

1 **A comparative study of machine learning algorithms in predicting severe complications**  
2 **after bariatric surgery**

3 Yang Cao<sup>1,2\*</sup>, Xin Fang<sup>2\*</sup>, Johan Ottosson<sup>3</sup>, Erik Näslund<sup>4</sup>, Erik Stenberg<sup>3</sup>

4

5 <sup>1</sup>Clinical Epidemiology and Biostatistics, School of Medical Sciences, Örebro University,  
6 Örebro, Sweden;

7 <sup>2</sup>Unit of Biostatistics, Institute of Environmental Medicine, Karolinska Institutet, Stockholm,  
8 Sweden;

9 <sup>3</sup>Department of Surgery, Faculty of Medicine and Health, Örebro University, Örebro,  
10 Sweden;

11 <sup>4</sup>Division of Surgery, Department of Clinical Sciences, Danderyd Hospital, Karolinska  
12 Institutet, Stockholm, Sweden

13

14 \***For correspondence:** [xin.fang@ki.se](mailto:xin.fang@ki.se); [yang.cao@ki.se](mailto:yang.cao@ki.se)

15 **Competing interest:** The authors declare that no competing interests exist.

16 **Funding:** This work was supported by grants from the Örebro Region County Council (E.S.),  
17 Örebro University (E.S.), Stockholm County Council (E.N.), SRP Diabetes (E.N.) and the  
18 NovoNordisk foundation (E.N.).

19

20 **Abstract**

21 Accurate models to predict severe postoperative complications could be of value in the  
22 preoperative assessment of potential candidates for bariatric surgery. Traditional statistical  
23 methods have so far failed to produce high accuracy. To find a useful algorithm to predict the  
24 risk for severe complication after bariatric surgery, we trained and compared 29 supervised  
25 machine learning (ML) algorithms using information from 37,811 patients operated with a  
26 bariatric surgical procedure between 2010 and 2014 in Sweden. The algorithms were then  
27 tested on 6,250 patients operated in 2015. Most ML algorithms showed high accuracy (>90%)  
28 and specificity (>0.9) in both the training and test data. However, none achieved an acceptable  
29 sensitivity in the test data. ML methods may improve accuracy of prediction but we did not  
30 yet identify one with a high enough sensitivity that can be used in clinical praxis in bariatric  
31 surgery. Further investigation on deeper neural network algorithms is needed.

32

## 33 **Introduction**

34 Morbid obesity is a global public health threat of growing proportions(Ng et al., 2014).

35 Bariatric surgery offers the best chance for long-term weight-loss and resolution of  
36 comorbidities(Sjostrom et al., 2004). Although modern bariatric surgery is considered to be  
37 safe, severe postoperative complications still occur(Finks et al., 2011; Stenberg et al., 2014).  
38 Accurate prediction models for severe postoperative complications could aid preoperative  
39 decision making for surgeons, anesthesiologists and patients. These models could also serve  
40 as basis for case-mix comparisons between different centers. Some prediction models based  
41 on linear regression of patient-specific data allow for relatively simple and interpretable  
42 inference; however, they have so far been proven inaccurate and can thus not be used in  
43 clinical practice(Geubbels et al., 2015; Stenberg et al., 2018).

44 In contrast, some machine learning (ML) methods have been shown to provide quite accurate  
45 predictions, and have increasingly been used in diagnosis and prognosis of different diseases  
46 and health conditions(Anderin et al., 2015; Kourou et al., 2015; Pan et al., 2017). ML  
47 methods are data-driven analytic approaches that specialize in the integration of multiple risk  
48 factors into a predictive algorithm(Passos et al., 2016). Over the past several decades, ML  
49 tools have become more and more popular for medical researchers. A variety of ML  
50 algorithms, including artificial neural networks, decision trees, Bayesian networks, and  
51 support vector machines (SVMs) have been widely applied with the aims to detect key  
52 features of the patient conditions and to model the disease progression after treatment from  
53 complex health information and medical datasets. The application of different ML methods  
54 for feature selection and classification in multidimensional heterogeneous data can provide  
55 promising tools for inference in medical practices. These highly nonlinear approaches have  
56 been utilized in medical research for the development of predictive models, resulting in  
57 effective and accurate decision making(Ali, 2017; Jiang et al., 2017).

58 Although new and improved software packages have significantly eased the implementation  
59 burden for many ML methods in recent years, few studies have used ML methods to examine  
60 the risk factors or predict the prognosis after bariatric surgery, including diabetes  
61 remission(Hayes et al., 2011; Pedersen et al., 2016), complication(Razzaghi et al., 2017),  
62 weight status(Piaggi et al., 2010; Thomas et al., 2017), and adverse events and death(Ehlers et  
63 al., 2017). Even though there is evidence that the use of ML methods can improve our  
64 understanding of postoperative progression of bariatric surgery, an appropriate level of  
65 validation is needed in order for these methods to be considered in the clinical practice.

66 In this study, we compared different conventional supervised ML algorithms in the modeling  
67 of severe postoperative complication after bariatric surgery. The study was based on the data  
68 from the Scandinavian Obesity Surgery Registry (SOReg). The SOReg is a national quality  
69 and research register, covering virtually all bariatric surgical procedures performed in Sweden  
70 since 2010. The register has been described in detail elsewhere(Hedenbro et al., 2015;  
71 Stenberg et al., 2014), and a prediction model based on logistic regression for the same group  
72 of patients has been described previously(Stenberg et al., 2018). The aim of the current study  
73 was to find an algorithm or algorithms that perform well not only on the training data but also  
74 on the test data that were not used to train the algorithms.

75

## 76 **Results**

77 Baseline characteristics of the patients in the training data and the test data are presented in  
78 Tables 1 and 2. The percentages of severe complication in the two data sets are 3.2% and  
79 3.0%, respectively. No statistically significant difference was found for percentages of severe  
80 complication between the two data sets (Pearson chi-square = 0.8283,  $p = 0.363$ ).

81 Univariable analyses indicate that differences of mean age, mean body mass index (BMI),  
82 median HbA1c, percentages of comorbidities for hypertension, diabetes, dyslipidaemia, and  
83 previous venous thromboembolism, and percentage of revisional surgery between the patients  
84 presenting and without severe complication are statistically significant in the training data  
85 (Table 1). In the test data, the statistically significant differences were found for age, waist  
86 circumference (WC), HbA1c, dyslipidaemia, and revisional surgery (Table 2).

87 Multivariable logistic regression analysis for the same data was published elsewhere (Stenberg  
88 et al., 2018). In brief, revisional surgery, age, low BMI, operation year, WC, and dyspepsia  
89 were associated with the an increased risk for severe postoperative complication, however, the  
90 performance of the multivariable logistic regression model for predicting the risk in individual  
91 patient case was poor. Validation of the model tested on patients operated in 2015 resulted in  
92 an area under the receiver operating characteristic (ROC) curve of only 0.53, a Hosmer-  
93 Lemshow goodness of fit 17.91 ( $p=0.056$ ) and Nagelkerke  $R^2$  0.013 (Stenberg et al., 2018).

94 In current study, 19 supervised machine learning algorithms were compared and ten of them  
95 were also trained using the synthetic minority oversampling technique (SMOTE), resulting in  
96 29 ML algorithms. Most of the machine learning algorithms shown high accuracy ( $>90\%$ ) and  
97 specificity ( $>0.9$ ) for both training data and test data (Table 3), except that bagging linear  
98 discriminant analysis (LDA), bagging quadratic discriminant analysis (QDA), adaptive  
99 boosting (AdaBoost) support vector machine (SVM), and multilayer perceptron (MLP) shown

100 low accuracy (<60%) for SMOTE training data, and oversampling-based bagging QDA  
101 shown low accuracy for test data (accuracy = 56.1%) (Table 3).

