

1 Running head: Predicting CVD

2

3

4 Predicting cardiovascular disease 5 (Supervised classification)

6

7 Esther Olabisi-Adeniyi

8

9 Department of Integrative Biology

10 University of Guelph, Guelph, Ontario, N1G 2W1, Canada

11

12

13

14

15

ABSTRACT

Cardiovascular disease (CVD) can be caused by several genetic and lifestyle factors. Nevertheless, machine learning (ML) algorithms can be used to predict CVD and assess the relevance of different risk factors. This project aims to predict CVD based on age, cholesterol, smoking status, alcohol consumption, and physical activity status using gradient boosting, k nearest neighbor, random forest, ridge, and decision tree classifiers. Feature selection algorithms (SelectKBest and RFE) indicate that age, cholesterol, and smoking are the most important features for predicting CVD. The gradient boosting and ridge classifiers were the best models based on F1 scores and learning plots. The ridge classifier had the lowest computational cost and is, therefore, the preferred option for scaling up to larger datasets. These findings could prompt early prevention and intervention of CVD.

Key words (maximum 6 keywords):

machine learning, cardiovascular disease, gradient boosting classification, ridge classification, feature selection

List of abbreviations:

ML: machine learning, CVD: cardiovascular disease, CV: cross-validation, RF: random forest, KNN: K nearest neighbor classification, GBC: gradient boosting classification, DTC: decision tree classification, WHO: world health organization, RFE: recursive feature elimination.

1. INTRODUCTION

Cardiovascular disease (CVD) is a broad term that includes all diseases of the heart and blood vessels. It is the leading cause of death globally, accounting for an estimated 17.9 million deaths per year according to the World Health Organization (WHO) (WHO, 2021). Risk factors for CVD include smoking, high blood pressure, high cholesterol, diabetes, obesity, physical inactivity, and hereditary factors (Benjamin et al., 2019). While some risk factors such as family history cannot be changed, many can be modified through lifestyle changes and medications. Moreover, over 80% of CVD diseases are preventable (Salah and Srinivas, 2022); therefore, an understanding of the relationship between CVD and its risk factors could prompt early prevention and intervention. A machine-learning algorithm that predicts the presence of CVD could effectively model the connection between risk factors and the disease enabling individuals to assess their CVD risk and take appropriate action by improving their lifestyle or seeking medical counsel. Several studies have employed machine learning algorithms to examine the relationship between CVD and its risk factors (Pal et al., 2022; (Salah and Srinivas, 2022)). Likewise, this project aims to adopt ML models in predicting CVD occurrence based on several risk factors: age, cholesterol, smoking status, alcohol consumption status, and physical activity status. The selected ML models are gradient boosting, k nearest neighbor, random forest, ridge, and decision tree classifiers. By building the aforementioned models to predict CVD, this study aims to answer the following questions:

- (i) Which of the classification models would have the best performance when predicting CVD?
- (ii) Which risk factors will emerge as the most relevant ones for predicting CVD?

49 2. MATERIALS AND METHODS

50 Data collection

51 This project uses a ‘cardiovascular disease dataset’ from *Kaggle.com* (The Devastator, 2023). The dataset records 12 different health
52 metrics of 70, 000 individuals, including their cardiovascular disease status (presence or absence of CVD). The following
53 measurements were taken: age, cholesterol, smoking status, alcohol consumption status, physical activity status, glucose level, gender,
54 height, weight, systolic blood pressure, diastolic blood pressure, and CVD status. Meanwhile, only the first five measurements of the
55 aforementioned list were selected to be features for predicting the latter measurement, CVD status. The reasoning behind this manual
56 feature selection was based on an assumption that we could have a scenario where these five measurements were the only recorded
57 features; and if so, “what would our results be?” However, it has now become clear that this manual feature selection step was not
58 ideal and that all features should have been included in the study.

