# Clinical decision support systems for 3-month mortality in elderly patients admitted to ICU with ischemic stroke using interpretable machine learning

Jian Huang[1,2,3,*], Xiaozhu Liu[4,*] and Wanlin Jin[1,2] (iD)

## Abstract

**Background:** Elderly patients are more likely to suffer from severe ischemic stroke (IS) and have worse outcomes, including death and disability. We aimed to develop and validate predictive models using novel machine learning algorithms for the 3-month mortality in elderly patients with IS admitted to the intensive care unit (ICU).

**Methods:** We conducted a retrospective cohort study. Data were extracted from Medical Information Mart for Intensive Care (MIMIC)-IV and International Stroke Perfusion Imaging Registry (INSPIRE) database. Ten machine learning algorithms including Categorical Boosting (CatBoost), Random Forest (RF), Support Vector Machine (SVM), Neural Network (NN), Gradient Boosting Machine (GBM), K-Nearest Neighbors (KNNs), Multi-Layer Perceptron (MLP), Naive Bayes (NB), eXtreme Gradient Boosting (XGBoost) and Logistic Regression (LR) were used to build the models. Performance was measured using area under the curve (AUC) and accuracy. Finally, interpretable machine learning (IML) models presenting as Shapley additive explanation (SHAP) values were applied for mortality risk prediction.

**Results:** A total of 1826 elderly patients with IS admitted to the ICU were included in the analysis, of whom 624 (34.2%) died, and endovascular treatment was performed in 244 patients. After feature selection, a total of eight variables, including minimum Glasgow Coma Scale values, albumin, lactate dehydrogenase, age, alkaline phosphatase, body mass index, platelets, and types of surgery, were finally used for model construction. The AUCs of the CatBoost model were 0.737 in the testing set and 0.709 in the external validation set. The Brier scores in the training set and testing set were 0.12 and 0.21, respectively. The IML of the CatBoost model was performed based on the SHAP value and the Local Interpretable Model-Agnostic Explanations method.

**Conclusion:** The CatBoost model had the best predictive performance for predicting mortality in elderly patients with IS admitted to the ICU. The IML model would further aid in clinical decision-making and timely healthcare services by the early identification of high-risk patients.

## Keywords

Prediction model, hospital mortality, machine learning, elderly patients, ischemic stroke

[1]Health Management Center, Hunan Provincial Clinical Medicine Research Center for Intelligent Management of Chronic Disease, The Second Xiangya Hospital of Central South University, Changsha, Hunan, China
[2]Department of General Medicine, The Second Xiangya Hospital of Central South University, Changsha, Hunan, China
[3]Department of Ultrasound, Sir Run Run Shaw Hospital, Zhejiang University College of Medicine, Hangzhou, China
[4]Department of Critical Care Medicine, Beijing Shijitan Hospital, Capital Medical University, Beijing, China

[*]These authors contributed equally to this work.

**Corresponding author:**
Wanlin Jin, Health Management Center, Hunan Provincial Clinical Medicine Research Center for Intelligent Management of Chronic Disease, The Second Xiangya Hospital of Central South University, Changsha, Hunan, China.
Email: wanlin.jin@csu.edu.cn

## Introduction

Ischemic stroke (IS) accounts for approximately 80% of strokes and has become the second leading cause of mortality.[1] Previous research has revealed that over 75% of strokes occur in the elderly, leading to a significant financial burden.[2,3] By 2050, the global population of elderly individuals will exceed that of individuals under 65 for the first time in recorded history.[4] With prolonged life expectancy, an aging population leads to significantly increased stroke rates.[5] By 2050, the number of stroke survivors is projected to exceed 200 million.[6] Therefore, it is important to focus on managing acute stroke in the elderly to decrease the occurrence and enhance outcomes in this vulnerable group.[7]

The risk factors for stroke and mechanisms of ischemic injury differ between young and elderly patients. Additionally, elderly patients with ischemic stroke frequently receive ineffective therapy and experience worse outcomes compared to younger individuals with the condition.[7] Because of the different profiles of risk factors and different frequencies of stroke etiologies and subtypes,[8] the models established from younger cohorts may lead to suboptimal care for elderly patients,[9] and special models are needed for the elderly. When stroke-related cerebral impairment impairs the function of other vital organs, patients may require intensive care unit (ICU) care. There were significant differences between those admitted to the ICU and those admitted to the neurological ward. The ICU group was characterized by higher neurological severity, measured using validated instruments (e.g. the National Institutes of Health Stroke Scale, NIHSS[10]); moderate to severe impairment of consciousness; need for mechanical ventilation in many cases[11]; and high in-hospital mortality. What's more, functional outcomes in survivors entering ICU appear to be poor, especially in elderly patients.[12] Except for stroke severity scoring systems, such as NIHSS, outcome assessment should include the clinical evolution and the quality of survival, using appropriate tools.[13] Therefore, special models should be constructed for managing elderly ICU individuals with IS.

Artificial intelligence is increasingly being used in medicine, and clinical decision aid systems that rely on artificial intelligence have become a research hotspot.[14] Recently, new machine learning algorithms have shown superior performance in various competitions, such as Categorical Boosting (CatBoost), and eXtreme Gradient Boosting (XGBoost).[15,16] What's more, the studies using machine learning (ML) algorithms to predict three-month ICU mortality in elderly stroke patients are limited.

Therefore, the aim of our study was to develop a prognostic model for elderly ICU individuals with IS that could reliably identify patients at a very high risk of death. We developed interpretable machine learning models for predicting three-month in-hospital mortality. We further analyzed the contribution of each variable of the interpretable machine learning (IML) Model outcome using the Shapley Additive exPlanation (SHAP) values.

## Methods

### Design and participants

Our study was conducted in accordance with the TRIPOD checklist,[17] and details are shown in Supplement Figure 6. We did a multicenter, retrospective study. For model training and testing, we used the data from the Medical Information Mart for Intensive Care (MIMIC). Data of elderly patients with IS were extracted from the publicly available critical care database, MIMIC-IV 2.0.[18] All data were extracted from MIMIC-IV 2.0 (certification ID: 43357625). Philips Healthcare provided the MIMIC-IV database in partnership with Massachusetts Institute of Technology (MIT) Laboratory for Computational Physiology. It included de-identified death data for 23,844 ICU patients admitted between 2008 and 2019. Based on the de-identified patient information, the database's official ethics committee approved the public distribution of these clinical data. Consent was waived because retrospective patient data were anonymized.

The inclusion criteria were as follows: (1) diagnosis of IS according to ICD-9-CM diagnoses: 434.91 (cerebral artery occlusion, unspecified with cerebral infarction) or ICD-10-CM diagnoses of IS: I63.50 (cerebral infarction due to unspecified occlusion or stenosis of unspecified cerebral artery); (2) first-time ICU visit; (3) age ≥65 years. The exclusion criteria were as follows: (1) ICU stay of less than 24 hours; (2) individuals with missing values of more than 30%.

The external validation set is from the International Stroke Perfusion Imaging Registry (INSPIRE) dataset Version 1.2, a publicly accessible research dataset dedicated to perioperative medicine. It encompasses around 130,000 patients, involving patients at a South Korean academic institution over the period 2011 to 2020. The inclusion criteria were as follows: (1) diagnosis of IS; (2) age ≥65 years. Patients with personal data missing of more than 30% were excluded. Finally, the external validation set of 515 patients came from the INSPIRE dataset.

