

1 Predictive Models for Secondary Epilepsy in Patients with Acute Ischemic Stroke Within One  
2 Year

3 Jinxin Liu<sup>1</sup>, Haoyue He<sup>1,2†</sup>, Yanglingxi Wang<sup>1</sup>, Jun Du<sup>6</sup>, Kaixin Liang<sup>7</sup>, Jun Xue<sup>8</sup>, Yidan Liang<sup>1</sup>,  
4 Peng Chen<sup>1</sup>, Shanshan Tian<sup>5</sup>, Yongbing Deng<sup>1,3,4</sup>,

5 1 Department of Neurosurgery, Chongqing Emergency Medical Center, Chongqing University  
6 Central Hospital, School of Medicine, Chongqing University, Chongqing, China

7 2 Bioengineering College of Chongqing University, Chongqing, China

8 3 Chongqing Key Laboratory of Emergency Medicine

9 4 Jinfeng Laboratory, Chongqing, China

10 5 Department of Prehospital Emergency, Chongqing University Central Hospital, Chongqing  
11 Emergency Medical Center, Chongqing, China

12 6 Department of Neurosurgery, Chongqing University Qianjiang Hospital, Chongqing, China

13 7 Department of Neurosurgery, Yubei District Hospital of Traditional Chinese Medicine,  
14 Chongqing, China

15 8 Department of Neurosurgery, Bishan hospital of Chongqing Medical University, Chongqing,  
16 China

17 †These authors have contributed equally to this work and share first authorship.

18

19 Corresponding author: Yongbing Deng Email: [dyb0913@cqu.edu.cn](mailto:dyb0913@cqu.edu.cn)

20 Shanshan Tian Email: [710836163@qq.com](mailto:710836163@qq.com)

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

### 38 **Data availability statement**

39 *The codes,models,analysis results was uploaded at <https://github.com/conanan/lasso-ml> . The*  
40 *full dataset can be provided for researchers if needed by the corresponding author.*

41

### 42 **Acknowledgements**

43 *The authors would like to thank the colleagues in the information and imaging departments for*  
44 *their hard work contributing to the final research results.*

### 45 **Ethics approval statement**

46 *We confirm that we have read the Journal's position on issues involved in ethical publication and*  
47 *affirm that this report is consistent with those guidelines.*

### 48 **Funding statement**

49 *The research is funded by Central University basic research young teachers and students research*  
50 *ability promotion sub-project(2023CDJYGRH-ZD06);by Emergency Medicine Chongqing Key*  
51 *Laboratory Talent Innovation and development joint fund project (2024RCCX10) .*

### 52 **Conflict of interests**

53 *The authors have no relevant conflicts of interest to disclose.*

54

### 55 **Patient consent statement**

56 *This study was a retrospective study and only deidentified patient data were collected,*  
57 *exempting the need for patient informed consent rights.*

58

### 59 **Permission to reproduce material from other sources**

60 *There are no reproduce material from other sources.*

61

### 62 **Clinical trial registration**

63 *The trail number is RS202406.*

64

65

## 66 Abstract

67 **Objective:** Post-stroke epilepsy (PSE) is a major complication that worsens both prognosis and  
68 quality of life in patients with ischemic stroke. This study aims to develop an interpretable  
69 machine learning model to predict PSE using medical records from four hospitals in Chongqing.

70 **Methods:** We collected and analyzed medical records, imaging reports, and laboratory test  
71 results from 21,459 patients diagnosed with ischemic stroke. Traditional univariable and  
72 multivariable statistical analyses were performed to identify key predictive factors. The dataset  
73 was divided into a 70% training set and a 30% testing set. To address class imbalance, the  
74 Synthetic Minority Oversampling Technique combined with Edited Nearest Neighbors was used.  
75 Nine widely applied machine learning algorithms were evaluated and compared using relevant  
76 prediction metrics. SHAP (SHapley Additive exPlanations) was used to interpret the model,  
77 assessing the contributions of different features.

78 **Results:** Regression analyses showed that complications such as hydrocephalus, cerebral hernia,  
79 and deep vein thrombosis, as well as brain regions (frontal, parietal, and temporal lobes),  
80 significantly contributed to PSE. Factors like age, gender, NIH Stroke Scale (NIHSS) scores, and  
81 laboratory results such as WBC count and D-dimer levels were associated with a higher risk of  
82 PSE. Among the machine learning models, tree-based methods such as Random Forest,  
83 XGBoost, and LightGBM demonstrated strong predictive performance, achieving an AUC of  
84 0.99.

85 **Conclusion:** Our model successfully predicts PSE risk, with tree-based models showing superior  
86 performance. The NIHSS score, WBC count, and D-dimer were identified as the most important  
87 predictors.

88

## 89 Introduction

90 Stroke is the second leading cause of death globally, with an annual mortality of  
91 approximately 5.5 million, and it is also the leading cause of disability, accounting for 50% of  
92 cases worldwide [1]. Ischemic stroke comprises about 80% of all stroke cases [2][3]. Post-stroke  
93 epilepsy (PSE) is a common complication, with studies reporting that 3-30% of stroke patients  
94 develop epilepsy, which adversely affects their prognosis and quality of life [4]. PSE can worsen  
95 cognitive, psychiatric, and physical impairments already caused by cerebrovascular disease and  
96 related conditions [5]. The highest incidence of PSE occurs within the first year after an acute  
97 stroke, accounting for nearly half of the cases [2]. Thus, early prediction and intervention for  
98 PSE, especially in ischemic strokes, are critical.

99           Currently, most studies rely on clinical data to build statistical models using survival  
100 analysis, Cox regression [2][6], and multiple linear regression [7] to create basic models for PSE  
101 prediction. Last year, Lin et al. developed a radiomics-based model that outperformed  
102 conventional clinical models in predicting PSE related to intracerebral hemorrhage (ICH). They  
103 suggested that a combined radiomics-clinical model could improve the assessment of individual  
104 PSE risk after the first occurrence of ICH, facilitating early diagnosis and treatment [8]. However,  
105 subsequent research raised concerns about the use of radiomics, indicating a need for further  
106 investigation [9]. Overall, research on PSE prediction remains limited, with most studies  
107 focusing on specific risk factors [10][11][8][12] and building simple models, without proposing  
108 more comprehensive and scientifically robust prediction models.

109           Machine learning has gained attention as a powerful tool for building medical models due to  
110 its ability to process large datasets and complex information. It has been increasingly applied in  
111 neuroscience and clinical prediction [13][14][15]. Previous studies have used machine learning  
112 to explore post-stroke cognitive impairments [16], predict stroke and myocardial infarction risks  
113 in large artery vasculitis patients [14], develop post-stroke depression models based on liver  
114 function tests [17], and predict hematoma expansion in traumatic brain injury (TBI) [18].  
115 Machine learning models can automatically manage both linear and complex nonlinear  
116 relationships between variables and offer insights into how different factors contribute to the  
117 prediction target—something that is difficult for traditional statistical models. However, machine  
118 learning requires substantial amounts of data and is prone to overfitting with small sample sizes.  
119 The quality and volume of input data are critical for the algorithm to detect underlying patterns  
120 and make accurate predictions.

121           This study aims to identify key risk factors from various features extracted from the clinical  
122 records and test data of ischemic stroke patients. Using these features, we will develop a machine  
123 learning-based prediction model for PSE. By leveraging early admission data, we seek to  
124 automatically predict the likelihood of PSE occurrence and provide guidance for clinical  
125 decision-making and patient care.

## 126 **Result**

### 127 **Filling of missing data**

128 Missing values were filled using a Random Forest (RF) model, handling one feature at a  
129 time. The imputed features were: Plt, WBC, RBC, HbA1c, CRP, TG, LDL, HDL, AST, ALT,  
130 bilirubin, albumin, urea, creatinine, BUA, PT, APTT, TT, INR, D-dimer, fibrinogen, CK, CK-  
131 MB, LDH, HBDH, IMA, lactate, anion gap, TCO2, and NIHSS.

### 132 **Characteristics of study participants**

133 A total of 21,459 patients were included in the study. The training set consisted of 15,021  
134 patients, with a PSE incidence of 4.3%. The test set contained 6,438 patients, also with a 4.3%  
135 incidence of PSE. The external validation cohort included 536 patients from three hospitals. The  
136 statistical details of the clinical characteristics are presented in Table 1.

137 Statistical analysis indicated that patients with a higher likelihood of developing PSE had  
138 complications such as uremia, a history of DVT, atrial fibrillation, hyperuricemia, cerebral hernia,  
139 and hydrocephalus. The affected brain regions included the frontal, parietal, occipital, and  
140 temporal lobes, as well as the cortex, subcortex, basal ganglia, and hypothalamus. General  
141 characteristics included age, gender, and NIHSS score. Laboratory indicators associated with a  
142 higher risk of PSE included WBC count, HbA1C, CRP, triglycerides, AST, ALT, bilirubin, urea,  
143 uric acid, APTT, PT, D-dimer, CK, CK-MB, LDH, HBDH, IMA, lactate, and anion gap.  
144 Additionally, significant p-values were found for fatty liver, coronary heart disease,  
145 hyperlipidemia, and HDL, with low or negative values of these indicators linked to a higher risk  
146 of secondary complications. The results of the statistical analyses, as well as the univariate and  
147 multivariate regression analyses, are detailed in Tables 1, 2, and 3.

### 148 **Performance of machine learning models**

149 The relevant performance indicators of the machine learning models are presented in Table  
150 4, while the ROC curves, calibration curve, and decision curve analysis (DCA) are shown in  
151 Figure 3. Among all models, tree-based models such as Random Forest (RF), XGBoost, and  
152 LightGBM had the highest AUC scores, outperforming other models. Notably, Random Forest  
153 had the highest positive predictive value (PPV) at 0.864, which was the most significant metric  
154 in our models. Complex machine learning algorithms performed better than traditional logistic  
155 regression. The Brier score of the calibration curve was 0.006, and the DCA demonstrated good  
156 clinical decision-making benefits, indicating strong practical value. In the external validation  
157 cohort, we used RF for predictions, achieving a sensitivity of 0.91 and a PPV of 0.95, confirming  
158 the model's strong predictive capability.

### 159 **Analysis of SHAP risk factors**

160 Figure 4 shows the SHAP (Shapley Additive Explanations) values, individual decision  
161 attempts, and overall decision curves. Among general characteristics, females had a higher rate  
162 of PSE. A higher NIHSS score was associated with a higher incidence of PSE. Additionally,  
163 elevated values of WBC count, D-dimer, CRP, AST, CK-MB, HbA1c, bilirubin, TCO<sub>2</sub>, and  
164 LDH at admission were linked to a greater likelihood of developing PSE. Conversely, lower  
165 levels of HBDH, PLT, and APTT were also associated with a higher probability of PSE. The  
166 specific brain regions affected did not have a significant individual effect on the overall outcome.  
167 Among complications, hypertension was more strongly associated with PSE development, while  
168 other conditions, such as coronary heart disease, diabetes, hyperlipidemia, and fatty liver, were  
169 less likely to be related to the outcome. We used the force plot of the first patient to illustrate  
170 how different features influenced the prediction. In this case, a prolonged APTT time contributed  
171 the most to PSE, followed by elevated AST levels, while a low NIHSS score contributed  
172 negatively to the final result. The decision plot aggregated model decisions to show how  
173 complex models arrived at their predictions.

## 174 Discussion

175 Our study used comprehensive clinical, imaging, and laboratory data from stroke patients to  
176 develop a predictive model using machine learning algorithms. This model achieved an AUC  
177 score above 0.95, demonstrating more accurate predictions compared to traditional statistical  
178 methods. Our research revealed that tree-based ensemble models provided superior predictive  
179 performance, especially when handling large datasets with high-dimensional features.

180 During the modeling process, due to the extreme imbalance between negative and positive  
181 samples, we applied the SMOTEENN technique to resample the dataset, improving the  
182 performance of the machine learning models. Through SHAP analysis, we conducted  
183 interpretability assessments of the model and identified the importance of different features.