59 According to the source of the dataset, the variable measurements were recorded from individuals during medical examinations
60 (Ulianova, 2019). Some of the variables were objectively considered (age, gender, height, and weight) while some were medically
61 assessed (cholesterol, glucose level, systolic blood pressure, and diastolic blood pressure). The rest were recorded as provided by the
62 patient. For the six variables used for this project, cholesterol was a categorical variable, age was numerical, and others were binary
63 variables.

64

65 Data processing

66 Quality control of the raw dataset involved a series of steps carried out using Python programming language and libraries. Although
67 the dataset appeared to be free of non-standard characters, I did not scan every single cell for such characters. Therefore, the most
68 practical option was to import the dataset through "utf-8" encoding. I also used the *Pandas* package to analyze the dataset for missing
69 values but none was found. Considering that, 6 out of 12 columns were selected for this project, the chances of having duplicate rows
70 and confounding records was relatively high. Consequently, 31,475 duplicate rows were found and removed; 38,525 remained. I
71 further inspected the remaining dataset for confounding records, that is, records with the same input values but different CVD statuses.
72 Confounding records were detected and removed, narrowing the dataset down to 29,886 records. The 'age' variable was the only
73 variable with a wide numerical range; therefore, it was examined for outliers using $\pm 2\text{IQR}$ as the cutoff. No outliers were detected. By
74 employing functions from the *Pandas* library, all binary variables were converted to dummy variables while 'cholesterol' was
75 converted to an ordinal variable on a scale of 1 to 3 indicating cholesterol level: 1= normal, 2=above normal, 3=well-above normal.

76 Problem definition

77 In order to predict the target variable, CVD occurrence, based on the input variables, the ML models will be designed to solve a
78 supervised classification problem. Each model will be trained with input variables and their corresponding CVD outputs. Afterwards,
79 the models will be tested with unseen data.

80

81 Statistical analyses

82 The refined dataset contained 29,886 complete records (CVD+ : 15,917, CVD- :13,969). *Scikit-learn*'s splitting function was used to
83 split the dataset into training and testing sets at a ratio of 80:20 ratio respectively. Afterwards, I applied two feature selection
84 algorithms on the training set. I expect the two algorithms to give the same results; otherwise, only one set will be used. SelectKBest
85 with *f_classif* was utilized in selecting the three most relevant features for output prediction. This option was used because it is
86 designed for classification problems. The function computes ANOVA F-values for each feature and high scores indicate the best
87 features. The second feature selection algorithm was Recursive Feature Elimination (RFE) with logistic regression. For this method,
88 the least relevant features are removed during iterations of the logistic regression model (Guyon et al., 2002). I requested the best three
89 features for RFE also.

90

91 Machine learning modelling

92 Ridge classification was chosen for this project because it is a regularization technique that efficiently handles the issue of
93 multicollinearity between input and output variables, as is the case for this study (Arashi et al., 2021). I utilized the ridge classification
94 function from *scikit-learn* which first converts output values from their binary state to $\{-1, 1\}$ before applying ridge regression
95 techniques (`sklearn.linear_model.RidgeClassifier`, 2023). The default settings for ridge classifier were used during training and testing.

96 Decision tree classification (DTC) is another powerful algorithm used for supervised learning problems. It works by
97 recursively splitting the data into subsets based on the most significant features that provide the best separation between classes

98 (Bahzad Taha Jijo and Adnan Mohsin Abdulazeez, 2021). DTC can be used to predict categorical, numerical, and binary variables.
99 *Scikit-learn*'s DTC function was applied to the training and testing sets of this project. Parameters were set as follows: 'gini impurity'
100 for quality criterion and 'best' for splitting strategy.

101 Gradient boosting classification (GBC) iteratively builds a series of decision trees, where each new tree is trained to correct the
102 errors of the previous ones (Alexey Natekin and Knoll, 2013). The final model is an ensemble of these weak learners, with each tree
103 weighted according to its performance on the training set. Gradient boosting is known for its ability to handle complex datasets and
104 provide relatively accurate estimates of the target variable (Alexey Natekin and Knoll, 2013). I applied *scikit-learn*'s GBC function
105 during training and testing with 100 as the number of boosting stages.

