



Benedetto, U., Sinha, S., Lyon, M. S., Dimagli, A., Gaunt, T. R., Angelini, G. D., & Sterne, J. A. C. (2020). Can machine learning improve mortality prediction following cardiac surgery? *European Journal of Cardio-Thoracic Surgery*, [ezaa229].
<https://doi.org/10.1093/ejcts/ezaa229>

Peer reviewed version

Link to published version (if available):
[10.1093/ejcts/ezaa229](https://doi.org/10.1093/ejcts/ezaa229)

[Link to publication record in Explore Bristol Research](#)
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Oxford University Press at <https://doi.org/10.1093/ejcts/ezaa229> . Please refer to any applicable terms of use of the publisher.

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:
<http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

1 **Can machine learning improve mortality prediction following cardiac surgery?**

2 Umberto Benedetto^{1,2}, Shubhra Sinha¹, Matt Lyon^{2,3,4}, Arnaldo Dimagli¹, Tom R
3 Gaunt^{2,3,4}, Gianni Angelini^{1,2}, Jonathan Sterne^{2,3,4}

4

5 1. Translational Health Sciences, Bristol Heart Institute, University of Bristol

6 2. NIHR Bristol Biomedical Research Centre, University of Bristol and University
7 Hospitals Bristol NHS Foundation Trust

8 3. Population Health Sciences, Bristol Medical School, University of Bristol

9 4. MRC Integrative Epidemiology Unit, University of Bristol

10

11

12

13 **Corresponding author**

14 Umberto Benedetto

15 Office Room 84 Level 7

16 Bristol Royal Infirmary, Upper Maudlin Street BS2 8HW

17 Tel. +44 (0) 117 3428854

18 umberto.benedetto@bristol.ac.uk

19

20 **Word count: 3828**

21 **Visual abstract**

22 Key question: Do modern machine learning models improve the prediction of in-
23 hospital mortality after cardiac surgery?

24 Key findings: machine learning models performed similarly to logistic regression
25 models.

26 Take-home message: prediction of in-hospital mortality is not improved by machine
27 learning relative to traditional methods based on logistic regression

28 **ABSTRACT**

29 **Objective(s):** Interest in the clinical usefulness of machine learning (ML) for risk
30 prediction has bloomed recently. Cardiac surgery patients are at high risk of
31 complications and therefore pre-surgical risk assessment is of crucial relevance. We
32 aimed to compare the performance of ML algorithms over traditional logistic
33 regression (LR) model to predict in-hospital mortality following cardiac surgery.

34 **Methods:** A single centre dataset of prospectively collected information from patients
35 undergoing adult cardiac surgery from 1996 to 2017 was split into 70% training set
36 and 30% testing set. Prediction models were developed using neural network, random
37 forest, naïve Bayes and retrained logistic regression based on features included in the
38 EuroSCORE. Discrimination was assessed using area under the receiver operating
39 characteristic curve (AUC) and calibration analysis was undertaken using calibration
40 belt method. Model calibration drift was assessed by comparing Goodness of fit chi-
41 squared statistics observed in two equal bins from the testing sample ordered by
42 procedure date.

43 **Results:** A total of 28,761 cardiac procedures were performed during the study period.
44 The in-hospital mortality rate was 2.7%. Retrained LR (AUC 0.80; 95% CI 0.77, 0.83)
45 and random forest model (0.80; 95%CI 0.76, 0.83) showed the best discrimination. All
46 models showed significant miscalibration. Retrained LR proved to have the weakest
47 calibration drift.

48 **Conclusions:** Our findings do not support the hypothesis that ML methods provide
49 advantage over LR model in predicting operative mortality after cardiac surgery.

50 **Keywords:** machine learning, mortality prediction, neural network, random forest,
51 naïve Bayes.

52 **ABBREVIATIONS**

53 AUC: Area Under the Receiver Operating Characteristic curve

54 LR: logistic regression

55 ML: machine learning

56 INTRODUCTION

57 Pre-operative assessment of surgical risk is of crucial importance in cardiac surgery
58 due to the high risk of intraoperative and postoperative complications. Risk models
59 can help health professionals to advise patients during the decision-making process,
60 as well as in monitoring surgical performance and cost-benefit analyses.