102 Although most of the algorithms shown low sensitivity for both the training data and the test  
103 data, some of them exhibited promising prediction ability in the training data. Sensitivities of  
104 oversampling-based bagging QDA, random forest, AdaBoost extremely randomized  
105 (AdaExtra) trees, AdaBoost gradient regression (AdaGradient) trees, bagging k-nearest  
106 neighbor (KNN), and deep learning neural network (NN) are 0.707, 0.965, 0.980, 0.968,  
107 0.996, and 0.757 for SMOTE training data, respectively (Table 3). Even for test data,  
108 oversampling-based bagging QDA and AdaBoost SVM show significant higher prediction  
109 ability than other algorithms. The sensitivities of the two algorithms are 0.417 and 0.364,  
110 respectively (Table 3). However, they still do not achieve an acceptable level for practical  
111 application.

112 When considering sensitivity and specificity together, most of the algorithms did not show  
113 better prediction ability than a random predictor, i.e. an area under ROC curve of 0.5. The  
114 areas under the ROC curves for all the algorithms, except for oversampling-based random  
115 forest, AdaExtra trees, and adaGradient trees, and KNN, are around 0.5 (Figures 1 - 4).

116 Although oversampling-based random forest, AdaExtra trees, AdaGradient trees, and KNN  
117 show outstanding prediction ability on the SMOTE training data (areas under ROC curves are  
118 above 0.9), their performance on the test data are not optimistic (Figures 2 and 3).

119 The performance of the three regression-based algorithms (logistic regression, LDA, QDA),  
120 SVM, and the two neural network-based algorithms (MLP and deep learning NN) was poor in  
121 any situation. However the bagging MLP and deep learning NN outperforms the tree-based  
122 algorithms (Figures 2 and 4) for test data, their areas under ROC curves for the test data are  
123 0.58 and 0.56, respectively (Figure 4) that are greatest among all the algorithms.

## 124 **Discussion**

125 Historically, laparoscopic gastric bypass has for a long time been the most common bariatric  
126 procedure in Sweden, although laparoscopic sleeve gastrectomy has increased in popularity  
127 over more recent years(Stenberg et al., 2014; The international federation for the surgery of  
128 obesity and metabolic disorders, 2017). The surgical technique is highly standardized with  
129 more than 99% of all gastric bypass procedure being the antecolic, antegastric, laparoscopic  
130 gastric bypass (so called Lönnroth technique)(Olbers et al., 2003). Virtually all patients  
131 receive pharmacologic prophylaxis for deep venous thrombosis and intraoperative antibiotic  
132 prophylaxis(Hedenbro et al., 2015; Stenberg et al., 2014). Patients who have bariatric surgery  
133 are exposed to the risk of having postoperative complications, which may increase the  
134 complexity of managing safety and healthcare costs.

135 Previous studies on postoperative complications of bariatric surgery have mainly used scoring  
136 for identifying patients who are more likely to have complications after surgery. However,  
137 these methods are not sensitive enough for clinical application(Geubbels et al., 2015;  
138 Stenberg et al., 2018). The potential of ML tools as clinical decision support in identifying  
139 risk factors and predicting health outcomes is therefore worth investigation on complications  
140 associated with bariatric surgery. To our knowledge, there is only one study that compared the  
141 performance of different ML algorithms in predicting the postoperative complications in  
142 imbalanced bariatric surgery data set(Razzaghi et al., 2017). Although the study indicates that  
143 the combination of a suitable feature selection method with ensemble ML algorithm equipped  
144 with SMOTE can achieve higher performance in predictive models for bariatric surgery risks,  
145 the ML algorithms were not validated using external test data. After all, for prediction  
146 purpose, we are not very interested in whether or not an algorithm accurately predicts severe  
147 complication for patients used to train the algorithm, since we already know which of those

148 patients have severe complications, but are interested in whether the algorithms may  
149 accurately predict the future patients based on their clinical measurements.

150 Our study compared in total 29 ML algorithms using real world data. Although the  
151 sensitivities of the algorithms were generally low, the study indicates that some ML  
152 algorithms were able to achieve higher accuracy than tradition logistic regression  
153 models(Geubbels et al., 2015; Stenberg et al., 2018). Four of 29 algorithms were able to  
154 achieve high sensitivity ( $>0.95$ ) and two achieved moderate sensitivity ( $>0.70$ ) in the training  
155 data, including three tree-based algorithms, bagging KNN, bagging QDA, and deep learning  
156 NN. We should notice that all the high or moderate sensitivities were obtained from SMOTE  
157 training data and/or using ensemble algorithms. Our findings support the previous study that  
158 ensemble ML algorithms equipped with SMOTE can achieve higher performance metrics for  
159 imbalanced data(Razzaghi et al., 2017).

160 Despite showing promising capability of prediction in training data, none of the 29 ML  
161 algorithms satisfactorily predicted severe postoperative complication after bariatric surgery in  
162 the test data. Why did the algorithms do a poor job of predicting the patients who had severe  
163 complication in test data? One potential explanation for this may be related to the limited  
164 number of severe postoperative complication in the current dataset, which cannot reveal the  
165 underlying relationship between risk factors and adverse health outcomes. Although there are  
166 several known risk factors, each of them only imposes a small increase in the risk for  
167 postoperative complication(Finks et al., 2011; Longitudinal Assessment of Bariatric Surgery  
168 et al., 2009; Maciejewski et al., 2012; Stenberg et al., 2014). Another likely explanation may  
169 be that preoperatively known variables are insufficient to predict postoperative complications.

170 In previous studies, the highest accuracy for prediction of postoperative complication has  
171 been models including operation data, mainly intraoperative complication and conversion to  
172 open surgery(Geubbels et al., 2015; Stenberg et al., 2014). Although including intraoperative



173 adverse events and conversion to open surgery may improve the accuracy of prediction  
174 models, such models would not be useful in the preoperative assessment for patients or for  
175 case mix comparisons. Furthermore, because the algorithms try to minimize the total error  
176 rate out of all classes, irrespective of which class the errors come from, they are not  
177 appropriate for imbalanced data such as what we used in our study(Maalouf et al., 2018).

178 Compared with traditional generalized linear predictive models, non-linear ML algorithms are  
179 more flexible and may achieve higher accuracy but at the expense of less interpretability.

180 Although there are interpretable models such as regression, Naïve Bayes, decision tree and  
181 random forests, several models are not designed to be interpretable(James et al., 2013). The  
182 aim of the methods is to extract information from the trained model to justify their prediction  
183 outcome, without knowing how the model works in details. The trade-off between prediction  
184 accuracy and model interpretability is always an issue when we have to consider in building a  
185 ML algorithm. A common quote on model interpretability is that with an increase in model  
186 complexity, model interpretability goes down at least as fast. Fully nonlinear methods such as  
187 bagging, boosting and support vector machines with nonlinear kernels are highly flexible  
188 approaches that are harder to interpret. Deep learning algorithms are notorious for their un-  
189 interpretability due to the sheer number of parameters and the complex approach to extracting  
190 and combining features. Feature importance is a basic (and often free) approach to  
191 interpreting the model. Although some nonlinear algorithms such as tree-based algorithms  
192 (e.g. random forest) may allow to obtain information on the feature importance, we cannot  
193 obtain such information from many ML algorithms.

194 Therefore, recent attempts have been made to improve interpretability for the black-box  
195 algorithms even such as deep learning. Local interpretable model-agnostic explanations  
196 (LIME) is one of them to make these complex models at least partly understandable. LIME is  
197 a more general framework that aims to make the predictions of ‘any’ ML model more

198 interpretable. In order to remain model-independent, LIME works by modifying the input to  
199 the model locally(Mishra et al., 2017; Ribeiro et al., 2016). So instead of trying to understand  
200 the entire model at the same time, a specific input instance is modified and the impact on the  
201 predictions are monitored.