106 The fourth classifier used in this project was k nearest neighbor classifier (KNN). KNN is popularly used for classification and
107 regression problems. This model first computes the distance between the new sample and training sample and then identifies the K
108 nearest data points for that new sample (Wang, 2019). Ultimately, a class label is assigned to the sample based on the majority class of
109 those neighboring points. The chosen K for training and testing with KNN was five.

110 Random forest classification (RF), like GBC, incorporates several models in building an ensemble. It works by creating
111 multiple decision trees on randomly sampled subsets of the data and features (Kulkarni and Sinha, 2013). The algorithm then
112 combines the predictions of all the decision trees to obtain the final prediction. Contrary to the weighted technique used in GBC, RF
113 uses a bagging method whereby the generated models are independent of each other (Kulkarni and Sinha, 2013). The number of trees
114 for the RF function from *scikit-learn* was set to 100.

115 Subsequently, the ML models were evaluated based on five metrics: k-fold cross-validation (CV), balanced accuracy, F1-
116 score, precision score, and learning curves. Firstly, a 10-fold CV of the training dataset was done with each classifier, after which the
117 average values for balanced accuracy, F1-score, and precision score were computed to compare the models. A similar process was
118 repeated for the testing dataset but with a *scikit-learn* function for predictions. Two types of learning curves were generated for each
119 model. Test-training learning curves were plotted to analyze fitting during training and CV while ‘training samples-fit time’ curves
120 (scalability of the model) were generated to analyze the computational costs required by each model.

121

122 3. RESULTS AND DISCUSSION

123 Feature selection

124 This project employed the use of SelectKBest and RFE in selecting the most important features for predicting CVD. Results for the
125 two algorithms demonstrate that age, cholesterol, and smoking, in the order listed, are the three most important features for predicting
126 CVD (Table 1). These results are consistent with previous findings of a ML study that highlights age and smoking as important
127 predictors for CVD (Salah and Srinivas, 2022).

128

129 ML performance based on evaluation scores

130 Each classification model was evaluated by tenfold CV averages of balanced accuracy, F1, and precision scores. All three metrics rank
131 the models equally during training (Table 2); however, since F1 score considers class imbalance, we will prioritize it during this

132 analysis. GBC was the top performer during training with an F1 score of 0.634 and ridge follows closely with 0.633 (Table 2). KNN,
133 RF, and DTC have 0.569, 0.556, and 0.554 respectively. Upon testing with unseen data, GBC and ridge maintained their positions
134 (0.639 and 0.638 respectively) (Table 2). The ML study by Salah and Srinivas also discovered that a gradient boosting model,
135 XGBoost had the best performance when predicting CVD based on risk factors (2022), thereby corroborating my findings. GBC likely
136 performed better than the other models because of its ensemble and boosting techniques, which are ideal for complex datasets like the
137 mixed-variable dataset for this study. Ridge also performed approximately as well as GBC, probably because of its efficiency when
138 multicollinearity is involved.

139 140 [ML performance based on learning plots](#)

141 In addition to analyzing F1, balanced accuracy, and precision for model performance, we can also examine learning and scalability
142 curves. The learning curves for DTC, RF and KNN indicate that a great amount of overfitting exist between their training and CV
143 scores (Fig 1). The low F1 score rankings of these models in addition to their high levels of overfitting suggest that the models require
144 some modifications such as hyperparamter tuning in order to fit this dataset better. GBC and ridge also had overfitting and underfitting
145 respectively, albeit relatively low. Moreover, GBC's curves show a rapid convergence of training and CV scores, further suggesting
146 that GBC may be the better model than ridge (Fig 1).