61 Several risk stratification models have been developed to predict in-hospital mortality
62 following cardiac surgery, for example as the European System for Cardiac Operative
63 Risk Evaluation, EuroSCORE [1,2] and the North American Society of Thoracic
64 Surgeons (STS) score [3]. However, a main limitation of these scores is overestimation
65 of risk in high risk patient subgroups [4,5]. This can potentially translate into risk-averse
66 practice, falsely reassuring conclusions about surgeon and centre performance, and
67 impaired decision-making.

68 Current risk scoring systems are based on logistic regression (LR). Development of
69 LR models requires input from the modeler to address complex interaction among
70 features and non-linear relationships of features with the outcome. For instance, the
71 contribution of advanced age to mortality risk may not be constant across the spectrum
72 of co-morbidities. If features interactions are overlooked in a LR model, its prediction
73 ability will be negatively affected. In contrast, machine learning (ML) algorithms require
74 less input from the modeler and interactions among features and non-linear
75 relationships can be learnt automatically from the data [6]. However, the extra flexibility
76 of ML algorithm requires larger sample to train the model.

77 Despite research on the utility of ML methods to improve prediction in healthcare has
78 exponentially increased, ML methods have not been widely adopted in the clinical
79 practice. Moreover, recent reports have challenged the additional value of ML in the
80 development of clinical prediction models in a variety of clinical conditions [6].

81 The objective of this study was to compare ML algorithms with LR model in the
82 prediction of in-hospital mortality after cardiac surgery, based on the set of features
83 included in the EuroSCORE [1].

84

85 **METHODS**

86 The present study was approved by Health Research Authority and Health and Care
87 Research Wales. Data were obtained from the National Adult Cardiac Surgery Audit
88 (NACSA) dataset which prospectively collects clinical information for all major heart
89 operations carried out in the United Kingdom. In the present analysis, we used a
90 subset of patients who underwent cardiac surgery at University Hospitals Bristol NHS
91 Trust between 1 April 1996 and 30 December 2017.

92 Missing or conflicting data for in-hospital mortality were obtained via record linkage to
93 the Office for National Statistics census database. For records where data required to
94 calculate a EuroSCORE variable was missing, it was assumed that the risk factor was
95 not present (equal to the reference level). Missing patient age at the time of surgery
96 was imputed as the median patient age for the corresponding financial year.

97 **Statistical analysis and models**

98 The primary endpoint was in-hospital mortality following cardiac surgery. Numerical
99 variables were summarised as mean and standard deviation or median and
100 interquartile range and compared using t-tests or Mann-Whitney tests. Categorical
101 variables were tabulated as frequencies and percentages and compared using chi-
102 squared **test**.

103 Procedures were ordered chronologically, the first 70% of records (01/04/96 -
104 27/09/11) were used for training and hyperparameter selection through five-fold cross-
105 validation. Final model performance was evaluated using the remaining 30%

106 (27/09/11 - 30/12/17). All prediction models were developed using the 17 features
107 included in the original EuroSCORE [1], which include information prior to surgery on
108 a range of patient, cardiac and operative factors. The features are age, gender, chronic
109 obstructive pulmonary disease, extracardiac arteriopathy, neurological dysfunction,
110 previous cardiac surgery, creatinine >200 $\mu\text{mol/l}$, active endocarditis, critical
111 preoperative state, unstable angina, left ventricular function, recent myocardial
112 infarction, pulmonary hypertension, emergency surgery, combined surgery other than
113 coronary artery bypass graft, surgery on thoracic aorta, post-infarct septal rupture.
114 We fitted a logistic regression (LR) ('retrained LR') model to the EuroSCORE risk
115 factors. We used the following ML approaches:

