contributions made by all covariates would be subject to conditioning on the outcome, which can have  
88 other adverse consequences on model interpretation [12-14]. In practice, the inclusion of covariates  
89 acting contemporaneously with or after the outcome' in prediction modelling is likely to be used only  
90 where the aim is to estimate the values for variables whose measurements are missing, imprecise or  
91 challenging to measure (i.e. in modelling that aims to achieve what might be called 'predictive  
92 interpolation' for diagnostic and related measurement/ascertainment purposes). These issues aside, it is  
93 important to stress that the coefficient estimates of all covariates (with the exception of the covariate  
94 closest in time to the outcome) cannot be causally interpreted, as they will be subject to inferential bias  
95 known as mediator bias [15], which undermines causal interpretation of their coefficients due to the so-  
96 called 'Table 2 Fallacy' [16].

97 Due to these caveats, predictions that are generated using GLMs cannot address the two key concerns  
98 of attending physicians, namely: "Which of the covariates (i.e. 'predictor' variables) are amenable to  
99 clinical intervention, so as to prevent or mitigate any adverse outcome (or promote and amplify any

100 favourable outcome) in each (or all) of these patients?” and “Which particular patients will experience  
101 an adverse (or favourable) outcome?”. To address the first of these questions, analysts need to switch  
102 their focus from predicting outcome values to estimating each of the relationships between covariates  
103 considered plausible targets for intervention and the outcome – an approach that can capitalise on recent  
104 advances in causal inference modelling techniques [17]. To address the second question, the best that  
105 can be achieved is to identify clinically meaningful *subgroups* of patients with shared characteristics  
106 that set them apart from other (*subgroups* of) patients – a relatively novel approach that involves  
107 multivariable ‘risk profiling’.

### 108 **Improving the acuity and clinical utility of predictive modelling for prognostication**

109 Multivariable risk profiling can be achieved using latent class analysis (LCA) in which the exploration  
110 of nonlinearity, and of important interactions amongst included covariates, forms an integral (albeit  
111 implicit) part of classifying patients into subgroups [18]. Despite these benefits, the clinical utility of  
112 the resulting latent classes ultimately depends upon the extent to which this approach optimally exploits  
113 the joint information amongst available covariates. This approach perhaps has greatest clinical utility  
114 where there are: (i) factors known to be *associated* with the outcome (which therefore facilitate  
115 prediction); but (ii) there are no known, modifiable causes of the outcome, or aetiological understanding  
116 is poor/contested (as is the case with many rare, novel or complex diseases). Indeed, providing that the  
117 specified outcome is excluded from the LCA process (to avoid conditioning on the outcome) [14],  
118 combining LCA class membership with candidate predictors provides increased complexity that can  
119 help exploit the joint covariate information in multivariable GLM prediction. That said, it is important  
120 to stress that causal interpretation of any covariate coefficients for latent class membership in such  
121 models remains deeply flawed for the very same reasons that causal interpretations of any covariate  
122 coefficient in prediction GLMs is flawed (as explained earlier). Ostensibly this consideration might  
123 appear to limit the clinical utility of LCA-generated class membership, and it is true that describing  
124 class membership as a ‘risk factor’ often generates, and commonly reflects, a lack of understanding.

125 Indeed, it risks conflating prediction and causal inference/determination just as it does when individual  
126 covariates are described in similar terms as ‘risk factors’ [7].

127 Thus, while classifying subgroups of individuals using LCA can improve analytical practice and  
128 strengthen consideration of nonlinear relationships (and important interactions amongst covariates), it  
129 does not address the clinical appetite for identifying so-called ‘modifiable risk factors’, or for  
130 individually tailored risk probabilities (the so-called ‘holy grail’ of personalised or precision medicine)  
131 [10]. This might explain why the use of latent variable methods in prediction modelling remains largely  
132 under-explored, even though more sophisticated approaches exist that incorporate such techniques  
133 within GLM and offer substantial advantages for clinicians through subgroup risk profiling. These  
134 approaches involve the construction of latent classes ‘across’ multivariable GLMs to: integrate  
135 consideration of nonlinear relationships and important interactions between covariates; and better  
136 capture (and exploit) the joint information amongst the available/included covariates. For example, in  
137 what is termed latent class regression (LCR) modelling, population data are partitioned into their  
138 constituent latent classes and a distinct GLM is simultaneously generated for each class. In the process,  
139 this approach accommodates any inherent population heterogeneity and thereby improves model  
140 precision.

141 In its simplest form, LCR models may be viewed as two distinct modelling concepts undertaken in a  
142 single estimation process: in the first, population data are probabilistically assigned to latent classes  
143 (population subgroups); while, in the second, separate GLMs are derived for each class/subgroup. The  
144 probability of an individual belonging to each class is based on similarities in the characteristics  
145 displayed by individuals attributed to different classes. Importantly, the assignment of individuals to  
146 classes is not limited to just those covariates available for analysis, since outcome differences  
147 attributable to unknown (i.e. latent) influences are also accommodated, *and* without inappropriate  
148 conditioning on the outcome. Individuals may thus belong to more than one class, with the sum of  
149 probabilities over all classes being one. Within each class, distinct GLMs are generated, with the  
150 selection of covariates acting as predictors (and their model coefficients) permitted to vary from one

151 class to the next. In this way, by ensuring that the consideration of potential nonlinearity and possible  
152 interactions is integral to the application of LCR models, these models should exploit the covariate joint  
153 information available in a more consistent fashion and thereby strengthen the acuity of the prediction  
154 achieved. An additional benefit of this approach is that the latent classes/subgroups identified using  
155 LCR may also strengthen the clinical utility of the prediction achieved because any variation in the  
156 risk of the outcome amongst different classes/subgroups can be used to target diagnostic, therapeutic  
157 or palliative resources more precisely and efficiently.

158 The aim of the present study was therefore to explore whether LCR models might improve the accuracy  
159 and precision of predictions at the population *and* individual level, by comparing LCR-generated  
160 predictions to standard GLM and LCA-informed GLM (including the use of LCA-generated class  
161 membership as either the *only* candidate predictor in univariable GLMs, or as an *additional* candidate  
162 predictor alongside all other available covariates in multivariable GLMs). We thus explore four models  
163 offering progressively more complex exploitation of the individual and joint information available from  
164 the covariates available for consideration as candidate predictors. To this end, the analyses that follow  
165 use data (on age, sex, haemoglobin level and diabetes) that are routinely available in a clinical context  
166 (cardiovascular medicine) in which Cox proportional hazards time-to-event analyses are commonly  
167 used in prognostic predictions of mortality, where survival and loss to follow-up are pertinent analytical  
168 endpoints.