147 The 'training samples-fit time' plots can help to make inference about the computational power required by each model (Fig
148 1). By understanding the computational costs of a model, we can infer how well the model would perform with similar cardiovascular

149 disease datasets of larger sizes. Fit time (measured in seconds) for each model is as follows (Fig. 1): ridge (0.09s), KNN (0.125s),
150 DTC (0.6s), GBC (10.8s), and RF (32s). Although all fit times in the plots continue to increase as more samples are added such that
151 none of them reach a plateau, the ridge classifier emerges with the lowest fit time and subsequently, the lowest computational cost.
152 KNN and DTC have low fit times as well but their fitting performance has been weak so far. GBC, compared to its close second,
153 ridge, requires more computational power, suggesting that in as much as GBC has had the best performance thus far, ridge may be the
154 preferred option for scaling up to larger datasets. RF's high fit time is most likely because the model builds a large amount of trees
155 when fitting samples. RF's computationally expensive model and weak performance makes for a bad CVD predictor in this study.

156 4. CONCLUSIONS

157 Despite the CVD still being a deadly illness with multiple risk factors, the potential for machine learning to predict this disease holds
158 great promise in reducing its future incidence and fatalities. According to the findings of this project, which are in line with previous
159 research outcomes, age, smoking habits, and cholesterol, are important predictors of CVD in ML (Salah and Srinivas, 2022). While
160 these important predictors do not explicitly represent a causative or associative relationship with CVD, their roles as influential factors
161 can contribute to other ML research and clinical studies for CVD. Furthermore, GBC was found to be the best performing
162 classification model for predicting CVD, with ridge classifier as its close second. Ridge classification model is presented as the
163 preferred option for larger datasets. Future work on this project should explore the applied models for potential improvements by
164 employing several modifications such as hyperparameter tuning and feature engineering.

165 ACKNOWLEDGEMENTS

166 The author would like to thank Dr. Dan Tulpan of the Department of Animal Biosciences for his instruction and support throughout
167 this machine-learning course.

168 AUTHORS' CONTRIBUTIONS

169 All authors improved and contributed to the editing of the manuscript. All authors read and approved the final manuscript.

170 DISCLOSURES

171 The authors declare no real or perceived conflicts of interest.

172

173

174

175

176

177 LITERATURE CITED

- 178 Arashi, M., Mahdi Roozbeh, N.R. Hamzah, and M. Gasparini. 2021. Ridge regression and its applications in genetic studies. PLOS
179 ONE. 16:e0245376–e0245376. doi:<https://doi.org/10.1371/journal.pone.0245376>. Available from:
180 <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0245376>
- 181 Alexey Natekin, and A. Knoll. 2013. Gradient Boosting Machines, A Tutorial. ResearchGate. Available from:
182 https://www.researchgate.net/publication/259653472_Gradient_Boosting_Machines_A_Tutorial
- 183 Bahzad Taha Jijo, and Adnan Mohsin Abdulazeez. 2021. Classification Based on Decision Tree Algorithm for Machine Learning.
184 ResearchGate. Available from:
185 https://www.researchgate.net/publication/350386944_Classification_Based_on_Decision_Tree_Algorithm_for_Machine_Learning
- 186 Benjamin, E. J., P. Muntner, A. Alonso, Márcio Sommer Bittencourt, C. W. Callaway, A. P. Carson, A. M. Chamberlain, A. R.
187 Chang, S. Cheng, S. R. Das, F. N. Dellling, Luc Djoussé, Mitchell S.V. Elkind, J. F. Ferguson, M. Fornage, L. C. Jordan, S. S. Khan,
188 B. M. Kissela, K. L. Knutson, and T. W. Kwan. 2019. Heart Disease and Stroke Statistics—2019 Update: A Report From the
189 American Heart Association. *Circulation*. 139. doi:<https://doi.org/10.1161/cir.0000000000000659>.
- 190 Guyon, I., J. Weston, S. D. Barnhill, and Vladimir Vapnik. 2002. *Machine Learning*. 46:389–422.
191 doi:<https://doi.org/10.1023/a:1012487302797>

192 Kulkarni, V., and P. Sinha. 2013. Random forest classifiers: A survey and future research directions. Available from:
 193 https://adiwijaya.staff.telkomuniversity.ac.id/files/2014/02/Random-Forest-Classifiers_A-Survey-and-Future.pdf

194 Pal, M., S. R. Parija, G. Panda, Kuldeep Dhama, and R. K. Mohapatra. 2022. Risk prediction of cardiovascular disease using machine
 195 learning classifiers. *Open Medicine*. 17:1100–1113. doi:<https://doi.org/10.1515/med-2022-0508>.

