



OPEN

Using explainable machine learning to identify patients at risk of reattendance at discharge from emergency departments

F. P. Chmiel¹✉, D. K. Burns¹, M. Azor², F. Borca^{2,3}, M. J. Boniface¹, Z. D. Zlatev¹, N. M. White¹, T. W. V. Daniels^{4,5} & M. Kiuber⁶

Short-term reattendances to emergency departments are a key quality of care indicator. Identifying patients at increased risk of early reattendance could help reduce the number of missed critical illnesses and could reduce avoidable utilization of emergency departments by enabling targeted post-discharge intervention. In this manuscript, we present a retrospective, single-centre study where we created and evaluated an extreme gradient boosting decision tree model trained to identify patients at risk of reattendance within 72 h of discharge from an emergency department (University Hospitals Southampton Foundation Trust, UK). Our model was trained using 35,447 attendances by 28,945 patients and evaluated on a hold-out test set featuring 8847 attendances by 7237 patients. The set of attendances from a given patient appeared exclusively in either the training or the test set. Our model was trained using both visit level variables (e.g., vital signs, arrival mode, and chief complaint) and a set of variables available in a patients electronic patient record, such as age and any recorded medical conditions. On the hold-out test set, our highest performing model obtained an AUROC of 0.747 (95% CI 0.722–0.773) and an average precision of 0.233 (95% CI 0.194–0.277). These results demonstrate that machine-learning models can be used to classify patients, with moderate performance, into low and high-risk groups for reattendance. We explained our models predictions using SHAP values, a concept developed from coalitional game theory, capable of explaining predictions at an attendance level. We demonstrated how clustering techniques (the UMAP algorithm) can be used to investigate the different sub-groups of explanations present in our patient cohort.

The use of emergency departments (EDs) has been growing steadily over the last decade^{1,2}, which in turn has contributed to increased overcrowding and extended waiting times. These factors have been linked to increased rates of adverse outcomes^{3,4}. Therefore, it is important to investigate the most efficient ways of using the available resources and minimise their unnecessary use. One way this can be achieved is by minimizing short-term reattendances, which describe a situation where a patient presents to an emergency department within 72 h of having been discharged. The number of short-term reattendances can be minimised by both delivering the highest levels of patient care, thereby reducing the chance of missed critical illness and injury at the initial attendance, and by mitigating reattendances for reasons at least partially unrelated to the initial ED attendance.

Research has shown there are several factors indicative of short-term reattendance risk, including social factors (e.g., living alone)⁵, depression⁶, initial diagnosis⁷, and historical emergency department usage⁸. Knowledge of these risk factors is important to clinical staff when planning discharge, but this is unlikely the most optimal way of determining those at risk of suffering from a significant illness following erroneous discharge or patients in need of additional support in the community following discharge. Predictive models, available as a decision support tool at the point of discharge, able to identify patients at increased risk of short-term reattendance may be able to significantly reduce the number of reattendances. Predictions, and any associated explanations, could

¹School of Electronics and Computer Science, University of Southampton, Southampton, UK. ²University Hospitals Southampton NHS Foundation Trust, Southampton, UK. ³Clinical Informatics Research Unit Faculty of Medicine, University of Southampton, Southampton, UK. ⁴Department of Respiratory Medicine, University Hospital NHS Foundation Trust, Southampton, UK. ⁵NIHR Southampton Biomedical Research Centre, Southampton Centre for Biomedical Research, Southampton General Hospital, Southampton, UK. ⁶Emergency Department, University Hospital Southampton NHS Foundation Trust, Southampton, UK. ✉email: francispeterchmiel@gmail.com

Variable	Data type	Encoding	Mean (std. dev.)	Mode	Missing fraction
Age (years)	Integer	–	46.5 (20.4)	21 years	0.0
Attendance complaint	Categorical	Target	–	–	0.0
Triage complaint	Categorical	Target	–	Limb problem	0.0
Diagnosis	Categorical	Target	–	No abnormality detected	0.0
Discriminator	Categorical	Target	–	Recent problem	0.0
Manchester Triage System score	Integer	–	3.3 (0.66)	3	0.0
Pain score	Integer	–	2.6 (2.2)	0	0.0
Arrival mode	Categorical	Target	–	Patient walk-in	0.0
Condition count	Integer	–	1.3 (2.2)	0	0.0
Condition indicators	Categorical	One-hot	–	–	–
30-day visit count	Integer	–	0.3 (1.3)	0	0.0
Temperature, vital signs (°C)	Continuous	–	36.7 (18.1)	–	66%
Systolic blood pressure, vital signs (mmHg)	Continuous	–	139.4 (25.6)	–	66%
Respiration rate, vital signs (Bpm)	Continuous	–	18.3 (4.2)	–	66%
Pulse rate, vital signs (Bpm)	Continuous	–	82.8 (18.2)	–	66%
Blood oxygen saturation, vital signs (%)	Continuous	–	97.2 (5.0)	–	66%
Hour of day (temporal)	Integer	–	13.4 (6.0)	11	0.0
Day of week (temporal)	Categorical	Nominal	–	Monday	0.0

Table 1. Variables used in our modelling and quantitative descriptors of them. For a full description of variables please see Supplementary Table 1. 66 % of attendances did not have associated vital signs recorded. This could of been because it was not deemed necessary by clinical staff, their condition was particularly severe, or because they were not recorded.

be used by clinical staff to help inform intervention (e.g., further diagnostic tests) or enrich discussions about a patients discharge plan.

Machine learning models are a class of predictive models which are well positioned to add value to emergency department processes. They are particularly powerful because they can ingest a large number of variables and learn complex interactions and statistical trends between these variables and associated outcomes, ultimately making highly accurate predictions of patient outcomes. The potential in using machine learning models to inform clinical care and provide high levels of personalization, is evident from the large amount of research into the applications of machine learning in healthcare⁹. Explainable machine learning, in particular, is attracting attention because of the strong assurance needs of healthcare^{10,11}. In the context of this work, research has shown that machine-learning models can use a large number of clinical and administrative variables to provide estimates of a patient's short-term reattendance risk^{12,13}. Explanations are also important in this context, as they can help to inform the patients care trajectory and guide post-discharge intervention plans.