### Outcome variables and predictors

The primary outcome was death of elderly stroke patients within three months after ICU admission, either in or out of the hospital. Data on deaths in discharged patients were collected during the follow-up. Clinical information was gathered within 24 hours of admission. A literature review was used to identify candidate predictor factors, with an emphasis on those available in the ICU. The 52

**Table 1.** The population demographics and clinical characteristics.

| Variables | Total (n = 1826) | Survival (n = 1202) | Death (n = 624) | P |
|---|---|---|---|---|
| Gender, n (%) | | | | 0.067 |
|    Female | 995 (54) | 636 (53) | 359 (58) | |
|    Male | 831 (46) | 566 (47) | 265 (42) | |
| Smoking, n (%) | | | | 0.741 |
|    No | 1283 (70) | 841 (70) | 442 (71) | |
|    Yes | 543 (30) | 361 (30) | 182 (29) | |
| Race, n (%) | | | | 0.095 |
|    Black | 210 (12) | 136 (11) | 74 (12) | |
|    White | 1119 (61) | 757 (63) | 362 (58) | |
|    Others | 497 (27) | 309 (26) | 188 (30) | |
| Age, median (Q1, Q3) | 78 (71, 85) | 76 (70, 83) | 81 (74, 87) | <0.001 |
| BMI, median (Q1, Q3) | 26.6 (22.8, 30.7) | 27.2 (23.3, 31.2) | 25.6 (21.7, 29.7) | <0.001 |
| Myocardial infarction, N (%) | | | | 0.002 |
|    No | 1432 (78) | 969 (81) | 463 (74) | |
|    Yes | 394 (22) | 233 (19) | 161 (26) | |
| Congestive heart failure, n (%) | | | | <0.001 |
|    No | 1203 (66) | 854 (71) | 349 (56) | |
|    Yes | 623 (34) | 348 (29) | 275 (44) | |
| Dementia, n (%) | | | | <0.001 |
|    No | 1643 (90) | 1105 (92) | 538 (86) | |
|    Yes | 183 (10) | 97 (8) | 86 (14) | |
| COPD, n (%) | | | | 0.077 |
|    No | 1426 (78) | 954 (79) | 472 (76) | |
|    Yes | 400 (22) | 248 (21) | 152 (24) | |
| DM, n (%) | | | | 1 |
|    No | 1173 (64) | 772 (64) | 401 (64) | |
|    Yes | 653 (36) | 430 (36) | 223 (36) | |

(continued)

**Table 1.** Continued.

| Variables | Total (n = 1826) | Survival (n = 1202) | Death (n = 624) | P |
|---|---|---|---|---|
| Renal disease, n (%) | | | | <0.001 |
| No | 1373 (75) | 934 (78) | 439 (70) | |
| Yes | 453 (25) | 268 (22) | 185 (30) | |
| Malignant cancer, n (%) | | | | <0.001 |
| No | 1670 (91) | 1125 (94) | 545 (87) | |
| Yes | 156 (9) | 77 (6) | 79 (13) | |
| Severe liver disease, n (%) | | | | 0.006 |
| No | 1800 (99) | 1192 (99) | 608 (97) | |
| Yes | 26 (1) | 10 (1) | 16 (3) | |
| APSIII, median (Q1, Q3) | 45 (33, 60) | 39 (30, 51) | 59 (44, 76) | <0.001 |
| LODS, median (Q1, Q3) | 4 (2, 7) | 3 (2, 5) | 6 (4, 9) | <0.001 |
| SOFA, median (Q1, Q3) | 4 (2, 6) | 3 (2, 5) | 5 (4, 8) | <0.001 |
| GCS min, median (Q1, Q3) | 12 (8, 14) | 13 (10, 14) | 9 (6, 13) | <0.001 |
| Heart rate, median (Q1, Q3) | 81 (70, 95) | 80 (69, 92) | 85.5 (73, 102) | <0.001 |
| SBP, median (Q1, Q3) | 138 (118, 155) | 139 (119, 156) | 135 (115, 155) | 0.018 |
| DBP, median (Q1, Q3) | 71 (60, 85) | 71 (60, 84) | 72 (59, 85) | 0.88 |
| MBP, median (Q1, Q3) | 90 (78, 102.75) | 90 (79, 103) | 89 (76, 102) | 0.291 |
| Respiratory rate, median (Q1, Q3) | 18 (16, 22) | 18 (16, 22) | 19 (16, 24) | <0.001 |
| Temperature, median (Q1, Q3) | 36.7 (36.4, 37.0) | 36.7 (36.5, 37.0) | 36.7 (36.4, 37.1) | 0.232 |
| $SpO_2$, median (Q1, Q3) | 98 (95, 100) | 98 (96, 100) | 98 (95, 100) | 0.028 |
| Glucose, median (Q1, Q3) | 128 (105, 167) | 124 (103, 160) | 136.5 (112, 187) | <0.001 |
| Hematocrit, median (Q1, Q3) | 35.1 (31.0, 39.2) | 35.7 (31.4, 39.5) | 34.1 (29.9, 38.5) | <0.001 |
| Hemoglobin, median (Q1, Q3) | 11.5 (9.9, 12.9) | 11.6 (10.2, 13.2) | 11.0 (9.4, 12.4) | <0.001 |
| Platelets, median (Q1, Q3) | 204 (158, 261) | 204 (160, 256) | 204 (156, 269) | 0.656 |
| WBC, median (Q1, Q3) | 9.6 (7.4, 12.8) | 9.2 (7.22, 12.1) | 10.5 (8, 14.12) | <0.001 |
| MCH, median (Q1, Q3) | 30.1 (28.4, 31.4) | 30.2 (28.6, 31.5) | 29.8 (28.1, 31.1) | 0.001 |
| MCHC, median (Q1, Q3) | 32.6 (31.6, 33.5) | 32.7 (31.8, 33.7) | 32.3 (31.3, 33.2) | <0.001 |

(continued)

**Table 1.** Continued.

| Variables | Total (n = 1826) | Survival (n = 1202) | Death (n = 624) | P |
|---|---|---|---|---|
| MCV, median (Q1, Q3) | 92 (88, 96) | 91 (88, 95) | 92 (88, 96) | 0.203 |
| RBC, mean ± SD | 3.84 ± 0.73 | 3.89 ± 0.73 | 3.73 ± 0.73 | <0.001 |
| RDW, Median (Q1, Q3) | 14.1 (13.3, 15.4) | 13.9 (13.2, 14.9) | 14.6 (13.6, 16.2) | <0.001 |
| Albumin, median (Q1, Q3) | 3.5 (3.2, 3.7) | 3.54 (3.3, 3.78) | 3.38 (3, 3.6) | <0.001 |
| Anion gap, median (Q1, Q3) | 15 (13, 17) | 14 (13, 16) | 15 (13, 18) | <0.001 |
| Bicarbonate, median (Q1, Q3) | 23 (21, 26) | 24 (21.25, 26) | 23 (20, 25) | <0.001 |
| BUN, median (Q1, Q3) | 20 (14, 29) | 19 (14, 26) | 22.3 (16, 34.3) | <0.001 |
| Calcium, median (Q1, Q3) | 8.7 (8.2, 9.1) | 8.7 (8.3, 9.1) | 8.6 (8.1, 9) | <0.001 |
| Chloride, median (Q1, Q3) | 104 (100, 107) | 104 (100, 107) | 104 (100, 107) | 0.594 |
| Creatinine, median (Q1, Q3) | 1 (0.8, 1.3) | 1 (0.8, 1.3) | 1 (0.8, 1.5) | <0.001 |
| Sodium, median (Q1, Q3) | 140 (137, 142) | 140 (137, 142) | 140 (137, 142) | 0.714 |
| Potassium, median (Q1, Q3) | 4.1 (3.8, 4.5) | 4.1 (3.7, 4.5) | 4.1 (3.8, 4.6) | 0.032 |
| ALT, median (Q1, Q3) | 19 (14, 27) | 18.7 (14.4, 25.2) | 20.7 (14.0, 30.6) | 0.004 |
| ALP, median (Q1, Q3) | 78 (64, 94) | 76 (63, 89) | 84 (65, 107) | <0.001 |
| AST, median (Q1, Q3) | 26.7 (21, 36) | 26.2 (21, 34) | 27.0 (21.6, 39.5) | 0.006 |
| Total bilirubin, median (Q1, Q3) | 0.6 (0.4, 0.71) | 0.6 (0.4, 0.7) | 0.6 (0.4, 0.8) | 0.336 |
| LDH, median (Q1, Q3) | 246.6 (218.3, 292.0) | 237 .0(210.7, 272.0) | 270.4 (234.3, 334.2) | <0.001 |
| INR, median (Q1, Q3) | 1.2 (1.1, 1.4) | 1.2 (1.1, 1.3) | 1.21 (1.1, 1.5) | <0.001 |
| PT, median (Q1, Q3) | 13 (11.83, 15.1) | 12.8 (11.7, 14.4) | 13.8 (12.2, 16.3) | <0.001 |
| PTT, median (Q1, Q3) | 29.4 (26.5, 34.3) | 29.1(26.5, 34.1) | 23.0 (26.5, 34.7) | 0.149 |
| Type, n (%) | | | | 0.086 |
| None | 1582 (86.6) | 1024 (85.2) | 558 (89.4) | |
| Cerebrovascular thrombectomy | 26 (1.4) | 19 (1.5) | 7 (1.1) | |
| Cerebrovascular thrombolysis | 196 (10.7) | 144 (12.0) | 52 (8.4) | |
| Cerebral artery stenting | 22 (1.3) | 15 (1.3) | 7 (1.1) | |