196 Salah, H., and S. Srinivas. 2022. Explainable machine learning framework for predicting long-term cardiovascular disease risk among
 197 adolescents. *Scientific Reports*. 12. doi:<https://doi.org/10.1038/s41598-022-25933-5>.

198 Wang, L.-S. 2019. Research and Implementation of Machine Learning Classifier Based on KNN. *IOP conference series*. 677:052038–
 199 052038. doi:<https://doi.org/10.1088/1757-899x/677/5/052038>. Available from: [https://iopscience.iop.org/article/10.1088/1757-](https://iopscience.iop.org/article/10.1088/1757-899X/677/5/052038)
 200 [899X/677/5/052038](https://iopscience.iop.org/article/10.1088/1757-899X/677/5/052038)

201 World Health Organization. 2021. Cardiovascular diseases (CVDs). Available from [https://www.who.int/news-room/fact-](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds))
 202 [sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds))

203 `sklearn.linear_model.RidgeClassifier`. 2023. *scikit-learn*. Available from: [https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.RidgeClassifier.html)
 204 [learn.org/stable/modules/generated/sklearn.linear_model.RidgeClassifier.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.RidgeClassifier.html)

205 The Devastator. 2023. Risk Factors for Cardiovascular Heart Disease. Kaggle.com. Available from:

206 <https://www.kaggle.com/datasets/thedevastator/exploring-risk-factors-for-cardiovascular-diseas>

207 Ulianova, S. 2019. Cardiovascular Disease dataset. Kaggle.com. Available from:

208 <https://www.kaggle.com/datasets/sulianova/cardiovascular-disease-dataset>

209

210

211

212

213

214

215

216 LIST OF FIGURES

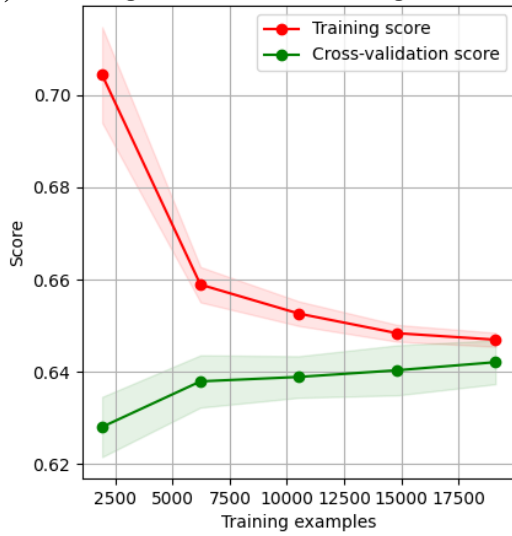
- 217 1. **a-e.** Learning curves showing how each model's accuracy score changes depending on the number samples in the training set.
218 A relatively high level of overfitting can be observed in **(b, c, d)** DTC, RF, KNN while **(a)** GBC overfits minimally. The ridge
219 classifier **(e)** is underfitting the training set. The scalability plots demonstrating how the time (seconds) for fitting the training
220 set changes as number of samples increase **(f-j)**. All the scalability plots show an increase in fit time as more samples are
221 added. However, **(j)** ridge has the lowest fit time at about 0.09s while **(h)** RF has the highest at about 32s.