In this manuscript, we present a machine-learning model, utilizing historical (coded, inpatient) discharge summaries, alongside contemporary clinical data recorded during emergency department attendances to identify patients at increased risk of short-term reattendance following an ED attendance. We explain our predictions by calculating SHAP values for our model output and cluster these explanations using the UMAP algorithm to explore the different sub-groups at risk of reattendance in our cohort.

Methods

Dataset curation. Our dataset featured an pseudonymized version of all attendances by adults to Southampton's Emergency Department (University Hospitals Southampton Foundation Trust) occurring between the 01/04/2019 and the 30/04/2020. The dataset included attendance level information such as the results of any near-patient observations and high-level information included in the standard UK Emergency Care Data Set (e.g., outcome, arrival mode, duration of visit, and chief complaint). In addition to attendance-level variables, our dataset included variables extracted from a patients electronic health record maintained by the University Hospitals Southampton Foundation Trust. These variables included any recorded medical complaints, which were obtained by extracting ICD10 coded conditions included in historical discharge summaries. An example of the most frequently observed medical conditions are presented in Table 2, alongside a breakdown of patient ages in our cohort. In total, our dataset included 14 variable sets which are displayed in Table 1 alongside key quantitative descriptors of them. Additional descriptions of variables are presented in Supplementary Table 1.

Our study cohort featured 54,015 attendances which resulted in discharge directly from the ED. Of these attendances we discarded those that resulted in planned reattendance (N = 328) and those occurring after 01/02/2020, removing attendances which occurred during the COVID-19 pandemic. Remaining attendances (N = 44,294) were randomly split at the patient level to create a hold-out test set (N = 8847) containing attendances from 20% of patients and the remaining attendances (N = 35,447) were used as the training set (Fig. 1).

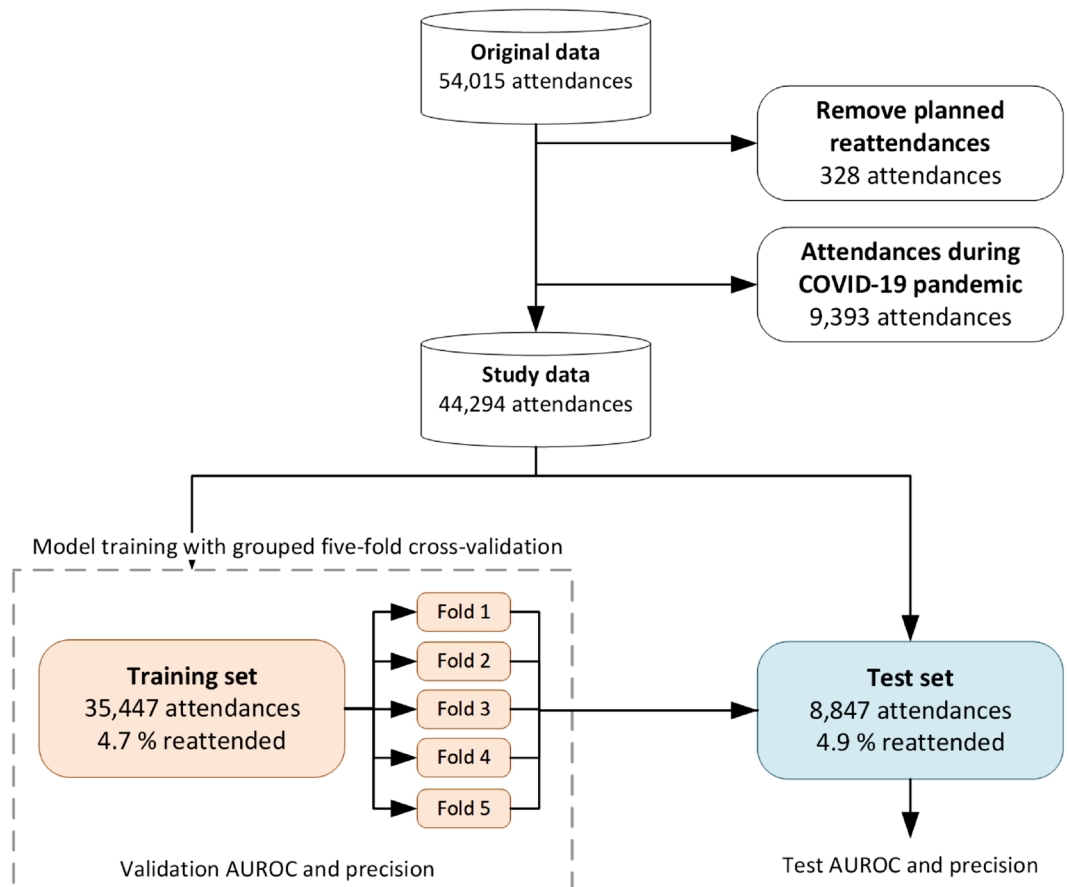


Figure 1. Partition of the study data into the training and hold-out test set. Attendances which resulted in planned reattendances (328 attendances) were removed. All attendances (9393 attendances) occurring after 01/02/2020 were also discarded to remove attendances which coincided with the COVID-19 pandemic. Reattendance rates (bottom row of rectangles representing training and test sets) display the observed 72-h reattendance rate for each cohort. Models were trained using (grouped) five-fold cross validation of the training set with patients in the train set not present in the validation set during the training phase. Model hyperparameters were selected such that the validation AUROC (mean out-of-fold AUROC across the five folds) was maximized. The final model was a mean ensemble of the five models (the average prediction of the five models) trained in the cross-validation process, which was then evaluated on the hold-out test set.

Ethics and data governance. Data was pseudonymized (and where appropriate linked) before being passed to the research team. The research team did not have access to the pseudonymisation key. Since the data was de-identified, Confidentiality Advisory Group approval was not needed and the study was approved by the NHS Health Research Authority (20/HRA/1102) without the need for informed consent because of the de-identified nature of the data. This study was approved by the University of Southampton's Ethics and Research governance committee (ERGO/FEPS/53164). Research was performed in the manner described in the protocol (v1) approved by the NHS Health Research Authority and conducted in accordance with relevant guidelines.

Reattendance identification. Reattendances were identified by using the patient pseudo identifier to calculate the time to a patient's next ED attendance. All reattendances (with the exception of planned reattendances, Fig. 1) were considered, even if the reattendance was for a different reason to the original attendance. We then dichotomized the time to next attendance (return within 72 h from the point of discharge), annotating each attendance with a binary indicator denoting whether it was followed by another attendance within 72 h. This formulation allowed us to frame the predictive task as a binary classification problem.