## 169 **Methods**

### 170 **Study design, data collection and ethics**

171 The analyses that follow used data from the United Kingdom Heart Failure Evaluation and Assessment  
172 of Risk Trial 2 (UK-HEART2) – a prospective cohort of ambulant patients with signs and symptoms of  
173 chronic heart failure (CHF) [19]. The study recruited 1,802 adult patients with CHF who attended  
174 specialist cardiology clinics in four UK hospitals between July 2006 and December 2014 [20]. Patients  
175 were eligible for recruitment if they: were aged 18 years or older; had had clinical signs and symptoms

176 of CHF for at least 3 months; and had a left ventricular ejection fraction that was less than or equal to  
177 45% [19,20]. Ethical approval was obtained from the research ethics committee at each participating  
178 hospital and eligible study participants were only recruited following informed consent [21]. Additional  
179 information regarding UK-HEART-2's study design, patient eligibility and inclusion criteria, together  
180 with a detailed description of the study cohort has been reported elsewhere [19-21].

## 181 Statistical methods

182 To simplify the methodological comparisons undertaken in the present study, the covariates selected as  
183 candidate predictors comprised two demographic variables (age, sex), a single physiological parameter  
184 (haemoglobin level), and a single clinical characteristic (type 2 diabetes). These four covariates were  
185 then used to generate prognostic predictions of survival amongst UK-HEART-2 participants using four  
186 separate statistical Procedures (for each of which the underlying principles and model building  
187 processes are described in detail in Part 1 of the Supplementary Materials):

- 188 • Procedure 1 (standard GLM) involved single step multivariable Cox proportional hazards models  
189 that considered all four covariates as candidate predictors of survival, with no consideration of  
190 nonlinear relationships or interactions between covariates.
- 191 • Procedure 2 (LCA-informed GLM *without* the inclusion of covariates) involved two sets of models,  
192 each involving two separate steps. First, LCA was used to identify any latent classes/subgroups of  
193 participants using the four selected covariates, with individual membership to each latent class  
194 allocated using modal (Procedure 2a) and probabilistic (Procedure 2b) assignment. Second,  
195 *univariable* Cox proportional hazards models examined latent class membership as the *sole*  
196 predictor of survival, with separate models generated using latent class membership derived using  
197 modal (Procedure 2a) or probabilistic (Procedure 2b) assignment.
- 198 • Procedure 3 (LCA-informed GLM *with* the inclusion of covariates) again involved two sets of  
199 models, each involving two separate steps. First, LCA was used to allocate latent class membership  
200 using modal (Procedure 3a) and probabilistic (Procedure 3b) assignment – as in the first step of



201 Procedure 2 (above). Second, *multivariable* Cox proportional hazards models considered all four  
202 covariates (as used in Procedure 1) *plus* latent class membership as *multiple* predictors of survival,  
203 with separate models generated using latent class membership derived using probabilistic  
204 (Procedure 3a) or modal (Procedure 3b) assignment.

205 • Procedure 4 (LCR) involved single step latent class regression (LCR) models that considered all  
206 four covariates as candidate predictors to simultaneously predict *both* latent class membership *and*  
207 survival within each latent class.

208 For all latent class models, entropy is reported which assesses the extent that individuals are aligned  
209 predominantly to a single class (i.e. having a large modal probability, leading to a greater entropy), as  
210 this facilitates clearer interpretation of each latent class as a near-complete collection of individuals.  
211 Model optimisation in terms of the number of latent classes may thus depend upon *both* the overall  
212 predictive acuity of the latent class structure (as evident from the model BIC) *and* the intended utility  
213 of the determined classes thereafter (as indicated by the model entropy). For this illustration, we  
214 prioritise overall predictive acuity.

215 All descriptive statistics and GLMs were generated using *R* (version 4.0.3) [22], as were the model  
216 specification, selection, validation and bootstrapping procedures (Part 2 of the Supplementary  
217 Materials). All latent class modelling was undertaken in *Mplus* (version 8.3) [23], using the *Mplus*  
218 automation package to run models in *Mplus* from within *R* [22].

## 219 **Results**

220 The first column of Table 1 summarises the distribution of each covariate amongst participants in the  
221 UK-HEART-2 cohort. These indicate that: the mean age of the cohort's participants was 70 years;  
222 around two thirds (69.7%) were male; over a quarter (28%) had type 2 diabetes; the mean level of  
223 circulating haemoglobin was 13.5 g/dl; and 60% died during the period of follow-up (equivalent to a  
224 median survival of 3.4 years).

225 In Procedure 1, the single step Cox proportional hazards models that considered all four covariates as  
226 candidate predictors of survival found that the model in which all four covariates were retained achieved  
227 the highest AUC (0.69) – a level of acuity considered ‘modest to poor’ [24].

228 In Procedure 2, the LCAs conducted during the first step found that the 5-class model which retained  
229 all four covariates had the most favourable BIC (Table 2). Applying this 5-class model during the second  
230 step as the sole predictor of survival in a Cox proportional hazards model, achieved an AUC of 0.65  
231 using modal assignment (Procedure 2a) and 0.66 using probabilistic assignment (Procedure 2b). These  
232 levels of acuity were both lower than that achieved using Procedure 1 (AUC=0.69).

233 In Procedure 3, the second step involved consideration not only of the four covariates as candidate  
234 predictors of survival in the Cox proportional hazards model (as in Procedure 1), but also membership  
235 of the same 5-class model developed in the first step of Procedure 2. These analyses found that: the best  
236 fitting GLMs did not retain class membership as a predictor; and forcibly retaining class membership  
237 in the model did not improve the AUC above that achieved in Procedure 1 or 2, regardless of how class  
238 membership was assigned (modal: AUC=0.65; probabilistic: AUC=0.66).

239 In Procedure 4, with all four covariates eligible for inclusion as candidate predictors of *both* latent class  
240 membership *and* the Cox Proportional Hazards models, some of the models were over-parameterised  
241 and failed to converge. Nonetheless, the most favourable of the models that successfully converged  
242 involved a latent class variable with just two classes and an AUC of 0.79 (Table 3). When compared to  
243 the best performing models in Procedures 1-3, these results suggest that Procedure 4 achieved a  
244 substantial improvement in predictive acuity of 10-14%.