222

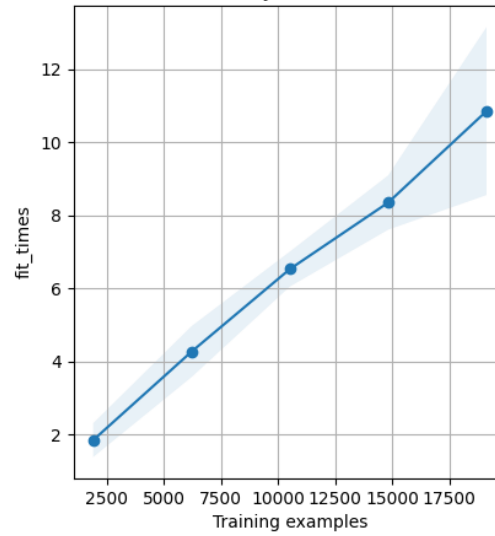
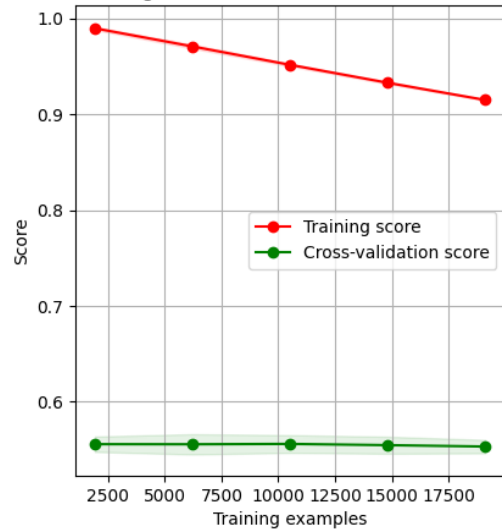
223

