

Identifying those at risk of reattendance at discharge from emergency departments using explainable machine learning

F. P. Chmiel^{1,*}, M. Azor², F. Borca^{2,3}, M. J. Boniface¹, D. K. Burns¹, Z. D. Zlatev¹, N. M. White¹, T. W. V. Daniels^{4,5}, and M. Kiuber⁶

¹School of Electronics and Computer Science, University of 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, Minerva House, University Hospital Southampton, UK

⁵School of Clinical and Experimental Sciences, Faculty of Medicine, University of Southampton, Southampton

General Hospital, LF13A, South Academic Block, Southampton, UK

⁶Emergency Department, University Hospital Southampton NHS Foundation Trust, Southampton, UK

*francispeterchmiel@gmail.com

ABSTRACT

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 create and evaluate a machine-learned classifier trained to identify patients at risk of reattendance within 72 hours of discharge from an emergency department. On a patient hold-out test set, our highest performing classifier obtained an AUROC of 0.749 and an average precision of 0.232; demonstrating that machine-learning algorithms can be used to classify patients, with moderate performance, into low and high-risk groups for reattendance. In parallel to our predictive model we train an explanation model, capable of explaining the predictions of the machine-learned classifier at an attendance level. These explanations can be used to help understand the decisions a model is making, evaluate biases present in its decisions and help inform the design of bespoke interventional strategies.

Introduction

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. Since delays in care and overcrowding have been linked to increased rates of adverse outcomes^{3,4}, it is important to investigate the most efficient ways of using the available resources and, importantly, minimise and mitigate their unnecessary use. One way this can be achieved is by the minimisation of short-term reattendances, which describe a situation where a patient presents to an emergency department within 72 hours 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 those in need of additional support in the community following discharge from an ED. Predictive models, available as a decision support tool at the point of discharge, able to reliably identify those at increased risk of short-term reattendance using known risk factors and attendance level information, may be able to significantly reduce the number of reattendances by appropriately quantifying and explaining a patient's risk of reattendance to clinical staff. Ultimately this would allow appropriate intervention (e.g., further diagnostic tests), more informed discussions about a patient's discharge plan, or support in the community for those recently discharged.

Machine learnt models are a class of predictive models which are particularly well positioned to add value to emergency department processes. By making use of large amounts of clinical and administrative data, these models can provide estimates of a patient's short-term reattendance risk^{9,10} and explain the reason for the patient's predicted risk. Explanation is particularly

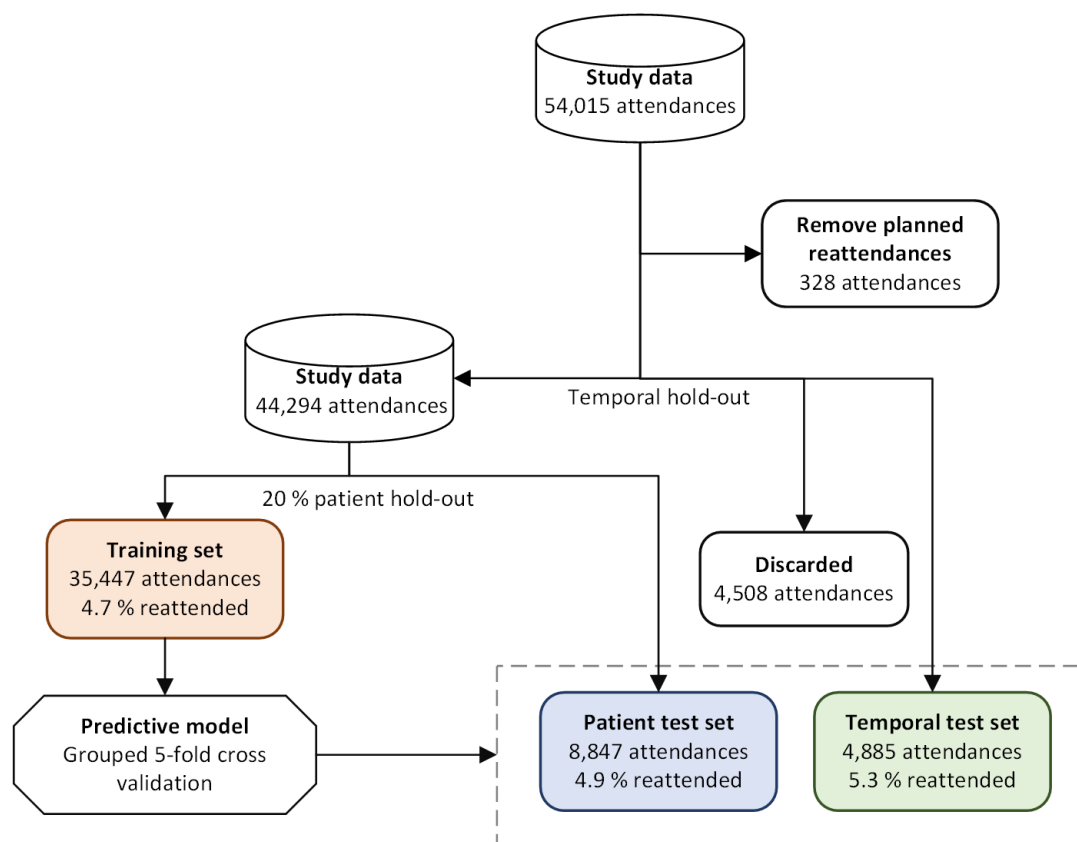


Figure 1. Segregation of the study data into training and the two hold-out test sets. Discarded attendances were those that occurred in either the first 30 days or last 72 hours of the temporal test, to avoid information leakage between the training and temporal test set and because the reattendance status could not be robustly calculated for attendances occurring in the last 72 hours of the dataset. Reattendance rates (bottom row of shaded boxes) display the observed 72-hour reattendance rate for each cohort.