184 In our study, age and NIHSS scores were treated as continuous variables. We found that  
185 female patients, older individuals, and those with higher NIHSS scores were more likely to  
186 develop PSE, consistent with recent studies. Higher NIHSS scores, indicating more severe  
187 strokes, significantly increased the risk of complications, second only to white blood cell (WBC)  
188 count and D-dimer in our model [5][19][10][20]. However, there are differing views on the  
189 effect of age. Some studies [5][21] suggest that age below 65 is a high-risk factor, which aligns  
190 with our findings, while other studies [22] have found that advanced age is the key factor.  
191 Yamada et al. [21] also agreed with our study, indicating that female patients have a higher risk  
192 of complications. On the other hand, Waafi et al. [10] reported that male patients are 3.325 times  
193 more likely to develop complications, which contradicts our findings.

194 Previous research has shown that patients with diabetes, dyslipidemia, hypertension,  
195 depression, or dementia are at higher risk of developing vascular epilepsy [12]. In our study,  
196 statistical analysis and multiple machine learning (ML) models examined the relationship  
197 between comorbidities and complications. We found that patients with coronary heart disease,  
198 diabetes, fatty liver, hyperlipidemia, or large artery stenosis or plaques (CCA and ICA) were less  
199 likely to develop epilepsy. According to the TOAST classification, ischemic stroke is divided  
200 into five categories: large artery atherosclerosis, cardioembolism, small vessel occlusion, other  
201 determined etiology, and undetermined etiology. Patients with multiple comorbidities often fall  
202 into the large artery atherosclerosis and cardioembolism categories, which are more clearly  
203 defined and easier to treat, resulting in a lower likelihood of epilepsy. In contrast, strokes of  
204 undetermined etiology tend to have worse prognoses and are more likely to lead to epilepsy.  
205 Among patients with diabetes, higher HbA1c levels indicate poor blood sugar control and a  
206 higher risk of complications. Patients with better control of their blood sugar have a lower  
207 overall risk of developing complications.

208 Alain et al. found that cortical infarction is more likely to lead to epilepsy in patients  
209 hospitalized with anterior circulation ischemic stroke [23]. Lin et al. found that factors such as  
210 cortical involvement and intracerebral hemorrhage volume increase the likelihood of PSE, which  
211 is consistent with our findings [8]. Al-Sahli et al. also suggested that cortical brain injury and  
212 large-area lesions raise the risk of PSE [5][21]. In our study, statistics showed that both cortical  
213 and subcortical involvement increased the likelihood of PSE, but these regions had less influence  
214 compared to other features and were not selected in the LASSO regression.

215 Previous studies have identified acute infection as a risk factor for ischemic stroke [24]. C-  
216 reactive protein (CRP) reflects inflammation levels and is an independent prognostic factor [25].  
217 In our study, both regression and SHAP analysis indicated that WBC count had a significant  
218 impact among routine blood test parameters, even surpassing the NIHSS score in SHAP analysis.  
219 A high WBC count may indicate severe inflammation or infection, as well as increased blood  
220 viscosity, making patients more prone to secondary complications. In general, a high red blood  
221 cell count and low platelet count also contributed to an increased risk of complications.

222 A large-scale study on Chinese individuals found a negative correlation between plasma  
223 high-density lipoprotein cholesterol (HDL-C) levels and the risk of ischemic stroke, a weak  
224 positive correlation between plasma triglyceride (TG) levels and stroke risk, and a strong  
225 correlation between plasma low-density lipoprotein cholesterol (LDL-C) and apolipoprotein B  
226 levels [26]. High HDL-C levels are linked to better prognosis [27]. Our study aligns with these  
227 findings, showing that high LDL-C, low HDL-C, and elevated TG levels are more likely to result  
228 in PSE. This can be understood as high cholesterol and triglyceride levels increase blood  
229 viscosity and contribute to vascular sclerosis, promoting clot formation [12][28][29]. Higher D-  
230 dimer levels indicate more significant brain tissue damage, increasing the likelihood of PSE. In  
231 general, lower activated partial thromboplastin time (APTT) and fibrinogen levels are associated  
232 with higher PSE risk, while INR, PT, and TT have a smaller impact. Among liver function  
233 indicators, aspartate aminotransferase (AST) had the greatest influence on PSE. High AST, low  
234 alanine aminotransferase (ALT), and low albumin levels also had some impact. Lingling Ding et  
235 al. found that liver enzyme subgroups defined by ALT and AST were linked to higher risks of  
236 adverse outcomes [30], which is consistent with our findings.

237 Studies have also shown that renal function biomarkers such as urinary microalbumin,  
238 cystatin C, and creatinine are associated with higher stroke recurrence rates and poorer prognosis  
239 [30]. In our study, low urea levels and high uric acid levels had a negative impact [31][32][33].  
240 Our research supports these conclusions. Elevated uric acid levels at admission were positively  
241 associated with PSE, although patients with a prior diagnosis of hyperuricemia were less likely  
242 to develop epilepsy. Since uric acid acts as a strong antioxidant and has neuroprotective  
243 properties [34], patients with normal liver and kidney function and mild hyperuricemia may have  
244 greater resilience in emergencies [35][36]. However, excessively high uric acid levels suggest  
245 metabolic disorders and poor liver and kidney function, which are linked to a poor prognosis.

246 When stroke patients are admitted, cardiac enzyme tests are often conducted to rule out  
247 myocardial ischemia. However, studies have shown that elevated CK-MB in stroke patients may  
248 not be solely heart-related [37]. Cardiac enzymes are important prognostic indicators [38][39]  
249 and have been incorporated into stroke scores [40]. Some studies have reported a higher  
250 incidence of abnormal serum cardiac enzyme levels in the acute phase of stroke. While the  
251 abnormalities are not related to the stroke type, they are associated with stroke severity, with  
252 patients exhibiting consciousness disorders having a significantly higher incidence of abnormal  
253 cardiac enzymes than those without such disorders [41]. In our study, CK, CK-MB, and IMA in  
254 the cardiac enzyme profile had a significant impact and high predictive value, though further  
255 research is required to understand the specific mechanisms involved [34].

256 Although our study incorporated extensive clinical, imaging, and laboratory data to build  
257 more accurate prediction models using machine learning algorithms, surpassing traditional  
258 statistical methods, there were still several limitations in the modeling process.

259 While the current study offers valuable insights, the data sample may not be fully  
260 representative, and the model's generalizability requires further evaluation. Although the data  
261 was collected from multiple tertiary hospitals and includes over 20,000 cases, earlier data was  
262 lost due to hospital system upgrades. The dataset mainly reflects patients diagnosed within the  
263 past five years and is predominantly from the Chongqing region, which may limit the model's  
264 applicability to other geographic areas.

265 Additionally, the retrospective nature of the study led to the absence of some important  
266 predictive indicators. Many potentially valuable features, such as hemorheology,  
267 thromboelastography, and hormone levels, were missing and had to be excluded. Including these  
268 features could potentially improve the model's accuracy.

269 To enhance the predictive power of the model, it would be beneficial to incorporate more  
270 data beyond baseline patient characteristics. The current analysis primarily used the results from  
271 the first examination upon admission, without fully utilizing information from subsequent exams.  
272 In future research, recurrent neural networks could be employed to extract features from the  
273 entire sequence of examinations more comprehensively.

274 To strengthen the study further, data standardization should be improved, and the number of  
275 cases and key indicators should continue to grow. Additionally, it would be advantageous to  
276 explore more advanced scientific methods, such as deep learning, and utilize all available data to  
277 improve prediction accuracy.



## 278 **Materials and methods**

### 279 **Research patients**

280 This study retrospectively included all stroke patients admitted to the Chongqing  
281 Emergency Center between June 2017 and June 2022 for the development of the prediction  
282 model. Data from three external validation centers—Qianjiang Central Hospital, Bishan District  
283 People's Hospital, and Yubei District Traditional Chinese Medicine Hospital—were collected  
284 between July 2022 and July 2023 to validate and evaluate the model externally. The external  
285 validation cohort emphasized collecting positive cases to test the model's ability to identify these  
286 cases accurately.

287 Inclusion criteria: (1) Age between 18 and 90 years at admission; (2) Diagnosed with acute  
288 ischemic stroke and hospitalized for treatment.

289 Exclusion criteria: (1) Patients with a history of stroke or transient ischemic attack (TIA); (2)  
290 Patients with a history of other conditions such as traumatic brain injury, intracranial tumors, or  
291 cerebral vascular malformations that may cause epilepsy; (3) Patients with a history of epilepsy  
292 or who have received antiseizure medications for the prevention of seizures or for other diseases  
293 (such as migraine or psychiatric disorders); (4) Patients who died within 72 hours after stroke  
294 onset.

295 This study collected de-identified data from relevant patients to build a multi-modal stroke  
296 patient database. The study protocol was approved by the Ethics Committees of Chongqing  
297 University Center Hospital, Chongqing University Qianjiang Central Hospital, Bishan District  
298 People's Hospital, and Yubei District Traditional Chinese Medicine Hospital.

299 The selection process is outlined in Figure 1. A total of 42,079 records were retrieved from  
300 the stroke database, and 24,733 patients were diagnosed with ischemic or lacunar stroke with  
301 new onset. Hemorrhagic strokes (4,565), a history of stroke (2,154), TIA (3,570), unclear cause  
302 strokes (561), and records with missing essential data (6,496) were excluded. Patients whose  
303 seizures might have been caused by other factors (such as brain tumors, intracranial vascular  
304 malformations, or traumatic brain injury) (865), those with a seizure history (152), and patients  
305 who died in the hospital (1,444) were also excluded. Additionally, patients lost to follow-up  
306 (those without outpatient records or unreachable by phone) or who died within three months of  
307 the stroke incident (813) were excluded. Finally, 21,459 cases were included in the study.

308

### 309 **Data collection**

310 We extracted all relevant records and data from the hospital databases. Using  
311 PostgreSQL, we wrote Structured Query Language (SQL) to manage the data as follows:

312 (1) General Information: This included gender, age, and NIH Stroke Scale (NIHSS) score at  
313 admission.

314 (2) Comorbidities and Complications: These included uremia, previous deep vein  
315 thrombosis (DVT), diabetes mellitus, hypertension, coronary atherosclerosis, atrial fibrillation,

316 cerebral hernia, hydrocephalus, hypoproteinemia, hyperuricemia, hyperlipidemia, internal carotid  
317 stenosis, and common carotid stenosis.

318 (3) Brain Involvement (CT or MRI records): We recorded involvement of the cortical lobes  
319 and subcortical areas, including the frontal, parietal, temporal, occipital, and insular lobes, as  
320 well as the basal ganglia, internal capsule, brain stem, cerebellum, periventricular area, centrum  
321 semiovale, and thalamus. The extent of cortical involvement (frontal, parietal, temporal, occipital,  
322 and insular lobes) was scored, with each lobe contributing 1 point. Similarly, subcortical  
323 involvement (basal ganglia, internal capsule, brain stem, periventricular area, thalamus, and  
324 cerebellum) was scored with each area contributing 1 point.

325 (4) Vascular Involvement (CTA, MRA, or DSA records): We recorded the presence of  
326 vascular stenosis or occlusion in the anterior cerebral artery (ACA), middle cerebral artery  
327 (MCA), posterior cerebral artery (PCA), vertebral artery (VA), and basilar artery (BA).

328 (5) Key Laboratory Indicators: These included blood lipids such as triglycerides (TG), high-  
329 density lipoprotein cholesterol (HDL), and low-density lipoprotein cholesterol (LDL); liver  
330 function indicators such as alanine transaminase (ALT), aspartate aminotransferase (AST),  
331 bilirubin, and albumin; renal function markers such as urea, blood uric acid (BUA), and  
332 creatinine; blood gas parameters such as lactate, anion gap, and total carbon dioxide (TCO<sub>2</sub>);  
333 coagulation markers such as international normalized ratio (INR), prothrombin time (PT),  
334 activated partial thromboplastin time (APTT), thrombin time (TT), D-dimer, and fibrinogen; and  
335 myocardial enzymes such as creatine kinase (CK), creatine kinase isoenzyme (CK-MB), lactate  
336 dehydrogenase (LDH), ischemic modified albumin (IMA), and  $\alpha$ -hydroxybutyrate  
337 dehydrogenase (HBDH).