224

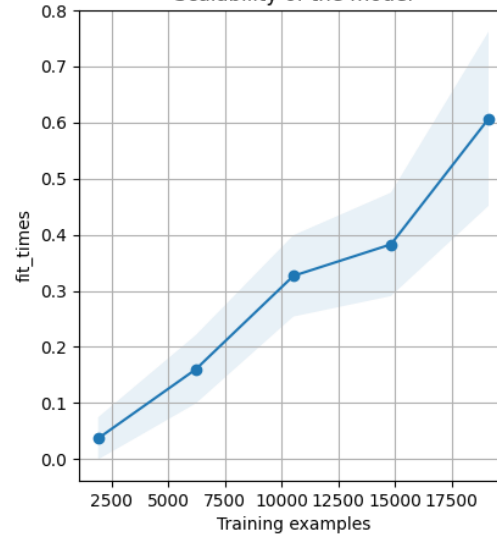
225

(a) Learning Curves Gradient Boosting Classification

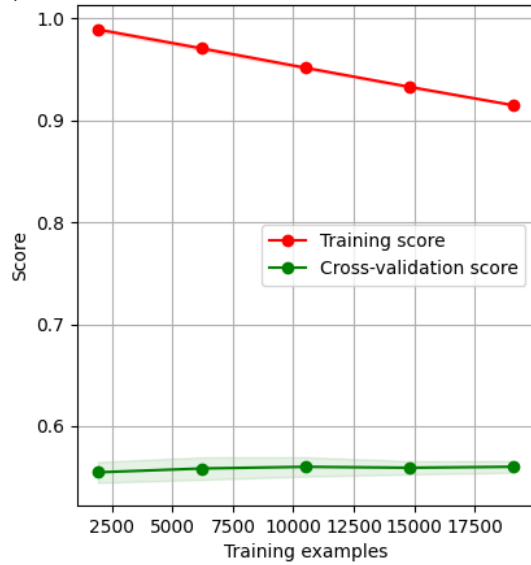
Scalability of the model

**(f)****(b)** Learning Curves Decision Tree Classification

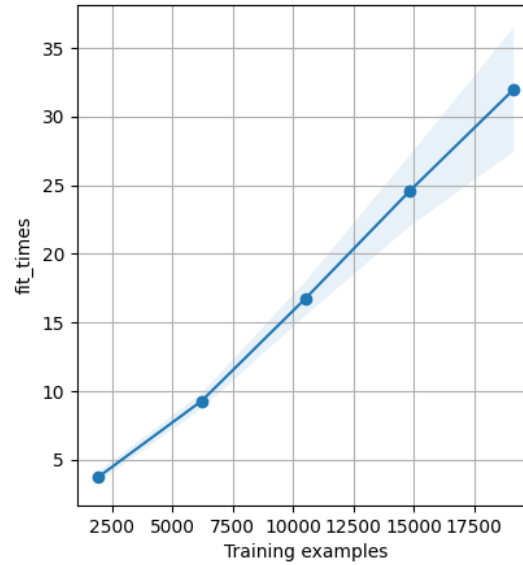
Scalability of the model

**(g)**

(c) Learning Curves Random Forest Classification

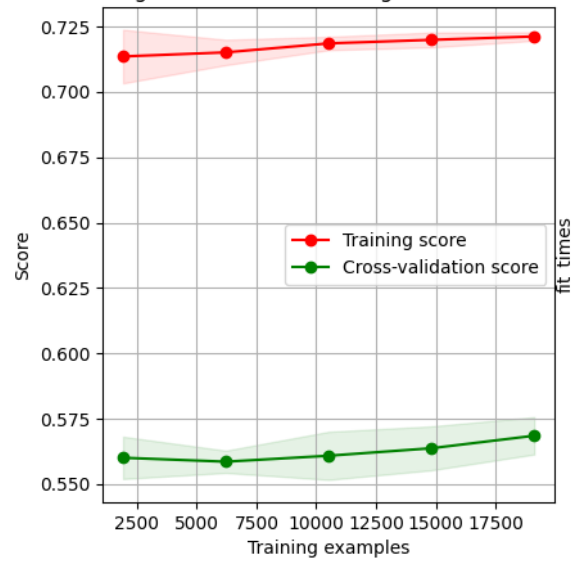


Scalability of the model

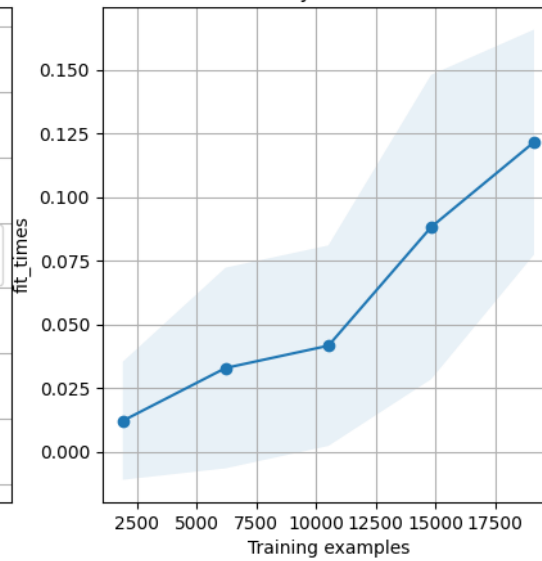


(h)

(d) Learning Curves K-Nearest Neighbour Classification



Scalability of the model



(i)