Predictive modelling. We used an extreme gradient boosting decision tree, as implemented in the XGBoost framework (version 1.3.3)¹⁶, as our machine-learning model. Models were trained using five-fold cross validation of the training set, with patients in the training folds not being present in the validation folds (Fig. 1). Final predictions were the mean prediction of the five models trained during the cross-validation process. Model hyperparameters were optimized by selecting the hyperparameters which maximized the validation AUROC (the mean performance of model on the five sets of out-of-fold samples). Hyperparameter search was performed using Bayesian optimization utilizing the Tree Parzen Estimator algorithm as implemented in the

	Number of patients	
	Train set (N = 28,945)	Test set (N = 7237)
Age groups		
18–24	5039	1244
25–34	5867	1500
35–44	4310	1108
45–54	3957	1024
55–64	3514	822
65–74	2876	701
75+	3382	838
Condition		
Hypertension	2911	681
History of smoking	1805	442
Asthma	1633	392
Depression	1544	393
Current Smoker	1436	359
Type 2 diabetes	1081	278
Hypercholesterolaemia	895	201
Osteoarthritis	775	189
COPD	706	157

Table 2. Breakdown of patient age and most frequently occurring conditions in the training and test sets. Conditions were extracted from their ICD10 codes. A given condition is only associated with a small fraction of attendances, but in total 40.0% of attendances resulting in discharge have at least one associated condition. Patient age and conditions are extracted from their last attendance in the respective data set.

hyperopt Python library (version 0.2.5)^{17,18}. Bayesian optimization was used because it is empirically more efficient than grid or randomized search methods¹⁹. Medical conditions associated with a patient were included as a one-hot-encoded feature vector, the day of the week encoded using ordinal encoding, and all other categorical variables were encoded using target encoding²⁰. The full encoding scheme is outlined in Table 1. Source code for our modelling is available online²¹.

To investigate the predictive ability of individual variables, under an independent variable assumption, we trained models using each distinct variable and evaluated them using five-fold cross validation on the training set. Models were trained (including hyperparameter tuning) using the previously described process. Next, we trained and evaluated a model using the three most predictive variables identified in this process and a model using all variables available at the point of discharge. During this analysis the standard error (95% confidence) of the models performance across the five-folds was calculated. The final model (including all variables available at discharge) was then evaluated on the hold-out test set, with 95% confidence intervals calculated using bootstrapping (n = 1000) of the hold-out test set. Models performance was evaluated using the Area Under the Receiving Operating Curve (AUROC) and the average precision under the precision-recall curve. Hyperparameters of our final model are presented in Supplementary Table 2.

We also evaluated the final model against patients *readmission* status (readmissions are a subset of reattendances for which the outcome was admission to hospital). This was performed without re-training the model and by evaluating the classifier with the ground truth outcome set to a patients readmission status.

Model explainability. To explain the predictions of our final model we made use of the TreeExplainer algorithm implemented in the SHAP Python library (version 0.37.0)^{22–24}. TreeExplainer is an algorithm which calculates SHAP values (i.e., Shapley values), in an optimized manner for decision trees. SHAP values are a concept from coalitional game theory which treats predictive variables as players in a game and distributes their contribution to the predicted probability. SHAP values are particularly powerful as they meet the four desirable theoretical conditions of an explanation algorithm and can provide instance (i.e., attendance) level explanations²⁴. Practically, for each attendance we have a scalar value for each variable used by the model which quantifies the contribution that variable had on the predicted reattendance risk. SHAP values of larger magnitudes indicate that a variable is of increased importance in determining the predicted reattendance risk. SHAP values can be negative (adding the variable reduces predicted reattendance risk) or positive (adding the variable increases the predicted reattendance risk).

To investigate the different explanations across the hold-out test set, we projected the SHAP values for all attendances in the hold-out test set into a two-dimensional ('explanation') space using Uniform Manifold Approximation and Projection (UMAP, version 0.5.1)²⁵. UMAP is a dimensionality reduction technique regularly used to visualise high-dimensional spaces in a low-dimensional embedding, such that global and local structure of the space can be explored^{26,27}. Attendances which are closer in proximity in this space share a more similar explanation for their predicted reattendance risk. To investigate the characteristics of different sub-groups in the explanation space, we assigned each attendance to a cluster in this space using the DBScan algorithm²⁸.

Model reference	Variables	Validation AUROC (95% CI)	Validation average precision (95% CI)
a	Day of week	0.501 (0.489–0.513)	0.047 (0.044–0.050)
b	Age	0.534 (0.516–0.552)	0.051 (0.046–0.057)
c	Manchester Triage System score	0.568 (0.558–0.578)	0.055 (0.050–0.059)
d	Vital signs	0.568 (0.565–0.572)	0.061 (0.057–0.065)
e	Hour of day	0.569 (0.563–0.575)	0.058 (0.054–0.062)
f	Pain score	0.572 (0.559–0.586)	0.058 (0.051–0.064)
g	Arrival mode	0.592 (0.576–0.608)	0.060 (0.055–0.064)
h	Triage discriminator	0.608 (0.583–0.633)	0.086 (0.062–0.110)
i	Triage complaint	0.640 (0.615–0.665)	0.091 (0.065–0.118)
j	Discharge diagnosis	0.641 (0.618–0.665)	0.090 (0.070–0.109)
k	Attendance complaint	0.645 (0.621–0.670)	0.092 (0.068–0.117)
l	30-day visit count	0.669 (0.653–0.684)	0.207 (0.149–0.264)
m	Condition count	0.692 (0.675–0.709)	0.107 (0.085–0.129)
n	Condition indicators	0.708 (0.694–0.723)	0.175 (0.134–0.216)
o	Top three features model (i, j, n)	0.736 (0.720–0.751)	0.242 (0.204–0.280)
p	Discharge model (b–n)	0.760 (0.746–0.774)	0.270 (0.207–0.332)

Table 3. Performance on the validation set for models using individual variables (models a–n) and sets of variables (models o and p). Bold text indicates models trained using multiple variable sets. Metrics are evaluated on the training set using grouped 5-fold CV at the patient level and we report the mean of the metric across the five validation folds. All models hyperparameters were tuned as described in the methods section to optimize the CV AUROC.

Assignment was chosen by visual inspection and finer grained cluster assignment can be achieved by tuning the hyperparameters of the DBScan algorithm.