245 Improvements in acuity aside, the most favourable of the LCR models had only three of the covariates  
246 (age, sex, and type 2 diabetes) retained in the Cox proportional hazards models for each membership  
247 class, and only one of these covariates (type 2 diabetes) and the remaining covariate (haemoglobin level)  
248 retained as covariates in the LCR class membership model (Table 4). Given that all four covariates were  
249 retained in the most favourable CPH models generated by Procedures 1 and 3, and in the LCA models

250 generated in the first step of Procedures 2 and 3, these findings suggest that Procedure 4's 10-14%  
251 improvement in AUC is likely to have been achieved by exploiting the available covariate information  
252 differently to each of the three other Procedures. An indication of what this entailed can be found in the  
253 distribution of covariate characteristics amongst the two classes of the most favourable LCR model  
254 (Table 5), which suggest that these classes might warrant post-hoc labelling as 'high risk' and 'low risk'  
255 subgroups and might thereby offer substantial additional clinical utility (in guiding the allocation of  
256 diagnostic, therapeutic and/or palliative resources).

257 A further key finding that emerges from closer examination of the Cox proportional hazards models  
258 generated for each of the two classes within the optimum LCR model (Table 4) is that the contribution  
259 made by each of the covariates therein varied by class, and was dissimilar to the contribution these  
260 covariates made in those Procedures where all covariates were available for inclusion as separate  
261 candidate predictors (i.e. Procedure 1 and 3a/b). While the coefficient estimates of covariates in each of  
262 these models cannot be interpreted as measures of causal effects [16], their contribution as candidate  
263 predictors is strikingly different and depends upon the choice of model(s) used in each Procedure (Table  
264 4). For example, the hazard of death associated with being male was 1.7 to 1.8 in Procedures 1 and 3,  
265 whilst for Procedure 4 being male was associated with a substantially higher hazard of death in one  
266 class (HR = 2.07; 1.58, 2.71) yet was unrelated to the hazard of death in the other class (HR = 1.01;  
267 0.64, 1.60). Likewise, Type 2 diabetes was consistently associated with an elevated hazard of death in  
268 models generated under Procedure 1 and 3, while in Procedure 4 this covariate was associated with both  
269 an *elevated* hazard of death in one class (HR = 1.26; 0.91, 1.75) and a *reduced* hazard of death in the  
270 other class (HR = 0.43; 0.23, 0.82).

271 Clearly, the joint information available amongst each of the candidate predictors is selected and utilised  
272 very differently by each of the Procedures examined in the present study (see Table 4). Nonetheless,  
273 what sets the LCR model in Procedure 4 apart from the models used in Procedures 1-3 is that LCR  
274 allows the predictive contribution from each covariate to be partitioned *across* any latent substructures  
275 existing within the study population, such that covariates are able to operate differently within each of

276 the latent subgroups – thereby capturing and reflecting population heterogeneity that is: unavailable to  
277 any of the other modelling Procedures; and, crucially, of substantial (additional) value when predicting  
278 the specified outcome.

## 279 **Discussion**

280 The present study provides proof of principle that LCR models can provide substantive improvements  
281 in predictive acuity *and* clinical utility over standard approaches using GLM (with or without LCA).  
282 Nonetheless, there are several potential limitations that warrant consideration and further investigation.  
283 In particular, it would be insightful to compare these alternative approaches to prediction using larger  
284 datasets and larger numbers of covariates than those chosen in this instance for illustration. This might  
285 involve comparing Procedures 1 through 4 using different numbers and sets of covariates from similar  
286 sized datasets; as well as extending the application of LCR modelling to more complex scenarios and  
287 much larger datasets. At the same time, it is important to point out that, in the context of the dataset  
288 used in the present study, the underlying ‘truth’ (and the data generating mechanisms involved) cannot  
289 be known with certainty, and exploring the potential strengths (and analytical limitations) of LCR would  
290 thus benefit from extensive simulations to evaluate a range of different circumstances for a range of  
291 different covariates *and* outcomes (including those comprising binary, continuous and count variables)  
292 to evaluate whether LCR continues to perform well (and better than GLM, LCA or both) under these  
293 circumstances. In the absence of subsequent research along these lines, the ‘proof of principle’ offered  
294 by the present study remains speculative; although it would also be worth exploring whether alternative  
295 approaches to prediction modelling might be incorporated into, or integrated with, the analytical  
296 principles underpinning LCR modelling, such as the inclusion of similar dual modelling structures  
297 within machine learning, to assess whether the apparent benefits of LCR models might be further  
298 enhanced.

299 These limitations, the present study successfully compared three different approaches for incorporating  
300 latent variable methods within prediction modelling and demonstrated that LCR models can outperform  
301 not only the standard approach using GLM (in which membership of latent classes is ignored –

302 Procedure 1), but also those that include latent class membership identified using LCA to generate an  
303 alternative (Procedure 2) or additional (Procedure 3) candidate predictor. This improvement in  
304 predictive acuity (which, as shown above, resulted in a 10-14 percentage point improvement in AUC,  
305 despite the modest number of participants and covariates involved) illustrates the potential benefits of  
306 LCR for prediction modelling which, in this instance, shifted the acuity of prediction from ‘modest to  
307 poor’ to ‘substantial’ [24].

308 The present study also demonstrated that the latent class/subgroup structure that is revealed through  
309 LCR may have potential clinical utility. This is because it might – as in the example examined here –  
310 facilitate the identification of discrete *subgroups* (i.e. latent classes) of populations with very different  
311 underlying risks of the outcome. While such subgroups may not necessarily be amenable to effective  
312 intervention (given that LCR models support prediction, not causal inference [16]), they should help to  
313 improve the efficient allocation/targeting of outcome-relevant diagnostic, therapeutic and/or palliative  
314 resources to those subgroups identified as more likely to require (and perhaps even benefit from) these.  
315 However, to maximise the clinical exploitation of latent subgroups identified using LCR (and similar  
316 techniques), model selection must focus on those achieving higher entropy – where the probability of  
317 class assignment is closer to one for most assignments – as this better aligns individuals/participants to  
318 a predominant single class (rather than aligning individuals/participants to multiple classes). For  
319 example, in Procedure 4, the 3-class model had lower predictive acuity but greater entropy than the 2-  
320 class model (see Table 3); and had the identification of clinically meaningful subgroups been the focus  
321 of these analyses (as opposed to overall predictive acuity), then it might have been appropriate to accept  
322 a modest reduction in predictive acuity in favour of enhanced clinical utility – i.e. recognising three  
323 (‘high’, ‘medium’ and ‘low’ risk) subgroups rather than just the two (‘high’ and ‘low’ risk) subgroups  
324 identified by the LCR model with the most favourable predictive acuity (Table 3). Indeed, when clinical  
325 resources are scarce, such an approach might prove a more reliable approach to resource allocation than  
326 one based upon a stringent interpretation of predictive acuity alone.