202 Regarding specific algorithm, though their motivations differ, the logistic regression and LDA  
203 or QDA methods are closely connected, therefore we were not surprised that LDA or QDA  
204 did not show significant improvement in prediction than logistic regression(Stenberg et al.,  
205 2018). KNN takes a complete different approach from classification which is completely non-  
206 parametric(James et al., 2013). Therefore, we can expect it to outperform parametric models  
207 such as logistic and LDA. However, KNN cannot tell us which predictor are of importance.  
208 QDA serves as a compromise between the non-parametric KNN and the LDA and logistic  
209 regression. Though not as flexible as KNN, QDA can perform better in the limited training  
210 data situation. MLP is a class of feedforward artificial neural network, which consists of at  
211 least three layers of nodes. Its multiple layers and non-linear activation can distinguish data  
212 that is not linearly separable. Deep learning NNs are high-level NNs including convolutional  
213 NN and recurrent NN et al. In our study, the deep learning NN with five hidden layers  
214 outperforms the conventional MLP with two hidden layers, especially on SMOTE training  
215 data (areas under ROC curves are 0.67 vs. 0.37), which deserves further investigation in the  
216 future.

217 Our study demonstrates that ensemble learning may improve predictions by combining  
218 several base algorithms. However, usually there are several ensemble methods available, such  
219 as bagging, boosting, and stacking(Zhou, 2012). A number of studies have shown that, when  
220 decomposing a classifier's error into bias and variance terms, AdaBoost is more effective at  
221 reducing bias, bagging is more effective at reducing variance, and stacking may improve  
222 predictions in general(Kotsiantis et al., 2007). There is no golden rule on which method works

223 best. The choice of specific ensemble methods is case by case and depends enormously on the  
224 data.

225 There are some limitations in our study. First, the study was limited to data registered within  
226 the SOReg. Cardiovascular and pulmonary comorbidities other than sleep apnea are not  
227 mandatory variables within the registry and could thus not be included in the model. Although  
228 these comorbidities are known risk factors for postoperative complications(Finks et al., 2011;  
229 Gupta et al., 2011; Maciejewski et al., 2012), they are not highly prevalent in European  
230 studies(Geubbels et al., 2015). Second, although we compared 29 ML algorithms investigated  
231 in our study, they are convenient and feasible methods for general medical researchers.  
232 Because of computational complexity and less interpretability, many complicated and  
233 advanced ML algorithms were not yet investigated in our study. However, our study at least  
234 points out a promising way for future investigations, i.e. deep learning NN equipped with  
235 SMOTE. Last but not the least, the exhaustive grid search was used in our hyperparameter  
236 optimization, which is extremely resource consuming and not optimal for complex ML  
237 algorithms, therefore other advanced methods such as gradient-based or evolutionary  
238 optimization would be considered in the future.

## 239 **Conclusion**

240 ML algorithms have the potential to improve the accuracy in predicting the severe  
241 postoperative complication among the 44,061 Swedish bariatric surgery patients during 2010  
242 - 2015. Because the imbalance nature of the data where the number of the interested outcome  
243 is relative small, oversampling technique needs to be adopted to balance the two outcomes  
244 (presenting or without severe complication). Ensemble algorithms outperform base  
245 algorithms. In general, deep learning NN results in better predictions for unseen patients.

246

## 247 **Materials and Methods**

### 248 **Patients and features**

249 Patients registered in the SOReg between 2010 and 2015 were included in the present study.  
250 All patients who underwent a bariatric procedure between 2010 and 2014 were used as  
251 training data in the ML. Data from patients who underwent a bariatric surgical procedure in  
252 2015 were used as test data to validate the algorithm's performance in predicting sever  
253 postoperative complication within 30 days after surgery. In total, 37,811 and 6,250 bariatric  
254 patients from SOReg were included in the training data and test data, respectively. In total 16  
255 features were included in ML, including five continuous features (age, HbA1c, body mass  
256 index [BMI], waist circumference [WC]), and operation year) and 11 binary features (sleep  
257 apnoea, hypertension, diabetes, dyslipidaemia, dyspepsia, depression, musculoskeletal pain,  
258 previous venous thromboembolism, revisional surgery, and severe postoperative  
259 complication). The last binary feature, i.e. severe postoperative complication, was used as  
260 output variable for the supervised ML classifiers. All the continuous features were  
261 standardized to have mean 0 and standard deviation 1 before they enter the classifier. HbA1c  
262 was log transformed before standardization because of its asymmetric distribution.

### 263 **Descriptive and inferential statistical methods**

264 Demographic and baseline characteristics of the patients were presented using descriptive  
265 statistical methods. Continuous variables were portrayed as mean and standard deviation  
266 (SD), or median and interquartile range where suitable, while categorical variables were  
267 outlined as counts and percentages. The difference between the patient presenting and without  
268 severe postoperative complication was tested using the Student's t-test or the Mann-Whitney  
269 U test for normally or asymmetrically distributed continuous variables, respectively; and  $\chi^2$   
270 test was used for binary variables.

## 271 **ML algorithms**

272 In current study, eight base ML algorithms, i.e. logistic regression, linear discriminant  
273 analysis (LDA), quadratic discriminant analysis (QDA), decision tree, k-nearest neighbor  
274 (KNN), support vector machine (SVM), multilayer perceptron (MLP) and deep learning  
275 neural network (NN), and 11 ensemble algorithms, i.e. adaptive boosting (AdaBoost) logistic  
276 regression, bagging LDA, bagging QDA, random forest, extremely randomized (Extra) trees,  
277 AdaBoost Extra trees, gradient regression tree, AdaBoost Gradient trees, bagging KNN,  
278 AdaBoost SVM, and bagging MLP, were implemented.

## 279 **Ensemble learning**

280 In order to improve generalizability and robustness over a single ML algorithm, we also used  
281 ensemble methods to combine multiple base or ensemble algorithms. Five ensemble methods  
282 were applied in our study:

- 283 • AdaBoost for logistic regression, Extra trees, gradient regression trees, and  
284 SVM(Schapire, 2003);
- 285 • bagging for LDA, QDA, KNN, and MLP(Kotsiantis et al., 2007);
- 286 • random forests for decision tree(Liaw & Wiener, 2002);
- 287 • Extra trees for decision tree(Geurts et al., 2006);
- 288 • gradient boosted regression trees for decision tree(Friedman, 2002).

## 289 **Initialization and optimization of hyperparameters**

290 ML algorithms involve a number of hyperparameters that have to be fixed before running the  
291 algorithms. In contrast to the parameters that are learned by training, hyperparameters  
292 determine the structure of a ML algorithm and how the algorithm is trained. The initial values  
293 of the hyperparameters for each ML algorithm used in our study are the default values

294 specified in the employed software packages based on recommendations or experience(Probst  
295 et al., 2018). In KNN algorithm, ten nearest neighbors were used. In MLP algorithm, two  
296 hidden layer were used with five and two neurons, respectively. In deep learning NN  
297 algorithm, the sequential linear stack of layers was used, with five hidden layers (three dense  
298 layers and two dropout regularization layers). For the detailed hyperparameterization of the  
299 algorithms, please refer the scikit-learn user manual at [http://scikit-](http://scikit-learn.org/stable/supervised_learning.html)  
300 [learn.org/stable/supervised\\_learning.html](http://scikit-learn.org/stable/supervised_learning.html)(Pedregosa et al., 2011) and the Keras  
301 Documentation at <https://keras.io/>.