Results

The results of our variable importance investigation are displayed in Table 3. The day of the week an attendance occurred was found not to be predictive of a patient's reattendance risk (model a, Table 3). Six variables (age, Manchester Triage System score, hour of day, vital signs, pain score, and arrival mode) were found to be weakly predictive of a patient's 72-h reattendance risk in isolation (AUROCs between 0.5 and 0.6, Table 3 models b–g). All other variables (Table 3 models h–n) were found to be moderately predictive (AUROC between 0.6 and 0.7) of 72-h reattendance risk in isolation.

Medical condition history was included in two representations. The count of the number of historical conditions (model m, Table 3) obtained a validation AUROC of 0.692 (95% CI: 0.675–0.709), reflecting that patients with a recorded medical history are more likely to reattend (8.3%, 95% CI 7.8–8.7%, reattendance rate) than those who do not (2.3%, 95% CI 2.1–2.5%, reattendance rate). When we included the full one-hot encoded matrix denoting whether the patient had a history of a specific condition, our model (model n, Table 3) obtained a validation AUROC of 0.708 (95% CI 0.694–0.723). The higher validation AUROC of the latter model suggests that different (medical) conditions are associated with differing degrees of reattendance risk.

The model that used the number of times a patient attended the emergency department in the 30-day prior to their current attendance (model l, Table 3) exhibited a validation AUROC of 0.669 (95% CI 0.653–0.684), agreeing with other studies that a patient's previous emergency department usage is an important consideration when considering their reattendance risk⁸. Three models (models i, j, and k, Table 3) make use of coded information describing the primary reason for the emergency department attendance, recorded at three distinct points in time and by potentially different members of clinical and non-clinical staff. Making use of the chief complaint collected at either the point of registration or Triage, respective validation AUROCs of 0.645 (95% CI 0.621–0.670) and 0.640 (95% CI 0.615–0.665) could be achieved. At the point of discharge, the recorded diagnosis obtained a validation AUROC of 0.641 (95% CI 0.618–0.665). This demonstrates that different diagnoses are associated with differing degrees of reattendance risk and indicates that a high-level, coded description of the patient's chief complaint is moderately predictive of reattendance risk, regardless of when it is recorded during the attendance.

Finally, models o and p in Table 3 present validation metrics for models trained using multiple variables. The model (model o in Table 3) using the three most predictive variables identified in our univariate importance study (Table 2) used the condition indicators, the condition count, and the number of times the patient visited the ED in the previous 30 days. We observed a validation AUROC of 0.736 (95% CI 0.720–0.751), demonstrating that models using multiple variables are more predictive of reattendance than a single variable. The model trained using all variables available at the point of discharge (model p in Table 3) increased the validation AUROC to 0.761 (95% CI 0.746–0.774).

The evaluation of our final model (model p) on the hold-out test set is presented in Fig. 2. This model obtained an AUROC and average precision of 0.747 (95% CI 0.722–0.773) and 0.233 (95% CI 0.194–0.277) respectively on the hold-out test set. This indicates that while the model retained its moderate performance, there was a reduction in model performance between the validation and the hold-out test set. This is likely the result of the

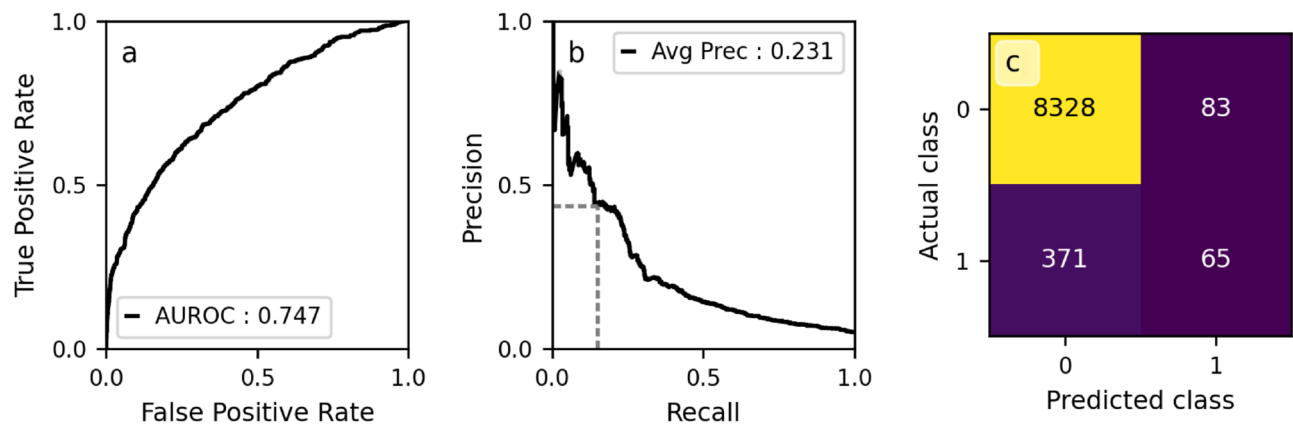


Figure 2. Performance of our model (model p in Table 3) evaluated on the hold-out test set. (a) Receiver operating characteristic curve for model's predictions. (b) Precision recall curve, the dashed grey line shows the configuration evaluated in the confusion matrix in panel (c). (c) Confusion matrix for predictions dichotomized using a threshold chosen such that the recall is equal to 0.15 (dashed grey line in panel (b)). A class of '1' indicates the patient reattended the emergency department within 72 h of discharge. Diagonal elements represent correct classifications and off-diagonal elements either False positives or negatives.

model overfitting the validation data (e.g., via the selection of hyperparameters to maximize model performance on the validation set), and is expected. In panel c of Fig. 2 we display a confusion matrix of our model at a single configuration, where the threshold for dichotomization of the continuous predictions output by our model was chosen such that a recall of 0.15 was obtained. This graphic presents the performance of our model on the hold-out test set when used as a binary decision support tool with a low-recall configuration.