BMI: body mass index; COPD: chronic obstructive pulmonary disease; DM: diabetes mellitus; APSIII: Acute Physiology Score III; LODS: Logistic Organ Dysfunction System; SOFA: Sequential Organ Failure Assessment; GCS min: minimum Glasgow Coma Scale values; SBP: systolic blood pressure; DBP: diastolic blood pressure; MBP: mean blood pressure; SpO$_2$: pulse oximetry; WBC: white blood cell count; MCH: mean corpuscular hemoglobin; MCHC: mean corpuscular hemoglobin concentration; MCV: mean corpuscular volume; RBC: red blood cell count; RDW: red blood cell volume distribution width; BUN: blood urea nitrogen; ALT: alanine aminotransferase; ALP: alkaline phosphatase; AST: aspartate aminotransferase; LDH: lactic dehydrogenase; INR: international normalized ratio; PT: prothrombin time; PTT: partial thromboplastin time; Type: types of endovascular surgery.

candidate variables were listed in Table 1. The variables include: (1) Demographic data: gender, age (years), race, and smoking. (2) Vital sign data: body mass index (BMI, kg/m$^2$), heart rate (bpm), systolic blood pressure (SBP, mmHg), diastolic blood pressure (DBP, mmHg), mean blood pressure (MBP, mmHg), pulse oximetry (SpO$_2$, %), respiratory rate (bpm), and temperature (degrees C). (3) Laboratory test data: glucose (mg/dL), hematocrit (%), hemoglobin (g/dL), platelets (K/µL), white blood cell count (WBC, K/µL), mean corpuscular hemoglobin (MCH, pg), mean corpuscular hemoglobin concentration (MCHC, %), mean corpuscular volume (MCV, fL), red blood cell count (RBC, m/µL), red blood cell volume distribution width (RDW, %), blood urea nitrogen (BUN, mg/dL), alanine aminotransferase (ALT, IU/L), alkaline phosphatase (ALP, IU/L), aspartate aminotransferase (AST, IU/L), LDH (lactic dehydrogenase, IU/L), international normalized ratio (INR), prothrombin time (PT, seconds), partial thromboplastin time (PTT, seconds). 4) Disease information: myocardial infarction, congestive heart failure, dementia, chronic obstructive pulmonary disease (COPD), diabetes mellitus (DM), renal disease, malignant cancer, severe liver disease, and types of endovascular surgery (Type). 5) Score data: Acute Physiology Score III (APSIII), Logistic Organ Dysfunction System (LODS), Sequential Organ Failure Assessment (SOFA), and minimum Glasgow Coma Scale values (GCS min). Instead of selecting the extreme value, we chose the average value of the features recorded multiple times in the electronic medical record system of a patient's hospital, thus mitigating the impact of data fluctuations on the outcome. Variables with more than 30% missing values were excluded from analysis to ensure accuracy. Missing values for the variables of the derivation cohort were listed in Supplementary Table 3. The K-Nearest Neighbor (KNN) algorithm was used to fill in the missing values for the variables with missing values less than 30%, which was fitted to the train data and applied to both the train and the test sets. Using the R package "DMwR" and the function "Knn Imputation", the remaining variables with missing values were imputed with default parameters (including k = 5).

### Model training and testing

The entire study population were randomly divided into training and testing sets at an 8:2 ratio.

### Selection procedure

Recursive feature elimination (RFE) was applied to select the most influential features for predicting outcome events using the training set. Area under the curve (AUC) values were then measured using the RFE algorithm.

### Machine learning model development

The study employed ten common machine learning algorithms, including CatBoost, Random Forest (RF), Support Vector Machine (SVM), Neural Network (NN), Gradient Boosting Machine (GBM), KNN, Multi-Layer Perceptron (MLP), Naive Bayes (NB), XGBoost, and Logistic Regression (LR), to predict the three-month mortality in elderly ICU patients with IS. We employed five-fold cross-validation on the training set to optimize the parameters. For each classifier, the hyperparameters that produced the highest average receiver operating characteristic curve (AUROC) in the five-fold cross-validations were chosen and adjusted before model testing. To evaluate the performance of these algorithms, we used a testing set that included 355 patients who were not part of the model training process. The AUROC, sensitivity, specificity, positive predictive value (PPV), negative predictive value (NPV), accuracy, and F1 score were used to assess the performance. Additionally, we drew calibration curves and calculated the Brier score. We also did a decision curve analysis (DCA) to evaluate the net clinical benefit.

### Statistical analysis

Data are expressed as mean ± standard deviation (SD) when normally distributed and as median and interquartile range (IQR) in the presence of skewed distribution. Categorical variables were expressed as frequencies and percentages and were compared using chi-squared analysis or Fisher's exact test. Outliers were identified and removed. For univariate analysis, the R packages "Nortest" and "CBCgrps" were used. The RFE feature selection was achieved using the rfe function of the "caret" package, within a cross-validation. The R package "caret" was used to propose the machine learning models. The receiver operating characteristic (ROC) analysis and area under the curve (AUC) calculations were performed using R package "pROC". Interpretability analysis based on SHAP. Local Interpretable Model-Agnostic Explanations (LIME) was performed by "modelstudio" package, and "lime" package in R (version 4.2.0). $P < 0.05$ was considered statistically significant.