338

## 339 Data processing and model building

340 Processing of Missing Data: We recorded all laboratory indicators from the first set of tests  
341 after stroke admission (every stroke patient undergoes routine blood tests, and liver and kidney  
342 function assessments). Indicators with more than 10% missing data were excluded. The  
343 remaining indicators with missing values were imputed using the random forest algorithm with  
344 default parameters. We processed the features in order of missing values, starting with those that  
345 had the least missing data (as this requires the least information for imputation). When imputing  
346 a feature, missing values in other features were temporarily replaced with 0. After each  
347 regression prediction, the predicted value was inserted into the original feature matrix before  
348 proceeding to the next feature. Once all features were processed, the dataset was complete.

349 Distribution of Characteristics: We used univariate analysis to compare the distribution of  
350 characteristics between the PSE-negative and PSE-positive groups. The data were then divided  
351 into a training set and a test set in a 7:3 ratio.

352 Processing of Unbalanced Data: Given the low incidence of PSE and the small proportion  
353 of positive cases, we augmented the positive data in the training set using the Synthetic Minority  
354 Over-sampling Technique combined with Edited Nearest Neighbors (SMOTEENN). The  
355 SMOTEENN method from the imblearn Python package was applied with default parameters,  
356 and a random seed of 42 was set to ensure reproducibility.

357

358 Processing of Categorical Data: For categorical variables, we used the one-hot encoding  
359 method for transformation. We then applied the LASSO method to the training set to identify the  
360 most important features.

361 Model Building: First, we used LASSO regression to select the 20 most important features.  
362 We then employed 9 commonly used machine learning methods, including Naive Bayes,  
363 Logistic Regression, Decision Tree, Random Forest, Gradient Boosting, Multi-Layer Perceptron,  
364 XGBoost, LightGBM, and K-Nearest Neighbors. Hyperparameters for each model were  
365 optimized through grid search to enhance performance. Model evaluation metrics included  
366 accuracy, sensitivity, specificity, F1-score, positive predictive value, and negative predictive  
367 value. We also generated ROC curves, calibration curves, and decision curves to further assess  
368 model performance. An independent external validation dataset was used to evaluate the  
369 generalization ability of the selected model. Lastly, we applied the SHAP algorithm to interpret  
370 the best-performing model, analyzing the contribution of each feature to the model's predictions  
371 and their clinical relevance. Through this process of model development, optimization, and  
372 interpretation, we constructed a machine learning model with strong predictive performance and  
373 interpretability, offering valuable support for clinical decision-making.

374

## 375 Statistical approach

376 PostgreSQL v15 (<http://www.postgresql.org/>) was used to search and extract data from the  
377 local database. The open-source statistical package "Scipy.stats" in Python was used for  
378 statistical analysis. The details of the univariate significance analysis for each feature are as  
379 follows:

380 The Shapiro-Wilk test was applied to assess the normality of each feature's distribution. For  
381 features that did not follow a normal distribution, the Mann-Whitney U test was used to evaluate  
382 their significance in relation to the target variable. For features with a normal distribution, the  
383 Levene test was performed to evaluate the homogeneity of variances. Features with  
384 homogeneous variances were analyzed using the Student's t-test for significance, while those  
385 with heterogeneous variances were analyzed using Welch's t-test.

386 Confidence intervals for AUC values and Brier scores were calculated using 1,000  
387 bootstrap resampling iterations on the datasets. Binary classification thresholds for the predicted  
388 probabilities from all models were established using the maximum Youden index derived from  
389 the training cohort.

390 Throughout the study, a two-tailed p-value of less than 0.05 was considered statistically  
391 significant.

392 All the code used in this study was uploaded to <https://github.com/conanan/lasso-ml>.

## 393 Conclusion

394 We developed an interpretable machine learning model to predict the risk of post-stroke  
395 epilepsy (PSE) in hospitalized patients with ischemic stroke. Using a large dataset of medical

396 records, our artificial intelligence model demonstrates strong predictive performance for PSE.  
397 The key predictors identified by the model include NIHSS score, D-dimer levels, lactate levels,  
398 and white blood cell count, along with liver function and cardiac enzyme profile indicators. The  
399 model's transparency and interpretability can build trust among clinicians and support decision-  
400 making. While the results are promising, further prospective studies are necessary to validate the  
401 clinical utility of this tool before it can be applied in real-world settings.

402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412

413 [1] Feigin V L, Krishnamurthi R V, Theadom A M, et al.. Global, Regional, and National  
414 Burden of Neurological Disorders during 1990–2015: A Systematic Analysis for the Global  
415 Burden of Disease Study 2015[J]. *The Lancet Neurology*, 2017, 16(11): 877–897.

416 [2] Galovic M, Döhler N, Erdélyi-Canavese B, et al.. Prediction of Late Seizures after  
417 Ischaemic Stroke with a Novel Prognostic Model (the SeLECT Score): A Multivariable  
418 Prediction Model Development and Validation Study[J]. *The Lancet Neurology*, 2018, 17(2):  
419 143.

420 [3] Krishnamurthi R V, Feigin V L, Forouzanfar M H, et al.. Global and Regional Burden  
421 of First-Ever Ischaemic and Haemorrhagic Stroke during 1990–2010: Findings from the  
422 Global Burden of Disease Study 2010[J]. *The Lancet Global Health*, 2013, 1(5): e259–e281.

423 [4] Zhao Y, Li X, Zhang K, et al.. The Progress of Epilepsy after Stroke[J]. *Curr*  
424 *Neuropharmacol*, 2018, 16(1): 71–78.

425 [5] Al-Sahli O a M, Tibekina L, Subbotina O P, et al.. Post-Stroke Epileptic Seizures: Risk  
426 Factors, Clinical Presentation, Principles of Diagnosis and Treatment[J]. *Epilepsy and*  
427 *paroxysmal conditions*, 2023, 15(2): 148–159.

428 [6] Chen Z, Churilov L, Chen Z, et al.. Association between Implementation of a Code  
429 Stroke System and Poststroke Epilepsy[J]. *Neurology*, 2018, 90(13): e1126–e1133.

430 [7] Merkler A E, Gialdini G, Lerario M P, et al.. Population-Based Assessment of the  
431 Long-Term Risk of Seizures in Survivors of Stroke[J]. *Stroke*, 2018, 49(6): 1319–1324.

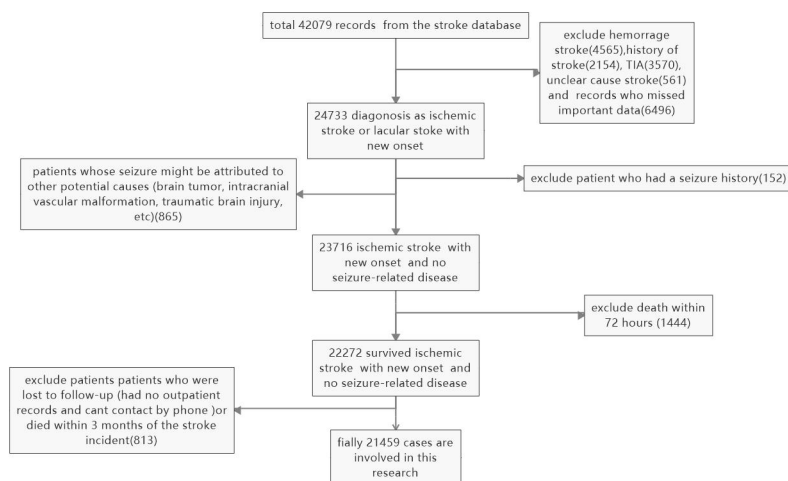
432 [8] Lin R, Lin J, Xu Y, et al.. Development and Validation of a Novel Radiomics-Clinical  
433 Model for Predicting PSE after First-Ever Intracerebral Haemorrhage[J]. *European*  
434 *Radiology*, 2023, 33(7): 4526–4536.

- 435 [9] Pszczolkowski S, Law Z K. Editorial Comment on «Development and Validation of a  
436 Novel Radiomics-Clinical Model for Predicting PSE after First-Ever Intracerebral  
437 Haemorrhage» [J]. *European Radiology*, 2023, 33(7): 4524–4525.
- 438 [10] Waafi A K, Husna M, Damayanti R, et al.. Clinical Risk Factors Related to PSE Patients  
439 in Indonesia: A Hospital-Based Study[J]. *Egyptian Journal of Neurology, Psychiatry and  
440 Neurosurgery*, 2023, 59(1).
- 441 [11] Herzig-Nichtweiß J, Salih F, Berning S, et al.. Prognosis and Management of Acute  
442 Symptomatic Seizures: A Prospective, Multicenter, Observational Study[J]. *Annals of  
443 Intensive Care*, 2023, 13(1).
- 444 [12] Pitkänen A, Roivainen R, Lukasiuk K. Development of Epilepsy after Ischaemic  
445 Stroke[J]. *The Lancet Neurology*, 2016, 15(2): 185–197.
- 446 [13] The Artificial Intelligence Revolution in Stroke Care: A Decade of Scientific Evidence  
447 in Review[J]. *World Neurosurgery*, Elsevier, 2024.
- 448 [14] Predicting Stroke and Myocardial Infarction Risk in Takayasu Arteritis with  
449 Automated Machine Learning Models[J]. *iScience*, Elsevier, 2023, 26(12): 108421.
- 450 [15] Daidone M, Ferrantelli S, Tuttolomondo A, et al.. Machine Learning Applications in  
451 Stroke Medicine: Advancements, Challenges, and Future Prospective[J]. *Neural  
452 Regeneration Research*, 2024, 19(4): 769–773.
- 453 [16] Lee M, Yeo N-Y, Ahn H-J, et al.. Prediction of Post-Stroke Cognitive Impairment after  
454 Acute Ischemic Stroke Using Machine Learning[J]. *Alzheimer's Research and Therapy*, 2023,  
455 15(1).
- 456 [17] Gong J, Zhang Y, Zhong X, et al.. Liver Function Test Indices-Based Prediction Model  
457 for Post-Stroke Depression: A Multicenter, Retrospective Study[J]. *BMC Medical Informatics  
458 and Decision Making*, 2023, 23(1).
- 459 [18] He H, Liu J, Li C, et al.. Predicting Hematoma Expansion and Prognosis in Cerebral  
460 Contusions: A Radiomics-Clinical Approach[J]. *Journal of Neurotrauma*, 2024:  
461 neu.2023.0410.
- 462 [19] Lin R, Yu Y, Wang Y, et al.. Risk of PSE Following Stroke-Associated Acute  
463 Symptomatic Seizures[J]. *Frontiers in Aging Neuroscience*, 2021, 13.
- 464 [20] Zöllner J P, Misselwitz B, Kaps M, et al.. National Institutes of Health Stroke Scale  
465 (NIHSS) on Admission Predicts Acute Symptomatic Seizure Risk in Ischemic Stroke: A  
466 Population-Based Study Involving 135,117 Cases[J]. *Scientific Reports*, 2020, 10(1).
- 467 [21] Yamada S, Nakagawa I, Tamura K, et al.. Investigation of Poststroke Epilepsy  
468 (INPOSE) Study: A Multicenter Prospective Study for Prediction of Poststroke Epilepsy[J]. *J  
469 Neurol*, 2020, 267(11): 3274–3281.