- 116 • Neural Network is a computational learning system that uses a network of functions
117 to understand and translate a data input of one form into a desired output. Machine
118 learning algorithms including neural networks generally do not need to be
119 programmed with specific rules that define what to expect from the input. The
120 neural net learning algorithm instead learns from processing many labelled
121 examples (i.e. data with "answers") that are supplied during training and using this
122 answer key to learn what characteristics of the input are needed to construct the
123 correct output. Once a sufficient number of examples have been processed, the
124 neural network can begin to process new, unseen inputs and successfully return
125 accurate results. The more examples and variety of inputs the program sees, the
126 more accurate the results typically become because the program learns with
127 experience. The basic unit of computation in a neural network is the **neuron**, often
128 called a node or unit. It receives input from some other nodes, or from an external
129 source and computes an output. Each input has an associated weight (w), which
130 is assigned on the basis of its relative importance to other inputs. The node applies

131 a function f to the weighted sum of its inputs (i.e. $f(w_1+w_2+w_3\dots)$ to introduce non-
132 linearity into the output of a neuron. Nodes are arranged in layers. Nodes from
133 adjacent layers have connections or edges between them. All these connections
134 have weights associated with them. Neural network consists of three types of
135 nodes: neural nets consist of 3 layers. 1) Input Layers: this entry point takes input
136 data (i.e. numbers, texts, etc); 2) Hidden Layers: are responsible for number
137 crunching i.e. mathematical operation, to detect patterns data. There can be a
138 minimum of one and many multiple hidden layers; 3) Output Layer: takes input
139 from the hidden layer to generate the desired output. [7,8]. As almost all ML
140 approaches, neural networks were not meant for time-related event, but as
141 research rapidly moved forward new methods have been introduced for this
142 purpose [9]. In our model, number of hidden layers and nodes per hidden layer
143 were configured manually in response to model discrimination (area under the
144 receiver operating characteristic curve [AUC]) evaluated with cross-validation. The
145 final model configuration used for evaluation was: input layer n=18 nodes, hidden-
146 layer one n=90 nodes, hidden-layer two n=36 nodes, output layer one node.

- 147 • Random Forest represents an ensemble of several decision trees. Decision tree
148 builds classification or regression models in the form of a tree structure. It breaks
149 down a dataset into smaller and smaller subsets while at the same time an
150 associated decision tree is incrementally developed. The final result is a tree
151 with decision nodes and leaf nodes. A decision node has two or more branches.
152 Leaf node represents a classification or decision. The topmost decision node in a
153 tree corresponds to the best predictor called root node, which splits the records
154 into mutually exclusive classes. After the root node, there are internal nodes which
155 lead to other internal nodes or to two or more terminal leaf nodes. An item is

156 classified according to which leaf node is reached. Each item can be trained using
157 resampling methods (i.e. bootstrapping) [10,11]. Random forest has several
158 parameters that have to be set by the user, e.g., number of trees in the forest
159 (estimator), max number of levels (depth) in each decision tree, min number of
160 data points placed in a node before the node is split and minimum samples leaf.
161 When new data are presented, each tree of the random forest votes for a class
162 and the final prediction is based on the class receiving the majority of the votes. In
163 our model, we manually tuned parameters in response to model discrimination
164 (AUC) evaluated with cross-validation. (estimators n=700, maximum depth n=10,
165 minimum samples split n=5, minimum samples leaf n=20).

- 166 • Naïve Bayes: is based on the Byes theorem. It is called “naïve” because it assumes
167 each feature contributes independently to the probability of classification. The final
168 prediction of the model is the a priori probability modified by the likelihood of each
169 predictor [12]. In our model, we used default parameters.

170 Full model configurations and discrimination are provided in the Supplementary Table
171 1. Models were developed and evaluate using scikit-learn v0.21.2 and TensorFlow
172 v1.14.0 through Anaconda Python 3 v2019.07.