## 327 **Online Supplementary Material**

### 328 **Part 1. Underlying principles and model-building processes involved in the GLM,** 329 **LCA and LCR techniques examined**

#### 330 GLM: the Cox proportional hazards model

331 A Cox proportional hazards model generates a (hazard) function which indicates the risk of the outcome  
332 occurring during the period of follow-up. Mathematically, a Cox regression model [25,26] is defined  
333 as:

$$334 \quad h(\mathbf{t} \mid \mathbf{x}, \beta_i) = h_0(\mathbf{t}) \exp(\mathbf{x} \beta^T) \quad (\text{Eq1})$$

335 where:  $\mathbf{t}$  is a non-negative random variable representing time to ‘death’, ‘loss to follow-up’ or ‘the end  
336 of the study’ for all participants (in this example, patients with CHF);  $h_0(\mathbf{t})$  is the baseline hazard function;  
337  $\mathbf{x}$  is the vector of predictors for the time-to-event outcome  $\mathbf{t}$ ; and  $\beta^T$  is the transpose of the vector of  
338 coefficients obtained from the Cox proportional hazards model. To make predictions using the Cox  
339 proportional hazards model, the survival function is defined as:

$$340 \quad S(\mathbf{t} \mid \mathbf{x}, \beta_i) = [S_0(\mathbf{t})]^{\exp(\mathbf{x} \beta^T)} \quad (\text{Eq2})$$

341 where, if the baseline hazard function  $h_0(\mathbf{t})$  is known, then:

$$342 \quad S_0(\mathbf{t}) = \exp \left\{ - \int_0^t h(u) du \right\} \quad (\text{Eq3})$$

#### 343 LCA: the general latent class (profile) model

344 Latent class (profile) models come from a family of finite mixture models that classify observations  
345 into classes associated with unobserved heterogeneity in a population. A population is partitioned into  
346  $g$  classes for the outcome  $\mathbf{y}$  with the mixture density function defined in relation to covariates  $\mathbf{x}$  as:

$$347 \quad f(\mathbf{y} \mid \mathbf{x}, \lambda) = \sum_{i=1}^g \pi_i f_i(\mathbf{y} \mid \mathbf{x}, \beta_i) \quad (\text{Eq4})$$

348 where  $f_i(\mathbf{y} | \mathbf{x}, \boldsymbol{\beta}_i)$  is the conditional probability density function for the observed response in the  $i^{\text{th}}$  class  
349 and  $\pi_i$  ( $i = 1 \dots g$ ) represent the class-membership probabilities that are estimated for each class such  
350 that:

$$351 \quad \sum_{i=1}^g \pi_i = 1. \quad (\text{Eq5})$$

352 For a class membership model, the structural part of the model is given by:

$$353 \quad \text{logit}(\pi_i(\mathbf{x} | \gamma_i, \boldsymbol{\delta}_i)) = \gamma_i + \mathbf{x} \boldsymbol{\delta}_i^T \quad (\text{Eq6})$$

354 hence

$$355 \quad \pi_i(\mathbf{x} | \gamma_i, \boldsymbol{\delta}_i) = \frac{\exp(\gamma_i + \mathbf{x} \boldsymbol{\delta}_i^T)}{\sum_{j=1}^g \exp(\gamma_j + \mathbf{x} \boldsymbol{\delta}_j^T)} \quad (\text{Eq7})$$

356 where:  $\mathbf{x}$  is a  $(p \times 1)$  covariate vector for the class-membership model; and  $\boldsymbol{\delta}_i^T$  is the transposition of the  
357 vector  $\boldsymbol{\delta}_i$  for the multinomial logistic class-membership model.

### 358 LCR: the latent class regression model

359 The latent class regression (LCR) model is an extended version of the generalised linear model where  
360 the concept of latent class mixtures is applied to the entire model specified, not just to a cluster of  
361 covariates. LCR survival analysis extends this to the time-to-event framework of Cox proportional hazards  
362 modelling to: (i) predict probabilistically assigned subgroups of participants with different futures (in this  
363 example, subgroups of patients with different prognoses of survival/death) based on the available  
364 covariates; and, (ii) *simultaneously* predict the survival distributions for each subgroup selecting from  
365 the same covariates acting as candidate predictors. The distribution of the survival time variable for each  
366 component in Eq4 can be:

- 367 • parametric – a scenario *with* distributional assumptions concerning the survival times;
- 368 • semi-parametric – a scenario with *relaxed* distributional assumptions; or

369 • non-parametric – a scenario *without* distribution assumptions concerning the survival times.

370 Assuming a parametric model for the specified outcome variable, the component's densities are  
371 assumed to be from the same family, so that a number of common distribution functions may be  
372 considered appropriate for survival times, such as: exponential; Gamma; and Weibull [38]. In a semi-  
373 parametric case, the Cox proportional hazards model is an example. Within a latent class framework, if  
374  $\mathbf{t}$  is the random variable representing time to event (e.g. 'death', 'loss to follow-up', or 'end of the  
375 study'), individuals are divided into  $g$  latent classes that are differentiated by the covariate vector  $\mathbf{z}$ , with  
376 individual survival in each class  $i$  predicted by covariate vector  $\mathbf{x}$ , and the survival model is defined as:

$$377 \quad S(\mathbf{t} \mid \mathbf{x}, \mathbf{z}, \theta) = \sum_{i=1}^g \pi_i(\mathbf{z} \mid \gamma_i, \delta_i) S_i(\mathbf{t} \mid \mathbf{x}_i, \beta_i) \quad (\text{Eq8})$$

378 where:  $\theta = (\gamma_i, \delta_i, \beta_i)$  is the collection of parameters to be estimated such that  $\pi_i(\mathbf{z} \mid \gamma_i, \delta_i)$  satisfies the  
379 constraints in Eq4. Vectors  $\mathbf{x}_i$  and  $\mathbf{z}$  are any available measures of participant characteristics, exposures  
380 and treatments etc., which may be the same or differ, just as the  $\mathbf{x}_i$  covariates may also differ for each  
381 class.