302 The hyperparameter optimization is defined as a tuple of hyperparameters that yields an  
303 optimal algorithm which minimizes a predefined loss function (i.e. cross entropy loss function  
304 in our study, see Annex 1) on a held-out validation set of the training data. The most wildly  
305 used however exhaustive grid search was used to perform hyperparameter optimization in our  
306 study, which specified subset of the hyperparameter space of a ML algorithm and was  
307 evaluated by cross-validation using the training data(Bergstra & Bengio, 2012).

### 308 **Cross validation**

309 For training data,  $k$ -fold ( $k = 5$  in our analyses) cross-validated predictions were used as  
310 predicted values. This approach involves randomly dividing the training data into  $k$  groups, or  
311 folds, of approximately equal size. Then an algorithm is trained on the  $k-1$  folds and the rest  
312 one fold is retained as the validation fold for testing the algorithm. The process is repeated  
313 until the algorithm is validated on all the  $k$  folds. For each patient in the training data, the  
314 predicted value that he/she obtained is the prediction when he/she was in the validation fold.  
315 Therefore, only cross-validation strategies that assign all patients to a validation fold exactly  
316 once can be used for the cross-validated prediction(James et al., 2013).

### 317 **SMOTE**

318 The bariatric surgery data is extreme imbalanced, i.e. only 1,408 of 44,061 (3.2%) patients  
319 experienced severe postoperative complication after bariatric surgery. The imbalance often  
320 results in serious bias in the performance metrics(Batista et al., 2004). Therefore, we  
321 performed synthetic minority oversampling technique (SMOTE) to tackle the  
322 imbalance(Chawla et al., 2002). SMOTE generates a synthetic instance by interpolating  $m$   
323 instances (for a given integer value  $m$ ) of the minority class that lies close enough to each  
324 other to achieve the desired ratio between the majority and minority classes. In our study, a  
325 1:1 ratio between the patients presenting severe postoperative complication and without  
326 severe postoperative complication was achieved in the training data, i.e. SMOTE training  
327 data. The aforementioned nine of the 11 ensemble ML algorithms and the deep learning NN  
328 were also implemented for the SMOTE training data.

### 329 **Performance metrics**

330 The performance of the in total 29 ML algorithms were evaluated using accuracy, sensitivity,  
331 specificity, and area under the receiver operating characteristic (ROC) curve. ROC curve  
332 shows the trade-off that the algorithms set the different threshold values for the posterior  
333 probability for the prediction.

334 Terminology and derivations of accuracy, sensitivity, specificity, and area under the ROC  
335 curve are given in Annex 1.

### 336 **Software and hardware**

337 The descriptive and inferential statistical analyses were performed using Stata 15.1  
338 (StataCorp, College Station). ML algorithms were achieved using packages scikit-learn 0.19.1  
339 (scikit-learn, <http://scikit-learn.org/>)(Pedregosa et al., 2011) and Keras 2.1.6 (Keras,  
340 <https://keras.io/>) in Python 3.6 (Python Software Foundation, <https://www.python.org/>).

341 All the computation was conducted in a computer with 64-bit Windows 7 Enterprise operation  
342 system (Service Pack 1), Intel ® Core™ i5-4210U CPU @ 2.40 GHz, and 16.0 GB installed  
343 random access memory.

344

#### 345 **Author contributions**

346 Yang Cao, Xin Fang, Conceptualization, Data curation, Software, Formal analysis, Writing-  
347 original draft; Yang Cao, Supervision, Investigation, Visualization, Methodology; Erik  
348 Stenberg Data curation, Conceptualization, Validation, Investigation, Writing-original draft;  
349 Johan Ottosson, Erik Näslund, Investigation, Writing-original draft.

350

#### 351 **Author ORCIDs**

352 Yang Cao, <https://orcid.org/0000-0002-3552-9153>

353 Xin Fang, <https://orcid.org/0000-0002-6846-7147>

354 Johan Ottosson, <https://orcid.org/0000-0002-9243-2390>

355 Erik Näslund, <https://orcid.org/0000-0002-0166-6344>

356 Erik Stenberg, <https://orcid.org/0000-0001-9189-0093>

357

#### 358 **Ethics**

359 The study was approved by the Regional Ethics Committee in Stockholm and was conducted  
360 in accordance with the ethical standards of the Helsinki Declaration (6th revision).

361



362

Table 1. Base line characteristics of the training patients

	All N=37,811	No serious complication N=36,591 (96.8%)	Having serious complication N=1,220 (3.2%)	p-value
Age in years, mean $\pm$ SD	41.2 $\pm$ 11.2	41.1 $\pm$ 11.2	42.9 $\pm$ 10.7	<0.001*
Sex, n (%)				
Female	28,682 (75.9%)	27,766 (75.9%)	916 (75.1%)	0.521 <sup>†</sup>
Male	9,129 (24.1%)	8,825 (24.1%)	304 (24.9%)	
BMI in kg/m <sup>2</sup> , mean $\pm$ SD	42.12 $\pm$ 5.66	42.13 $\pm$ 5.66	41.79 $\pm$ 5.58	0.0355*
WC in cm, mean $\pm$ SD	126.0 $\pm$ 14.0	126.0 $\pm$ 14.0	126.2 $\pm$ 13.8	0.6018*
HbA1c, median (P25, P75)	38 (35, 42)	38 (38, 32)	38 (35, 43)	0.0090 <sup>‡</sup>
Comorbidity, n (%)				
Sleep apnoea	3,792 (10.0%)	3,656 (10.0%)	136 (11.2%)	0.186 <sup>†</sup>
Hypertension	9,760 (25.8%)	9,404 (25.7%)	356 (29.2%)	0.006 <sup>†</sup>
Diabetes	5,407 (14.3%)	5,204 (14.2%)	203 (16.6%)	0.018 <sup>†</sup>
Dyslipidaemia	3,802 (10.1%)	3,667 (10.0%)	135(11.1%)	0.233 <sup>†</sup>
Dyspepsia	3,970 (10.5%)	3,803 (10.4%)	167 (13.7%)	<0.001 <sup>†</sup>
Depression	5,609 (14.8%)	5,409 (14.8%)	200 (16.4%)	0.119 <sup>†</sup>
Musculoskeletal pain	4,905 (13.0%)	4,754 (13.0%)	151 (12.4%)	0.529 <sup>†</sup>
Previous venous thromboembolism	918 (2.4%)	875 (2.39%)	43 (3.52%)	0.011 <sup>†</sup>
Revisional surgery	1,367 (3.6%)	1,261 (3.5%)	106 (8.7%)	<0.001 <sup>†</sup>

363 SD: standard deviation; BMI, body mass index; WC, waist circumference; P25, the 25th percentile;

364 P75, the 75th percentile.

365 \*t-test was used; <sup>†</sup> $\chi^2$  test was used; <sup>‡</sup>Mann-Whitney U test was used.