Figure 3 summarizes the SHAP values of all attendances in the hold-out test set. In Fig. 3a the SHAP values for the ten most important variables (as determined by mean absolute SHAP values) are shown for each attendance. Visual inspection of this graphic provides insight into the trends our model has learned: the model associates anyone with a recorded medical condition as being at increased risk of reattendance and learned that some medical conditions represent a greater reattendance risk than others. For example, the model generally associates living alone with a higher reattendance risk than having a history of depression (the mean SHAP value is greater for those who live alone). In Fig. 3b,c, we plot the SHAP values across all attendances for two variables, the hour of day the attendance occurred and 30 day visit count. Visual inspection of panel b shows that the model associated the patients risk of reattendance with a periodic dependence on the hour of day the attendance occurred. By inspection of Fig. 3c, we can see that the model learned an approximately linear dependence between a patients reattendance risk and the number of times they have attended the emergency department in the last 30 days. It is important to note that these insights do not necessarily reflect actual risk factors for reattendance (since the model is an imperfect classifier) but only reflect the trends the model has learned to make its decisions.

The projection of all SHAP values for attendances in the hold-out test set into the two dimensional explanation space (see Methods) is displayed in Fig. 4. In Fig. 4a we colour attendances by the reattendance risk predicted by our machine-learning model. The reattendance risk (colour) is relatively uniform, reflecting that the majority of attendances do not result in reattendance (the observed reattendance rate across the test set is 4.9%). However, there are clear subgroups of the population identified to be at heightened risk of reattendance. In Fig. 4b we colour the attendances by the high-level cluster assignment obtained using the DBScan clustering algorithm. The characteristics of attendances in each cluster is displayed in Table 4. Visual inspection of Table 4 provides high-level insight into the logic of the machine-learning model. For example, for attendances assigned to cluster seven, on average, the most important variable for determining a patients reattendance risk is their 30-day visit count.

When evaluated on the readmission status (whether a patient reattended within 72 h *and* was subsequently admitted to hospital) of attendances in the hold-out test set our model achieved an AUROC of 0.804 (95% CI 0.774–0.832) and an average precision of 0.044 (95% CI 0.030–0.063); reflecting that the model demonstrated moderate discriminative ability in identifying the subset of reattendances which resulted in admission.

Discussion

Our final 72-h reattendance risk model achieved an AUROC of 0.747 (95% CI 0.722–0.773) and an average precision of 0.233 (95% CI 0.194–0.277) on the hold-out test set. Qualitatively, our model can use variables describing an emergency department attendance and a view of a patient's medical history to predict their reattendance risk with moderate performance. We calculated SHAP values to provide an explanation of our predictions at an attendance level (Fig. 3 and Supplementary Fig. 2) and projected these explanations into a two-dimensional space (Fig. 4). We used this low-dimensional embeddings to identify different patient sub-groups at risk of reattendance.

Our final model (model p in Table 3) makes use of all variables in our dataset available at the point of discharge. This set of variables is not necessarily the optimal set of variables for a reattendance predictor, it is plausible the variable set contains obsolete information because of correlations between variables. High correlation

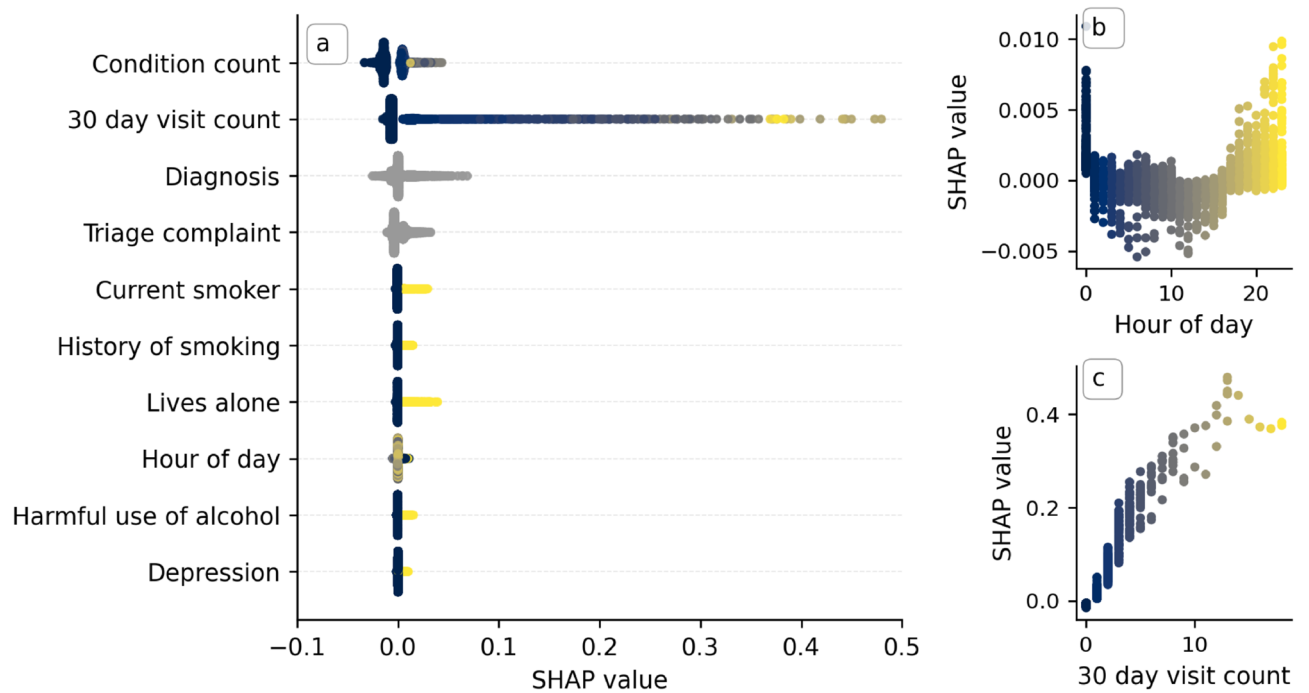


Figure 3. SHAP values for each attendance in the hold-out test set. **(a)** Plot summarizing the SHAP values for the ten most important variables (by mean absolute SHAP value) for each attendance in the test set. They are ordered by the global impact the feature has on the explanation (equal to the mean absolute SHAP value of the feature across all attendances). For the binary variables (i.e., the condition indicators) this favours variables with a high number of occurrences (i.e., more common conditions), not necessarily those associated with the high reattendance risk. **(b)** SHAP value against recorded hour of day of attendance registration (dots). **(c)** SHAP value against number of emergency department visits in the 30 days prior to the given attendance (dots). In panels **(b)** and **(c)** vertical dispersion is the result of interaction with other variables. All panels are coloured by the magnitude of the respective variable for the given data point, with lighter colours indicating higher values (e.g., inspect panels **(b)** and **(c)**). Grey data points correspond to non-binary, nominal categorical variables and therefore have no natural ordering which could be used to colour the data points.