173 Discrimination was assessed by calculating model AUC with its relative 95%
174 confidence interval using bootstrapping (2000 repetitions) (pROC R-package v1.15.3).
175 The assessment of calibration, i.e., the model's ability to provide reliable predictions,
176 is crucial to test risk models. Statistical techniques such, the Hosmer–Lemeshow
177 statistics and the Cox calibration test, are all non-informative with respect to calibration
178 across risk classes. To better characterise the calibration of new models we used the
179 calibration belt model [13]. In this new approach, the relation between the logits of the
180 probability predicted by a model and of the event rates observed in a sample is

181 represented by a polynomial function, whose coefficients are fitted and its degree is
182 fixed by a series of likelihood-ratio tests. This method also enables confidence
183 intervals to be computed for the curve, which can be plotted [13] (R-package givitiR
184 v1.3) (R-package ResourceSelection v0.3.5). The calibration belt produces a trend
185 with the 95% confidence interval containing the line of equality. Open-source code is
186 available from: <https://github.com/MRCIEU/cvd-mortality-ml>.

187 We also reported the performance of the original EuroSCORE I and EuroSCORE II
188 for completeness. We were able to calculate the EuroSCORE II [2] only in 1889
189 (21.9%) patients for whom exact values of serum creatinine were available.

190

191 **RESULTS**

192 **Participants**

193 A total of 28,761 cardiac procedures were included in the final dataset (Supplementary
194 Figure 1). Patients younger than 18 years at time of surgery were excluded (n=41) to
195 avoid inclusion of congenital abnormalities. The outcome and full set of features were
196 available for all records after imputation. The overall percentage of missing data in the
197 EuroSCORE variables was very low (1.7%) and records of age were missing in 86
198 patients. Patient characteristics are presented in Table 1. All features included in
199 EuroSCORE I were robustly associated with the outcome in univariable analyses,
200 except of elevated systolic pulmonary pressure. In-hospital mortality rate was 2.7%
201 (n=786).

202 **Model discrimination**

203 Results of model selection and hyperparameter tuning using the training set are
204 reported in Supplementary Table 1. Discrimination ability of models selected in the
205 testing set is presented in Figure 1. Retrained LR showed good discrimination (AUC

206 0.80; 95% CI 0.77, 0.83). Among the ML classifiers, random forest showed the best
207 discrimination ability (0.80; 95% CI 0.76, 0.83) which was comparable to retrained LR
208 model. Neural network and naïve Bayes AUC were 0.77 (95% CI 0.73, 0.80) and 0.77
209 (95% CI 0.74, 0.80) respectively. Original EuroSCORE I and II AUC were 0.76 (95%
210 CI 0.73, 0.79) and 0.77 (0.70, 0.84) respectively.

211 **Probability calibration**

212 Retrained LR had strong evidence against the null hypothesis of well calibrated
213 probabilities when applied to our data ($P < 0.001$; Figure 2 Panel A). Among the
214 contemporary classifiers, neural network and random forest also showed poor
215 calibration ($P < 0.001$; Figure 2 Panel B and C) although the latter produced probabilities
216 that did not depart far from the line of equality. Naïve Bayes produced probabilities
217 that suggest very poor calibration. EuroSCORE I showed poor calibration ($P < 0.001$
218 Figure 3 Panel A) while EuroSCORE II was well calibrated although the sample size
219 and event number was smaller increasing the possibility of a type II error ($P = 0.64$
220 Figure 3 Panel B). To evaluate calibration drift in the retrained LR and ML models, the
221 test dataset was divided into two equal bins ordered by procedure date with
222 approximately equal number of events ($n = 102$ vs $n = 105$). Hosmer-Lemeshow
223 goodness of fit chi-squared statistics were calculated for first and second quantiles
224 (Table 2). Retrained LR had the weakest change in test statistic between quantiles
225 (+15.9%) and therefore weakest calibration drift. Random forest had the second
226 smallest effect (+21.2%). EuroSCORE II had too few events and could not be reliably
227 evaluated.

228

229 **DISCUSSION**

230 The main finding of the present study is that when trained on the same set of variables,
231 ML algorithms do not improve prediction over LR model. Both LR and random forest
232 model proved to be associated with good discrimination ability but substantial
233 miscalibration. However, these two models showed the least calibration drift.