382 If the effects of the  $\mathbf{x}_i$  covariates on the hazards (i.e. the instantaneous risk of event) in each class is  
383 constant during the duration of follow-up, then the hazard function can be specified as:

$$384 \quad h_i(\mathbf{t} \mid \mathbf{x}_i, \beta_i) = h_{0i}(\mathbf{t}) \exp(\mathbf{x}_i \beta_i^T) \quad (\text{Eq9})$$

385 where:  $h_{0i}(\mathbf{t})$  is the baseline hazard for class  $i$ ; and  $\exp(\mathbf{x}_i \beta_i^T)$  is the relative risk associated with a  
386 vector of the  $\mathbf{x}_i$  covariates acting as candidate predictors. We can then derive a survival function from  
387 equation Eq9 as follows:

$$388 \quad S(\mathbf{t} \mid \mathbf{x}_i, \beta_i) = [S_{0i}(\mathbf{t})]^{\exp(\mathbf{x}_i \beta_i^T)} \quad (\text{Eq10})$$

389 where:

$$390 \quad S_{0i}(\mathbf{t}) = \exp \left\{ - \int_0^t h_i(\mathbf{t} \mid \mathbf{x}_i, \beta_i) du \right\} \quad (\text{Eq11})$$



391 is the baseline survival for individuals at times  $\mathbf{t}$ , given a vector of candidate predictors  $\mathbf{x}_i$  for class  $i$ .  
392 The baseline hazard  $h_{0i}(\mathbf{t})$  in Eq9 is assumed to be an unknown, arbitrary and non-negative function of  
393 time and the only parametric part of the model in Eq10 is  $\exp(\mathbf{x}_i\boldsymbol{\beta}_i^T)$  [25]. The maximum likelihood  
394 procedure fails to estimate parameters for the likelihood function of Eq9 accurately because the baseline  
395 hazard function is not assumed to take any particular form. Instead, these parameters can be estimated  
396 using the partial-likelihood approach [26]. This is derived by taking the product of the conditional risk  
397 at time  $t_i$  given the set of individuals not yet dead, lost to follow-up or censored by that time.

398

## 399 **Part 2. Model specification, selection, validation and software used in** 400 **Procedures 1, 2a/b, 3a/b and 4.**

### 401 Model specification, selection and validation

402 All subsets regression was deployed [27], along with  $k$ -fold cross-validation as recommended by Grimm  
403 et al. [28], to find the best-fitting model for Procedures 1-4, with four covariates considered for both  
404 Cox proportional hazards models *and* (where applicable) the latent class models. The area under the  
405 receiver operating characteristic (ROC) curve (AUC) was used to evaluate all models generated – an  
406 approach that has been widely used in medical research to assess the diagnostic acuity of biomarkers to  
407 discriminate between diseased and healthy subjects [29-31]. In this way the AUC was used in the present  
408 study to quantify the extent to which each modelling Procedure was able to discriminate between  
409 individuals/participants and classes at risk of mortality. AUC values range from 0.5 to 1, where 0.5  
410 indicates that the discrimination achieved is equivalent to (and no better than that that could be achieved)  
411 by chance; a value of 1 indicates perfect discrimination; and a value  $>0.8$  is interpreted as evidence of  
412 good discrimination.  $k$ -fold cross-validation involved randomly dividing the dataset into  $k$  partitions of  
413 approximately equal size, where  $k - 1$  partitions were used as a training set and the model was evaluated  
414 and validated using the remaining  $k^{\text{th}}$  partition, repeated  $k$  times. The value  $k = 10$  was chosen based on  
415 established (and evaluated) best practice [32], with  $k = 10$  favoured for less biased model parameters,

416 according to experimentation [33]. The AUC was calculated for each of the 10 test samples, with  
417 subsequent confirmation of the results obtained from 10 iterations assessed using a bootstrap re-sampling  
418 procedure 1000 times (creating datasets from the original data without making further assumptions) to  
419 provide empirical 95% confidence intervals [34].

420 Covariate selection was guided by the desire to achieve parsimonious models according to the Bayesian  
421 Information Criterion (BIC) – the statistic preferred as the most parsimonious penalised likelihood  
422 statistic to minimise the risk of overfitting [35]. In choosing the optimum number of latent classes for  
423 the latent variable models (i.e. LCA and LCR), BIC was again the preferred statistic as simulations have  
424 demonstrated it outperforms other model fit statistics [36]. Strategies for determining the optimal  
425 number of classes may also be influence by interpretability (such as clinical salience and/or utility  
426 [33,37]), which is often reflected in ‘entropy’ [38] – a measure of the consistency between the modal  
427 (highest probability) and probabilistic (exact probabilities) assignment of individuals to latent classes.  
428 A high entropy indicates that individuals are more aligned to a single class (large modal probability),  
429 which leads to clearer interpretation of each latent class. A low entropy does not preclude latent classes  
430 having utility and substantive meaning, but individuals may not be as clearly aligned to just one class,  
431 making modal assignment a poor representation of the latent class structure.

#### 432 Software

433 An important challenge with latent class modelling is its sensitivity to starting values, because these are  
434 used to maximize the likelihood function when estimating model parameters. Where the starting values  
435 are far from the optimum solution, the likelihood function takes longer to converge and may even fail to  
436 do so. Occasionally, up to 50% of the random starts chosen will generate meaningful solutions when the  
437 likelihood function is maximized. For a solution to be meaningful, the highest likelihood value is  
438 expected to be replicated many times. When this does not occur, it signifies that either: no solution has  
439 been achieved and the number of random starts needs to be increased to converge on a global optimum  
440 solution; or the specified model structure is unsuitable for the given dataset. While this can add to the

441 time required to explore optimum solutions, once the target values are estimated they can be used as  
442 initial values for the final models derived, thereby reducing the duration of the final search process [23].