37 important, as this could help either inform the patients care trajectory or guide the post-discharge intervention plan. In
38 this manuscript we discuss a machine-learnt model, utilizing historical (coded, inpatient) discharge summaries, alongside
39 contemporary clinical data recorded during emergency department attendances such as observations and the results of standard
40 triage processes, to identify patients at increased risk of short-term reattendance following an emergency department attendance.
41 In addition to our predictive model, we construct an explanation model which enables us to evaluate the trends our model has
42 learned and explain our model's prediction at an attendance level.

43 Methods

44 Dataset curation

45 The dataset features a pseudonymized version of all attendances by adults to Southampton's Emergency Department (University
46 Hospitals Southampton Foundation Trust) occurring between the 1st April 2019 and the 30th of April 2020. For our study
47 cohort, we take only attendances which resulted in discharge directly from the ED, of which there were 54,015. The core
48 dataset includes patients' year of birth, results of any near-patient observations, and high-level information about the attendance
49 included in the standard UK Emergency Care Data Set (e.g., outcome, arrival mode, duration of visit, and chief complaint). To
50 provide the machine learning classifier with a view of patients' medical history we make use of historical discharge summaries
51 associated with the patient, both from the emergency department and from the patients electronic health record maintained by
52 the University Hospitals Southampton Foundation Trust. For a given patient, from any discharge summary occurring prior to a
53 given emergency department attendance, we make use of (ICD10) coded conditions (e.g., type 2 diabetes, current smoker) and
54 create a binary indicator which indicates whether a patient has a given condition coded in their electronic health record prior to
55 a given ED attendance. The electronic health records used by our models are available to review by clinicians and are used in
56 regular practice. Our model does not have access to any free text fields in the electronic health record, where as a clinician

Condition	Attendances with condition	Fraction of attendances (%)
Hypertension	3,801	10.7
Depression	2,725	7.7
History of smoking	2,680	7.6
Asthma	2,608	7.4
Current Smoker	2,590	7.3
Type 2 diabetes	1,528	4.3
Hypercholesterolaemia	1,156	3.3
Harmful use of alcohol	1,072	3.0

Table 1. Most frequently occurring conditions for attendances in the training set. The left column denotes the noted conditions (as specified by ICD10 codes) and the right column the number of attendances in the training set noted to have this condition. 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.

would. Previous studies have shown that (free text) clinical notes can be predictive of patient outcomes across the broader hospital network^{11,12}, but including these notes was beyond the scope of our study as they would limit the explainability of our algorithms. An example of the most frequently observed conditions are presented in Table 1.

Reattendance identification

Patient reattendances are identified by using the patient pseudo identifier to calculate the time to their next ED attendance. All reattendances (with the exception of planned reattendances, Figure 1) are considered, even if the second attendance is for a different condition to the original attendance. This is then dichotomized (return within 72 hours from the point of discharge) to annotate the reattendance state of each attendance. This formulation allowed us to frame the predictive task as a binary classification problem.

Predictive modelling

We separated our data into a training set and two independent test sets (Figure 1). The last 3 months of attendances (01/02/2020 to 30/04/2020, inclusive of the COVID-19 pandemic) were segregated as a temporal test set, excluding any visit which took part in either the first 30 days or the final 72 hours. These exclusions remove information leakage between the training and temporal test set and remove attendances where reattendance to the emergency department could not be calculated reliably (i.e., attendances which occur in the last 72 hours of the data extract). The remaining attendances (N=44,294) were randomly split at the patient level to create a patient-level hold test set containing attendances from 20 % of the remaining patients. Remaining attendances (N=35,447) were used as the training and validation set. The number of patients in each respective dataset was 4,464, 7,237, and 28,945. The relation between patients in each dataset is displayed in Supplementary Figure 1, demonstrating patient exclusivity between the training set and the patient hold-out test set. Results of our predictive model, evaluated on the temporal test set is discussed in the Supplementary Information.

As our machine-learned classifier, we used a gradient boosted decision tree as implemented in the XGBoost framework¹⁶. Features used in modelling include : patient age (estimated from year of birth), number of emergency department attendances in the 30 days prior to the attendance, the chief complaint of the attendance (e.g., ‘abdominal pain’), the patients mode of arrival, previously described medical condition indicators, the count of the number of medical conditions a patient has, vital signs (temperature, pulse and respiration rate, systolic blood pressure, and blood oxygen saturation levels), the Manchester Triage System score, triage pain score, (coded) discharge diagnosis, and the hour of day and day of the week the attendance occurred. A full data schema is presented in Supplementary Table 1. Medical conditions associated with the patient at a given attendance 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 encoding¹³. Model hyperparameters were optimized using five-fold cross-validation (CV) of the training set at the patient level (the set of attendances from a unique patient appear exclusively in the validation or training set for each fold) using Bayesian optimization utilizing the Tree Parzen Estimator algorithm as implemented in the hyperopt Python library^{14,15}.

We evaluate our final model performance (the mean output of the five models trained during cross-validation) on the two hold-out test sets. Models performance is evaluated using the Area Under the Receiving Operating Curve (AUROC) and the average precision under the precision-recall curve.