234 Interest in risk-prediction models has bloomed in clinical use to aid in multidisciplinary
235 shared-decision making. They are also used for benchmarking outcomes and both
236 monitoring innovations. All this applies especially in an era of expanding multimodal
237 therapy for coronary artery and valve disease where risk prediction plays an important
238 role in determining which patients would benefit most from surgery or percutaneous
239 therapy. Moreover, national cardiac surgical registries have been established in many
240 countries and they are used to develop risk prediction model with improved
241 performance for local populations. Two of the most used risk stratification models in
242 cardiac surgery the European System for Cardiac Operative Risk Evaluation version
243 (EuroSCORE and EuroSCORE II) [1,2] and the STS-PROM Score [3] were both
244 developed based on LR. The EuroSCORE I and II have been extensively criticized
245 [14] including poor performance in external validation particularly for high-risk
246 subgroup [15,16]. This has been partially attributed to the small proportion (10%) of
247 patients aged 75 years and above in the reference dataset [17]. On the other hand,
248 STS provides superior discrimination when compared to EuroSCORE II, but it shows
249 suboptimal calibration, especially in the high-risk subgroup [18, 19].

250 It is possible that poor calibration of EuroSCORE II and STS score can be partially
251 attributed to the fact these LR-based models overlook complex interactions among
252 features and non-linear relationship. ML methods can capture interaction among
253 features and non-linearity without input from the modeller and this can potentially result
254 in improved prediction. A recent systematic review [20] on the application of ML

255 methods in cardiovascular diseases acknowledged the potential premise of ML in
256 certain applications such as automated imaging interpretation. However, the
257 advantage of ML methods over traditional risk stratification tools remains unclear.
258 Mendes et al. [21] found that neural networks did not outperform LR when predicting
259 mortality in patients after coronary artery bypass grafting. Other studies have
260 suggested an advantage from ML methods over LR. Random forest has been shown
261 to provided better discrimination when compared to LR, EuroSCORE and
262 EuroSCORE II [22,23]. Ghavidel et al. [24] found that decision trees achieved better
263 discrimination power when compared to EuroSCORE and retrained LR. Nilsson et al.
264 found that neural networks using 34 features determined a small improvement in
265 accuracy in mortality risk prediction when compared to LR and EuroSCORE [25].
266 Recently, Kilic et al. [26] reported that a new ML method (i.e. extreme gradient
267 boosting) may improve prediction in cardiac surgery when compared to the STS risk
268 models. These discordant results can partially be explained by the fact that ML
269 methods and in particular neural network need far more events per variable to be
270 trained and therefore their application should only be considered if very large data sets
271 are available [27]. An important limitation of available studies is that they focused on
272 model discrimination while calibration has been inconsistently reported. Discrimination
273 does not assess the model accuracy in individual risk predictions (calibration), which
274 is crucial when using a predictive model to inform decisions about individual patient.
275 Thus, a model might perform well based on discrimination measures while suffering
276 substantial miscalibration [28].
277 The present study was designed to get insights into the usefulness of ML methods to
278 improve individual risk prediction in cardiac surgery. We used a large dataset
279 collecting information on the set of features included in the EuroSCORE and we

280 assessed both model discrimination and calibration. We failed to show any significant
281 advantage from ML methods over traditional LR model based on the same set of
282 features included in the original EuroSCORE.

283 There are possible explanations for the lack of advantage from ML model over LR
284 observed in the present study. We had a limited number of events (hospital deaths) to
285 train and test prediction models despite the large original sample. This may have
286 limited our ability to exploit the superiority of ML methods in identifying patterns of
287 features related to the outcome. Moreover, automatic ML model hyper-tuning could
288 not be performed as dedicated technology required was not available. Age at the time
289 of surgery was the only continuous variable included in the models and this may have
290 limited the ability of ML models to capture non-linear interaction for continuous
291 variables. We did not train models using features included in the EuroSCORE II
292 because preoperative creatinine value was reported as dichotomous variable (<200
293 or ≥ 200 mmol/l) while the actual value, which is part of the EuroSCORE II, was
294 available only for a minority of patients. Similarly, we could not use the set of features
295 of the STS-PROM score because our dataset did not include some of the items
296 needed for its calculation. The present analysis aimed to compare the performance of
297 different algorithms based on the same set of features. Therefore, data-driven variable
298 selection to improve model performance was not performed. Finally, we limited our
299 analysis to in-hospital mortality to be consistent with current prediction models [2,3]
300 but we cannot exclude that ML algorithms can improve prediction of long-term
301 outcomes [29].