- 470 [22] Lidetu T, Zewdu D. Incidence and Predictors of Post Stroke Seizure among Adult  
471 Stroke Patients Admitted at Felege Hiwot Compressive Specialized Hospital, Bahir Dar,  
472 North West Ethiopia, 2021: A Retrospective Follow up Study[J]. BMC Neurology, 2023,  
473 23(1).
- 474 [23] Lekoubou A, Ssentongo P, Maffie J, et al.. Associations of Small Vessel Disease and  
475 Acute Symptomatic Seizures in Ischemic Stroke Patients[J]. Epilepsy & Behavior, 2023, 145:  
476 109233.
- 477 [24] Bova I Y, Bornstein N M, Korczyn. Acute Infection as a Risk Factor for Ischemic  
478 Stroke[J]. Stroke, 1996, 27(12): 2204–2206.
- 479 [25] Di Napoli M, Papa F, Bocola V. C-Reactive Protein in Ischemic Stroke an Independent  
480 Prognostic Factor[J]. Stroke, 2001, 32(4): 917–924.
- 481 [26] Sun L, Clarke R, Bennett D, et al.. Causal Associations of Blood Lipids with Risk of  
482 Ischemic Stroke and Intracerebral Hemorrhage in Chinese Adults[J]. Nat Med, Nature  
483 Publishing Group, 2019, 25(4): 569–574.
- 484 [27] Bandedali S, Farmer J. High-Density Lipoprotein and Atherosclerosis: The Role of  
485 Antioxidant Activity[J]. Current Atherosclerosis Reports, 2012, 14(2): 101–107.
- 486 [28] Gasparini S, Neri S, Brigo F, et al.. Late Epileptic Seizures Following Cerebral Venous  
487 Thrombosis: A Systematic Review and Meta-Analysis[J]. Neurol Sci, 2022, 43(9): 5229–  
488 5236.
- 489 [29] Abraira L, Giannini N, Santamarina E, et al.. Correlation of Blood Biomarkers with  
490 Early-Onset Seizures after an Acute Stroke Event[J]. Epilepsy & Behavior, 2020, 104:  
491 106549.
- 492 [30] Ding L, Liu Y, Meng X, et al.. Biomarker and Genomic Analyses Reveal Molecular  
493 Signatures of Non-Cardioembolic Ischemic Stroke[J]. Sig Transduct Target Ther, Nature  
494 Publishing Group, 2023, 8(1): 1–16.
- 495 [31] Zhang W, Cheng Z, Fu F, et al.. Serum Uric Acid and Prognosis in Acute Ischemic  
496 Stroke: A Dose–Response Meta-Analysis of Cohort Studies[J]. Frontiers in Aging  
497 Neuroscience, 2023, 15.
- 498 [32] Wang D, Hu B, Dai Y, et al.. Serum Uric Acid Is Highly Associated with Epilepsy  
499 Secondary to Cerebral Infarction[J]. Neurotox Res, 2019, 35(1): 63–70.
- 500 [33] Wang C, Cui T, Wang L, et al.. Prognostic Significance of Uric Acid Change in Acute  
501 Ischemic Stroke Patients with Reperfusion Therapy[J]. Eur J Neurol, 2021, 28(4): 1218–  
502 1224.
- 503 [34] Ng G J L, Quek A M L, Cheung C, et al.. Stroke Biomarkers in Clinical Practice: A  
504 Critical Appraisal[J]. Neurochemistry International, 2017, 107: 11–22.

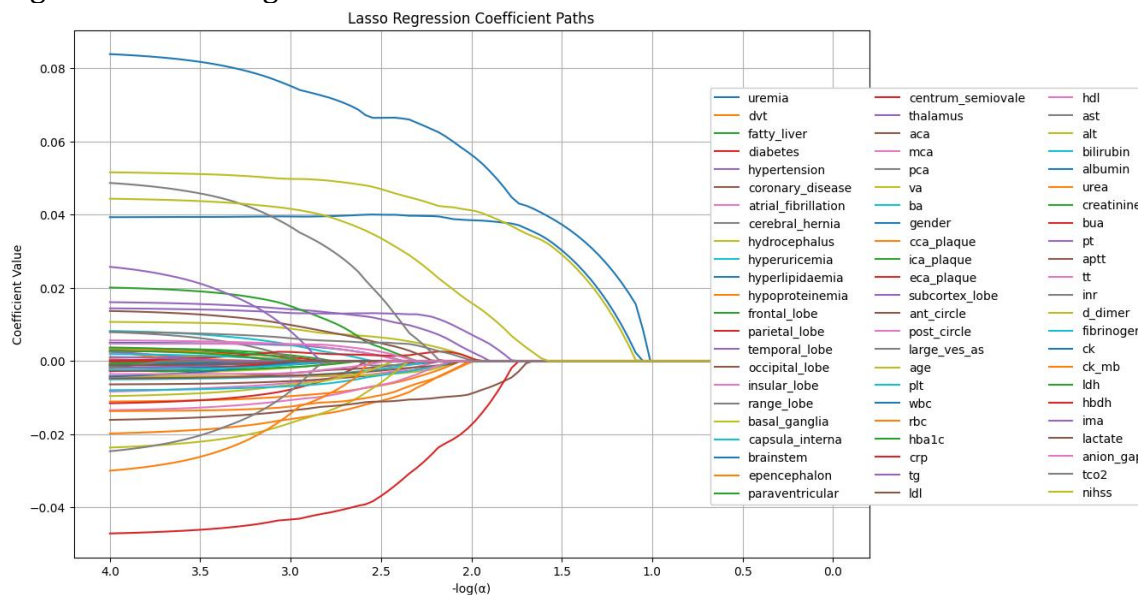
- 505 [35] Amaro S, Urrea X, Gómez-Choco M, et al.. Uric Acid Levels Are Relevant in Patients  
506 With Stroke Treated With Thrombolysis[J]. Stroke, American Heart Association, 2011,  
507 42(1\_suppl\_1): S28–S32.
- 508 [36] Amaro S, Urrea X, Gómez-Choco M, et al.. Uric Acid Levels Are Relevant in Patients  
509 with Stroke Treated with Thrombolysis[J]. Stroke, 2011, 42(SUPPL. 1): S28–S32.
- 510 [37] Ay H, Arsava E M, Sarba O. Creatine Kinase-MB Elevation after Stroke Is Not Cardiac  
511 in Origin Comparison with Troponin T Levels[J]. Stroke, 2002, 33(1): 286–289.
- 512 [38] Liu X, Chen X, Wang H, et al.. Prognostic Significance of Admission Levels of Cardiac  
513 Indicators in Patients with Acute Ischaemic Stroke: Prospective Observational Study[J]. J  
514 Int Med Res, SAGE Publications Ltd, 2014, 42(6): 1301–1310.
- 515 [39] Zeng Y-Y, Zhang W-B, Cheng L, et al.. Cardiac Parameters Affect Prognosis in Patients  
516 with Non-Large Atherosclerotic Infarction[J]. Molecular Medicine, 2021, 27(1): 2.
- 517 [40] Hijazi Z, Lindbäck J, Alexander J H, et al.. The ABC (Age, Biomarkers, Clinical History)  
518 Stroke Risk Score: A Biomarker-Based Risk Score for Predicting Stroke in Atrial  
519 Fibrillation[J]. European Heart Journal, 2016, 37(20): 1582–1590.
- 520 [41] Zheng Yuan-Hui, ZHENG Jin-Yi, ZHANG Jian. Changes of serum myocardial enzyme  
521 profile in acute stage of stroke [J]. Chinese Journal of Advanced Medical Doctors, China  
522 Medical Journal, 2009, 32(07): 46 -- 47.

Figure 1. Selection and Exclusion Procedure of Patients



A total of 42,079 records were retrieved from the stroke database, and 24,733 patients were diagnosed with ischemic or lacunar stroke with new onset. Hemorrhagic strokes (4,565), a history of stroke (2,154), TIA (3,570), unclear cause strokes (561), and records with missing essential data (6,496) were excluded. Patients whose seizures might have been caused by other factors (such as brain tumors, intracranial vascular malformations, or traumatic brain injury) (865), those with a seizure history (152), and patients who died in the hospital (1,444) were also excluded. Additionally, patients lost to follow-up (those without outpatient records or unreachable by phone) or who died within three months of the stroke incident (813) were excluded. Finally, 21,459 cases were included in the study.

Figure 2. LASSO Regression Coefficient Paths

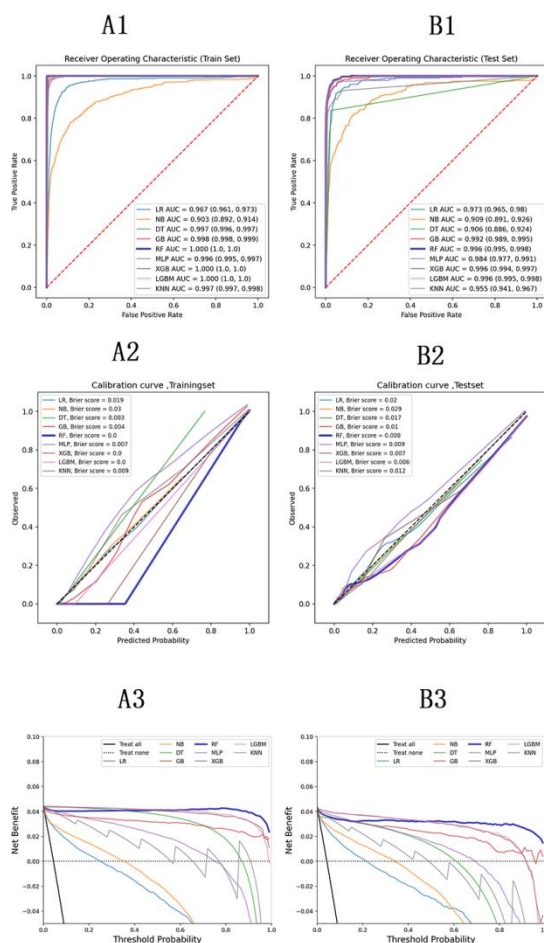




The image shows the LASSO regression coefficient paths for various features related to a medical or research study. The x-axis represents the log of the regularization parameter alpha, and the y-axis shows the regression coefficient values.

The lines in the plot represent the coefficient paths for different features as the regularization parameter changes. The features are labeled on the right side of the plot, and the most important features selected by the LASSO model are shown at the bottom of the image.

Figure 3. Model Evaluation Metrics and Curves



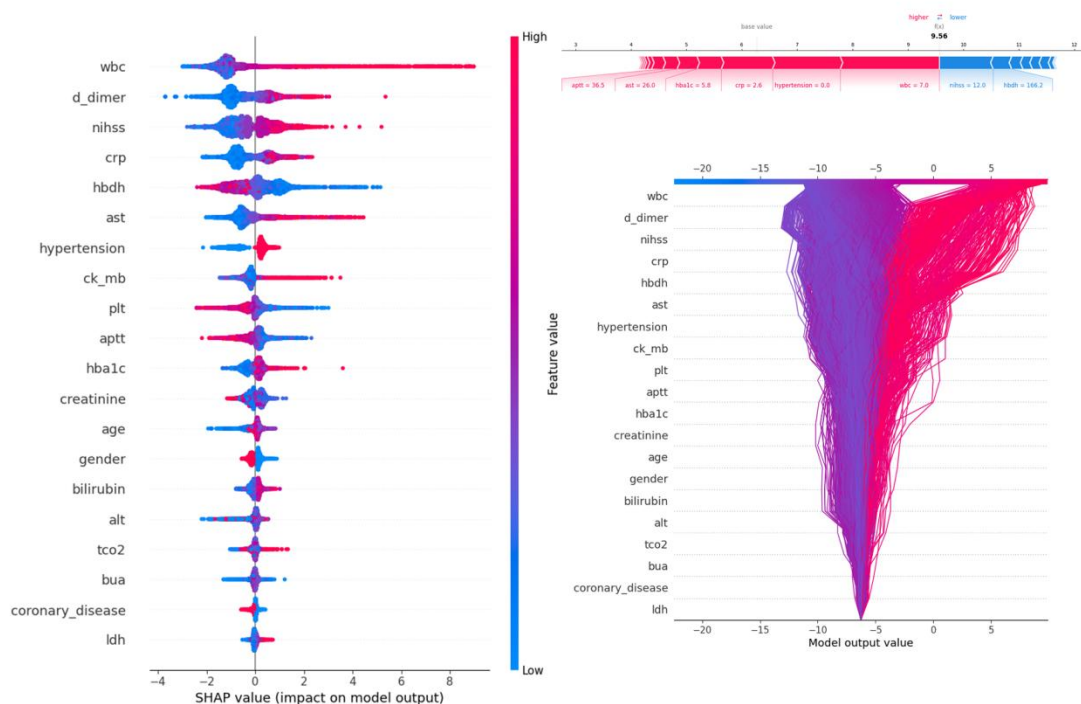
The figure displays model performance curves across six sections (A1, A2, A3 on the left; B1, B2, B3 on the right) for training and test sets.