## Results

### Population demographics

There were 3115 admissions for IS. The 1826 IS patients were eligible for further analysis according to the inclusion and exclusion criteria (Figure 1), of which 54% (995) were women and 46% (831) were men, with a median age of 78 years (IQR, 71–85 years). A total of 1202 patients with IS survived for three months in the hospital, and 624 died. Table 1 shows the characteristics of the survival and mortality groups.
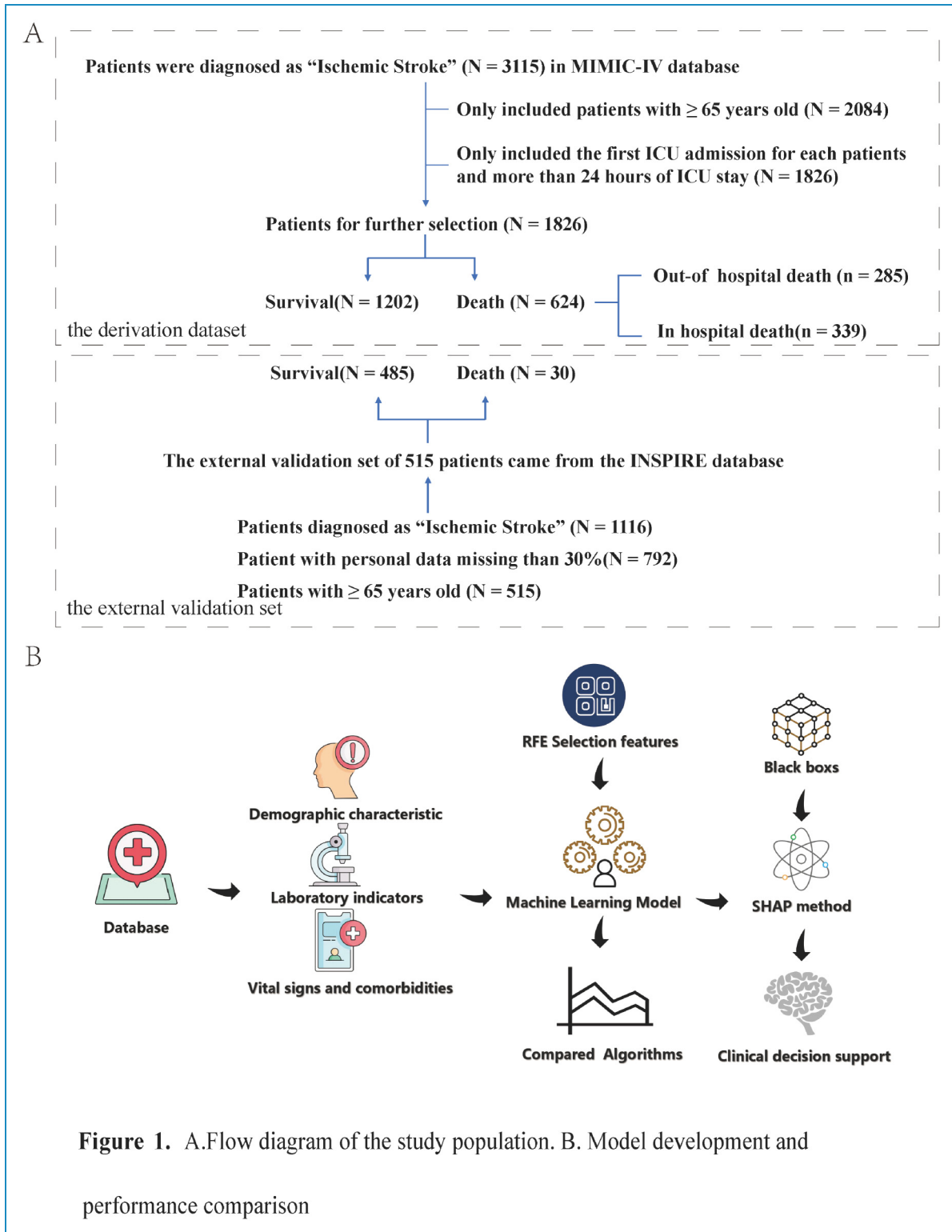
**Figure 1.** A.Flow diagram of the study population. B. Model development and

performance comparison

**Figure 1.** (a) Flow diagram of the study population. (b) Model development and performance comparison.

Compared to the death group, the BMI, SBP, temperature, $SpO_2$, hematocrit, hemoglobin, RBC, MCH, MCHC, albumin, blood calcium, and bicarbonate were significantly higher in the survival group ($P < 0.05$). Nevertheless, the levels of bilirubin, age, heart rate, GCS min, LODS, SOFA score, respiratory rate, glucose, WBC, RDW,

anion gap, BUN, creatinine, potassium, ALT, AST, ALP, LDH, INR, and PT were higher in the death group. The proportion of patients with mild liver disease was higher in the death group ($P < 0.05$). In addition, we compared the baseline data of patients who died in-hospital and those who died out-of-hospital (Supplementary Table 1). The differences in age, dementia, malignant cancer, glucose, WBC, bicarbonate, calcium, PT, BMI, GCS, LODS, SOFA, and APSIII in the two groups were statistically significant ($P < 0.05$).

The external validation set baseline table is shown in Supplementary Table 2. Compared to the death group, the BMI, albumin, and GCS min were higher in the survival group ($P < 0.05$).

## Feature selection

The RFE method was applied for feature selection (Figure 1), after which the most important eight features were selected. The prediction models could be built using features such as GCS min, albumin, LDH, age, ALP, BMI, platelets, and types of surgery.

## Model evaluation and comparison

The prediction models were constructed using several widely used machine learning algorithms (CatBoost, RF, SVM, NN, GBM, KNN, MLP, NB, and XGBoost and LR). The ROC curve, cutoff value, Youden index, F1 score, accuracy, specificity, sensitivity, PPV, and NPV were used to evaluate the prediction model. As shown in Figure 2, CatBoost had the best predictive performance in the testing set, with an AUC

of 0.737, which is better than that of the other models. In terms of the F1 index and Youden index, as well as the accuracy, the CatBoost model also exhibited excellent predictive performance in the testing set. Comparison of the AUC values of the CatBoost model with other models by DeLong test is shown in Supplementary Table 4. In both the training set and the test set, the $P$-values are less than 0.05, suggesting that the difference in AUC of CatBoost compared with other models is statistically significant.

The accuracy of the CatBoost mode in the training set was 0.821, which was higher than the accuracy of the LR model. Table 2 displays all parameters of the models developed using different algorithms. The AUC value of the model in the external validation set was 0.709 (Supplementary Figure 3).

In the training set, the CatBoost model had a Brier score of 0.12, and in the testing set, it had a Brier score of 0.21. When the Brier score <=0.25, the model was considered to have favorable calibration.[19] In both training and testing sets, the DCA curve indicated a net benefit and threshold probability for the CatBoost model (Figure 3(a) and (b)). According to the calibration plot (Figure 3(c) and (d)), the CatBoost model adequately predicted mortality in both training and testing sets. The DCA curves of the other models in the training sets and testing sets are shown in Supplementary Figures 7 and 8, respectively.

## Model interpretation

As shown in Figure 4, CatBoost analyzes an independent testing set using the Tree-Explainer class imported from Shapley additive explanation (SHAP).[20] Among the characteristics associated
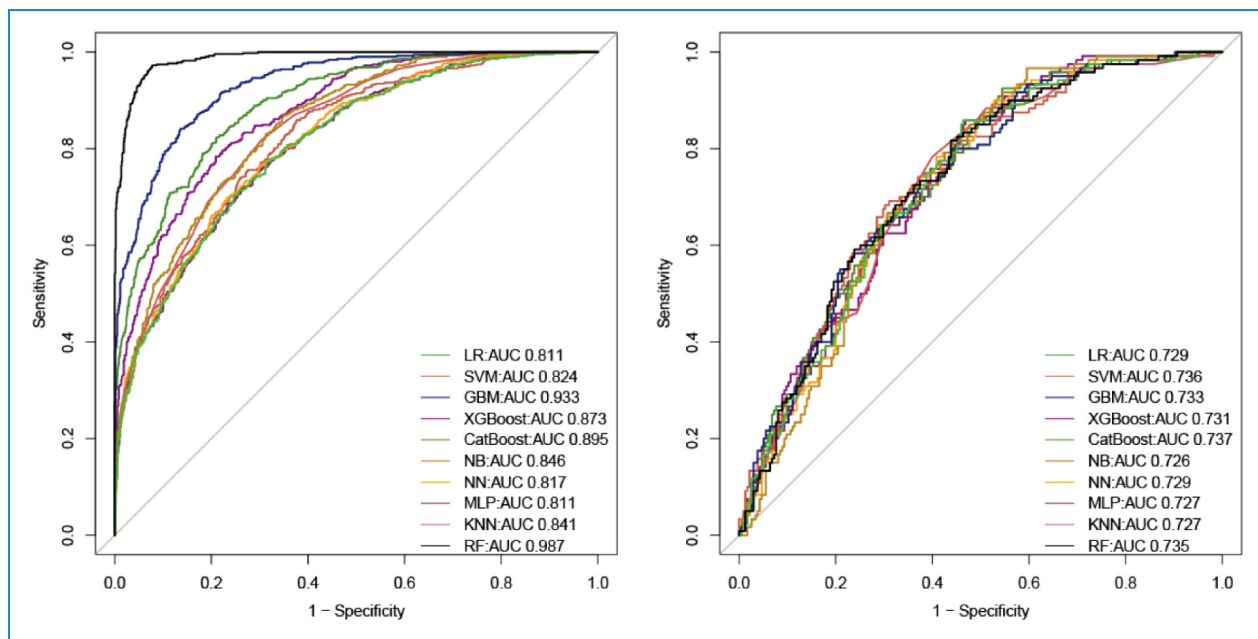


**Figure 2.** ROC curve of model. (a) AUC values of all models in the training set. (b) AUC values for all models in the testing set.