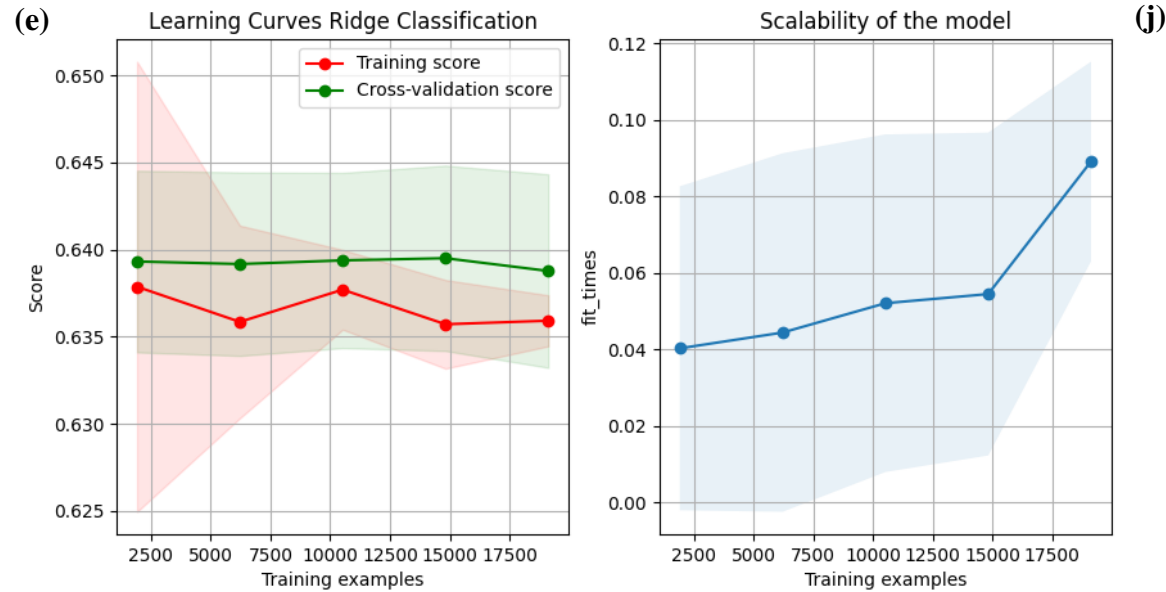


Fig. 1 a-e. Learning curves showing how each model's accuracy score changes depending on the number samples in the training set. A relatively high level of overfitting can be observed in (b, c, d) DTC, RF, KNN while (a) GBC overfits minimally. The ridge classifier (e) is underfitting the training set. The scalability plots demonstrate how the time (seconds) for fitting the training set changes as number of samples increase (f-j). All the scalability plots show an increase in fit time as more samples are added. However, (j) ridge has the lowest fit time at about 0.09s while (h) RF has the highest at about 32s.

232 LIST OF TABLES

- 233 1. Feature selection results for **(a)** RFE and **(b)** SelectKBest algorithms. The two algorithms rank the features similarly. Age,
234 cholesterol, and smoking are the most important features for predicting CVD.
- 235 2. Evaluation scores for each model's fitting of the training set. Highest to lowest performance (1 - 5). **(b)**. Evaluation metrics for
236 each model's performance when predicting testing dataset. Highest to lowest performance (1 - 5).
- 237
- 238
- 239
- 240
- 241
- 242
- 243
- 244
- 245
- 246
- 247
- 248
- 249
- 250
- 251
- 252
- 253
- 254
- 255
- 256
- 257
- 258
- 259
- 260
- 261

RFE			SelectKBest	
Rank	Feature	ANOVA F-values	Rank	Feature
1	Age	1371.4424	1	Age
2	Cholesterol	1567.2153	1	Cholesterol
3	Smoking	61.1316	1	Smoking
4	Alcohol	31.0252	2	Alcohol
5	Physical Activity	17.9860	3	Physical Activity
(a)			(b)	

Table 1. Feature selection rankings based on (a) RFE and (b) SelectKBest algorithms. The two algorithms rank the features similarly. Age, cholesterol, and smoking are the most important features for predicting CVD.

Training scores				
Rank	Algorithm	Balanced accuracy	F1 Score	Precision
1	GBC	0.635	0.634	0.639
2	Ridge	0.633	0.633	0.634
3	KNN	0.569	0.569	0.569
4	RandomForest	0.556	0.556	0.556
5	DTC	0.555	0.554	0.554

(a)

Testing scores				
Rank	Algorithm	Balanced accuracy	F1 Score	Precision
1	GBC	0.639	0.639	0.642
2	Ridge	0.638	0.638	0.638
3	RandomForest	0.563	0.564	0.564
4	DTC	0.562	0.562	0.562
5	KNN	0.543	0.543	0.543

(b)

Table 2 (a). Evaluation scores for each model's fitting of the training set. Highest to lowest performance (1 - 5). **(b).** Evaluation metrics for each model's performance when predicting testing dataset. Highest to lowest performance (1 - 5).