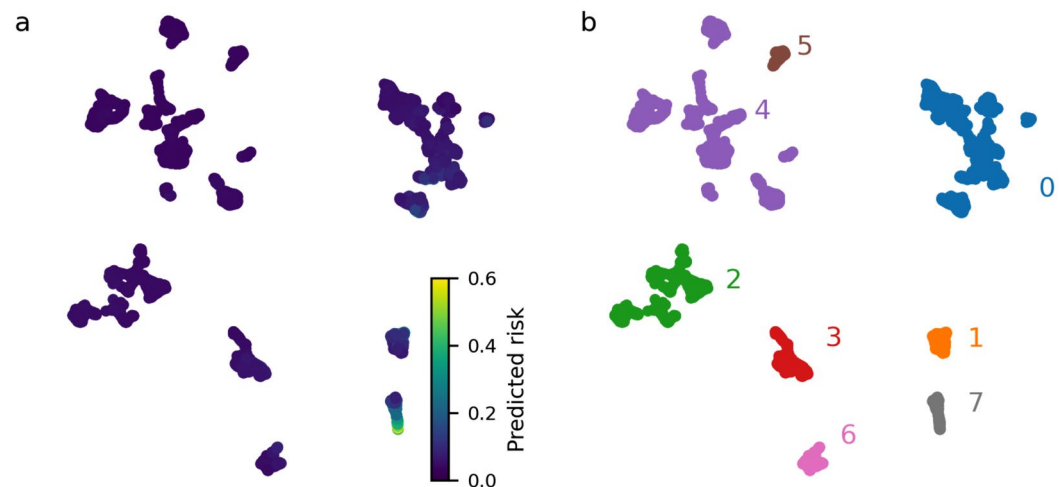


Figure 4. Embedding of the hold-out test set into a two-dimensional ‘explanation’ space using the UMAP algorithm. **(a)** All attendances in the hold-out test set visualised in the explanation space, colour indicates the predicted reattendance risk for the respective attendance. **(b)** Attendances in the explanation space coloured by cluster assignment. Cluster assignment was performed with the DBScan algorithm²⁸. Assignment was chosen by visual inspection and finer grained cluster assignment can be achieved by tuning the hyperparameters of the DBScan algorithm. The explanation embedding was created by clustering the SHAP values (explanations) for each attendance using the UMAP algorithm. Closer data points share a more similar explanation for their predicted reattendance risk. Descriptive properties of each cluster are displayed in Table 4.

Cluster	Count	Age (SD)	Reattendance rate (95% CI)	Predicted reattendance rate (95% CI)	Condition count	30-day visit count	Most important variables (mean absolute SHAP values)		
							1st	2nd	3rd
0	2723	55.4 (22.6)	5.9 (5.1–6.9)	6.8 (5.0–8.7)	2.9	0.0	30-day visit count	Triage complaint	Condition count
1	446	52.8 (21.4)	8.5 (6.3–11.5)	10.6 (7.1–14.0)	3.5	1.0	30-day visit count	Medical history	Triage complaint
2	1363	42.5 (19.2)	2.6 (1.9–3.6)	4.6 (4.4–5.0)	0.0	0.0	Condition count	30-day visit count	Triage complaint
3	671	40.9 (18.5)	4.0 (2.8–5.79)	5.6 (5.0–6.3)	0.0	0.0	Condition count	Diagnosis	30-day visit count
4	2761	41.5 (16.9)	1.8 (1.3–2.3)	3.8 (3.6–4.0)	0.0	0.0	Condition count	30 day visit count	Triage complaint
5	221	40.5 (15.9)	0.0 (0.0–1.7)	3.6 (3.4–3.7)	0.0	0.0	Condition count	30-day visit count	Triage complaint
6	321	39.7 (18.1)	5.6 (3.6–0.9)	5.9 (4.7–7.0)	0.0	1.0	Condition count	30-day visit count	Triage complaint
7	341	42.3 (17.4)	31.1 (26.4–36.2)	27.2 (13.6–40.7)	3.3	3.75	30-day visit count	Medical history	Condition count

Table 4. Properties of the attendances assigned to each of the explanation clusters (Fig. 4). The count column displays the number of attendances in a given cluster. The age (in years), reattendance rate, predicted reattendance rate, condition count, and 30-day visit count column display the mean of the respective variable for all attendances in the cluster. Values in brackets in the age column display the standard deviation patient age in the respective cluster. The final three columns display the most important variables in making a decision, averaged across all attendances in a given cluster.

between variables is expected for clinical data. For example, it is likely that patient age, arrival mode, and vital signs all latently encode patient frailty, which is known to be related to a patient's reattendance risk²⁹. Despite this, our predictive algorithm of choice, extreme gradient boosting decision trees, is relatively robust against correlations between variables and high variable redundancy is unlikely to be significantly detrimental to model performance.

In our exploratory analysis, we found that the hour of day the attendance began correlates to the reattendance rate, with higher reattendance rates observed during the night (Supplementary Fig. 1). By inspection of the observed SHAP values for the hour of day (Fig. 3b), we observe that our model has learned a similar trend, associating attendance registration during the night with an increased (between zero and one percent increase) reattendance risk. This trend could have several different origins. Firstly, we have found that the hour of day displays correlation with the reason for attendance, with either complaints associated with a higher (lower) risk of 72-h reattendance more (less) likely to present during the night (not shown). Secondly, it is plausible that staff fatigue and lower staffing levels may contribute to the increased reattendance rate for attendances occurring during the night, although we have no way of testing this hypothesis in our dataset.

Our analysis (e.g., Fig. 2 and model e in Table 3) shows that certain complaints are associated with a higher risk of 72-h reattendance. For example, attendances whose chief complaint at registration was 'abdominal pain' had a mean 72-h reattendance rate of 6.8% (95% CI 6.0–7.7%), compared to the mean reattendance rate of 4.8% (95% CI 4.6–5.0%) in the training set. Coding compliance was not evaluated in our dataset, which may affect this observation. For example, the most common chief complaint at registration was 'unwell adult' (with 20.5% of attendances listing this as the chief complaint in the training set) which is utilized when either the chief complaint is not clear at registration, when the patient presents with multiple complaints or as a result of inappropriate coding.