**Table 2.** Model performance metrics in the training set and in the validation set.

| Database | Model | Cutoff Value | F1 | Accuracy | Sensitivity | Specificity | PPV | NPV |
|---|---|---|---|---|---|---|---|---|
| Training set (*n* = 1471) | | | | | | | | |
| | GBM | 0.383 | 0.769 | 0.853 | 0.714 | 0.925 | 0.833 | 0.862 |
| | KNN | 0.291 | 0.567 | 0.764 | 0.450 | 0.928 | 0.764 | 0.764 |
| | LR | 0.305 | 0.591 | 0.751 | 0.526 | 0.868 | 0.674 | 0.778 |
| | MLP | 0.427 | 0.621 | 0.744 | 0.613 | 0.813 | 0.631 | 0.801 |
| | NB | 0.179 | 0.619 | 0.782 | 0.518 | 0.919 | 0.770 | 0.785 |
| | NN | 0.321 | 0.608 | 0.757 | 0.550 | 0.866 | 0.681 | 0.787 |
| | RF | 0.363 | 0.886 | 0.927 | 0.827 | 0.979 | 0.954 | 0.916 |
| | SVM | 0.220 | 0.618 | 0.760 | 0.565 | 0.861 | 0.680 | 0.792 |
| | XGBoost | 0.474 | 0.487 | 0.760 | 0.330 | 0.981 | 0.904 | 0.739 |
| | CatBoost | 0.346 | 0.722 | 0.821 | 0.679 | 0.895 | 0.770 | 0.843 |
| Testing set (*n* = 355) | | | | | | | | |
| | GBM | 0.280 | 0.560 | 0.699 | 0.567 | 0.766 | 0.553 | 0.776 |
| | KNN | 0.241 | 0.394 | 0.679 | 0.308 | 0.868 | 0.544 | 0.711 |
| | LR | 0.181 | 0.477 | 0.673 | 0.442 | 0.791 | 0.519 | 0.735 |
| | MLP | 0.249 | 0.532 | 0.687 | 0.525 | 0.770 | 0.538 | 0.761 |
| | NB | 0.189 | 0.471 | 0.665 | 0.442 | 0.779 | 0.504 | 0.732 |
| | NN | 0.284 | 0.477 | 0.673 | 0.442 | 0.791 | 0.519 | 0.735 |
| | RF | 0.268 | 0.466 | 0.690 | 0.400 | 0.838 | 0.558 | 0.732 |
| | SVM | 0.343 | 0.533 | 0.699 | 0.508 | 0.796 | 0.560 | 0.760 |
| | XGBoost | 0.466 | 0.300 | 0.685 | 0.200 | 0.932 | 0.600 | 0.695 |
| | CatBoost | 0.133 | 0.528 | 0.687 | 0.517 | 0.774 | 0.539 | 0.758 |

with the three-month mortality in elderly patients with IS admitted to the ICU, as shown in SHAP summary plots, GCS min, LDH, type, albumin, age, ALP, platelets, and BMI had the highest importance scores.

## ML explainability results for two patients

The SHAP force plot visualizes the Shapley value, which indicates whether a prediction increases or decreases from its baseline.[21]

*Patient 1.* Patient 1 was an elderly individual admitted to the ICU for IS. Indeed, the patient passed away on the 90th day after the admission. The factors identified by the model that contributed to the higher mortality prediction for this patient included age, LDH level, platelet count, ALP, albumin level, and GCS min score. The IML model predicting that the patient had a high risk of death, the patient's actual outcome during the 3 months of follow-up after ICU admission was death (Figure 4(c)). The SHAP plot indicates that this patient is at a high risk of poor
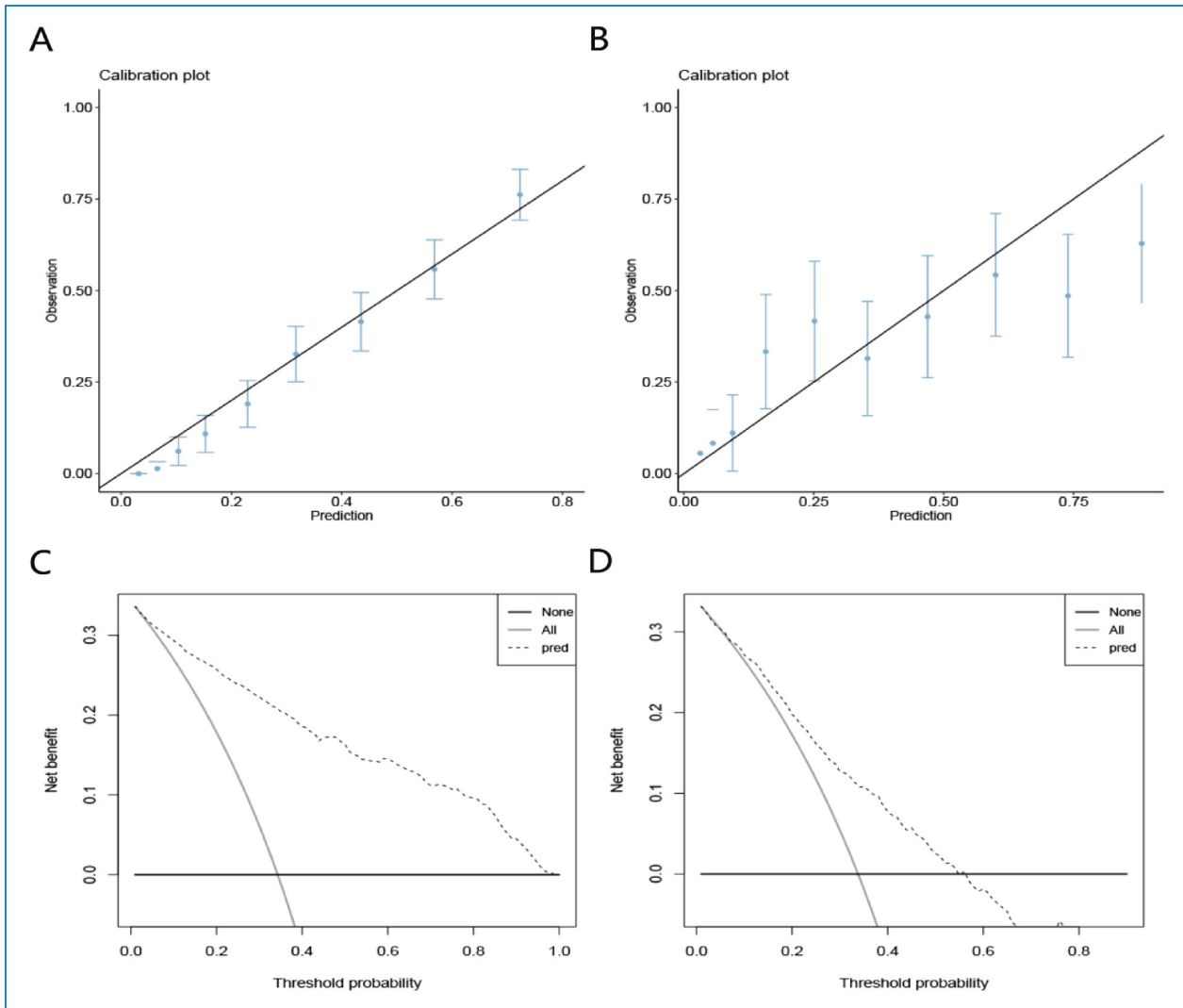
**Figure 3.** Assessment of CatBoost model. (a) The calibration plot of model in the training set. (b) The calibration plot of model in the testing set. (c) The DCA curve of model in the training set. (d) The DCA curve of model in the testing set.

prognosis after three months, which means a higher short-term risk of mortality. This is due to the patient's higher age, LDH levels, and lower GCS min score at the time of admission to the ICU. Specifically, the elderly the patient, the higher the risk of death (as shown in the dependent plot, Figure 5(e)). LDH levels greater than 300 were followed by a SHAP value greater than 0, indicating an increased risk of death (Figure 5(d)). A GCS min score of less than 8 was followed by a SHAP value greater than 0, also indicating an increased risk of death (Figure 5(a)).