302 In conclusion, the present findings suggest that the application of ML algorithms alone,
303 is unlikely to determine a substantial gain in prediction of in-hospital mortality following
304 cardiac surgery if a small set of structured clinical data is available. A precise

305 estimation of individual risk is likely to be achieved only by the identification of new
306 powerful predictors that can explain more of the variance observed.

307 **Funding:** The present study was funded by the Bristol Biomedical Research Centre
308 (Bristol BRC).

309 **Conflict of interest:** none declared.

310 **FIGURE LEGENDS**

311 **Figure 1.** Receiver operating characteristic curve of EuroSCORE I & II, logistic
312 regression and machine learning classifiers: neural network, naïve Bayes and random
313 forest using EuroSCORE I features. The axes are true positive rate against 1 – false
314 positive rate. The area under the curve provides a measure of discrimination accuracy.
315 The dashed line represents no classification discrimination ability.

316 **Figure 2.** External probability calibration of logistic regression (Panel A), Neural
317 Network (Panel B), and Random Forest (Panel C) using the calibration belt method.
318 The method regresses true mortality on classifier probability of mortality (via logit
319 function) using polynomial logistic regression. All models showed significant
320 miscalibration ($P < 0.001$).

321 **Figure 3.** External probability calibration of EuroSCORE I (Panel A) and EuroSCORE
322 II (Panel B) using the calibration belt method. The method regresses true mortality on
323 classifier probability of mortality (via logit function) using polynomial logistic
324 regression. EuroSCORE I ($P < 0.001$) but not EuroSCORE II ($P = 0.64$) showed
325 significant model miscalibration.

326 **Supplementary Figure 1.** Flow of participants in the study

327 **Table 1.** Distribution of features included in the EuroSCORE stratified for in-hospital
 328 mortality in patients who underwent adult cardiac surgery from 1996 to 2017. (SD,
 329 standard deviation. LVEF, left-ventricle ejection fraction).

	Alive		Dead		P
	N= 27934		N=786		
Age (yrs mean, SD)	65.29	12.10	69.38	11.85	<0.001
Female	7149	25.59%	286	36.39%	<0.001
Serum creatinine \geq200 μmol/l	332	1.19%	56	7.12%	<0.001
Extracardiac arteriopathy	2346	8.40%	131	16.67%	<0.001
Pulmonary disease	3370	12.06%	146	18.58%	<0.001
Neurological dysfunction	593	2.12%	27	3.44%	0.018
Previous cardiac surgery	1734	6.21%	128	16.28%	<0.001
Recent myocardial infarct	6665	23.86%	226	28.75%	0.002
LVEF 30-50%	5539	19.83%	226	28.75%	<0.001
LVEF <30%	1391	4.98%	129	16.41%	<0.001
Systolic pulmonary pressure >60 mmHg	836	2.99%	28	3.56%	0.414
Active endocarditis	285	1.02%	23	2.93%	<0.001
Unstable angina	2554	9.14%	155	19.72%	<0.001
Emergency operation	884	3.16%	208	26.46%	<0.001
Critical preoperative state	417	1.49%	128	16.28%	<0.001
Ventricular septal rupture	53	0.19%	32	4.07%	<0.001

Other than isolated coronary surgery	10461	37.45%	464	59.03%	<0.001
Thoracic aortic surgery	1363	4.88%	148	18.83%	<0.001

330

331 **Table 2.** Evaluation of calibration drift. The test dataset was divided into two equal
 332 bins ordered by procedure date with approximately equal number of events (n=102 vs
 333 n=105). Goodness of fit chi-squared statistics were calculated for first (G1) and second
 334 (G2) group.