366

Table 2. Base line characteristics of the test patients

	All N=6,250	No serious complication N=6,062 (97.0%)	Having serious complication N=188 (3.0%)	p-value
Age in years, mean $\pm$ SD	41.2 $\pm$ 11.5	41.2 $\pm$ 11.5	42.9 $\pm$ 11.8	0.0423*
Sex, n (%)				
Female	4,832 (77.3%)	4,682 (77.2%)	150 (79.8%)	0.411 <sup>†</sup>
Male	1,418 (22.7%)	1,380 (22.8%)	38 (20.2%)	
BMI in kg/m <sup>2</sup> , mean $\pm$ SD	41.22 $\pm$ 5.87	41.20 $\pm$ 5.89	41.95 $\pm$ 5.40	0.0848*
WC in cm, mean $\pm$ SD	123.3 $\pm$ 14.1	123.2 $\pm$ 14.0	126.2 $\pm$ 14.7	0.0086*
HbA1c, median (P25, P75)	37 (34, 41)	37 (34, 41)	38 (35, 44)	0.0017 <sup>‡</sup>
Comorbidity, n (%)				
Sleep apnoea	622 (10.0%)	607 (10.0%)	15 (8.0%)	0.359 <sup>†</sup>
Hypertension	1,563 (25.0%)	1,506 (24.8%)	57 (30.3%)	0.088 <sup>†</sup>
Diabetes	761 (12.2%)	734 (12.1%)	27 (14.4%)	0.352 <sup>†</sup>
Dyslipidaemia	518 (8.3%)	493 (8.13%)	25 (13.3%)	0.011 <sup>†</sup>
Dyspepsia	645 (10.3%)	620 (10.2%)	25 (13.3%)	0.173 <sup>†</sup>
Depression	1,096 (17.5%)	1,053 (17.4%)	43 (22.9%)	0.051 <sup>†</sup>
Musculoskeletal pain	1,315 (21.0%)	1,268 (20.9%)	47 (25.0%)	0.176 <sup>†</sup>
Previous venous thromboembolism	182 (2.9%)	177 (2.99%)	5 (2.7%)	0.834 <sup>†</sup>
Revisional surgery	61 (1.0%)	54 (0.9%)	7 (3.7%)	<0.001 <sup>†</sup>

367 SD: standard deviation; BMI, body mass index; WC, waist circumference; P25, the 25th percentile;

368 P75, the 75th percentile.

369 \*t-test was used; <sup>†</sup> $\chi^2$  test was used; <sup>‡</sup>Mann-Whitney U test was used.

370

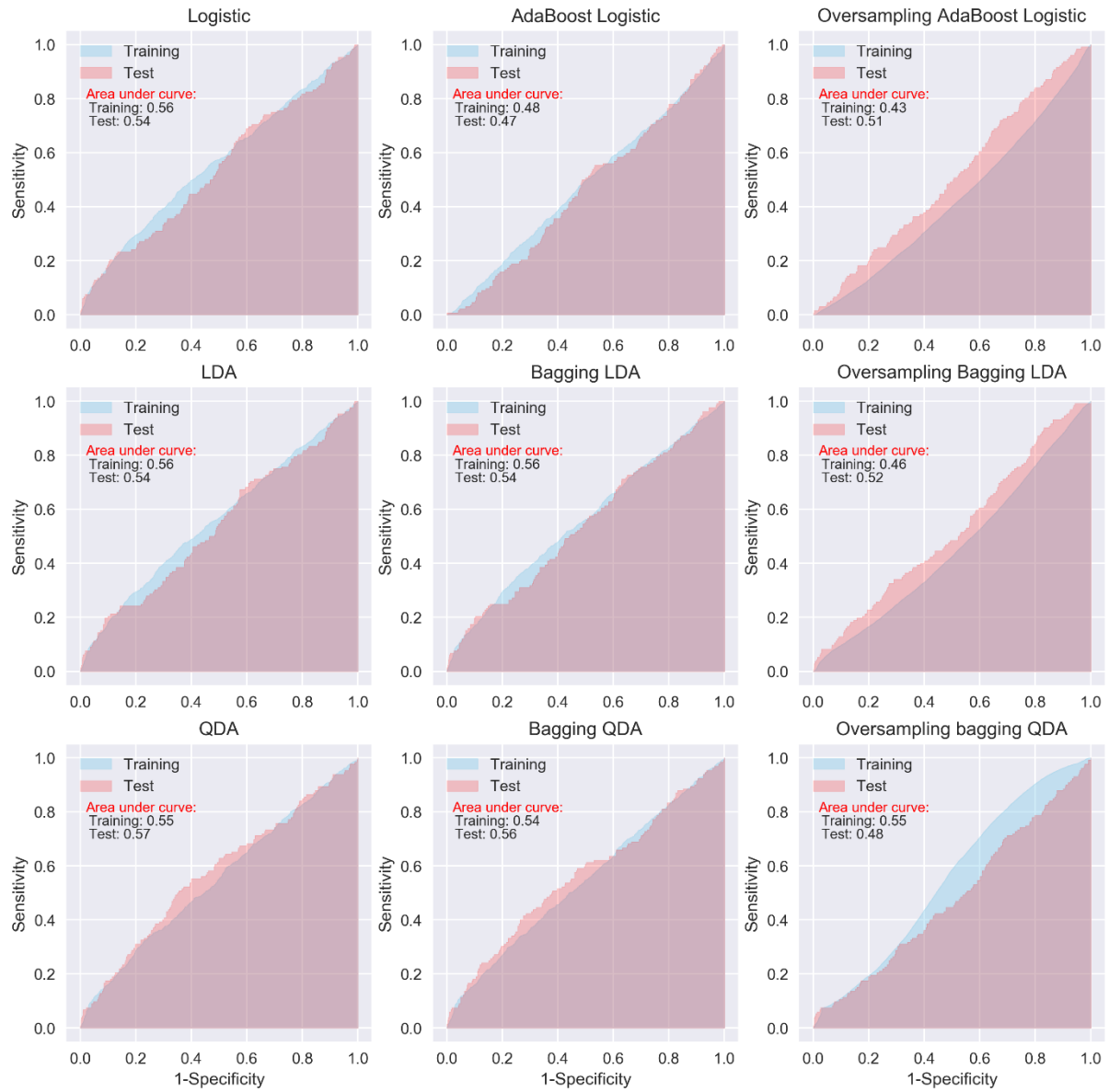
371

Table 3. Performance of the algorithms

Algorithm	Training data			Test data		
	Accuracy (%)	Specificity	Sensitivity	Accuracy (%)	Specificity <sub>y</sub>	Sensitivity <sub>y</sub>
Logistic	96.9	1.000	0.000	97.1	1.000	0.000
AdaBoost Logistic	96.9	1.000	0.000	97.1	1.000	0.000
Oversampling AdaBoost Logistic	46.5	0.382	0.547	76.9	0.786	0.227
LDA	96.9	1.000	0.000	97.1	1.000	0.000
Bagging LDA	96.9	1.000	0.000	97.1	1.000	0.000
Oversampling Bagging LDA	46.3	0.370	0.556	79.0	0.807	0.212
QDA	92.8	0.954	0.107	94.7	0.973	0.076
Bagging QDA	93.2	0.958	0.103	95.5	0.982	0.068
Oversampling bagging QDA	55.4	0.401	0.707	56.1	0.566	0.417
Decision tree	93.5	0.963	0.038	93.1	0.958	0.045
Random Forest	96.9	1.000	0.000	97.0	1.000	0.000
Oversampling Random Forest	94.5	0.925	0.965	96.6	0.995	0.008
ExtRa Trees	96.6	0.997	0.006	96.7	0.996	0.015
AdaBoost ExtRa Trees	96.6	0.996	0.004	96.6	0.995	0.008
Oversampling AdaExtra Trees	93.0	0.881	0.980	95.3	0.982	0.015
Gradient regression trees	96.9	1.000	0.000	97.1	1.000	0.008
AdaBoost Gradient trees	96.8	0.998	0.000	97.0	0.999	0.000
Oversampling AdaGradient trees	97.0	0.972	0.968	97.0	0.999	0.000
KNN	96.9	1.000	0.000	97.1	1.000	0.000
Bagging KNN	96.9	1.000	0.000	97.1	1.000	0.000
Oversampling Bagging KNN	79.4	0.592	0.996	82.3	0.841	0.235
SVM	96.9	1.000	0.000	97.1	1.000	0.000
AdaBoost SVM	96.9	1.000	0.000	97.1	1.000	0.000
Oversampling AdaBoost SVM	53.6	0.397	0.675	60.6	0.614	0.364
MLP	96.9	1.000	0.000	97.1	1.000	0.000
Bagging MLP	96.9	1.000	0.000	97.1	1.000	0.000
Oversampling bagging MLP	45.7	0.226	0.687	96.6	0.994	0.015
Deep learning NN	96.9	1.000	0.000	97.1	1.000	0.000
Oversampling deep learning NN	62.1	0.484	0.757	93.3	0.959	0.068