*Patient 2.* This was an elderly patient admitted to the ICU because of IS. After 90 days, the patient is still alive. The ML model predicted that the patient had a low risk of death, which was correct (Figure 4(d)). Given that patient's indicators were relatively stable in the ICU, observing this result of the IS patient in the ICU is reasonable. What's

more, the patient has lower levels of LDH, platelet, and age, and higher levels of GCS min score and BMI. These instructions indicate that the patient has a clearer level of consciousness, a better physical recovery state, a better nutritional status of the body, lower level of bodily nerve damage, and a lighter inflammatory state.

## The contribution of the feature to the outcome

Decreased GCS min scores, BMI and albumin, increased age, LDH, ALP, platelet levels, would have a positive effect on the occurrence of outcome events (Figure 5). Figure 5 shows that the minimum GCS score of the patient was negatively correlated with the SHAP value, which is consistent with clinical practice. The clearer the consciousness of the patient, the better is condition of the patient. Higher ALP, LDH, and age result in more positive
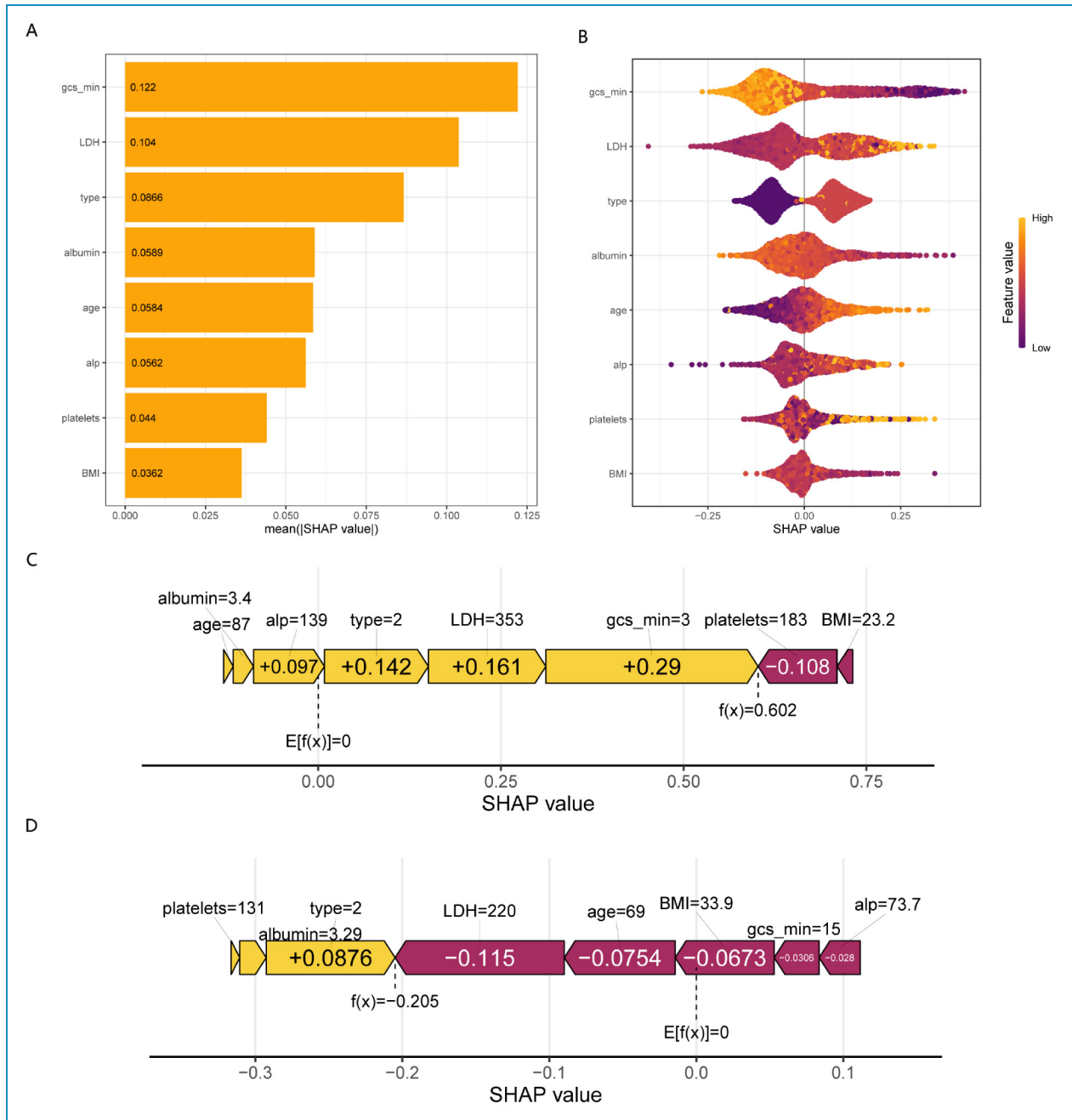
**Figure 4.** (a/b) SHAP summary plot for the eight clinical features contributing to model prediction for mortality, GCS min, LDH, type, albumin, age, ALP, platelets, and BMI. (c/d) SHAP explanation force plot for two patients from the held-out testing set of the ML model.

SHAP values, suggesting that higher values of ALP, LDH, or the patient's age are associated with a higher SHAP value and a greater likelihood of mortality. Higher albumin levels are associated with a smaller SHAP value, which is also consistent with clinical practice. The contribution of the features to the outcomes is shown in Supplementary Figure 1. The plots of feature importance, feature distribution, Shapley values, accumulated dependence, ceteris paribus, target versus feature, and break down are shown

in Supplementary Figure 2 (example for age (A) and GCS min (B) of one positive patient) and Supplementary Web attachment.

## Model performance in different subgroups

Subgroup analyses provide more insight into the diagnostic performance of models in specific patient populations. Therefore, model performance was studied in different

**Figure 5.** The level of the feature corresponds to the SHAP value. (a) GCS min, (b) ALP, (c) albumin, (d) LDH, (e) age.

subgroups (age, gender). Previous studies have shown that age (≥80 years old) is a significant independent predictor of the three-month mortality,[22] so the age subgroups were divided at the 80-year mark. In the training set, the AUC values of the CatBoost model in the subgroups aged greater than 80 years old and aged between 65 and 80 years old were 0.869 and 0.886, respectively (Supplementary Figure 4(A) and (B)). In the test set, The AUC values of the CatBoost model in the subgroups aged greater than 80 years old and aged between 65 and 80 years old were 0.752 and 0.746, respectively (Supplementary Figure 4(C) and (D)).

In the training set, the AUC values of the CatBoost model in the subgroups male and female were 0.890 and 0.882, respectively (Supplementary Figure 5(A) and (B)). In the test set, The AUC values of the CatBoost model in the subgroups male and female were 0.793 and 0.743 respectively (Supplementary Figure 5(C) and (D)).

## Discussion

Ischemic stroke is the most common type of stroke, predominantly affecting elderly individuals. Short-term mortality rates after stroke hospitalization are particularly high. To stratify elderly patients with IS admitted to the ICU, a death prediction model is essential and helps improve healthcare quality and clinical decision-making by early identification of high-risk patients.

This is the first interpretable machine learning prediction model for the three-month mortality in patients with IS in an ICU cohort in the United States. In addition, the influence of endovascular treatment on the occurrence of outcome

events was analyzed, and the contribution of thrombectomy and thrombolysis to the occurrence of outcome events is presented in the figure.