In addition to identifying complaints associated with a heightened short-term reattendance risk, our model also makes use of ICD10 coded conditions extracted from a patient's electronic health record. These variables allow the model to identify medical conditions, comorbidities, and risks which are associated with increased reattendance risk and enables models to achieve moderate predictive performance (Table 3). Excluding the medical condition indicators, the most important feature is the 30-day visit count which, in part, reflects the disproportionate use of EDs by frequent users³⁰. In the visualisation of the attendances in the patient hold-out test set in the two-dimensional explanation space (Fig. 4), the most frequent attenders (30 day visit count of two or more) are clearly segregated (cluster 7 in Fig. 4b and Table 4). Visual inspection of the properties of each cluster in Fig. 4 could be used to help guide the design of interventional strategies. For example, the reattendance risk of patients within cluster 7 (i.e., frequent attenders) could be reduced by increased provision of community support. Conversely, such support is unlikely to be the most appropriate intervention strategy for patients suffering from an acute injury (such as a severe burn) associated with increased reattendance risk.

From a clinical perspective it is interesting to investigate the subset of reattendances which are also readmissions (i.e., reattendances to the emergency department which result in subsequent admission to an inpatient ward). In these cases, there is increased risk that there was missed critical illness or injury at the initial attendance and these readmissions are important to evaluate for clinical assurance purposes. Overall, 37.1% of reattendances end in readmission, resulting in a 72-h readmission rate of 2.0%. Evaluating our models predictions with a target equal to whether the patient readmitted in 72 h, we find it has an AUROC of 0.804 (95% CI 0.774–0.832) and an average precision of 0.044 (95% CI 0.030–0.063) on the hold-out test set. The high AUROC means the model displays high discernibility between attendances which result in readmission and those that do not. The low average precision reflects that readmissions only make up a minority fraction of reattendances and the false positive rate increases as a result of the large class imbalance. Overall, these results demonstrate our model can identify the subset of reattendances which are also readmissions with a similar predictive performance as reattendances

which do not result in admission, an important result since these two different outcomes will require different interventional strategies to reduce the risk of reattendance/readmission.

A limitation of our study, shared with other investigations of machine-learning use in EDs³¹, is that its primary data source is structured past medical histories. A medical history is not available for all patients and this could lead to our model discriminating against these patients. An example of this bias can be observed for cluster two (Table 3), where the two most important variables for determining a patient's reattendance rate is the absence of any visit to the emergency department in the last 30 days and the absence of any recorded medical conditions. Unless implemented with consideration, using such a model could have adverse impact on patient care, because the model may not have sufficient descriptive information about a patient to make a reliable prediction of their reattendance risk. Patients who have several long-term health conditions, but have no past medical history with the hospital are particularly likely to be disadvantaged by this. We partially mitigate this bias by providing our model with visit-level information. In the future this bias could be reduced further by linking to community datasets (e.g., GP records) to obtain a view of a patient's medical history, by providing confidence intervals alongside risk scores, and by educating clinicians on the limitations of the machine-learning model. In a deployment scenario, additional improvements could be achieved by using the model as an alert tool. In this setting, model predictions would only be presented to clinical staff for a small subset of patients; those the model predicts to be at particularly high risk of reattendance and for whom a decision to discharge has been made. Otherwise, the model would be invisible to clinical staff who would be free to carry out standard clinical practice in cases where an alarm is not raised.

Since our model only uses information available to clinicians at the time of the emergency department attendances it has a relatively low barrier to implementation. Despite this, it will be essential to perform prospective, randomized clinical trials of any implementation, investigating the efficacy of the model and the associated interventions. Ultimately, deployment of a machine-learning model could eventually invalidate the model by changing the behaviours and descriptors of reattendances by altering the clinical decisions made. In the short term, a relatively low-risk implementation of a model trained to identify patients at risk of reattendance would be in the implementation of a low-recall and high-precision alert system (for example, the configuration presented in Fig. 2c). This would only raise alarms for the cases the model believes are at the highest risk of reattendance and highlight the need for additional clinical review. On average, using the configuration displayed in Fig. 2c, this would have raised an alarm for only 1.7% of attendances in which a decision to discharge was made (approximately 2 times per day) and would expect to be correct approximately 44% of the time—mitigating the risk of alarm fatigue. A model deployed in this manner would be of limited impact (because of its low recall), but the configuration could be re-evaluated if model performance improves. Performance could be improved by the inclusion of free-text notes recorded by clinical staff, which previous studies have shown can be predictive of patient outcomes across the broader hospital network^{14,15}.

It is important to discuss the context in which our model could be prospectively deployed. Our model was trained and retrospectively evaluated using data available to clinicians during standard clinical practice at the emergency department in Southampton. This is an advantage if the model was to be used at this location because the characteristics of future attendances will likely reflect the attendances the model was trained on. Conversely, this means that the model will not necessarily generalize to other EDs without first training on their local data. Poor model generalization will be particularly prominent in EDs with a catchment zone with different demographics to Southampton and, therefore, different disease and illness prevalences. Because our model contains variables either in the standard UK emergency care dataset or regularly available to EDs nationally, it would be possible to evaluate this model directly in other EDs. External validation of our model is essential before prospective deployment beyond the department at which the training data was sourced.

Conclusion

In conclusion, we have constructed and retrospectively evaluated an extreme gradient boosting decision tree model capable of predicting the 72-h reattendance risk for a patient at the point of discharge from an emergency department. The highest performing model achieved an AUROC of 0.747 (95% CI 0.722–0.773) and an average precision of 0.233 (95% CI 0.194–0.277) on a hold-out test set. We investigated the variables most indicative of risk and showed these were patient level factors (medical history) rather than visit level variables such as recorded vital signs. We calculated SHAP values to explain our predictions (Fig. 3) and used these to investigate the trends our model learned. We projected explanations into a (2D) explanation space using the UMAP algorithm and used this to explore the different sub-groups at risk of reattendance. We suggested an implementation of the algorithm in a low-recall, high-precision configuration such that alarms are only raised for patients predicted to be at a (clinically defined) heightened risk of reattendance. External validation and prospective clinical trials of our models is essential, with considerable consideration given to the planned interventions and the impact this would have on clinical decisions.

Data availability

Due to patient privacy concerns the dataset used in this study is not publicly available. However, it will be made available upon reasonable request.