ROC Curve: Illustrates the trade-off between sensitivity and specificity, with the AUC indicating overall model performance.

Calibration Curve: Compares predicted probabilities to actual outcomes, assessing the model's confidence accuracy.

Precision-Recall Curve: Analyzes the balance between precision and recall at various thresholds, particularly useful for imbalanced datasets.

Figure 4. Description of the SHAP Values and Feature Importance



SHAP Value (Left): Displays the impact of each feature on the model's predictions, with features sorted by importance. The color gradient indicates the range of feature values, from low (blue) to high (red).

Force Plot (Upper Right): Illustrates the contribution of individual features of the first sample to the final model output, highlighting how each feature value pushes the prediction away from the baseline value.

Decision Plot (Lower Right): Visualizes the cumulative impact of features on the model output for each sample, showing how the feature values combine to produce the final prediction.

<b>Feature</b>	<b>positive, N=954</b>	<b>negative, N = 20789</b>	<b>method</b>	<b>P</b>	<b>stats</b>
eca_plaque	-	-	Chi-Square	0.438971	0.59897
—0	942 (98.742%)	20591 (99.048%)	-	-	-
—1	12 (1.258%)	198 (0.952%)	-	-	-
subcortex_lobe	-	-	Chi-Square	0.001273	10.381551
—0	814 (85.325%)	18454 (88.768%)	-	-	-
—1	140 (14.675%)	2335 (11.232%)	-	-	-
ba	-	-	Chi-Square	0.991017	0.000127
—0	945 (99.057%)	20605 (99.115%)	-	-	-
—1	9 (0.943%)	184 (0.885%)	-	-	-
hypertension	-	-	Chi-Square	0.602539	0.271184
—0	290 (30.398%)	6497 (31.252%)	-	-	-
—1	664 (69.602%)	14292 (68.748%)	-	-	-
ica_plaque	-	-	Chi-Square	0.152086	2.051203
—0	878 (92.034%)	19392 (93.28%)	-	-	-
—1	76 (7.966%)	1397 (6.72%)	-	-	-
frontal_lobe	-	-	Chi-Square	0	53.171781
—0	868 (90.985%)	19943 (95.931%)	-	-	-
—1	86 (9.015%)	846 (4.069%)	-	-	-
cerebral_hernia	-	-	Chi-Square	0.000032	17.284355

	—0	934 (97.904%)	20626 (99.216%)	-	-	-
	—1	20 (2.096%)	163 (0.784%)	-	-	-
thalamus		-	-	Chi-Square	0.060918	3.512207
	—0	937 (98.218%)	20565 (98.923%)	-	-	-
	—1	17 (1.782%)	224 (1.077%)	-	-	-
occipital_lobe		-	-	Chi-Square	0.000034	17.17679
	—0	919 (96.331%)	20422 (98.235%)	-	-	-
	—1	35 (3.669%)	367 (1.765%)	-	-	-
pca		-	-	Chi-Square	0.891182	0.018717
	—0	952 (99.79%)	20729 (99.711%)	-	-	-
	—1	2 (0.21%)	60 (0.289%)	-	-	-
paraventricular		-	-	Chi-Square	0.213759	1.545786
	—0	899 (94.235%)	19786 (95.175%)	-	-	-
	—1	55 (5.765%)	1003 (4.825%)	-	-	-
mca		-	-	Chi-Square	0.393066	0.729435
	—0	912 (95.597%)	19998 (96.195%)	-	-	-
	—1	42 (4.403%)	791 (3.805%)	-	-	-
coronary_disease		-	-	Chi-Square	0	26.19087
	—0	599 (62.788%)	11288 (54.298%)	-	-	-
	—1	355 (37.212%)	9501 (45.702%)	-	-	-
hypoproteinemia		-	-	Chi-Square	0	53.351931
	—0	774	18479	-	-	-

	(81.132%)	(88.888%)			
—1	180 (18.868%)	2310 (11.112%)	-	-	-
parietal_lobe	-	-	Chi-Square	0	57.137771
—0	884 (92.662%)	20180 (97.071%)	-	-	-
—1	70 (7.338%)	609 (2.929%)	-	-	-
aca	-	-	Chi-Square	0.928981	0.007944
—0	941 (98.637%)	20524 (98.725%)	-	-	-
—1	13 (1.363%)	265 (1.275%)	-	-	-
brainstem	-	-	Chi-Square	0.294979	1.096759
—0	938 (98.323%)	20532 (98.764%)	-	-	-
—1	16 (1.677%)	257 (1.236%)	-	-	-
hyperuricemia	-	-	Chi-Square	0.000001	25.147468
—0	801 (83.962%)	18547 (89.215%)	-	-	-
—1	153 (16.038%)	2242 (10.785%)	-	-	-
temporal_lobe	-	-	Chi-Square	0	57.872112
—0	886 (92.872%)	20209 (97.21%)	-	-	-
—1	68 (7.128%)	580 (2.79%)	-	-	-
diabetes	-	-	Chi-Square	0.389926	0.739172
—0	617 (64.675%)	13737 (66.078%)	-	-	-
—1	337 (35.325%)	7052 (33.922%)	-	-	-
range_lobe	-	-	Chi-Square	0	85.377485

	—0	830 (87.002%)	19559 (94.083%)	-	-	-
	—1	43 (4.507%)	467 (2.246%)	-	-	-
	—2	32 (3.354%)	329 (1.583%)	-	-	-
	—3	31 (3.249%)	224 (1.077%)	-	-	-
	—4	15 (1.572%)	175 (0.842%)	-	-	-
	—5	3 (0.314%)	35 (0.168%)	-	-	-
epencephalon		-	-	Chi-Square	1	0
	—0	934 (97.904%)	20362 (97.946%)	-	-	-
	—1	20 (2.096%)	427 (2.054%)	-	-	-
hydrocephalus		-	-	Chi-Square	0	181.23517
	—0	895 (93.816%)	20565 (98.923%)	-	-	-
	—1	59 (6.184%)	224 (1.077%)	-	-	-
insular_lobe		-	-	Chi-Square	0.391042	0.735699
	—0	938 (98.323%)	20519 (98.701%)	-	-	-
	—1	16 (1.677%)	270 (1.299%)	-	-	-
gender		-	-	Chi-Square	0	44.244052
	—0	372 (38.994%)	10407 (50.06%)	-	-	-
	—1	582 (61.006%)	10382 (49.94%)	-	-	-
uremia		-	-	Chi-Square	0.00008	15.568169
	—0	934 (97.904%)	20618 (99.177%)	-	-	-
	—1	20 (2.096%)	171 (0.823%)	-	-	-
atrial_fibrillation		-	-	Chi-Square	0.008017	7.029734

	—0	838 (87.841%)	18811 (90.485%)	-	-	-
	—1	116 (12.159%)	1978 (9.515%)	-	-	-
centrum_semiovale		-	-	Chi-Square	0.36206	0.830735
	—0	922 (96.646%)	20207 (97.2%)	-	-	-
	—1	32 (3.354%)	582 (2.8%)	-	-	-
basal_ganglia		-	-	Chi-Square	0.005355	7.755329
	—0	893 (93.606%)	19869 (95.575%)	-	-	-
	—1	61 (6.394%)	920 (4.425%)	-	-	-
dvt		-	-	Chi-Square	0	40.790867
	—0	847 (88.784%)	19534 (93.963%)	-	-	-
	—1	107 (11.216%)	1255 (6.037%)	-	-	-
fatty_liver		-	-	Chi-Square	0.000171	14.123893
	—0	812 (85.115%)	16655 (80.114%)	-	-	-
	—1	142 (14.885%)	4134 (19.886%)	-	-	-
hyperlipidaemia		-	-	Chi-Square	0.000317	12.969155
	—0	801 (83.962%)	16439 (79.075%)	-	-	-
	—1	153 (16.038%)	4350 (20.925%)	-	-	-
cca_plaque		-	-	Chi-Square	0.376965	0.780577
	—0	751 (78.721%)	16100 (77.445%)	-	-	-
	—1	203 (21.279%)	4689 (22.555%)	-	-	-

va	-	-	Chi-Square	0.797483	0.065847
—0	927 (97.17%)	20159 (96.97%)	-	-	-
—1	27 (2.83%)	630 (3.03%)	-	-	-
fibrinogen	3.518 ± 0.663	3.602 ± 0.464	Mann-Whitney U	0.434584	10064078.5
d_dimer	4.362 ± 4.398	1.198 ± 0.98	Mann-Whitney U	0	3555180.5
bu	342.521 ± 74.651	344.132 ± 58.336	Mann-Whitney U	0.000037	10698805.5
tco2	22.739 ± 1.025	22.781 ± 1.225	Mann-Whitney U	0.166751	10178363
hbdh	209.295 ± 57.826	175.906 ± 48.18	Mann-Whitney U	0	6107843
anion_gap	13.026 ± 1.456	12.345 ± 1.368	Mann-Whitney U	0	6496800
ldl	2.686 ± 0.372	2.685 ± 0.361	Mann-Whitney U	0.23394	10140916.5
tt	16.636 ± 0.809	16.432 ± 0.615	Mann-Whitney U	0	7950954.5
nihss	11.529 ± 2.564	7.886 ± 2.871	Mann-Whitney U	0	2984725.5
albumin	40.734 ± 2.37	40.886 ± 2.257	Mann-Whitney U	0.025821	10338834.5
inr	1.068 ± 0.072	1.076 ± 0.149	Mann-Whitney U	0	9016933.5
tg	1.662 ± 0.484	1.536 ± 0.433	Mann-Whitney U	0	7582690.5
bilirubin	16.516 ± 4.009	15.197 ± 3.981	Mann-Whitney U	0	7522775
ima	81.624 ± 8.559	75.458 ± 12.891	Mann-Whitney U	0	4487861



pt	13.822 ± 0.627	13.843 ± 1.151	Mann-Whitney U	0	8374380.5
crp	55.681 ± 48.823	15.314 ± 18.865	Mann-Whitney U	0	3060302
wbc	11.79 ± 3.084	8.316 ± 1.286	Mann-Whitney U	0	2667973
age	65.335 ± 13.909	66.806 ± 12.597	Mann-Whitney U	0.013188	10386092
hdl	1.246 ± 0.146	1.249 ± 0.149	Mann-Whitney U	0.619502	10008026
lactate	2.825 ± 0.376	2.505 ± 0.411	Mann-Whitney U	0	4480425
rbc	4.408 ± 0.274	4.304 ± 0.324	Mann-Whitney U	0	7811417
ast	38.25 ± 18.205	26.05 ± 12.823	Mann-Whitney U	0	3814876
plt	180.251 ± 36.939	190.132 ± 26.424	Mann-Whitney U	0	11826502.5
alt	26.827 ± 10.349	24.193 ± 10.108	Mann-Whitney U	0	7632233.5
aptt	35.045 ± 1.881	35.702 ± 2.313	Mann-Whitney U	0	11737054.5
ldh	296.455 ± 111.282	215.357 ± 75.036	Mann-Whitney U	0	5261997.5
creatinine	83.837 ± 24.574	85.199 ± 52.439	Mann-Whitney U	0	8567930.5
hba1c	6.759 ± 1.048	6.662 ± 0.916	Mann-Whitney U	0.000035	9132523
urea	6.33 ± 1.354	6.419 ± 1.438	Mann-Whitney U	0.001566	10515532
ck	1029.594 ± 872.8	195.007 ± 273.212	Mann-Whitney U	0	3469376

---

Table 1. Single factor significant analysis results

This table presents the results of Chi-Square and Mann-Whitney U tests used to evaluate the association of various features with positive and negative samples.