Model reference	Variables	Validation AUROC	Validation average precision
a	Day of week	0.500	0.047
b	Age	0.534	0.051
c	Manchester Triage System score	0.567	0.055
d	Hour of day	0.569	0.058
e	Vital signs	0.571	0.061
f	Pain score	0.571	0.058
g	Arrival mode	0.592	0.060
h	Triage discriminator	0.606	0.085
i	Triage complaint	0.640	0.092
j	Discharge diagnosis	0.642	0.090
k	Attendance complaint	0.647	0.089
l	30-day visit count	0.669	0.209
m	Condition count	0.692	0.107
n	Condition indicators	0.706	0.160
o	Top three feature model (i, j, n)	0.741	0.261
p	Discharge model (b-n)	0.761	0.260

Table 2. Performance on the validation set for models using individual features (models a-n) and sets of features (models o and p). 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.

Model explainability

To explain the predictions of our model we make use of the TreeExplainer algorithm implemented in the SHAP Python library¹⁷⁻¹⁹. TreeExplainer calculates SHAP values (i.e. Shapley values), a concept from coalitional game theory which treats predictive variables as players in a game and distributes their contribution to the predicted probability. To calculate the SHAP value for a given feature, one trains a model for each possible feature set (with and without the given feature) and calculates the mean change in the predicted probability when the feature was added to a feature set for all possible sets of features. This mean change is the SHAP value and can be negative (adding the feature predicted reduces reattendance risk) or positive (adding the feature increases the predicted reattendance risk). 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 explanations¹⁹. Practically, for each attendance we will have a scalar value for each variable used by the model which quantifies the contribution that variable had on the predicted reattendance risk for the given attendance, with SHAP values of larger magnitudes indicating that the relevant variable was of greater importance in determining the predicted reattendance risk.

To investigate the different explanations across the whole dataset, we project the SHAP values for all attendances into a two-dimensional ('explanation') space using Uniform Manifold Approximation and Projection (UMAP)²⁰. 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^{21,22}. Attendances which are closer in proximity in this two-dimensional space share a more similar explanation for their predicted reattendance risk.

Ethics and data governance

This study was approved by the University of Southampton's Ethics and Research governance committee (ERGO/FEPS/53164) and approval was obtained from the Health Research Authority (20/HRA/1102). 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.

Results

To investigate the potential of individual and sets of variables at predicting 72-hour reattendance we constructed a series of XGBoost models, evaluating their performance on the training set using five fold cross-validation as described in the Methods section, the results of this experiment are displayed in Table 2.

Six variables (age, Manchester Triage System score, hour of day, vital signs, pain score, and arrival model) were found to be weakly predictive of a patient's 72-hour reattendance risk in isolation (AUROCs between 0.5 and 0.6, Table 2 models b-g). All other variables (Table 2 models h-n) were found to be moderately predictive (AUROC between 0.6 and 0.70) of 72-hour

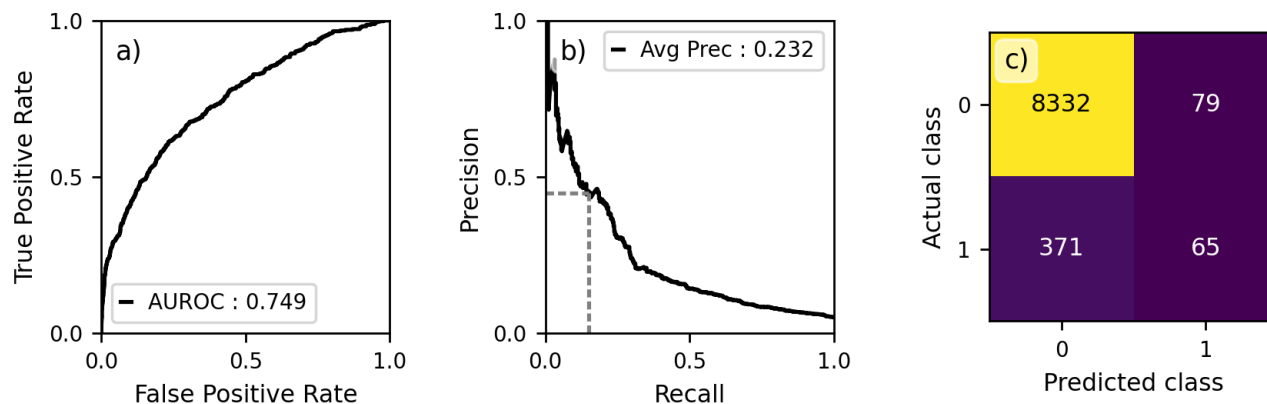


Figure 2. Performance of our classifier (model p in Table 2) evaluated on the patient hold-out test set. a) Receiver operating characteristic curve for model’s prediction. b) Precision recall curve for predictions on test set, the dashed grey line shows the configuration evaluated in the confusion matrix 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 hours of discharge. Diagonal elements represent correct classifications and off-diagonal elements either False positives or negatives.

121 reattendance risk in isolation, with the exception of the day of the week the attendance occurred which was not predictive of a
122 patients reattendance risk (model a, Table 2).