Among the ten models used in this study, the CatBoost model demonstrated the best overall performance, followed by SVM. RF, XGBoost, LR, NN, GBM, KNN, MLP, and NB showed lower performance levels. LR is the most commonly used model to describe the relationship between a dependent variable and one or more explanatory variables. Our research shows that the LR model does not perform as well as the CatBoost model. This may be due to the imbalance in our data outcome variables and the non-linear nature of the prediction results and features. Therefore, we resampled the data and used cross-validation to adjust for overfitting. KNN can be extremely valuable in improving prediction accuracy when outcomes and interrelationships between variables may be non-linear or unknown. In this study, the performance of the KNN model was slightly worse, which may be attributed to the poor performance of the Euclidean distance metric used by KNN. The SVM model can be affected by sample imbalance, which may cause the model to be biased toward the class with a large sample number, thereby reducing the model's prediction efficiency. In this study, when the outcome event is a binary classification problem, GBM typically utilizes the Logistic Loss function and applies the sigmoid function to convert the continuous predicted value into a probability value at the final output. We dedicated a significant amount of time to the GBM model, which is a limitation of this algorithm. The performance of the RF model and the CatBoost model varies considerably, despite both being tree-based models. Although the specificity of the

RF model is high, its sensitivity is low, which leads to unsatisfactory prediction results. The CatBoost model is significantly better than the RF model, demonstrating the effectiveness of the CatBoost model in identifying stroke-related deaths in the elderly. CatBoost is a comprehensive new algorithm based on gradient boosting of decision trees. Previously, the two mainstream algorithms in the Boosting family were XGBoost and LightGBM. According to official evaluations, the CatBoost model, a new addition to the Boosting family, outperforms the two algorithms mentioned above. In this study, the CatBoost model outperforms the XGBoost model. In addition, as an emerging algorithm, CatBoost has unique advantages. It can automatically handle categorical features, requires minimal hyperparameter adjustment, improves model stability, and reduces the risk of overfitting.

By analyzing demographic data, biochemical tests, ICU scores, and surgical treatment of patients, machine learning algorithms were used to fill in the missing values and extract the most important features. A total of eight features (GCS min, albumin, LDH, age, ALP, BMI, platelets, and types of surgery) were finally used for model construction. Machine learning models were used for the first time to predict the three-month mortality in elderly patients with IS admitted to the ICU, allowing the identification of critically ill elderly patients with IS earlier. We evaluated different supervised machine learning algorithms and compared them with classic LR approaches to identify the best model for predicting short-term death in elderly patients with IS. In the testing set, the prediction performance of each model was similar to its training set performance, indicating that the model was robust and generalizable. Finally, the AUC of CatBoost model used to identify patients who died was 0.737 in the testing set, better than ten models (such as the LR model with an AUC of 0.729). The CatBoost model achieved better accuracy than the traditional LR method and other machine learning methods.

Li et al.[23] analyzed 30-day stroke mortality using the MIMIC database. However, since outcomes are typically more severe in elderly individuals, models based on individuals aged 18 and older may not be directly applicable to elderly populations,[8,9,22] potentially resulting in suboptimal care. What's more we included types of surgery, which was also the feature in the final model. Additionally, we have an external dataset. Someeh et al.[24] predicted mortality in brain stroke patients. We are focusing on elderly stroke patients in the ICU and have utilized the neural network model, but NN does not perform as well as the CatBoost model. This may be because functional outcomes in survivors entering the ICU appear to be poor, especially in elderly patients.[12] The lack of interpretability of machine learning models for IS patients[25] is a major barrier and has limited clinical applications. We used the SHAP method to enhance the interpretability of the CatBoost model. The SHAP summary plot provides a global explanation of the dataset's prediction results, whereas the SHAP force plot provides an explanation of each individual patient's prediction results.

We developed a predictive model for assessing the in-hospital mortality risk for elderly IS patients admitted to the ICU and presented a user-friendly interface to improve healthcare quality and clinical decision-making by early identification of high-risk patients. Our model was better than the THRIVE score in predicting 90-day outcomes among stroke patients undergoing endovascular treatment, with an AUC of 0.709.[26] Compared with other studies, this study included the GCS score, which reflects the degree of coma in patients and can improve the effect of laboratory examination indicators in predicting the outcome of stroke patients.[27,28]

We found that GCS min, LDH, types of surgery, albumin, age, ALP, platelets, and BMI were tested and selected for the three-month mortality prediction model. The Glasgow Coma Scale was used to predict the 90-day mortality in patients with IS.[29] Some studies have also reported the application of GCS in predicting 30-day mortality[30] and 10-year stroke mortality.[31] Wang et al. reported that in individuals with acute ischemic stroke or transient ischemic attack, increased lactate levels are associated with adverse outcomes.[32] Albumin increases the risk of death after a stroke.[33] You et al. reported that a low platelet count upon admission was independently related to the three-month mortality and pneumonia in patients with acute IS.[34] According to Uehara et al., patients with transient ischemic attacks caused by intracranial atherosclerosis have higher serum ALP levels at admission.[35] In patients with symptomatic intracranial atherosclerosis, elevated serum ALP levels can predict early neurological decline.[35] ALP level was an independent predictor of all-cause and vascular death after ischemic or hemorrhagic stroke.[36] Among patients with IS treated with intravenous or endovascular therapy, elderly age is associated with poor outcomes.[37] Age ≥80 years is a significant independent predictor of 90-day mortality.[22] BMI has been reported to be associated with ICU.[36,38] Our machine learning algorithm also underlines the significance of BMI. A higher BMI implies a better nutritional status of the body, which is more advantageous in resisting disease invasion.

This study also analyzed the correlation between each variable (GCS min, platelet, age, and albumin) and the outcome (Supplementary Figure 1). This study also analyzed the local interpretability of the CatBoost model. Taking three patients as examples, the results of feature importance, feature SHAP value, feature accumulated dependence, ceteris paribus, break down, and so on were analyzed. Interpretable machine-learning models can be viewed by clinicians through the web, making it easier for clinicians to use these models to develop treatment strategies. See Supplementary Appendix.

In routine clinical practice, when elderly ischemic stroke patients are admitted to the ICU on the first day, the model can be used to predict the mortality rate of elderly ischemic stroke patients in the ICU within three months. This model may have a positive impact on improving patient outcomes in the following ways: First, it provides individual three-month mortality risk assessments. Interpretable machine learning models can provide an intuitive understanding of patient mortality risk, including factors associated with poor prognosis within three months. The interpretability of the model makes it easier for doctors to understand the causes of mortality risk, which helps improve doctors' understanding of the disease and enables better prevention and treatment strategies. Second, it assists in clinical decision-making and optimizing patient management. It helps provide early warning of the risk of death for elderly ischemic stroke patients with more severe conditions, enabling timely intervention measures to reduce mortality. For instance, if the model predicts a high risk of mortality for a patient, doctors can opt for a more proactive treatment plan. This plan may include more frequent monitoring, comprehensive care, or other interventions that could impact patient outcomes. Third, it improves coordination among healthcare teams. By predicting the level of risk using the model, the healthcare team can enhance their collaboration to ensure that the patient receives optimal medical and nursing care. This can help healthcare professionals rationalize the allocation of resources, allocate healthcare staff and equipment inputs according to patients' priorities, and improve the efficiency and quality of healthcare services.

Despite their relative newness, machine learning methods outperform the currently available tools for healthcare applications. The outcomes of acute ischemic stroke have been predicted using machine learning methods with a better AUC for deep neural network models.[39] Using artificial neural networks, stroke mimics have been distinguished from strokes, and patients at high risk for TIAs and minor strokes have been identified.[40] Prediction modeling employing machine learning shows promise but requires further investigation for different applications.

There are some shortcomings in this study. Due to its retrospective observational design, selection bias could not be eliminated. The selection bias can be caused by non-participation or exclusion because of missing values in our main variables. These biases may affect the results in the following ways. Firstly, selection bias resulting from a retrospective observational design may lead to incomplete and unrepresentative samples. This implies that our sample may not be entirely representative of the entire target population, as certain groups might be more prone to exclusion, or the representativeness of the sample could be compromised due to missing data. Secondly, if certain specific types of participants are more likely to participate in the study, our estimates of these factors may favor the former, leading to biased results. This incompleteness or bias can affect the ability to generalize our findings. Validation in various populations or settings evaluates the generalizability of the findings.