Sample Sizes: Positive samples (N=954) and negative samples (N=20,789).

Statistical Methods: The Chi-Square test assesses the relationship between categorical variables, while the Mann-Whitney U test compares differences between independent groups for continuous data.

P-values: Indicate the significance of the associations, with lower values suggesting stronger evidence against the null hypothesis.

Statistical Values: Include counts and percentages of samples for each feature in both groups, along with the calculated statistics for each test.

Feature	0 (N=20 789)	1 (N=95 4)	OR (univariable)	coef	std err	z	P >  z	[0. 02 5	0. 97 5]	Label_ 1	Label_0
age	66.806 ± 12.597	65.335 ± 13.909	0.991 (0.986- 0.996, p=0.0)	- 0.09 0	0.03	- 3.5 08	0. 0 0	- 0. 0	- 0. 0	-	-
plt	190.13 2 ± 26.424	180.25 1 ± 36.939	0.986 (0.983- 0.988, p=0.0)	- 0.01 1	0.01	- 11. 32 0	0. 0 0	- 0. 0	- 0. 0	-	-
wbc	8.316 ± 1.286	11.79 ± 3.084	2.23 (2.149- 2.314, p=0.0)	0.8 02 2	0.019	42. 30 6	0. 0 0	0. 0. 0	0. 76 83 5 9	-	-
rbc	4.304 ± 0.324	4.408 ± 0.274	2.622 (2.162- 3.177, p=0.0)	0.9 63 8	0.098	9.8 05	0. 0 0	0. 0. 0	1. 77 15 6	-	-
hba1c	6.662 ± 0.916	6.759 ± 1.048	1.112 (1.042- 1.186, p=0.001)	0.1 05 9	0.033	3.1 76	0. 0 1	0. 0. 0	0. 04 17 1	-	-

			1.033				0.					
	15.314	55.681	(1.031-	0.0		36.	0	0.	0.			
	±	±	1.035,	32	0.0	79	0	03	03			
crp	18.865	48.823	p=0.0)	6	01	2	0	1	4	-	-	
			1.617				0.					
	1.536	1.662	(1.441-	0.4			0	0.	0.			
	±	±	1.815,	80	0.0	8.1	0	36	59			
tg	0.433	0.484	p=0.0)	7	59	70	0	5	6	-	-	
			1.009				0.	-				
	2.685	2.686	(0.843-	0.0			9	0.	0.			
	±	±	1.207,	08	0.0	0.0	2	17	18			
ldl	0.361	0.372	p=0.924)	7	91	95	4	1	8	-	-	
			0.87	-			0.	-				
	1.249	1.246	(0.562-	0.1		-	5	0.	0.			
	±	±	1.349,	38	0.2	0.6	3	57	29			
hdl	0.149	0.146	p=0.534)	9	23	22	4	7	9	-	-	
			1.028				0.					
	26.05	38.25	(1.024-	0.0		17.	0	0.	0.			
	±	±	1.031,	27	0.0	00	0	02	03			
ast	12.823	18.205	p=0.0)	7	02	7	0	4	1	-	-	
			1.017				0.					
	24.193	26.827	(1.012-	0.0			0	0.	0.			
	±	±	1.021,	16	0.0	7.5	0	01	02			
alt	10.108	10.349	p=0.0)	9	02	07	0	2	1	-	-	
			1.068				0.					
	15.197	16.516	(1.054-	0.0			0	0.	0.			
	±	±	1.082,	66	0.0	9.8	0	05	07			
bilirubin	3.981	4.009	p=0.0)	2	07	26	0	3	9	-	-	
			0.971	-			0.	-	-			
	40.886		(0.945-	0.0		-	0	0.	0.			
	±	40.734	0.999,	29	0.0	2.0	4	05	00			
albumin	2.257	± 2.37	p=0.042)	1	14	36	2	7	1	-	-	
			0.955	-			0.	-				
	6.419		(0.91-	0.0		-	0	0.	0.			
	±	6.33 ±	1.002,	45	0.0	1.8	6	09	00			
urea	1.438	1.354	p=0.063)	9	25	62	3	4	2	-	-	
			0.999	-			-	-				
	85.199	83.837	(0.998-	0.0		-	0.	0.	0.			
creatinine	±	±	1.001,	00	0.0	0.7	4	00	00	-	-	

	52.439	24.574	p=0.425)	6	01	98	2	2	1		
							5				
				-			0.	-			
	344.13	342.52	1.0 (0.998-	0.0		-	4	0.	0.		
	2 ±	1 ±	1.001,	00	0.0	0.8	1	00	00		
bua	58.336	74.651	p=0.411)	5	01	22	1	2	1	-	-
			0.982	-			0.	-			
	13.843	13.822	(0.925-	0.0		-	5	0.	0.		
	±	±	1.043,	17	0.0	0.5	6	07	04		
pt	1.151	0.627	p=0.564)	7	31	77	4	8	2	-	-
			0.863	-			0.	-	-		
	35.702	35.045	(0.835-	0.1		-	0	0.	0.		
	±	±	0.891,	47	0.0	8.9	0	18	11		
aptt	2.313	1.881	p=0.0)	3	17	17	0	0	5	-	-
			1.411				0.				
	16.432	16.636	(1.287-	0.3			0	0.	0.		
	±	±	1.547,	44	0.0	7.3	0	25	43		
tt	0.615	0.809	p=0.0)	2	47	28	0	2	6	-	-
			0.643	-			0.	-			
	1.076	1.068	(0.385-	0.4		-	0	0.	0.		
	±	±	1.074,	42	0.2	1.6	9	95	07		
inr	0.149	0.072	p=0.091)	1	62	89	1	5	1	-	-
			1.717				0.				
	1.198	4.362	(1.662-	0.5		32.	0	0.	0.		
	± 0.98	±	1.774,	40	0.0	72	0	50	57		
d_dimer		4.398	p=0.0)	5	17	4	0	8	3	-	-
			0.675	-			0.	-	-		
	3.602	3.518	(0.585-	0.3		-	0	0.	0.		
	±	±	0.778,	93	0.0	5.4	0	53	25		
fibrinogen	0.464	0.663	p=0.0)	1	73	08	0	6	1	-	-
			1.002				0.				
	195.00	1029.5	(1.002-	0.0	6.1	38.	0	0.	0.		
	7 ±	94 ±	1.002,	02	5e-	32	0	00	00		
ck	273.21	872.8	p=0.0)	4	05	6	0	2	2	-	-
			1.005				0.				
	215.35	5 ±	(1.005-	0.0		21.	0	0.	0.		
	7 ±	111.28	1.006,	05	0.0	42	0	00	00		
ldh	75.036	2	p=0.0)	3	00	4	0	5	6	-	-

			1.006				0.					
	175.90	209.29	(1.005-	0.0		15.	0	0.	0.			
	6 ±	5 ±	1.007,	06	0.0	63	0	00	00			
hbdh	48.18	57.826	p=0.0)	2	00	7	0	5	7	-	-	
			1.015				0.					
	75.458	81.624	(1.012-	0.0		10.	0	0.	0.			
	±	±	1.017,	14	0.0	70	0	01	01			
ima	12.891	8.559	p=0.0)	7	01	7	0	2	7	-	-	
			3.12				0.					
	2.505	2.825	(2.784-	1.1		19.	0	1.	1.			
	±	±	3.494,	37	0.0	58	0	02	25			
lactate	0.411	0.376	p=0.0)	7	58	7	0	4	1	-	-	
			1.344				0.					
	12.345	13.026	(1.29-	0.2		14.	0	0.	0.			
	±	±	1.399,	95	0.0	36	0	25	33			
anion_g ap	1.368	1.456	p=0.0)	3	21	8	0	5	6	-	-	
			0.972	-			0.	-				
	22.781	22.739	(0.921-	0.0		-	2	0.	0.			
	±	±	1.025,	28	0.0	1.0	9	08	02			
tco2	1.225	1.025	p=0.293)	7	27	51	3	2	5	-	-	
			1.342				0.					
	7.886	11.529	(1.318-	0.2		30.	0	0.	0.			
	±	±	1.368,	94	0.0	95	0	27	31			
nihss	2.871	2.564	p=0.0)	2	10	7	0	6	3	-	-	
uremia_ 0	20618 (99.17 7%)	934 (97.90 4%)		-	-	-	-	-	-	4.334% (934 / 21552)	95.666% (20618 / 21552)	
			2.582				0.					
	171	20	(1.618-	0.9			0	0.	1.	10.471	89.529%	
	(0.823 %)	(2.096 %)	4.121,	48	0.2	3.9	0	48	41	% (20 /	(171 /	
uremia_ 1			p=0.0)	5	39	74	0	1	6	191)	191)	
	19534	847		-	-	-	-	-	-	4.156% (847 / 20381)	95.844% (19534 / 20381)	
dvt_0	(93.96 3%)	(88.78 4%)										
			1.966				0.					
	1255	107	(1.595-	0.6			0	0.	0.	7.856%	92.144%	
	(6.037 %)	(11.21 6%)	2.423,	76	0.1	6.3	0	46	88	(107 /	(1255 /	
dvt_1			p=0.0)	1	07	40	0	7	5	1362)	1362)	

fatty_liver_0	16655 (80.11 4%)	812 (85.11 5%)	-	-	-	-	-	-	-	4.649% (812 / 17467)	95.351% (16655 / 17467)
			0.705 (0.587- 0.845, p=0.0)	- 0.3 50 2	- 0.0 93	- 3.7 82	0 0 0	- 0 53 2	- 0 16 9	3.321% (142 / 4276)	96.679% (4134 / 4276)
diabetes_0	13737 (66.07 8%)	617 (64.67 5%)	-	-	-	-	-	-	-	4.298% (617 / 14354)	95.702% (13737 / 14354)
			1.064 (0.929- 1.219, p=0.371)	0.0 62 0	0.0 69	0.8 95	0 7 1	- 0 7 4	- 0 07 8 19	4.561% (337 / 7389)	95.439% (7052 / 7389)
hypertension_0	6497 (31.25 2%)	290 (30.39 8%)	-	-	-	-	-	-	-	4.273% (290 / 6787)	95.727% (6497 / 6787)
			1.041 (0.904- 1.198, p=0.578)	0.0 40 0	0.0 72	0.5 56	0 5 8	- 0 7 8	- 0 10 1 18	4.44% (664 / 14956)	95.56% (14292 / 14956)
coronary_disease_0	11288 (54.29 8%)	599 (62.78 8%)	-	-	-	-	-	-	-	5.039% (599 / 11887)	94.961% (11288 / 11887)
			0.704 (0.616- 0.805, p=0.0)	- 0.3 50 8	0.0 68	- 5.1 28	0 0 0	- 0 48 5	- 0 21 7	3.602% (355 / 9856)	96.398% (9501 / 9856)
atrial_fibrillation_0	18811 (90.48 5%)	838 (87.84 1%)	-	-	-	-	-	-	-	4.265% (838 / 19649)	95.735% (18811 / 19649)
			1.316 (1.078- 1.608, p=0.007)	0.2 74 9	0.1 02	2.6 99	0 7	0 07	0 47 5	5.54% (116 / 2094)	94.46% (1978 / 2094)
hyperuricemia_0	18547 (89.21 5%)	801 (83.96 2%)	-	-	-	-	-	-	-	4.14% (801 / 19348)	95.86% (18547 / 19348)
			1.58 (1.322- 1.889,	0.4 57	0.0 0.0	5.0 5.0	0 0	0 27	0 63	6.388% (153 / 2242)	93.612% (2242 / 2242)