123 Patients conditions were included in two representations. The count of the number of historical conditions (model m, Table
124 2) obtained a validation AUROC of 0.692, reflecting that those with a recorded medical history with the emergency department
125 or the associated hospital are more likely to reattend (8.3 % (95 % CI: 7.8-8.7 %) reattendance rate) than those who do not
126 (2.3% (95 % CI: 2.1-2.5%) reattendance rate). When we included the full one-hot encoded matrix denoting whether the patient
127 had a history of a specific condition, our model (model n, Table 2) obtained a validation AUROC of 0.706 – higher than when
128 our model used just the number of historical conditions. This suggests that different (medical) conditions are associated with a
129 differing degrees of reattendance risk.

130 The model that used the number of times a patient attended the emergency department in the 30-day prior to their current
131 attendance (Table 2, model l) exhibited a validation AUROC of 0.669, agreeing with other studies that a patients previous
132 emergency department usage is an important consideration when considering their reattendance risk⁸. Three models (models
133 i, j, and k, Table 2) make use of coded information describing the primary reason for the emergency department attendance,
134 recorded at three distinct timepoints and by potentially different members of clinical and non-clinical staff. Making use of
135 the chief complaint collected at either the point of registration or Triage, respective validation AUROCs of 0.640 and 0.647
136 could be achieved. At the point of discharge, the recorded coded diagnosis obtained a validation AUROC of 0.642. This
137 demonstrates that different diagnoses are associated with differing degrees of reattendance risk and indicates that a high-level,
138 coded description of the patients chief complaint is moderately predictive of reattendance risk, regardless of when it is recorded
139 during the attendance.

140 Finally, we investigated models which utilize multiple variables (models o and p in Table 2). Firstly, we trained a model
141 using only the three variables which were most predictive as determined by our univariate feature importance investigation
142 (Table 2). This model (model o in Table 2) used the condition indicators, the chief complaint recorded at triage, and the number
143 of times the patient visited the ED in the previous 30 days. Ultimately, it obtained a validation AUROC of 0.741, demonstrating
144 that using multiple variables is more predictive of reattendance than a single variable. We also trained and evaluated a model
145 using all features available at discharge, which used all variables available at the time of discharge and increased the validation
146 AUROC to 0.761, but with a small decrease in the validation average precision.

147 Next, we applied our final model (model p, Table 2) on the patient wise hold-out test set, the evaluation of which is presented
148 in Figure 2. The AUROC and average precision was 0.749 and 0.232; while the model obtained its moderate performance there
149 was a reduction in performance between the validation and patient wise hold-out test set. This could arise because of a degree
150 of overfitting of the model to the validation data (via the selection of hyperparameters which optimize model performance on
151 the validation set) or because of a different set of patients in the hold-out test set. To demonstrate the evaluation of our model
152 as a binary decision support tool, we display a confusion matrix for our classifier at a single configuration in Figure 2 c. The
153 threshold for dichotomization of the predictions was chosen such that a recall of 0.15 was obtained.

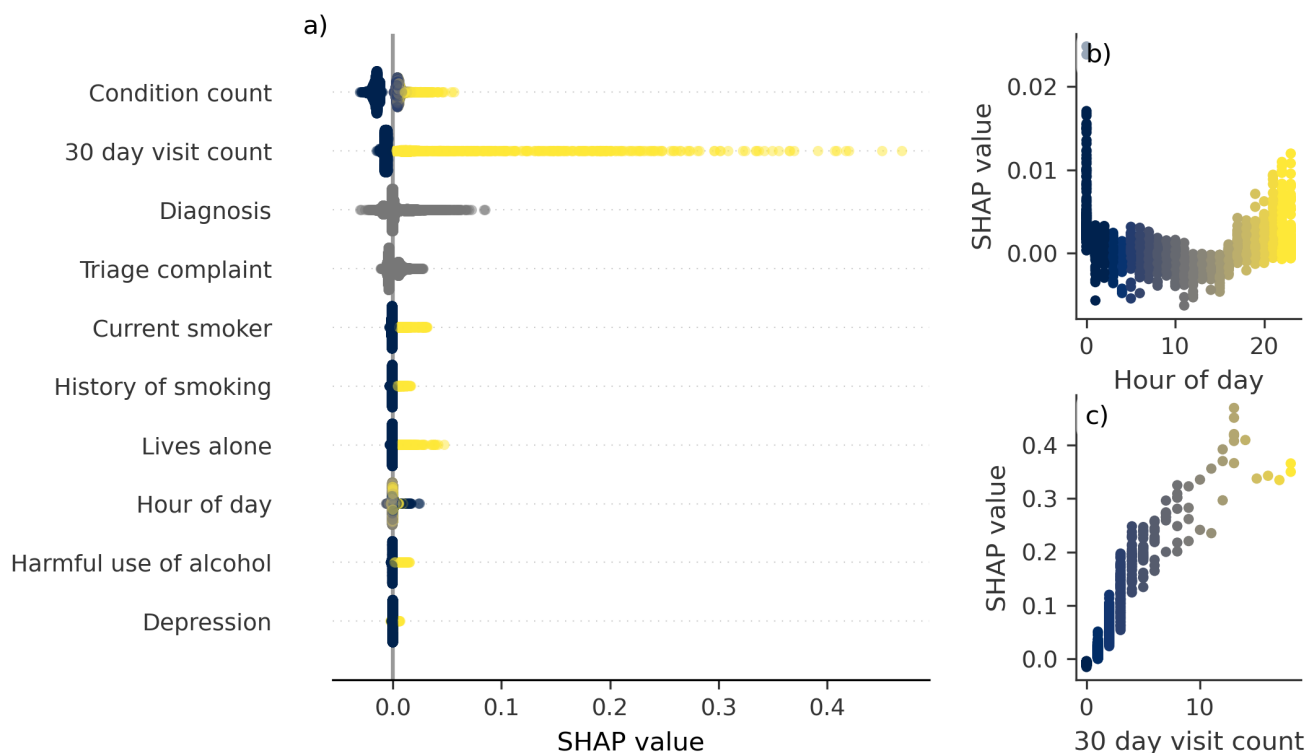


Figure 3. Explanations of model predictions using SHAP for each attendance in the patient hold-out test set. a) Plot summarizing the SHAP values for ten variables for each patient in the patient test set. They are ordered by the global impact the feature has on the explanation (practically, equal to 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 which have the highest reattendance risk. b) SHAP value against recorded hour of day for time of registration for a given attendance (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 in the feature set. 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 categorical variables.