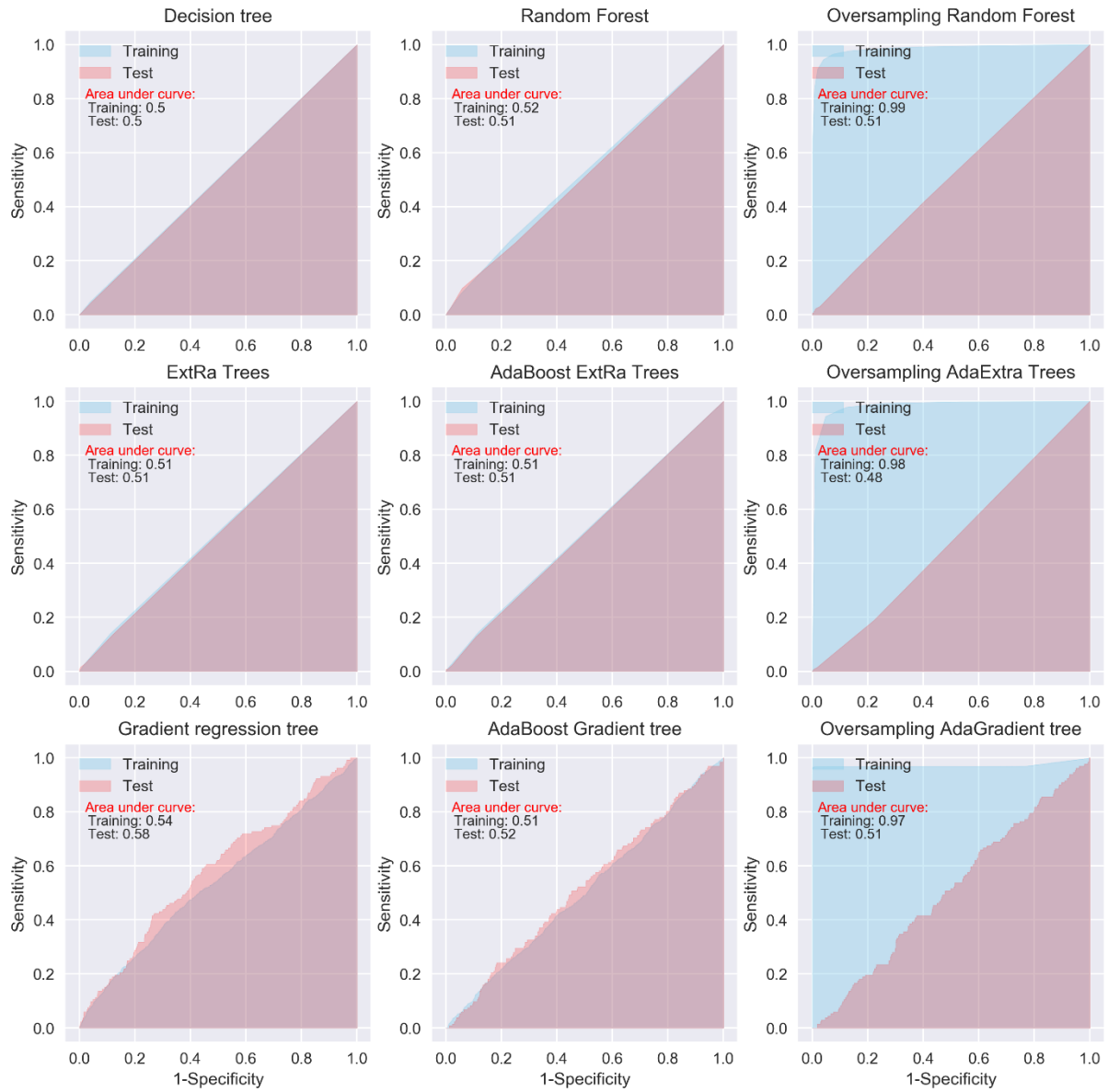
372 AdaBoost, adaptive boosting; LDA, linear discriminant analysis; QDA, quadratic discriminant  
 373 analysis; ExtRa, extremely randomized; AdaExtra, adaptive boosting extremely randomized;  
 374 AdaGradient, adaptive boosting gradient; KNN, k-nearest neighbor; SVM, support vector machine;  
 375 MLP, multilayer perceptron; NN, neural network

376



377  
378

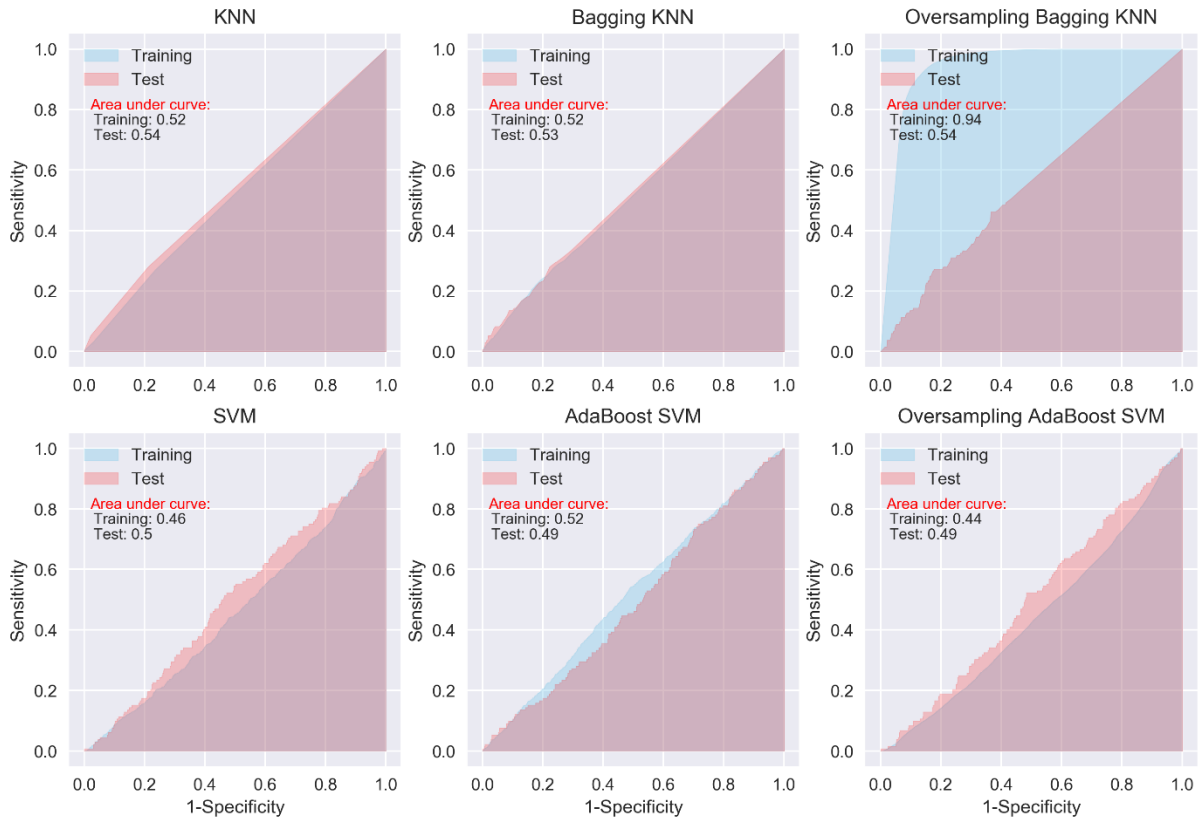
Figure 1. ROC curves of logistic regression, LDA and QDA