1	5%)	8%)	p=0.0)	5	91	27	0	9	6	2395)	2395)
							0				
hyperlipidaemia_0	16439 (79.07 5%)	801 (83.96 2%)	-	-	-	-	-	-	-	4.646% (801 / 17240)	95.354% (16439 / 17240)
hyperlipidaemia_1	4350 (20.92 5%)	153 (16.03 8%)	0.722 (0.605-0.861, p=0.0)	0.3 25 9	0.0 90	3.6 27	0 0	0 50 2	0 15 0	3.398% (153 / 4503)	96.602% (4350 / 4503)
hypoproteinemia_0	18479 (88.88 8%)	774 (81.13 2%)	-	-	-	-	-	-	-	4.02% (774 / 19253)	95.98% (18479 / 19253)
hypoproteinemia_1	2310 (11.11 2%)	180 (18.86 8%)	1.86 (1.573-2.201, p=0.0)	0.6 20 8	0.0 86	7.2 48	0 0	0 45 3	0 78 9	7.229% (180 / 2490)	92.771% (2310 / 2490)
cerebral_hernia_0	20626 (99.21 6%)	934 (97.90 4%)	-	-	-	-	-	-	-	4.332% (934 / 21560)	95.668% (20626 / 21560)
cerebral_hernia_1	163 (0.784 %)	20 (2.096 %)	2.71 (1.696-4.332, p=0.0)	0.9 96 8	0.2 39	4.1 66	0 0	0 52 8	1 46 6	10.929% (20 / 183)	89.071% (163 / 183)
hydrocephalus_0	20565 (98.92 3%)	895 (93.81 6%)	-	-	-	-	-	-	-	4.171% (895 / 21460)	95.829% (20565 / 21460)
hydrocephalus_1	224 (1.077 %)	59 (6.184 %)	6.052 (4.509-8.125, p=0.0)	1.8 00 4	0.1 50	11. 98 2	0 0	1 50 6	2 09 5	20.848% (59 / 283)	79.152% (224 / 283)
frontal_lobe_0	19943 (95.93 1%)	868 (90.98 5%)	-	-	-	-	-	-	-	4.171% (868 / 20811)	95.829% (19943 / 20811)
frontal_lobe_1	846 (4.069 %)	86 (9.015 %)	2.336 (1.852-2.945, p=0.0)	0.8 48 3	0.1 18	7.1 66	0 0	0 61 6	1 08 0	9.227% (86 / 932)	90.773% (846 / 932)
parietal_lobe_0	20180 (97.07 1%)	884 (92.66 2%)	-	-	-	-	-	-	-	4.197% (884 / 21064)	95.803% (20180 / 21064)

			2.624					0.				
parietal_lobe_1	609 (2.929 %)	70 (7.338 %)	(2.03- 3.391, p=0.0)	0.9 64 7	0.1 31	7.3 75	0 0	0 70 8	1. 22 1	10.309 % (70 / 679)	89.691% (609 / 679)	
temporal_lobe_0	20209 (97.21 %)	886 (92.87 2%)	-	-	-	-	-	-	-	4.2% (886 / 21095)	95.8% (20209 / 21095)	
			2.674					0.				
temporal_lobe_1	580 (2.79 %)	68 (7.128 %)	(2.063- 3.469, p=0.0)	0.9 83 6	0.1 33	7.4 13	0 0	0 72 4	1. 24 4	10.494 % (68 / 648)	89.506% (580 / 648)	
occipital_lobe_0	20422 (98.23 5%)	919 (96.33 1%)	-	-	-	-	-	-	-	4.306% (919 / 21341)	95.694% (20422 / 21341)	
			2.119					0.				
occipital_lobe_1	367 (1.765 %)	35 (3.669 %)	(1.489- 3.016, p=0.0)	0.7 51 1	0.1 80	4.1 70	0 0	0 39 8	1. 10 4	8.706% (35 / 402)	91.294% (367 / 402)	
insular_lobe_0	20519 (98.70 1%)	938 (98.32 3%)	-	-	-	-	-	-	-	4.372% (938 / 21457)	95.628% (20519 / 21457)	
			1.296					0.	-			
insular_lobe_1	270 (1.299 %)	16 (1.677 %)	(0.78- 2.155, p=0.317)	0.2 59 5	0.2 59	1.0 00	1 7	0 24 9	0. 76 8	5.594% (16 / 286)	94.406% (270 / 286)	
range_lobe_0	19559 (94.08 3%)	830 (87.00 2%)	-	-	-	-	-	-	-	4.071% (830 / 20389)	95.929% (19559 / 20389)	
			2.17					0.				
range_lobe_1	467 (2.246 %)	43 (4.507 %)	(1.576- 2.989, p=0.0)	0.7 74 6	0.1 63	4.7 45	0 0	0 45 5	1. 09 5	8.431% (43 / 510)	91.569% (467 / 510)	
			2.292					0.				
range_lobe_2	329 (1.583 %)	32 (3.354 %)	(1.584- 3.317, p=0.0)	0.8 29 4	0.1 89	4.3 99	0 0	0 46 0	1. 19 9	8.864% (32 / 361)	91.136% (329 / 361)	
			3.261					0.				
range_lobe_3	224 (1.077 %)	31 (3.249 %)	(2.226- 4.778, p=0.0)	1.1 82 1	0.1 95	6.0 66	0 0	0 80 0	1. 56 4	12.157 % (31 / 255)	87.843% (224 / 255)	





vale_0	%)	6%)								21129)	21129)
centrum _semio vale_1	582 (2.8%)	32 (3.354 %)	1.205 (0.839- 1.73, p=0.313)	0.1 86 5	0.1 85	1.0 10	0. 3 3	- 0. 5	0. 17 8	5.212% (32 / 614)	94.788% (582 / 614)
thalamu s_0	20565 (98.92 3%)	937 (98.21 8%)	-	-	-	-	-	-	-	4.358% (937 / 21502)	95.642% (20565 / 21502)
thalamu s_1	224 (1.077 %)	17 (1.782 %)	1.666 (1.013- 2.74, p=0.044)	0.5 10 2	0.2 54	2.0 11	0. 0 4	- 0. 4	1. 01 3	7.054% (17 / 241)	92.946% (224 / 241)
aca_0	20524 (98.72 5%)	941 (98.63 7%)	-	-	-	-	-	-	-	4.384% (941 / 21465)	95.616% (20524 / 21465)
aca_1	265 (1.275 %)	13 (1.363 %)	1.07 (0.611- 1.874, p=0.813)	0.0 67 6	0.2 86	0.2 36	0. 8 3	- 0. 3	0. 49 8	4.676% (13 / 278)	95.324% (265 / 278)
mca_0	19998 (96.19 5%)	912 (95.59 7%)	-	-	-	-	-	-	-	4.362% (912 / 20910)	95.638% (19998 / 20910)
mca_1	791 (3.805 %)	42 (4.403 %)	1.164 (0.848- 1.598, p=0.348)	0.1 52 1	0.1 62	0.9 39	0. 3 8	- 0. 5	0. 16 9	5.042% (42 / 833)	94.958% (791 / 833)
pca_0	20729 (99.71 1%)	952 (99.79 %)	-	-	-	-	-	-	-	4.391% (952 / 21681)	95.609% (20729 / 21681)
pca_1	60 (0.289 %)	2 (0.21 %)	0.726 (0.177- 2.974, p=0.656)	0.3 20 5	0.7 20	0.4 45	0. 6 6	- 1. 1	1. 73 0	3.226% (2 / 62)	96.774% (60 / 62)
va_0	20159 (96.97 %)	927 (97.17 %)	-	-	-	-	-	-	-	4.396% (927 / 21086)	95.604% (20159 / 21086)
va_1	630 (3.03 %)	27 (2.83 %)	0.932 (0.631- 1.377, p=0.724)	0.0 70 4	0.1 99	0.3 53	0. 0. 7 2	- 0. 1	0. 46 0	4.11% (27 / 657)	95.89% (630 / 657)

4

ba_0	20605 (99.11 5%)	945 (99.05 7%)	-	-	-	-	-	-	-	4.385% (945 / 21550)	95.615% (20605 / 21550)
ba_1	184 (0.885 %)	9 (0.943 %)	1.067 (0.544- 2.09, p=0.851)	0.0 64 4	0.3 43	0.1 88	0. 8 1	- 0. 60 8	0. 0. 73	4.663% (9 / 193)	95.337% (184 / 193)
gender_0	10407 (50.06 %)	372 (38.99 4%)	-	-	-	-	-	-	-	3.451% (372 / 10779)	96.549% (10407 / 10779)
gender_1	10382 (49.94 %)	582 (61.00 6%)	1.568 (1.373- 1.791, p=0.0)	0.4 50 0	0.0 68	6.6 35	0. 0 0	0. 31 7	0. 58 3	5.308% (582 / 10964)	94.692% (10382 / 10964)
cca_plaque_0	16100 (77.44 5%)	751 (78.72 1%)	-	-	-	-	-	-	-	4.457% (751 / 16851)	95.543% (16100 / 16851)
cca_plaque_1	4689 (22.55 5%)	203 (21.27 9%)	0.928 (0.792- 1.088, p=0.356)	0.0 74 6	0.0 81	0.9 23	0. 3 6	- 0. 23 3	0. 0. 08 4	4.15% (203 / 4892)	95.85% (4689 / 4892)
ica_plaque_0	19392 (93.28 %)	878 (92.03 4%)	-	-	-	-	-	-	-	4.332% (878 / 20270)	95.668% (19392 / 20270)
ica_plaque_1	1397 (6.72 %)	76 (7.966 %)	1.202 (0.945- 1.528, p=0.135)	0.1 83 6	0.1 23	1.4 96	0. 1 3	- 0. 05 7	0. 0. 42 4	5.16% (76 / 1473)	94.84% (1397 / 1473)
eca_plaque_0	20591 (99.04 8%)	942 (98.74 2%)	-	-	-	-	-	-	-	4.375% (942 / 21533)	95.625% (20591 / 21533)
eca_plaque_1	198 (0.952 %)	12 (1.258 %)	1.325 (0.737- 2.382, p=0.347)	0.2 81 2	0.2 99	0.9 40	0. 3 4	- 0. 30 5	0. 0. 86 8	5.714% (12 / 210)	94.286% (198 / 210)
subcort_ex_lobe_0	18454 (88.76 8%)	814 (85.32 5%)	-	-	-	-	-	-	-	4.225% (814 / 19268)	95.775% (18454 / 19268)
subcort	2335	140	1.359							5.657%	94.343%

ex_lobe _1	(11.23 2%)	(14.67 5%)	(1.131- 1.634, p=0.001)	0.3 07 0	0.0 94 0	3.2 62 0	0. 0 0	0. 12 3	0. 49 1	(140 / 2475)	(2335 / 2475)
---------------	---------------	---------------	-------------------------------	----------------	----------------	----------------	--------------	---------------	---------------	-----------------	------------------

Table 2. Single Factor Significant Analysis Results

This table presents the results of a single factor significance analysis for various features across two groups of samples: negative samples (0) and positive samples (1).

Sample Size:

Group 0 (Negative): N = 20,789

Group 1 (Positive): N = 954

Feature Analysis: For each feature, the table includes the mean and standard deviation ( $\pm$ ) for both groups, odds ratios (OR) from univariable analysis, coefficients (coef), standard errors (std err), z-scores (z), p-values ( $P > |z|$ ), and 95% confidence intervals ([0.025, 0.975]).

Significance Levels: Features with statistically significant differences are indicated by p-values less than 0.05. An odds ratio greater than 1 suggests an increased risk associated with the feature in the positive group, while an odds ratio less than 1 suggests a decreased risk.

Labels: The last two columns provide the proportions of the positive and negative samples for selected features.