154 To investigate the trends our model has learned we made use of the TreeExplainer algorithm¹⁸; a demonstration of the global
155 explanation of our reattendance model is presented in Figure 3. In Figure 3 a the SHAP values (quantifying, at an instance
156 level, the impact a given variable has on the model's prediction) for 10 variables are shown for each attendance. Looking at this
157 plot for a large number of variables allows a high level understanding of the model to be obtained: the model associates anyone
158 with a recorded medical condition as being at increased risk of reattendance and learns that some medical conditions represent
159 a greater reattendance risk than others, for example living alone is generally associated with a higher reattendance risk than
160 having a history of depression (the mean SHAP value is greater for those who live alone). In panels b and c of Figure 3, we
161 plot the same information for two features (the hour of day attendance occurred and 30 day visit count respectively) but in
162 a 2D plane which allows greater insight into the associations our model has learned between these variables and a patients
163 reattendance risk. The model has learned that the patients risk of reattendance displays a periodic dependence (Figure 3 b) with
164 the hour of day the attendance occurred (vertical dispersion is the result of interactions with other variables in the dataset) and it
165 has also learned an approximately linear dependence between a patients reattendance risk and the number of times they have
166 attended the emergency department in the last 30 days (Figure 3 c). It is important to note that these insights do not necessarily
167 reflect the actual risk factors for reattendance (since the model is an imperfect classifier) but only reflect the trends the model
168 has learned to make its decisions.

169 We then project the explanations for all attendances into a lower-dimensional (2D) 'explanation' space using the UMAP
170 algorithm (see Methods for details), this projection provides insight into the different high-level groups of explanations provided
171 by our model. The two-dimensional embedding of the attendances in the patient hold-out test set into the explanation space is
172 visualised in Figure 4, attendances close in this space share more similar explanations for their predicted reattendance risk.

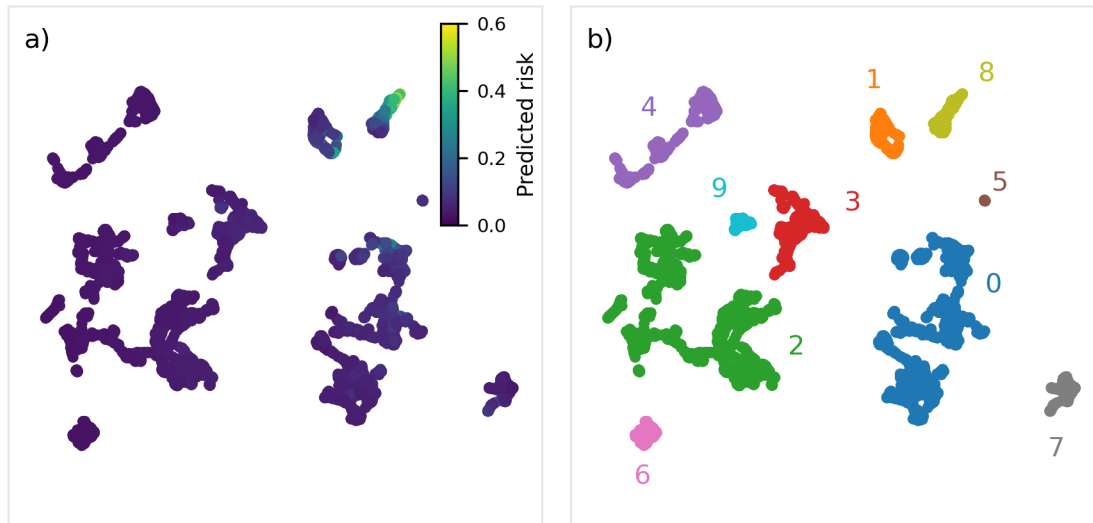


Figure 4. Embedding of the patient hold-out test set into a two-dimensional ‘explanation’ space using the UMAP algorithm. a) All attendances in the patient hold-out test set visualised in the explanation space, colour indicates the predicted reattendance risk for the respective attendance. b) Attendances in the explanation space coloured by cluster assignment. Cluster assignment was performed with the DBScan algorithm²³. 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 explanations (generated using the TreeExplainer algorithm) for each emergency department 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 3.

173 There are clear regions (colour) in the explanation space associated with increased reattendance risk. In Figures 4 b we colour
174 the attendances by a high-level cluster assignment obtained using the DBScan algorithm²³), where hyperparameters were
175 selected by visual inspection. Key descriptors of each cluster are displayed in Table 3 and visual inspection of this table
176 provides high-level insight into the reasoning of the machine-learnt model. For example, for attendances assigned to cluster
177 five, on average, the most important variable for determining a patients reattendance risk is their medical history.