## 443 **References**

444

- 445 1. Wang Y, Kung L, Byrd TA. Big data analytics: Understanding its capabilities and potential  
446 benefits for healthcare organizations. *Technol Forecast Soc Change*. 2018; 126: 3-13.  
447 <http://dx.doi.org/10.1016/j.techfore.2015.12.019>
- 448 2. Diebold FX. On the origin(s) and development of the term ‘Big Data’. Penn Institute for  
449 Economic Research Working Paper. 2012; 12-037: 1-7.  
450 <https://economics.sas.upenn.edu/sites/default/files/filevault/12-037.pdf>
- 451
- 452 3. Harrell Jr FE, Lee KL, Califf RM, Pryor DB, Rosati RA. Regression modelling strategies for  
453 improved prognostic prediction. *Stat Med*. 1984; 3: 143-52.  
454 <http://dx.doi.org/10.1002/sim.4780030207>
- 455
- 456 4. DeFilippis AP, Young R, Carrubba CJ, McEvoy JW, Budoff MJ, Blumenthal RS et al. An  
457 analysis of calibration and discrimination among multiple cardiovascular risk scores in a  
458 modern multiethnic cohort. *Ann Intern Med*. 2015; 162: 266-75.  
459 <https://dx.doi.org/10.7326%2FM14-1281> PMID: 25686167
- 460
- 461 5. O’Donnell CJ. Opportunities and challenges for polygenic risk scores in prognostication and  
462 prevention of cardiovascular disease. *JAMA Cardiol*. 2020; 5: 399-400.  
463 <https://dx.doi.org/10.1001/jamacardio.2019.6232>
- 464
- 465 6. Holmberg C, Parascandola M. Individualised risk estimation and the nature of prevention.

- 466 Health Risk Soc. 2010; 12: 441-52. <https://doi.org/10.1080/13698575.2010.508835>
- 467
- 468 7. Huitfeldt A. Is caviar a risk factor for being a millionaire? *Br Med J.* 2016; 355: i6536.
- 469 <https://doi.org/10.1136/bmj.i6536>
- 470
- 471 8. Killu AM, Granger CB, Gersh BJ. Risk stratification for stroke in atrial fibrillation: a critique.
- 472 *Eur Heart J.* 2019; 40: 1294-1302. <https://dx.doi.org/10.1093/eurheartj/ehy731> PMID:
- 473 30508086
- 474
- 475 9. Arnold KF, Davies V, de Kamps M, Tennant PW, Mbotwa J, Gilthorpe MS. Reflections on
- 476 modern methods: generalized linear models for prognosis and intervention—theory, practice
- 477 and implications for machine learning. *Int J Epidemiol.* 2020; May 7: dyaa049.
- 478 <https://dx.doi.org/10.1093/ije/dyaa049> PMID: 32380551
- 479
- 480 10. Rockhill B, Kawachi I, Colditz GA. Individual risk prediction and population-wide disease
- 481 prevention. *Epidemiologic Reviews* 2000; 22: 176-80.
- 482 <https://dx.doi.org/10.1093/oxfordjournals.epirev.a018017> PMID: 10939025
- 483
- 484 11. Wilkinson J, Arnold KF, Murray EJ, van Smeden M, Carr K, Sippy R et al. Time to reality
- 485 check the promises of machine learning-powered precision medicine. *Lancet Digit Health.*
- 486 2020; Sep 16. [https://doi.org/10.1016/S2589-7500\(20\)30200-4](https://doi.org/10.1016/S2589-7500(20)30200-4)
- 487
- 488 12. Cole SR, Platt RW, Schisterman EF, Chu H, Westreich D, Richardson D et al. Illustrating bias
- 489 due to conditioning on a collider. *Int J Epidemiol.* 2010; 39: 417-20.
- 490 <https://dx.doi.org/10.1093/ije/dyp334> PMID: 19926667
- 491
- 492 13. Elwert F, Winship C. Endogenous selection bias: The problem of conditioning on a collider

- 493 variable. *Annu Rev Sociol.* 2014; 40: 31-53. <https://doi.org/10.1146/annurev-soc-071913->  
494 043455
- 495
- 496 14. Gadd SC, Tennant PWG, Heppenstall AJ, Boehnke JR, Gilthorpe MS. Analysing trajectories  
497 of a longitudinal exposure: a causal perspective on common methods in Lifecourse  
498 research. *Plos One.* 2019;14(12)
- 499 15. Richiardi L, Bellocco R, Zugna D. Mediation analysis in epidemiology: methods,  
500 interpretation and bias. *Int J Epidemiol.* 2013; 42: 1511-19.  
501 <https://dx.doi.org/10.1093/ije/dyt127> PMID: 24019424
- 502
- 503 16. Westreich D, Greenland S. The table 2 fallacy: presenting and interpreting confounder and  
504 modifier coefficients. *Am J Epidemiol.* 2013; 177: 292-8.  
505 <https://dx.doi.org/10.1093/aje/kws412> PMID: 23371353
- 506
- 507 17. Tennant PWG, Harrison WJ, Murray EJ, Arnold KF, Berrie L, Fox MP et al. Use of directed  
508 acyclic graphs (DAGs) in applied health research: review and recommendations. *Int J*  
509 *Epidemiol.* 2020; in press. <https://dx.doi.org/10.1101/2019.12.20.19015511>
- 510
- 511 18. Dean N, Raftery AE. Latent class analysis variable selection. *Ann Inst Stat.* 2010; 62: 11-35.  
512 <https://dx.doi.org/10.1007/s10463-009-0258-9> PMID: 20827439
- 513
- 514 19. Witte KK, Drozd M, Walker AM, Patel PA, Kearney JC, Chapman S et al. Mortality reduction  
515 associated with  $\beta$ -adrenoceptor inhibition in chronic heart failure is greater in patients with  
516 diabetes. *Diabetes Care* 2018; 41: 136-42. <https://doi.org/10.2337/dc17-1406>
- 517
- 518 20. Witte KK, Patel PA, Walker AM, Schechter CB, Drozd M, Sengupta A et al. Socioeconomic  
519 deprivation and mode-specific outcomes in patients with chronic heart failure. *Heart* 2018;