## Conclusions

We offer a user-friendly predictive model for the three-month mortality in elderly patients with IS admitted to the ICU. This model helps improve healthcare quality and clinical decision-making by early identification of high-risk patients. The early warning of the high risk of death for elderly ischemic stroke patients with more severe conditions provides the opportunity for early intervention to reduce mortality.

**ORCID iD:** Wanlin Jin 🔟 https://orcid.org/0000-0001-6622-9631

## References

1. Katan M and Luft A. Global burden of stroke. *Semin Neurol* 2018; 38: 208–211.
2. Feigin VL, Lawes CM, Bennett DA, et al. Stroke epidemiology: a review of population-based studies of incidence, prevalence, and case-fatality in the late 20th century. *Lancet Neurol* 2003; 2: 43–53.
3. Rosamond W, Flegal K, Furie K, et al. Heart disease and stroke statistics–2008 update: a report from the American Heart Association statistics committee and stroke statistics subcommittee. *Circulation* 2008; 117: e25–146.
4. Powell JL and Cook IG. Global ageing in comparative perspective: a critical discussion. *Int J Sociol Soc Policy* 2009; 29: 388–400.
5. Béjot Y. Forty years of descriptive epidemiology of stroke. *Neuroepidemiology* 2022; 56: 157–162.
6. Global, regional, and national burden of stroke and its risk factors, 1990-2019: a systematic analysis for the global burden of disease study 2019. *Lancet Neurol* 2021; 20: 795–820.
7. Chen RL, Balami JS, Esiri MM, et al. Ischemic stroke in the elderly: an overview of evidence. *Nat Rev Neurol* 2010; 6: 256–265.
8. Denti L, Scoditti U, Tonelli C, et al. The poor outcome of ischemic stroke in very old people: a cohort study of its determinants. *J Am Geriatr Soc* 2010; 58: 12–17.
9. Kauffmann J, Grün D, Yilmaz U, et al. Acute stroke treatment and outcome in the oldest old (90 years and older) at a tertiary care medical centre in Germany-a retrospective study showing safety and efficacy in this particular patient population. *BMC Geriatr* 2021; 21: 611.
10. Lyden P, Brott T, Tilley B, et al. Improved reliability of the NIH stroke scale using video training. NINDS TPA stroke study group. *Stroke* 1994; 25: 2220–2226.
11. de Montmollin E, Terzi N, Dupuis C, et al. One-year survival in acute stroke patients requiring mechanical ventilation: a multicenter cohort study. *Ann Intensive Care* 2020; 10: 53.
12. Alonso A, Ebert AD, Kern R, et al. Outcome predictors of acute stroke patients in need of intensive care treatment. *Cerebrovasc Dis* 2015; 40: 10–17.
13. Sonneville R, Gimenez L, Labreuche J, et al. What is the prognosis of acute stroke patients requiring ICU admission? *Intensive Care Med* 2017; 43: 271–272.
14. Haug CJ and Drazen JM. Artificial intelligence and machine learning in clinical medicine, 2023. *N Engl J Med* 2023; 388: 1201–1208.
15. Sagi O and Rokach L. Ensemble learning: a survey. *WIREs Data Min Knowl Discov* 2018; 8: e1249.
16. Giacinto G, Perdisci R, Del Rio M, et al. Intrusion detection in computer networks by a modular ensemble of one-class classifiers. *Inf Fusion* 2008; 9: 69–82.
17. Collins GS, Reitsma JB, Altman DG, et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *J Clin Epidemiol* 2015; 68: 134–143.
18. Pollard TJ, Johnson AEW, Raffa JD, et al. The eICU collaborative research database, a freely available multi-center database for critical care research. *Sci Data* 2018; 5: 180178.
19. Lin J, Yin M, Liu L, et al. The development of a prediction model based on random survival forest for the postoperative prognosis of pancreatic cancer: A SEER-based study. *Cancers (Basel)* 2022; 14: 4667.
20. Lundberg SM, Erion G, Chen H, et al. From local explanations to global understanding with explainable AI for trees. *Nat Mach Intell* 2020; 2: 56–67.
21. Lundberg SM, Nair B, Vavilala MS, et al. Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. *Nat Biomed Eng* 2018; 2: 749–760.
22. Zhang MY, Mlynash M, Sainani KL, et al. Ordinal prediction model of 90-day modified Rankin scale in ischemic stroke. *Front Neurol* 2021; 12: 727171.
23. Li XD and Li MM. A novel nomogram to predict mortality in patients with stroke: a survival analysis based on the MIMIC-III clinical database. *BMC Med Inform Decis Mak* 2022; 22: 92.
24. Someeh N, Mirfeizi M, Asghari-Jafarabadi M, et al. Predicting mortality in brain stroke patients using neural networks: outcomes analysis in a longitudinal study. *Sci Rep* 2023; 13: 18530.
25. Wang W, Rudd AG, Wang Y, et al. Risk prediction of 30-day mortality after stroke using machine learning: a nationwide registry-based cohort study. *BMC Neurol* 2022; 22: 195.
26. Flint AC, Cullen SP, Faigeles BS, et al. Predicting long-term outcome after endovascular stroke treatment: the totaled health risks in vascular events score. *AJNR Am J Neuroradiol* 2010; 31: 1192–1196.
27. Smith EE, Shobha N, Dai D, et al. Risk score for in-hospital ischemic stroke mortality derived and validated within the get with the guidelines-stroke program. *Circulation* 2010; 122: 1496–1504.
28. Zhu H and Hill MD. Stroke: the Elixhauser Index for comorbidity adjustment of in-hospital case fatality. *Neurology* 2008; 71: 283–287.
29. Namale G, Kamacooko O, Makhoba A, et al. Predictors of 30-day and 90-day mortality among hemorrhagic and ischemic stroke patients in urban Uganda: a prospective hospital-based cohort study. *BMC Cardiovasc Disord* 2020; 20: 442.
30. Birkner MD, Kalantri S, Solao V, et al. Creating diagnostic scores using data-adaptive regression: an application to prediction of 30-day mortality among stroke victims in a rural hospital in India. *Ther Clin Risk Manag* 2007; 3: 475–484.
31. Szlachetka WA, Pana TA, Mamas MA, et al. Predicting 10-year stroke mortality: development and validation of a nomogram. *Acta Neurol Belg* 2022; 122: 685–693.
32. Wang A, Tian X, Zuo Y, et al. High lactate dehydrogenase was associated with adverse outcomes in patients with acute ischemic stroke or transient ischemic attack. *Ann Palliat Med* 2021; 10: 10185–10195.
33. Carter AM, Catto AJ, Mansfield MW, et al. Predictive variables for mortality after acute ischemic stroke. *Stroke* 2007; 38: 1873–1880.
34. You S, Sun X, Zhou Y, et al. The prognostic significance of white blood cell and platelet count for inhospital mortality and pneumonia in acute ischemic stroke. *Curr Neurovasc Res* 2021; 18: 427–434.

35. Uehara T, Ohara T, Minematsu K, et al. Predictors of stroke events in patients with transient ischemic attack attributable to intracranial stenotic lesions. *Intern Med* 2018; 57: 295–300.

36. Zhu HJ, Sun X, Guo ZN, et al. Prognostic values of serum alkaline phosphatase and globulin levels in patients undergoing intravenous thrombolysis. *Front Mol Neurosci* 2022; 15: 932075.

37. Gattringer T, Posekany A, Niederkorn K, et al. Predicting early mortality of acute ischemic stroke. *Stroke* 2019; 50: 349–356.

38. Pepper DJ, Sun J, Welsh J, et al. Increased body mass index and adjusted mortality in ICU patients with sepsis or septic shock: a systematic review and meta-analysis. *Crit Care* 2016; 20: 181.

39. Heo J, Yoon JG, Park H, et al. Machine learning-based model for prediction of outcomes in acute stroke. *Stroke* 2019; 50: 1263–1265.

40. Chan KL, Leng X, Zhang W, et al. Early identification of high-risk TIA or Minor stroke using artificial neural network. *Front Neurol* 2019; 10: 171.