178 Discussion

179 Our final 72-hour reattendance risk model achieved an AUROC of 0.749 and an average precision of 0.232 on a set of
180 attendances independent to the training set. Qualitatively, our model can use a patient’s (local) medical history and attendance
181 level information to predict their reattendance risk with moderate performance. In parallel, we trained an explanation model,
182 which can explain the model’s predictions at an attendance level (Figure 3 and Supplementary Figure 3) level. These
183 explanations were projected into a two-dimensional space (Figure 4). Such a visualisation can facilitate improved understanding
184 of the model’s high-level reasoning and can be used as a tool to understand the different sub-groups at risk of reattendance,
185 which could be used by the clinical care team to design interventions based on where a given attendance resides within the
186 explanation space. Ultimately, this facilitates the deployment of the machine-learnt models in a more informed manner.

187 Our final model (model p in Table 2) makes use of all variables found to be predictive of short-term reattendance. This is
188 not necessarily the optimal set of variables for a reattendance predictor and the variable likely contain obsolete information
189 because of correlations between variables. High correlation between variables is expected for clinical data, for example, it is
190 likely that patient age, arrival mode, and vital signs all to latently encode the patient’s frailty, which is known to be related to a
191 patients reattendance risk²⁴.

192 In our exploratory analysis, we found that the hour of day the attendance began correlates to the reattendance rate, with
193 higher reattendance rates observed during the night (Supplementary Figure 2). By evaluating the observed SHAP values for
194 the hour of day (Figure 3 b) we can observe that our model has learned a similar trend, associating attendance registration
195 during the night with an increased (between zero and two percent increase) reattendance risk. This trend could have several
196 different origins. Firstly, we have found (not shown) that the hour of day displays correlation with the reason for attendance with
197 either complaints associated with a higher risk of 72-hour reattendance more likely to present during the night or complaints
198 associated with a lower risk of 72-hour reattendance less likely to present during the night. Secondly, it is possible that the staff

Cluster	Count	Age (years)	Reattendance rate (95 % CI)	Predicted reattendance rate (95 % CI)	Condition count	30-day visit count	Most important variables (mean absolute SHAP values)		
							1 st	2 nd	3 rd
0	2,701	55.2	5.9(5.1-6.9)	6.8 (5.9-7.8)	2.8	0.0	Medical history	30-day visit count	Diagnosis
1	445	52.8	8.5 (6.3-11.5)	10.5 (8.0-13.7)	3.5	1.0	30-day visit count	Medical history	Diagnosis
2	2,885	42.8	2.2 (1.7-2.8)	4.0 (3.3-4.7)	0.0	0.0	Condition count	30-day visit count	Triage complaint
3	801	40.6	3.6 (2.5-5.2)	5.6 (4.2-7.4)	0.0	0.0	Condition count	Diagnosis	30-day visit count
4	895	41.5	1.2 (0.7-2.2)	3.4 (2.4-4.8)	0.0	0.0	Condition count	30 day visit count	Triage complaint
5	23	79.4	8.7 (2.4-26.8)	9.6 (2.8-28.0)	3.2	0.3	Medical history	Condition count	Diagnosis
6	296	32.9	2.0 (0.9-4.4)	3.0 (1.6-5.6)	0.0	0.0	Condition count	Diagnosis	30-day visit count
7	321	39.6	5.6 (3.6-8.7)	5.8 (3.8-9.0)	0.0	1.0	Condition count	30-day visit count	Diagnosis
8	341	42.3	31.1 (26.4-36.2)	26.6 (22.2-31.5)	3.3	3.8	30-day visit count	Medical history	Age
9	139	43.5	1.4 (0.4-5.1)	4.2 (1.9-8.9)	0.0	0.0	Condition count	30-day visit count	Diagnosis

Table 3. Properties of the attendances assigned to each of the explanation clusters (Figure 4). The count column displays the number of attendances in a given cluster. The age, 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. The final three columns display the most important variables in making a decision, averaged across all attendances in a given cluster.

199 fatigue and lower staffing levels may contribute to the increased reattendance rate for attendances occurring during the night,
200 although we have no way of testing this hypothesis in our dataset.

201 Our analysis (e.g., Figure 2 and model e in Table 2) shows that certain complaints are associated with a higher risk of
202 72-hour reattendance. For example, attendances whose chief complaint at registration is ‘Abdominal pain’ had a mean 72-hour
203 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 %).
204 Coding compliance was not evaluated in our dataset, which may effect this observation. For example, the most common chief
205 complaint at registration was ‘Unwell Adult’ (with 20.5 % of attendances listing this as the chief complaint in the training set)
206 which is utilized when either the chief complaint is not clear at registration, when the patient presents with multiple complaints
207 or as a result of inappropriate coding.

208 In addition to identifying complaints associated with a heightened short-term reattendance risk, our model also makes use
209 of ICD10 coded conditions (e.g., type 1 diabetes, lives alone) extracted from a patients electronic health record. These variables
210 allow the model to identify medical conditions, comorbidities, and risks which are associated with increased reattendance rates
211 and enables models to achieve moderate predictive performance (Table 2). Excluding the medical condition indicators, the
212 most important feature is the 30-day visit count which, in part, reflects the disproportionate use of EDs by frequent users²⁵. In
213 the visualisation of the attendances in the patient hold-out test set in the two-dimensional explanation space (Figure 4), the
214 most frequent attenders (30 day visit count of two or more) are clearly segregated (cluster 8 in Figure 4 b and Table 3). Visual
215 inspection of the properties of each cluster in Figure 4 could be used to help guide the design of interventional strategies. For
216 example, while attendances within cluster 8 (i.e., frequent attenders) may benefit from support in the community to mitigate
217 their reattendance risk, this will not necessarily be appropriate patients with a heightened reattendance risk but are suffering from
218 an acute injury associated with increased reattendance risk, such as a severe burn.