379

380

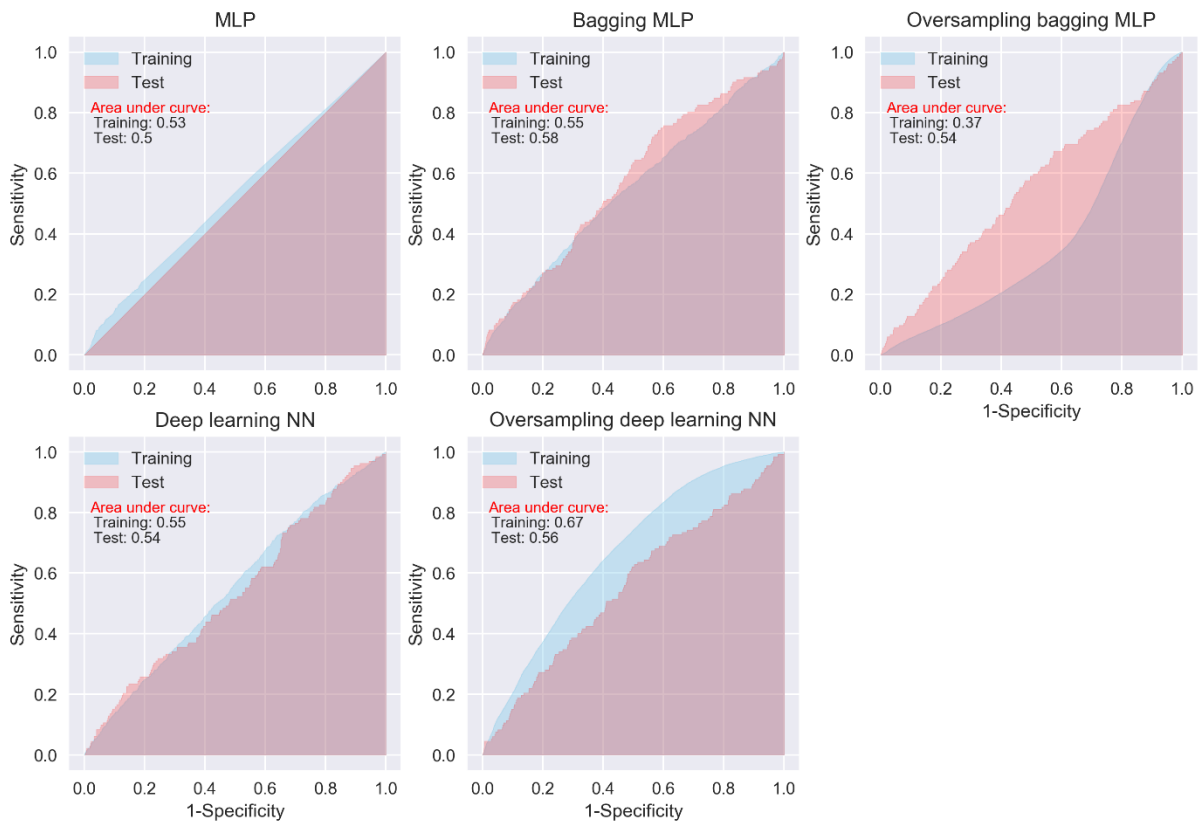
Figure 2. ROC curves of tree-based algorithms



381

382

Figure 3. ROC curves of KNN and SVM



383

384

Figure 4. ROC curves of neural network algorithms

## 385 References

- 386 Ali, A.-R. (2017). Deep Learning in Oncology—Applications in Fighting Cancer. Retrieved from  
387 <https://www.techemergence.com/deep-learning-in-oncology/>
- 388 Anderin, C., Gustafsson, U. O., Heijbel, N., & Thorell, A. (2015). Weight loss before bariatric surgery  
389 and postoperative complications: data from the Scandinavian Obesity Registry (SOReg). *Ann*  
390 *Surg*, 261(5), 909-913.
- 391 Batista, G. E., Prati, R. C., & Monard, M. C. (2004). A study of the behavior of several methods for  
392 balancing machine learning training data. *ACM SIGKDD explorations newsletter*, 6(1), 20-29.
- 393 Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of*  
394 *Machine Learning Research*, 13(Feb), 281-305.
- 395 Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority  
396 over-sampling technique. *Journal of Artificial Intelligence Research*, 16, 321-357.
- 397 Ehlers, A. P., Roy, S. B., Khor, S., Mandagani, P., Maria, M., Alfonso-Cristancho, R., & Flum, D. R.  
398 (2017). Improved Risk Prediction Following Surgery Using Machine Learning Algorithms.  
399 *eGEMs (Generating Evidence & Methods to improve patient outcomes)*, 5(2), 3.
- 400 Finks, J. F., Kole, K. L., Yenumula, P. R., English, W. J., Krause, K. R., Carlin, A. M., ... & Michigan  
401 Bariatric Surgery Collaborative. (2011). Predicting risk for serious complications with bariatric  
402 surgery: results from the Michigan Bariatric Surgery Collaborative. *Ann Surg*, 254(4), 633-640.
- 403 Friedman, J. H. (2002). Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38(4),  
404 367-378.
- 405 Geubbels, N., de Brauw, L. M., Acherman, Y. I. Z., van de Laar, A. W. J. M., & Bruin, S. C. (2015). Risk  
406 Stratification Models: How Well do They Predict Adverse Outcomes in a Large Dutch Bariatric  
407 Cohort? *Obesity Surgery*, 25(12), 2290-2301.
- 408 Geurts, P., Ernst, D., & Wehenkel, L. (2006). Extremely randomized trees. *Machine learning*, 63(1), 3-  
409 42.
- 410 Gupta, P. K., Franck, C., Miller, W. J., Gupta, H., & Forse, R. A. (2011). Development and Validation of  
411 a Bariatric Surgery Morbidity Risk Calculator Using the Prospective, Multicenter NSQIP  
412 Dataset. *Journal of the American College of Surgeons*, 212(3), 301-309.
- 413 Hayes, M. T., Hunt, L. A., Foo, J., Tychinskaya, Y., & Stubbs, R. S. (2011). A model for predicting the  
414 resolution of type 2 diabetes in severely obese subjects following Roux-en Y gastric bypass  
415 surgery. *Obes Surg*, 21(7), 910-916.
- 416 Hedenbro, J. L., Naslund, E., Boman, L., Lundegardh, G., Bylund, A., Ekelund, M., ... & Naslund, I.  
417 (2015). Formation of the Scandinavian Obesity Surgery Registry, SOReg. *Obes Surg*, 25(10),  
418 1893-1900.
- 419 James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning with*  
420 *Application in R* (Vol. 112). New York: Springer.
- 421 Jiang, F., Jiang, Y., Zhi, H., Dong, Y., Li, H., Ma, S., ... & Wang, Y. (2017). Artificial intelligence in  
422 healthcare: past, present and future. *Stroke Vasc Neurol*, 2(4), 230-243.
- 423 Kotsiantis, S. B., Zaharakis, I., & Pintelas, P. (2007). Supervised machine learning: A review of  
424 classification techniques. *Emerging artificial intelligence applications in computer*  
425 *engineering*, 160, 3-24.
- 426 Kourou, K., Exarchos, T. P., Exarchos, K. P., Karamouzis, M. V., & Fotiadis, D. I. (2015). Machine  
427 learning applications in cancer prognosis and prediction. *Comput Struct Biotechnol J*, 13, 8-  
428 17.
- 429 Liaw, A., & Wiener, M. (2002). Classification and regression by randomForest. *R news*, 2(3), 18-22.
- 430 Longitudinal Assessment of Bariatric Surgery, C., Flum, D. R., Belle, S. H., King, W. C., Wahed, A. S.,  
431 Berk, P., ... & Wolfe, B. (2009). Perioperative safety in the longitudinal assessment of bariatric  
432 surgery. *N Engl J Med*, 361(5), 445-454.