- 520 104: 993-8. <https://dx.doi.org/10.1136/heartjnl-2017-312539> PMID: 29386325
- 521
- 522 21. Cubbon RM, Gale CP, Kearney LC, Schechter CB, Brooksby WP, Nolan J et al. Changing  
523 characteristics and mode of death associated with chronic heart failure caused by left  
524 ventricular systolic dysfunction: a study across therapeutic eras. *Circ Heart Fail.* 2011; 4: 396-  
525 403. <https://dx.doi.org/10.1161/CIRCHEARTFAILURE.111.963066> PMID: 21772014
- 526
- 527 22. Hallquist M, Wiley J. Mplus Automation: Automating Mplus model estimation and  
528 interpretation. 2014 <https://cran.r-project.org/web/packages/MplusAutomation/index.html>.
- 529
- 530 23. Mplus. Los Angeles, CA: Muthén & Muthén, 2014.
- 531
- 532 24. Mandrekar JN, Receiver operating characteristic curve in Diagnostic test assessment. *J.*  
533 *Thorac. Oncol.*, 5 (2010), pp. 1315-1316
- 534
- 535 25. Cox DR. Regression models and life tables (with discussion). *J R Stat Soc Series B Stat*  
536 *Methodol.* 1972; 34: 187-220.
- 537
- 538 26. Cox DR. Partial likelihood. *Biometrika.* 1975; 62: 269-76.
- 539
- 540 27. Kuk AY. All subsets regression in a proportional hazards model. *Biometrika* 1984; 71: 587-92.
- 541
- 542 28. Grimm KJ, Mazza GL, Davoudzadeh P. Model selection in finite mixture models: A k-fold  
543 cross-validation approach. *Struct Equ Modeling.* 2017; 24: 246-56.  
544 <https://doi.org/10.1080/10705511.2016.1250638>
- 545
- 546 29. Metz CE. Basic principles of ROC analysis. *Semin Nucl Med.* 1978; 8: 283-98.

- 547 [https://doi.org/10.1016/S0001-2998\(78\)80013-0](https://doi.org/10.1016/S0001-2998(78)80013-0)
- 548
- 549 30. Hajian-Tilaki K. Receiver operating characteristic (ROC) curve analysis for medical  
550 diagnostic test evaluation. *Caspian J Intern Med.* 2013; 4: 627-35. PMID: 24009950
- 551
- 552 31. Heagerty PJ, Lumley T, Pepe MS. Time-dependent ROC curves for censored survival data and  
553 a diagnostic marker. *Biometrics.* 2000; 56: 337-44. [https://dx.doi.org/10.1111/j.0006-](https://dx.doi.org/10.1111/j.0006-341x.2000.00337.x)  
554 [341x.2000.00337.x](https://dx.doi.org/10.1111/j.0006-341x.2000.00337.x). PMID: 10877287
- 555
- 556 32. Kuhn M, Johnson K. *Applied Predictive Modeling*. Springer-Verlag, New York; 2013.  
557 <https://doi.org/10.1007/978-1-4614-6849-3>
- 558
- 559 33. Harrison WJ, Gilthorpe MS, Downing A, Baxter PD. Multilevel latent class modelling of  
560 colorectal cancer survival status at three years and socioeconomic background whilst  
561 incorporating stage of disease. *Int J Stat Prob.* 2013; 2: 85-95.  
562 <https://doi.org/10.5539/ijsp.v2n3p85>
- 563
- 564 34. Bland JM, Altman DG. Statistics notes: bootstrap resampling methods. *Br Med J.* 2015; 350:  
565 h2622. <https://doi.org/10.1136/bmj.h2622>
- 566
- 567 35. Hitchcock C, Sober E. Prediction versus accommodation and the risk of overfitting. *The Br J*  
568 *Philosoph Sci.* 2004; 55: 1-34. <https://doi.org/10.1093/BJPS/55.1.1>
- 569
- 570 36. Nylund KL, Asparoutiov T, Muthen BO. Deciding on the number of classes in latent class  
571 analysis and growth mixture modeling: A Monte Carlo simulation study. *Struct Equ Modeling.*  
572 2007; 14: 535-69. <https://doi.org/10.1080/10705510701575396>
- 573

574 37. Gilthorpe MS, Dahly DL, Tu YK, Kubzansky LD, Goodman E. Challenges in modelling the  
575 random structure correctly in growth mixture models and the impact this has on model  
576 mixtures. *J Dev Orig Health Dis.* 2014; 5: 197-  
577 205. <https://doi.org/10.1017/S2040174414000130> PMID: 24901659

578  
579 38. Celeux G, Soromenho G. An entropy criterion for assessing the number of clusters in a  
580 mixture model. *J Classif.* 1996; 13: 195-212. <https://doi.org/10.1007/BF01246098>

581  
582

583

584

585



**Table 1. Descriptive characteristics of the study cohort.**

<b>Study Cohort</b>	
	<b>N (%)</b>
<b>Participants</b>	1,796 (100.0)
<b>Deaths</b>	1,061 (59.1)
<b>Male</b>	1,313 (73.1)
<b>Type 2 Diabetes</b>	504 (28.1)
<hr/>	
	<b>Median (IRQ)</b>
<b>Survival Time (years)</b>	3.40 (2.11, 5.78)
<hr/>	
	<b>Mean (95% CI)</b>
<b>Age (years)</b>	69.7 (69.1, 70.2)
<b>Haemoglobin (g/dl)</b>	13.46 (13.38, 13.54)

N = number; % = percentage; IQR = interquartile range; CI = confidence interval.

**Table 2. Latent class analysis (LCA) model summaries – the preferred model from this step was used in Procedures 2 and 3.**

Number of classes	Number of parameters	BIC	Entropy	Class	Modal N (%)	Probabilistic N (%)
1	6	19,818.53	-		1,796 (100.0)	-
2	11	19,537.79	0.75	Class 1	1,452 (80.8)	1425.3 (79.4)
				Class 2	344 (19.2)	370.7 (20.6)
3	16	19,445.74	0.74	Class 1	1,203 (67.0)	11744.0 (65.4)
				Class 2	480 (26.7)	500.7 (27.9)
				Class 3	113 (6.3)	120.3 (6.7)
4	21	19,422.35	0.80	Class 1	811 (45.2)	797.0 (44.4)
				Class 2	486 (27.1)	504.4 (28.1)
				Class 3	381 (21.2)	371.4 (20.7)
				Class 4	118 (6.6)	123.2 (6.9)
5	26	19,421.44	0.67	<b>Class 1</b>	<b>586 (32.6)</b>	<b>566.7 (31.6)</b>
				<b>Class 2</b>	<b>470 (26.2)</b>	<b>459.7 (25.6)</b>
				<b>Class 3</b>	<b>324 (18.0)</b>	<b>296.9 (16.5)</b>
				<b>Class 4</b>	<b>317 (17.7)</b>	<b>368.6 (20.5)</b>
				<b>Class 5</b>	<b>99 (5.5)</b>	<b>104.1 (5.8)</b>
6	31	19,422.87	0.63	Class 1	527 (29.3)	517.7 (28.8)
				Class 2	474 (26.4)	470.5 (26.2)
				Class 3	276 (15.4)	247.7 (13.8)
				Class 4	234 (13.0)	232.6 (13.0)
				Class 5	186 (10.4)	229.8 (12.8)
				Class 6	99 (5.5)	97.6 (5.4)

BIC = Bayesian information criterion; N = number; % = percentage; the optimal LCA model according to the BIC is emboldened.