219 From a clinical perspective it is important to investigate the subset of reattendances which are also readmissions (i.e.,
220 reattendances to the emergency department which result in the patient being admitted to an inpatient ward). In these cases,
221 there is increased risk that there was missed critical illness or injury at the initial attendance, and they are important to evaluate
222 for clinical assurance purposes. Overall, 37.1 % of reattendances end in readmission, resulting in a 72-hour readmission rate
223 of 2.0 %. Evaluating our models predictions with a target equal to whether the patient readmitted in 72 hours, we find it has
224 an AUROC of 0.805 and an average precision of 0.040 on the hold-out patient test set. The high AUROC means the model
225 displays high discernibility between attendances which result in readmission and those that do not. The low average precision
226 reflects that readmissions only make up a minority fraction of reattendances and the false positive rate increases as a result
227 of the large class imbalance. Overall, these results demonstrate our classifier can identify the subset of reattendances which
228 are also readmissions with a similar predictive performance as reattendances which do not result in admission – a particularly
229 important result since these two different outcomes will likely merit different interventional strategies to reduce the risk of
230 reattendance/readmission.

231 A limitation of our study, shared with other investigations of machine-learning use in EDs²⁶, is that its primary data source
232 is structured past medical histories, which is unavailable for many patients. This could lead to our model discriminating against
233 people without a clinical history at the emergency department and associated hospital. An example of this bias can be observed
234 for cluster two (Table 3), where the two most important variables for determining a patients reattendance rate is the absence
235 of any visit to the emergency department in the last 30 days and the absence of any recorded medical conditions; which may
236 not be desirable behaviour. We mitigate this through the use of visit-level information and this bias could be further reduced

237 by linking to community datasets (e.g., GP records) to get a view of a patients medical history. However, in a deployment
238 scenario this bias could be minimized further by using the model as an alert tool – with results only being displayed for patients
239 it predicts to be at high risk of reattendance and otherwise being entirely invisible to clinical staff who would be free to carry
240 out standard clinical practice in cases where the alarm is not raised.

241 Practically, since the model uses only information available to clinicians at the time of the emergency department visit the
242 model has a relatively low barrier to implementation. Despite this, it will be essential to perform prospective, randomized
243 clinical trials of any implementation, investigating the efficacy of these predictive risk models, the associated interventions
244 and, importantly, analysing how they impact decision making. Ultimately, deployment of a machine-learning model could
245 eventually invalidate the model by changing the behaviours and descriptors of reattendances by altering the clinical decisions
246 made. In the short term, a relatively low-risk implementation of a machine-learned model trained to identify patients at risk of
247 reattendance would be in the implementation of a low-recall and high-precision alert system (for example, the configuration
248 presented in Figure 2 c). This would only raise alarms for the cases the model believes are at the highest risk of reattendance
249 and suggest appropriate clinically-validated intervention or additional clinical review. On average, using the configuration
250 displayed in Figure 2 c, this would have raised an alarm for only 1.6 % of attendances in which a decision to discharge was
251 made (approximately 2 times per day) and would expect to be correct approximately 50 % of the time – mitigating the risk of
252 alarm fatigue. While such a model would only be of limited impact because of the models low recall, as model performance
253 improves, this configuration could be re-evaluated and changed to increase the impact of the model.

254 When considering deployment it is important to discuss the context in which these predictive models could be prospectively
255 deployed. Our model was trained and retrospectively evaluated using data obtained local to the emergency department in
256 Southampton, using data available to clinicians during standard clinical practice. This is clearly an advantage if the model was
257 to be used at this location – the biases in the data and attendance characteristics will likely reflect what the model will encounter
258 in production. Conversely, this does mean that the model will not necessarily generalize to different EDs without first training
259 on their local data, this will be particularly prominent in EDs with a catchment zone with very different demographics to
260 Southampton, which would have a differing disease prevalence and characteristics at presentation to the emergency department.
261 Despite this since our model contains variables either in the standard UK emergency care dataset or regularly available to EDs
262 nationally, it would be possible to evaluate this model directly in other EDs with little alteration. External validation of our
263 model using data from different EDs is essential before prospective deployment beyond the department at which the training
264 data was sourced.

265 **Conclusion**

266 In conclusion, we have constructed and retrospectively evaluated a gradient boosted decision tree classifier capable of predicting
267 the 72-hour reattendance risk for a patient at the point of discharge from an emergency department. The highest performing
268 model achieved an AUROC of 0.749 and an average precision of 0.232 on a set of attendances independent to the training set.
269 We investigated the variables most indicative of risk and showed these were patient level factors (medical history) rather than
270 visit level variables such as recorded vital signs. We demonstrated how explainable machine learning can be used to investigate
271 the decisions a model is making. We suggested an implementation of the algorithm in a low-recall high-precision configuration
272 such that alarms are only raised if the model deems the patient to be at a (clinically defined) heightened risk of reattendance.
273 External validation and prospective clinical trials of these models are essential, with considerable consideration given to the
274 planned intervention resulting from the model's recommendation and the impact this would have on clinical decisions.