Model	χ^2 (G1)	χ^2 (G2)	Change
Logistic regression (retrained)	12.45	14.81	15.9%
naïve Bayes	1242.96	2126.79	41.6%
neural network	2.51	7.00	64.2%
random forest	15.53	19.70	21.2%
EuroSCORE I	15.94	26.93	40.8%

335 .

336 **REFERENCES**

- 337 1. Nashef SA, Roques F, Michel P, Gauducheau E, Lemeshow S, Salamon R.
338 European system for cardiac operative risk evaluation (EuroSCORE). *Eur J*
339 *Cardiothorac Surg.* 1999;16:9–13.
- 340 2. Nashef, S. A. M., Roques, F., Sharples, L. D., Nilsson, J., Smith, C., Goldstone, A.
341 R., et al. EuroSCORE II. *European Journal of Cardio-Thoracic Surgery.* 2012; 41:
342 734–745.
- 343 3. Edwards FH, Clark RE, Schwartz M. Coronary artery bypass grafting: the Society
344 of Thoracic Surgeons National Database experience. *Ann Thorac Surg.* 1994
345 ;57:12–9.
- 346 4. Provenchère S, Chevalier A, Ghodbane W, Bouleti C, Montravers P, Longrois D,
347 et al. Is the EuroSCORE II reliable to estimate operative mortality among
348 octogenarians? *PLoS One.* 2017;12:e0187056.
- 349 5. Guida P, Mastro F, Scрасcia G, Whitlock R, Paparella D. Performance of the
350 European System for Cardiac Operative Risk Evaluation II: a meta-analysis of 22
351 studies involving 145,592 cardiac surgery procedures. *J Thorac Cardiovasc Surg.*
352 2014;148:3049–57.e1
- 353 6. Christodoulou E, Ma J, Collins GS, Steyerberg EW, Verbakel JY, Van Calster B. A
354 systematic review shows no performance benefit of machine learning over logistic
355 regression for clinical prediction models. *J Clin Epidemiol.* 2019;110:12–22.
- 356 7. Drew PJ, Monson JRT. Artificial neural networks. *Surgery.* 2000;127:3–11.
- 357 8. Kingma, D.P., & Ba, J. (2014). Adam: A Method for Stochastic Optimization. *CoRR*,
358 abs/1412.6980.
- 359 9. Kvamme, H., Borgan, Ø., Scheel, I. Time-to-Event Prediction with Neural Networks
360 and Cox Regression. *ArXiv.* 2019:abs/1907.00825.

- 361 10. Kingsford, C., Salzberg, S. L. What are decision trees? *Nature biotechnology*.
362 2008; 26:1011–1013.
- 363 11. Sarica, A., Cerasa, A., Quattrone, A. Random Forest Algorithm for the
364 Classification of Neuroimaging Data in Alzheimer's Disease: A Systematic Review.
365 *Frontiers in aging neuroscience*. 2017; 9:329
- 366 12. Zhang Z. Naïve Bayes classification in R. *Annals of translational medicine*,
367 2016;4:241.
- 368 13. Nattino, G., Finazzi, S., & Bertolini, G. A new calibration test and a reappraisal of
369 the calibration belt for the assessment of prediction models based on dichotomous
370 outcomes. *Statistics in Medicine*. 2014; 33:2390–2407.
- 371 14. Sergeant P, Meuris B, Pettinari M. EuroSCORE II, illum qui est gravitates magni
372 observe*. *Eur J Cardio-Thoracic Surg*. 2012;41:729–31.
- 373 15. Gummert JF, Funkat A, Osswald B, Beckmann A, Schiller W, Krian A, et al.
374 EuroSCORE overestimates the risk of cardiac surgery: results from the national
375 registry of the German Society of Thoracic and Cardiovascular Surgery. *Clin Res*
376 *Cardiol*. 2009;98:363–9.
- 377 16. Ad N, Holmes SD, Patel J, Pritchard G, Shuman DJ, Halpin L. Comparison of
378 EuroSCORE II, Original EuroSCORE, and The Society of Thoracic Surgeons Risk
379 Score in Cardiac Surgery Patients. *Ann Thorac Surg*. 2016;102:573–9.
- 380 17. Shanmugam, G.; West, M.; Berg, G. Additive and logistic EuroSCORE
381 performance in high risk patients. *Interact. Cardiovasc. Thorac. Surg*. 2005, 4,
382 299–303.
- 383 18. Osnabrugge RL, Speir AM, Head SJ, et al. Performance of EuroSCORE II in a
384 large US database: implications for transcatheter aortic valve implantation. *Eur J*
385 *Cardiothorac Surg*. 2014;46(3):400-408.