Received: 4 February 2021; Accepted: 20 October 2021

Published online: 02 November 2021

References

1. Berchet, C. *Emergency Care Services: Trends, Drivers and Interventions to Manage the demand.* (2015).

2. Baier, N. *et al.* Emergency and urgent care systems in Australia, Denmark, England, France, Germany and the Netherlands—Analyzing organization, payment and reforms. *Health Policy* **123**, 1–10 (2019).
3. Bernstein, S. L. *et al.* The effect of emergency department crowding on clinically oriented outcomes. *Acad. Emerg. Med.* **16**, 1–10 (2009).
4. Guttman, A. *et al.* Association between waiting times and short term mortality and hospital admission after departure from emergency department: population based cohort study from Ontario, Canada. *BMJ* **342**, d2983 (2011).
5. Besga, A. *et al.* Risk factors for emergency department short time readmission in stratified population. *BioMed Res. Int.* <https://doi.org/10.1155/2015/685067> (2015).
6. Deschodt, M., *et al.* Characteristics of older adults admitted to the emergency department (ED) and their risk factors for ED reattendance based on comprehensive geriatric assessment: A prospective cohort study. *BMC Geriatr.* **15**, 1 (2015).
7. Martin-Gill, C. & Reiser, R. C. Risk factors for 72-hour admission to the ED. *Am. J. Emerg. Med.* **22**(6), 448–453 (2004).
8. Arendts, G. *et al.* Derivation of a nomogram to estimate probability of revisit in at-risk older adults discharged from the emergency department. *Intern. Emerg. Med.* **8**(3), 249–254 (2013).
9. Kun-Hsing, Y., Beam, A. L. & Kohane, I. S. Artificial intelligence in healthcare. *Nat. Biomed. Eng.* **2**(10), 719–731 (2018).
10. Moncada-Torres, A., van Maaren, M. C., Hendriks, M. P., Siesling, S. & Geleijnse, G. Explainable machine learning can outperform Cox regression predictions and provide insights in breast cancer survival. *Sci. Rep.* **11**(1), 1–13 (2021).
11. Li, Y. *et al.* An Interpretable Machine Learning Survival Model for Predicting Long-term Kidney Outcomes in IgA Nephropathy. *In AMIA Annual Symposium Proceedings* 737 (2020).
12. Hao, S. *et al.* Risk prediction of emergency department revisit 30 days post discharge: A prospective study. *PLoS ONE* **9**(11), e112944 (2014).
13. Hong, W. S., Haimovich, A. D. & Taylor, R. A. Predicting 72-hour and 9-day return to the emergency department using machine learning. *JAMIA Open* **2**(3), 346–352 (2019).
14. Huang, K., Altsaer J. and Ranganath R. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *arXiv*, [arXiv:1904.05342](https://arxiv.org/abs/1904.05342) (2019).
15. Sterling, N. W., Patzer, R. E., Di, M. & Schragger, J. D. Prediction of emergency department patient disposition based on natural language processing of triage notes. *Int. J. Med. Inform.* **129**, 184–188 (2019).
16. Chen, T. & Guestrin, C. XGBoost: A scalable tree boosting system. *In Proceedings of the 22nd ACM sigkdd international conference on knowledge discovery and data mining* 785–794 (2016).
17. Bergstra, J. S., Bardenet, R., Bengio, Y. & Kégl, B. Algorithms for hyper-parameter optimization. *In Advances in neural information processing systems* 2546–2554 (2011).
18. Bergstra, J., Yamins, D. & Cox, D. D. Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures. *In International Conference on Machine Learning* 115–123 (2013).
19. Snoek, J., Larochelle, H. & Adams, R. P. Practical Bayesian optimization of machine learning algorithms. *Adv. Neural Inf. Process. Syst.* **25**, 1–9 (2012).
20. Micci-Barreca, D. A preprocessing scheme for high-cardinality categorical attributes in classification and prediction problems. *ACM SIGKDD Explor. Newsl.* **3**(1), 27–32 (2001).
21. Chmiel, F. P., Burns, D. K. <https://git.soton.ac.uk/it-innovation-public-projects/triaged-readmission-model> (2021).
22. Lundberg, S. M. & Su-In, L. A unified approach to interpreting model predictions. *In Advances in neural information processing systems* 4765–4774 (2017).
23. Lundberg, S. M. *et al.* From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.* **2**(1), 2522–5839 (2020).
24. Lundberg, S. M. *et al.* Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. *Nat. Biomed. Eng.* **2**(10), 749–760 (2018).
25. McInnes, L., Healy, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, *arXiv* [2arXiv:1802.03426](https://arxiv.org/abs/1802.03426) (2018).
26. Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* **37**(1), 38–44 (2019).
27. Diaz-Papkovich, A., Anderson-Trocmé, L., & Gravel, S. UMAP reveals cryptic population structure and phenotype heterogeneity in large genomic cohorts. *PLoS Genet.* **15** 11, e1008432 (2019).
28. Ester, M., Kriegl, H. P., Sander, J. & Xu, X. A density-based algorithm for discovering clusters in large spatial databases with noise. *In Kdd* **96**(34), 226–231 (1996).
29. Kahlon, S. *et al.* Association between frailty and 30-day outcomes after discharge from hospital. *Cmaj* **187**(11), 799–804 (2015).
30. LaCalle, E. & Rabin, E. Frequent users of emergency departments: The myths, the data, and the policy implications. *Ann. Emerg. Med.* **56**(1), 42–48 (2010).
31. Joseph, J. W. *et al.* Deep-learning approaches to identify critically ill patients at emergency department triage using limited information. *JACEP Open* **1**, 773–781 (2020).

Acknowledgements

This work was supported by The Alan Turing Institute under the EPSRC Grant EP/N510129/1. We acknowledge support from the NIHR Wessex ARC.

Author contributions

F.P.C. and D.K.B. performed the data analysis and modelling. Z.D.Z. discussed and commented on the analysis with F.P.C. and D.K.B. N.W. and F.P.C. obtained governance and ethical approval. M.A. and F.B. created the data extract. F.P.C., M.J.B. and N.W. managed the study at the UoS. M.K. managed the study at UHS. F.P.C. and M.K. designed the study with assistance from T.W.V.D. M.K. and T.W.V.D. provided clinical guidance and insight. F.P.C. wrote the first draft of the manuscript with assistance from M.K. and T.W.V.D. All authors frequently discussed the work and commented and contributed to future drafts of the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-00937-9>.

Correspondence and requests for materials should be addressed to F.P.C.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021