- 433 Maalouf, M., Homouz, D., & Trafalis, T. B. (2018). Logistic regression in large rare events and  
434 imbalanced data: A performance comparison of prior correction and weighting methods.  
435 *Computational Intelligence*, 34(1), 161-174.
- 436 Maciejewski, M. L., Winegar, D. A., Farley, J. F., Wolfe, B. M., & DeMaria, E. J. (2012). Risk  
437 stratification of serious adverse events after gastric bypass in the Bariatric Outcomes  
438 Longitudinal Database. *Surg Obes Relat Dis*, 8(6), 671-677.
- 439 Mishra, S., Sturm, B. L., & Dixon, S. (2017). *LOCAL INTERPRETABLE MODEL-AGNOSTIC EXPLANATIONS*  
440 *FOR MUSIC CONTENT ANALYSIS*. Paper presented at the The 18th ISMIR Conference, Suzhou,  
441 China.
- 442 Ng, M., Fleming, T., Robinson, M., Thomson, B., Graetz, N., Margono, C., ... & Gakidou, E. (2014).  
443 Global, regional, and national prevalence of overweight and obesity in children and adults  
444 during 1980-2013: a systematic analysis for the Global Burden of Disease Study 2013. *Lancet*,  
445 384(9945), 766-781.
- 446 Olbers, T., Lonroth, H., Fagevik-Olsen, M., & Lundell, L. (2003). Laparoscopic gastric bypass:  
447 development of technique, respiratory function, and long-term outcome. *Obes Surg*, 13(3),  
448 364-370.
- 449 Pan, L., Liu, G., Lin, F., Zhong, S., Xia, H., Sun, X., & Liang, H. (2017). Machine learning applications for  
450 prediction of relapse in childhood acute lymphoblastic leukemia. *Sci Rep*, 7(1), 7402.
- 451 Passos, I. C., Mwangi, B., & Kapczynski, F. (2016). Big data analytics and machine learning: 2015 and  
452 beyond. *Lancet Psychiatry*, 3(1), 13-15.
- 453 Pedersen, H. K., Gudmundsdottir, V., Pedersen, M. K., Brorsson, C., Brunak, S., & Gupta, R. (2016).  
454 Ranking factors involved in diabetes remission after bariatric surgery using machine-learning  
455 integrating clinical and genomic biomarkers. *Npj Genomic Medicine*, 1.
- 456 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E.  
457 (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12,  
458 2825-2830.
- 459 Piaggi, P., Lippi, C., Fierabracci, P., Maffei, M., Calderone, A., Mauri, M., ... & Santini, F. (2010).  
460 Artificial neural networks in the outcome prediction of adjustable gastric banding in obese  
461 women. *PLoS One*, 5(10), e13624.
- 462 Probst, P., Bischl, B., & Boulesteix, A.-L. (2018). Tunability: Importance of hyperparameters of  
463 machine learning algorithms. *arXiv preprint arXiv:1802.09596*. Retrieved from  
464 <https://arxiv.org/abs/1802.09596>
- 465 Razzaghi, T., Safro, I., Ewing, J., Sadrfaridpour, E., & Scott, J. D. (2017). Predictive Models for Bariatric  
466 Surgery Risks with Imbalanced Medical Datasets. Clemson, South Carolina: TigerPrints.
- 467 Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). *Why should i trust you?: Explaining the predictions of*  
468 *any classifier*. Paper presented at the Proceedings of the 22nd ACM SIGKDD International  
469 Conference on Knowledge Discovery and Data Mining, San Francisco, California, USA
- 470 Schapire, R. E. (2003). The boosting approach to machine learning: An overview *Nonlinear estimation*  
471 *and classification* (pp. 149-171): Springer.
- 472 Sjoström, L., Lindroos, A. K., Peltonen, M., Torgerson, J., Bouchard, C., Carlsson, B., ... & Swedish  
473 Obese Subjects Study Scientific, G. (2004). Lifestyle, diabetes, and cardiovascular risk factors  
474 10 years after bariatric surgery. *N Engl J Med*, 351(26), 2683-2693.
- 475 Stenberg, E., Cao, Y., Szabo, E., Naslund, E., Naslund, I., & Ottosson, J. (2018). Risk Prediction Model  
476 for Severe Postoperative Complication in Bariatric Surgery. *Obes Surg*.
- 477 Stenberg, E., Szabo, E., Agren, G., Naslund, E., Boman, L., Bylund, A., ... & Scandinavian Obesity  
478 Surgery Registry Study, G. (2014). Early complications after laparoscopic gastric bypass  
479 surgery: results from the Scandinavian Obesity Surgery Registry. *Ann Surg*, 260(6), 1040-  
480 1047.
- 481 The international federation for the surgery of obesity and metabolic disorders. (2017). *Third IFSO*  
482 *Global Registry Report 2017*. Retrieved from Oxfordshire, United Kingdom:  
483 [http://www.ifso.com/wp-content/themes/ypo-theme/pdfs/3rd-ifso-report-21-august-](http://www.ifso.com/wp-content/themes/ypo-theme/pdfs/3rd-ifso-report-21-august-2017.pdf)  
484 [2017.pdf](http://www.ifso.com/wp-content/themes/ypo-theme/pdfs/3rd-ifso-report-21-august-2017.pdf)



485 Thomas, D. M., Kuiper, P., Zaveri, H., Surve, A., & Cottam, D. R. (2017). Neural Networks to Predict  
486 Long-term Bariatric Surgery Outcomes. *Bariatric Times*, 14(12), 14-17.  
487 Zhou, Z.-H. (2012). *Ensemble methods: foundations and algorithms*. New York: CRC press.

488

489

490 Annex 1: Terminology and derivations

491 Cross Entropy loss (Log loss):

492 
$$V(f(\vec{x}), y) = -y\ln(f(\vec{x})) - (1 - y)\ln(1 - f(\vec{x}))$$

493 where  $y$  is true classifier  $\in \{0, 1\}$  and  $f(\vec{x})$  is predicted value.

494

495 True or false refers to the predicted outcome being correct or wrong, while positive or negative refers  
496 to presenting severe complication or no severe complication.

497 ACC: accuracy

498 SEN: sensitivity

499 SPE: specificity

500 TP: number of true positives, i.e. patient presenting severe complication correctly predicted as positive

501 TN: number of true negatives, i.e. patient without severe complication correctly predicted as negative

502 FP: number of false positives, i.e. patient without severe complication wrongly predicted as positive

503 FN: number of false negatives, i.e. patient presenting severe complication wrongly predicted as  
504 negative

505 Total: total number of the patients, i.e. TP+TN+FP+FN

506 P: number of patients presenting severe complication, i.e. TP+FN

507 N: number of patients without severe complication, i.e. TN+FP

508 AUC: area under the receiver operating characteristic (ROC) curve for binary outcome

509 T: threshold for a patient is classified as presenting severe complication if  $X > T$ , where  $X$  is predicted  
510 probability of a patients presenting severe complication by an algorithm.

$$511 \quad ACC = \frac{TP + TN}{Total} \times 100\%$$

$$512 \quad SEN = \frac{TP}{P}$$

$$513 \quad SPE = \frac{TN}{N}$$

$$514 \quad AUC = \int_0^1 SEN_T(1 - SPE)_T dT$$