**Table 3. Latent class regression (LCR) model summaries for Procedure 4**

<b>Number of classes</b>	<b>Number of parameters</b>	<b>BIC</b>	<b>Entropy</b>
1	3	3695.06	----
<b>2</b>	<b>10</b>	<b>3659.49</b>	<b>0.68</b>
3	17	3682.44	0.91
4	24	3722.89	0.94

BIC = Bayesian information criterion; the optimal LCA model according to the BIC is emboldened.

**Table 4. Covariate coefficients for each preferred model (Procedures 1-4) executed on the complete data, along with median AUC and empirical 95% empirical confidence intervals generated through 10-fold cross-validation.**

<b>Model (AUC: 95% CI)</b>		<b>HR (95% CI)</b>
<b>Procedure 1 - CPH</b> (AUC = 0.69: 0.67, 0.72)		
	Type 2 Diabetic vs. not	1.35 (1.16, 1.59)
	Male vs. Female	1.76 (1.47, 2.11)
	Age (per 5 years)	1.24 (1.20, 1.29)
	Haemoglobin (per g/dl)	0.82 (0.78, 0.86)
<b>Procedure 2a - LCA (modal) / CPH</b> (AUC = 0.65: 0.62, 0.67)		
†Class 1 (N = 586) vs:	Class 2 (470)	0.35 (0.30, 0.44)
	Class 3 (324)	1.33 (1.10, 1.60)
	Class 4 (317)	0.71 (0.57, 0.87)
	Class 5 (99)	0.17 (0.10, 0.29)
<b>Procedure 2b - LCA (probabilistic) / CPH:</b> (AUC = 0.66: 0.64, 0.68)		
‡Class 1 (32.0%) vs:	Class 2 (26.0%)	0.26 (0.19, 0.34)
	Class 3 (18.0%)	1.00 (0.71, 1.39)
	Class 4 (18.0%)	1.58 (1.27, 1.97)
	Class 5 (6.0%)	0.17 (0.09, 0.32)
<b>Procedure 3a - LCA (modal) / CPH</b> (AUC = 0.69: 0.66, 0.72)		
	Type 2 Diabetic vs. not	1.51 (1.13, 2.01)
	Male vs. Female	1.80 (1.49, 2.17)
	Age (per 5 years)	1.21 (1.13, 1.29)
	Haemoglobin (per g/dl)	0.82 (0.79, 0.86)
†Class 1 (N = 586) vs:	Class 2 (470)	0.77 (0.53, 1.10)
	Class 3 (324)	0.84 (0.59, 1.19)
	Class 4 (317)	0.92 (0.71, 1.20)
	Class 5 (99)	0.79 (0.38, 1.67)
<b>Procedure 3b - LCA (probabilistic) / CPH</b> (AUC = 0.69: 0.66, 0.72)		
	Type 2 Diabetic vs. not	1.44 (1.01, 2.06)
	Male vs. Female	1.70 (1.31, 2.21)
	Age (per 5 years)	1.21 (1.11, 1.32)
	Haemoglobin (per g/dl)	0.81 (0.76, 0.88)
‡Class 1 (32.0%) vs:	Class 2 (26.0%)	0.78 (0.41, 1.49)
	Class 3 (18.0%)	0.90 (0.55, 1.48)
	Class 4 (18.0%)	1.15 (0.56, 2.36)
	Class 5 (6.0%)	0.99 (0.35, 2.78)
<b>Procedure 4 – LCR</b> (AUC = 0.79: 0.74, 0.85)		
<i>Cox proportional hazards model</i>		
Class 1 ('High risk'):	Type 2 Diabetic vs. not	1.26 (0.91, 1.75)
	Male vs. Female	2.07 (1.58, 2.71)
	Age (per 5 years)	1.36 (1.28, 1.44)
Class 2 ('Low risk'):	Type 2 Diabetic vs. not	0.44 (0.23, 0.82)
	Male vs. Female	1.01 (0.64, 1.60)
	Age (per 5 years)	1.17 (1.06, 1.29)
<i>Class membership model</i>		
‘High’ vs. ‘Low’ risk:	Type 2 Diabetic vs. not	0.27 (0.09, 0.76)
	Haemoglobin (per g/dl)	2.16 (1.64, 2.84)

AUC = area under the curve; CI = empirical confidence interval obtained from the 2.5% to 97.5% centiles of bootstrapped samples following 10-fold cross-validation; HR = hazards ratio; OR = odds ratio; CPH = Cox proportional hazards; LCA = latent class analysis (modal assignment or probabilistic assignment); LCR = latent class regression.

**Table 5. Descriptive characteristics for the 2-class Cox proportional hazards latent class regression model.**

	<b>Latent Class Regression Model</b>			
	<b>Class 1 ('High risk')</b>		<b>Class 2 ('Low risk')</b>	
	<b>Modal N (%)</b>	<b>Probabilistic N (%)</b>	<b>Modal N (%)</b>	<b>Probabilistic N (%)</b>
<b>Participants</b>	1,566 (87.2)	1507.8 (84.0)	230 (22.8)	288.2 (16.0)
<b>Deaths</b>	1,046 (66.8)	1014.7 (67.3)	15 (6.5)	45.8 (15.9)
<b>Male</b>	1,160 (74.1)	1112.8 (73.8)	153 (66.5)	200.9 (69.7)
<b>Type 2 Diabetes</b>	368 (23.5)	342.3 (22.7)	136 (59.1)	162.5 (56.4)
	<b>Median (IRQ)</b>		<b>Median (IRQ)</b>	
<b>Survival Time (years)</b>	3.86 (2.41, 5.89)		1.13 (0.50, 2.27)	
	<b>Mean (95% CI)</b>		<b>Mean (95% CI)</b>	
<b>Age (years)</b>	69.2 (68.6, 69.9)		72.5 (71.1, 73.9)	
<b>Haemoglobin (g/dl)</b>	13.80 (13.72, 13.88)		11.14 (10.99, 11.30)	

N = number; % = percentage; IQR = interquartile range; CI = confidence interval.