Feature	0 (N=20789)	1 (N=954)	OR (multivariable)	Coef.	Std. Err.	z	P> z	[0.025	0.975]
tg	1.536 $\pm$ 0.433	1.662 $\pm$ 0.484	2.458 (2.069-2.92, p=0.0)	0.899	0.088	10.23	0	0.727	1.071
rbc	4.304 $\pm$ 0.324	4.408 $\pm$ 0.274	4.731 (3.274- 6.837, p=0.0)	1.554	0.188	8.275	0	1.186	1.922
age	66.806 $\pm$ 12.597	65.335 $\pm$ 13.909	1.012 (1.004- 1.021, p=0.003)	0.012	0.004	2.971	0.003	0.004	0.021
ast	26.05 $\pm$ 12.823	38.25 $\pm$ 18.205	1.048 (1.04- 1.055, p=0.0)	0.046	0.004	12.413	0	0.039	0.054
plt	190.13 2 $\pm$ 26.424	180.25 1 $\pm$ 36.939	0.977 (0.973-0.98, p=0.0)	0.024	0.002	13.375	0	0.027	0.027

alt	24.193 ± 10.108	26.827 ± 10.349	0.953 (0.942- 0.964, p=0.0)	- 0.0 48	0.006	8.17 7	0	0.05 9	0.03 6
ima	75.458 ± 12.891	81.624 ± 8.559	1.006 (1.001- 1.012, p=0.014)	0.0 06	0.003	2.45 3	0.0 14	0.00 1	0.01 2
ldh	215.35 7 ± 75.036	296.45 5 ± 111.28 2	0.984 (0.982- 0.987, p=0.0)	- 0.0 16	0.001	12.9 92	0	0.01 8	0.01 4
tt	16.432 ± 0.615	16.636 ± 0.809	1.13 (1.009- 1.265, p=0.034)	0.1 22	0.058	2.11 6	0.0 34	0.00 9	0.23 5
crp	15.314 ± 18.865	55.681 ± 48.823	1.032 (1.028- 1.036, p=0.0)	0.0 31	0.002	15.5 85	0	0.02 7	0.03 5
wbc	8.316 ± 1.286	11.79 ± 3.084	2.091 (1.985- 2.204, p=0.0)	0.7 38	0.027	27.5 83	0	0.68 5	0.79
ck	195.00 7 ± 273.21 2	1029.5 94 ± 872.8	1.001 (1.001- 1.001, p=0.0)	0.0 01	0	7.86	0	0.00 1	0.00 1
subcortex_lobe_0	18454 (88.768 %)	814 (85.325 %)	-	-	-	-	-	-	-
subcortex_lobe_1	2335 (11.232 %)	140 (14.675 %)	1.188 (0.827- 1.707 ,p=0.3 52)	0.1 72	0.185	0.93	0.3 52	0.19 1	0.53 5
frontal_lobe_0	19943 (95.931 %)	868 (90.985 %)	-	-	-	-	-	-	-
frontal_lobe_1	846 (4.069 %)	86 (9.015 %)	4.577 (1.381- 15.17 ,p=0.0 13)	1.5 21	0.611	2.48 8	0.0 13	0.32 3	2.71 9
cerebral_hernia_0	20626 (99.216 %)	934 (97.904 %)	-	-	-	-	-	-	-

cerebral_hernia_1	163 (0.784 %)	20 (2.096 %)	0.846 (0.387- 1.85 ,p=0.67 6)	- 0.1 67	0.399	- 0.41 8	0.6 76	- 0.94 9	0.61 5
thalamus_0	20565 (98.923 %)	937 (98.218 %)	-	-	-	-	-	-	-
thalamus_1	224 (1.077 %)	17 (1.782 %)	0.669 (0.327- 1.373 ,p=0.2 73)	- 0.4 01	0.366	- 1.09 5	0.2 73	- 1.11 9	0.31 7
occipital_lobe_0	20422 (98.235 %)	919 (96.331 %)	-	-	-	-	-	-	-
occipital_lobe_1	367 (1.765 %)	35 (3.669 %)	2.172 (0.741- 6.368 ,p=0.1 57)	0.7 76	0.549	1.41 4	0.1 57	-0.3	1.85 1
coronary_disease_0	11288 (54.298 %)	599 (62.788 %)	-	-	-	-	-	-	-
coronary_disease_1	9501 (45.702 %)	355 (37.212 %)	1.408 (1.151- 1.724 ,p=0.0 01)	0.3 42	0.103	3.32 2	0.0 01	0.14	0.54 5
hypoproteinemia_0	18479 (88.888 %)	774 (81.132 %)	-	-	-	-	-	-	-
hypoproteinemia_1	2310 (11.112 %)	180 (18.868 %)	1.183 (0.9- 1.554 ,p=0.2 28)	0.1 68	0.139	1.20 6	0.2 28	- 0.10 5	0.44 1
parietal_lobe_0	20180 (97.071 %)	884 (92.662 %)	-	-	-	-	-	-	-
parietal_lobe_1	609 (2.929 %)	70 (7.338 %)	6.939 (2.253- 21.375 ,p=0. 001)	1.9 37	0.574	3.37 5	0.0 01	0.81 2	3.06 2
hyperuricemia_0	18547 (89.215 %)	801 (83.962 %)	-	-	-	-	-	-	-

hyperuricemia _1	2242 (10.785 %)	153 (16.038 %)	0.938 (0.691- 1.275 ,p=0.6 84)	- 0.0 64	0.156	- 0.40 7	0.6 84	-0.37	0.24 3
temporal_lobe _0	20209 (97.21 %)	886 (92.872 %)	-	-	-	-	-	-	-
temporal_lobe _1	580 (2.79%)	68 (7.128 %)	5.242 (1.548- 17.752 ,p=0. 008)	1.6 57	0.622	2.66 2	0.0 08	0.43 7	2.87 6
range_lobe_0	19559 (94.083 %)	830 (87.002 %)	-	-	-	-	-	-	-
range_lobe_1	467 (2.246 %)	43 (4.507 %)	0.359 (0.111- 1.159 ,p=0.0 87)	- 1.0 25	0.598	- 1.71 3	0.0 87	- 2.19 7	0.14 7
range_lobe_2	329 (1.583 %)	32 (3.354 %)	0.084 (0.01- 0.703 ,p=0.0 22)	- 2.4 79	1.085	- 2.28 5	0.0 22	- 4.60 5	- 0.35 2
range_lobe_3	224 (1.077 %)	31 (3.249 %)	0.011 (0.0- 0.231 ,p=0.0 04)	- 4.5 42	1.569	- 2.89 5	0.0 04	- 7.61 7	- 1.46 7
range_lobe_4	175 (0.842 %)	15 (1.572 %)	0.001 (0.0- 0.057 ,p=0.0 01)	- 6.6 66	1.943	- -3.43	0.0 01	- 10.4 75	- 2.85 7
range_lobe_5	35 (0.168 %)	3 (0.314 %)	0.001 (0.0- 0.115 ,p=0.0 04)	- 6.5 86	2.259	- 2.91 5	0.0 04	- 11.0 14	- 2.15 9
hydrocephalus _0	20565 (98.923 %)	895 (93.816 %)	-	-	-	-	-	-	-
hydrocephalus _1	224 (1.077 %)	59 (6.184 %)	3.251 (1.939- 5.451 ,p=0.0 )	1.1 79	0.264	4.47 1	0	0.66 2	1.69 6
gender_0	10407 (50.06	372 (38.994	-	-	-	-	-	-	-

		%)	%)						
gender_1	10382 (49.94 %)	582 (61.006 %)	0.572 (0.454- 0.72 ,p=0.0)	- 0.5 58	0.117	- 4.75 3	0	- 0.78 9	- 0.32 8
uremia_0	20618 (99.177 %)	934 (97.904 %)	-	-	-	-	-	-	-
uremia_1	171 (0.823 %)	20 (2.096 %)	1.979 (1.098- 3.564 ,p=0.0 23)	0.6 82	0.3	2.27 3	0.0 23	0.09 4	1.27 1
atrial_fibrillatio n_0	18811 (90.485 %)	838 (87.841 %)	-	-	-	-	-	-	-
atrial_fibrillatio n_1	1978 (9.515 %)	116 (12.159 %)	1.446 (1.087- 1.923 ,p=0.0 11)	0.3 69	0.145	2.53 4	0.0 11	0.08 4	0.65 4
basal_ganglia_ 0	19869 (95.575 %)	893 (93.606 %)	-	-	-	-	-	-	-
basal_ganglia_ 1	920 (4.425 %)	61 (6.394 %)	1.024 (0.642- 1.633 ,p=0.9 21)	0.0 24	0.238	0.09 9	0.9 21	- 0.44 3	0.49
dvt_0	19534 (93.963 %)	847 (88.784 %)	-	-	-	-	-	-	-
dvt_1	1255 (6.037 %)	107 (11.216 %)	1.254 (0.922- 1.706 ,p=0.1 49)	0.2 27	0.157	1.44 3	0.1 49	- 0.08 1	0.53 4
hyperlipidaemi a_0	16439 (79.075 %)	801 (83.962 %)	-	-	-	-	-	-	-
hyperlipidaemi a_1	4350 (20.925 %)	153 (16.038 %)	0.825 (0.646- 1.052 ,p=0.1 21)	- 0.1 93	0.124	- 1.55 2	0.1 21	- 0.43 7	0.05 1
fatty_liver_0	16655 (80.114 %)	812 (85.115 %)	-	-	-	-	-	-	-



	%)	%)							
fatty_liver_1	4134 (19.886 %)	142 (14.885 %)	0.759 (0.59- 0.978 ,p=0.0 33)	- 0.2 75	0.129	2.13 5	0.0 33	- 0.52 8	- 0.02 3

### Table 3. Multivariable Analysis Results

This table summarizes the results of a multivariable analysis for various features across two groups of samples: negative samples (0) and positive samples (1).

Sample Size:

Group 0 (Negative): N = 20,789

Group 1 (Positive): N = 954

Feature Analysis: For each feature, the table includes the mean and standard deviation ( $\pm$ ) for both groups, odds ratios (OR) from multivariable analysis, coefficients (Coef.), standard errors (Std. Err.), z-scores (z), p-values ( $P > |z|$ ), and 95% confidence intervals ([0.025, 0.975]).

Significance Levels: Features with statistically significant differences are indicated by p-values less than 0.05. An odds ratio greater than 1 indicates an increased risk associated with the feature in the positive group, while an odds ratio less than 1 suggests a decreased risk.

Labels: The last column presents the proportions of the positive and negative samples for selected features.

Model	AUC	Accuracy	Sensitivity/Recall	Specificity	F1-score	PPV/precision
LR	0.967   0.973	0.928   0.927	0.920   0.929	0.928   0.927	0.530   0.524	0.373   0.365
NB	0.903   0.909	0.938   0.936	0.634   0.662	0.952   0.949	0.474   0.472	0.378   0.367
DT	0.997   0.906	0.993   0.970	1.000   0.836	0.993   0.976	0.930   0.706	0.870   0.610
GB	0.998   0.992	0.987   0.980	0.976   0.900	0.988   0.983	0.871   0.794	0.786   0.711
RF	1.000   0.996	0.997   0.989	1.000   0.883	0.997   0.994	0.967   0.873	0.936   0.864
MLP	0.996   0.984	0.977   0.972	0.975   0.932	0.977   0.974	0.790   0.744	0.664   0.619
XGB	1.000   0.996	0.996   0.988	1.000   0.897	0.996   0.992	0.961   0.867	0.926   0.840

LGBM	1.000   0.996	0.997   0.989	1.000   0.886	0.997   0.993	0.970   0.869	0.941   0.853
KNN	0.997   0.955	0.965   0.955	0.999   0.890	0.964   0.958	0.717   0.631	0.560   0.489

---

#### Table 4. Model Performance Evaluation Results

This table presents the performance evaluation metrics for various machine learning models, including AUC, Accuracy, Sensitivity (Recall), Specificity, F1-score, Positive Predictive Value (PPV/Precision), and Negative Predictive Value (NPV)

**AUC:** Area Under the Curve, indicating the model's ability to distinguish between positive and negative samples; values closer to 1 indicate better performance.

**Accuracy:** The proportion of correctly classified samples among the total samples.

**Sensitivity/Recall:** The proportion of correctly identified positive samples out of all actual positive samples.

**Specificity:** The proportion of correctly identified negative samples out of all actual negative samples.

**F1-score:** The harmonic mean of precision and recall, considering both the accuracy and completeness of the model.

**Positive Predictive Value (PPV/Precision):** The proportion of correctly identified positive samples among all samples predicted as positive.

**Negative Predictive Value (NPV):** The proportion of correctly identified negative samples among all samples predicted as negative.