- 386 19. Kirmani BH, Mazhar K, Fabri BM, Pullan DM. Comparison of the EuroSCORE II
387 and Society of Thoracic Surgeons 2008 risk tools. *Eur J Cardiothorac Surg.*
388 2013;44(6):999-1005.
- 389 20. Kilic A. Artificial Intelligence and Machine Learning in Cardiovascular Health Care.
390 *Ann Thorac Surg.* 2020;109:1323–9.
- 391 21. Mendes, R. G., de Souza, C. R., Machado, M. N., Correa, P. R., Di Thommazo-
392 Luporini, L., et al. Predicting reintubation, prolonged mechanical ventilation and
393 death in post-coronary artery bypass graft surgery: a comparison between artificial
394 neural networks and logistic regression models. *Archives of medical science AMS.*
395 2015; 11:756–763.
- 396 22. Allyn J, Allou N, Augustin P, Philip I, Martinet O, Belghiti M, et al. A Comparison of
397 a Machine Learning Model with EuroSCORE II in Predicting Mortality after Elective
398 Cardiac Surgery: A Decision Curve Analysis. *PLoS One.* 2017;12:e0169772.
- 399 23. Mejia OAV, Antunes MJ, Goncharov M, Dallon LRP, Veronese E, Lapenna GA, et
400 al. Predictive performance of six mortality risk scores and the development of a
401 novel model in a prospective cohort of patients undergoing valve surgery
402 secondary to rheumatic fever. *PLoS One.* 2018;13:e0199277.
- 403 24. Ghavidel AA, Javadikasgari H, Maleki M, Karbassi A, Omrani G, Noohi F. Two new
404 mathematical models for prediction of early mortality risk in coronary artery bypass
405 graft surgery. *J Thorac Cardiovasc Surg.* 2014;148:1291–1298.e1.
- 406 25. Nilsson J, Ohlsson M, Thulin L, Höglund P, Nashef SA, Brandt J. Risk factor
407 identification and mortality prediction in cardiac surgery using artificial neural
408 networks. *J Thorac Cardiovasc Surg.* 2006;132:12–19.

- 409 26. Kilic, A., Goyal, A., Miller, J. K., Gjekmarkaj, E., Tam, W. L., Gleason, T. G., et al.
410 Predictive Utility of a Machine Learning Algorithm in Estimating Mortality Risk in
411 Cardiac Surgery. *The Annals of thoracic surgery*. 2019;S0003-4975(19)31620-0.
- 412 27. van der Ploeg T, Austin PC, Steyerberg EW. Modern modelling techniques are
413 data hungry: a simulation study for predicting dichotomous endpoints. *BMC Med*
414 *Res Methodol*. 2014;14:137.
- 415 28. Omar RZ, Ambler G, Royston P, Eliahoo J, Taylor KM. Cardiac surgery risk
416 modeling for mortality: a review of current practice and suggestions for
417 improvement. *Ann Thorac Surg*. 2004;77:2232–2237.
- 418 29. Linden A, Yarnold PR. Modeling time-to-event (survival) data using classification
419 tree analysis. *J Eval Clin Pract*. 2017;23:1299-1308.