275 **Acknowledgements**

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

278 **Author contributions statement**

279 FPC performed the data analysis and modelling. DKB and ZDZ discussed and commented on the analysis with FPC. NW and
280 FPC obtained governance and ethical approval. MA and FB created the data extract. FPC, MJB and NW managed the study
281 at the UoS. MK managed the study at UHS. FPC and MK designed the study with assistance from TWVD. MK and TWVD
282 provided clinical guidance and insight. FPC wrote the first draft of the manuscript with assistance from MK and TWVD. All
283 authors frequently discussed the work and commented and contributed to future drafts of the manuscript.

284 Additional information

285 **Competing interests** The authors declare no competing interests.

286

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

289 References

- 290 1. Berchet, C., Emergency care services: trends, drivers and interventions to manage the demand. (2015)
- 291 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)
- 292 3. Bernstein, S. L., *et al.* The effect of emergency department crowding on clinically oriented outcomes. *Academic Emergency Medicine* **16**, 1-10 (2009)
- 293 4. Guttman, A., *et al.* Association between waiting times and short term mortality and hospital admission after departure
294 from emergency department: population based cohort study from Ontario, Canada. *BMJ* **342**, d2983 (2011)
- 295 5. Besga, A., *et al.* Risk factors for emergency department short time readmission in stratified population. *BioMed research international* (2015)
- 296 6. Deschodt, M., *et al.* Characteristics of older adults admitted to the emergency department (ED) and their risk factors for
297 ED reattendance based on comprehensive geriatric assessment: a prospective cohort study. *BMC geriatrics* **15**, 1 (2015)
- 298 7. Martin-Gill, C., Reiser, R.C. Risk factors for 72-hour admission to the ED. *The American journal of emergency medicine*
299 **22** 6, 448-453 (2004)
- 300 8. Arendts, G., Fitzhardinge, S., Pronk, K., Hutton, M., Nagree, Y. and Donaldson, M. Derivation of a nomogram to estimate
301 probability of revisit in at-risk older adults discharged from the emergency department. *Internal and emergency medicine*
302 **8** 3, 249-254 (2013)
- 303 9. Hao, S., *et al.* Risk prediction of emergency department revisit 30 days post discharge: a prospective study. *PLoS one* **9**
304 **11**, e112944 (2014)
- 305 10. Hong, W.S., Haimovich, A.D. and Taylor, R.A. Predicting 72-hour and 9-day return to the emergency department using
306 machine learning. *JAMIA open* **2** 3, 346-352 (2019)
- 307 11. Huang, K., Altosaar J. and Ranganath R. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *arXiv*,
308 arXiv:1904.05342 (2019)
- 309 12. Sterling, N. W., Patzer, R. E., Di, M., and Schrager, J. D. Prediction of emergency department patient disposition based on
310 natural language processing of triage notes. *International journal of medical informatics*, **129**, 184-188 (2019)
- 311 13. Micci-Barreca, D. A preprocessing scheme for high-cardinality categorical attributes in classification and prediction
312 problems. *ACM SIGKDD Explorations Newsletter*, **3** 1, 27-32 (2001)
- 313 14. Bergstra, J.S., Bardenet, R., Bengio, Y. and Kégl, B. Algorithms for hyper-parameter optimization. *In Advances in neural
314 information processing systems*, 2546-2554 (2011)
- 315 15. Bergstra, J., Yamins, D., and Cox, D. D. Making a Science of Model Search: Hyperparameter Optimization in Hundreds
316 of Dimensions for Vision Architectures. *In International Conference on Machine Learning*, 115-123 (2013)
- 317 16. Chen, T. and Guestrin, C. XGBoost: A scalable tree boosting system. *In Proceedings of the 22nd ACM sigkdd international
318 conference on knowledge discovery and data mining*, 785-794 (2016)
- 319 17. Lundberg, S. M., and Su-In L. A unified approach to interpreting model predictions. *In Advances in neural information
320 processing systems*, 4765-4774 (2017)
- 321 18. Lundberg, S. M., *et al.* From local explanations to global understanding with explainable AI for trees. *Nature machine
322 intelligence* **2** 1, 2522-5839 (2020)
- 323 19. Lundberg, S. M., *et al.* Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. *Nature
324 biomedical engineering* **2** 10, 749-760 (2018)
- 325 20. McInnes, L., Healy, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, *arXiv* **2**
326 arXiv:1802.03426 (2018)
- 327
- 328
- 329

- 330 **21.** Becht, E., *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nature biotechnology* **37** 1,38-44
331 (2019)
- 332 **22.** Diaz-Papkovich, A., Anderson-Trocmé, L., Gravel, S. UMAP reveals cryptic population structure and phenotype
333 heterogeneity in large genomic cohorts. *PLoS genetics* **15** 11, e1008432 (2019)
- 334 **23.** Ester, M., Kriegel, H.P., Sander, J. and Xu, X. A density-based algorithm for discovering clusters in large spatial databases
335 with noise. *In Kdd* **96** 34, 226-231 (1996)
- 336 **24.** Kahlon, S., *et al.* Association between frailty and 30-day outcomes after discharge from hospital. *Cmaj* **187** 11, 799-804
337 (2015)
- 338 **25.** LaCalle, E. and Rabin, E. Frequent users of emergency departments: the myths, the data, and the policy implications.
339 *Annals of emergency medicine* **56** 1, 42-48 (2010)
- 340 **26.** Joseph, J.W., *et al.* Deep-Learning Approaches to Identify Critically Ill Patients at Emergency Department Triage Using
341 Limited Information. *JACEP Open* **1**, 773-781 (2020)