

# Development and validation of a machine learning-based stroke suicidal ideation prediction model : A retrospective study

**Seung Il Song**

Gumi University Faculty of Health

**Hyeon Taek Hong**

Daegu University

**Changwoo Lee**

Seoul National University Hospital

**Seung Bo Lee** (✉ [koreateam23@gmail.com](mailto:koreateam23@gmail.com))

Keimyung University School of Medicine

---

## Article

**Keywords:** anxiety, depression, machine learning, prediction model, stroke, suicidal ideation

**Posted Date:** May 19th, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1512786/v2>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# Abstract

**Background:** Currently, the identification of stroke patients with an increased suicide risk is mainly based on self-report questionnaires, and this method suffers from a lack of objectivity. This study developed and validated a suicide ideation (SI) prediction model using clinical data and identified SI predictors.

**Method:** Significant variables were selected through traditional statistical analysis based on retrospective data of 385 stroke patients; the data was collected from October 2012 to March 2014. The data was then applied to three boosting models (Xgboost, CatBoost, LGBM) to identify the comparative and best performing models.

**Results:** Demographic variables that showed significant differences between the two groups were age, onset, type, socioeconomic, and education level. Additionally, functional variables also showed a significant difference with regard to ADL and emotion ( $p < 0.05$ ). The CatBoost model (0.900) showed higher performance than the other two models; and depression, anxiety, self-efficacy, and rehabilitation motivation (RM) were found to have high importance. Negative emotions such as depression and anxiety showed a positive relationship with SI and rehabilitation motivation and self-efficacy displayed an inverse relationship with SI.

**Conclusion:** Machine learning-based SI models could augment SI prevention by helping rehabilitation and medical professionals identify high-risk stroke patients in need of SI prevention intervention.

## Background

Stroke is a cerebrovascular disease characterized by neurological deficits, including hemiplegia, sensory dysfunction, aphasia, neglect, and intellectual and mental disabilities [1]. Post-stroke depression (PSD) is considered the most frequent and important sequela of stroke [2], and is the largest indicator of the occurrence of suicidal ideation (SI) [3]. Recently, due to the spread of COVID-19, mental health issues such as depression, anxiety, and suicide have increased. Such mental health issues may increase SI risk, especially in patients with PSD [4–6].

SI precedes suicidal attempts or suicidal behaviors, and understanding the effect of SI contributes to understanding and preventing the risk of suicidal behavior [7]. SI is more prevalent among those with persistent physical and cognitive impairments resulting from stroke [8]. Previous studies have reported that the occurrence of depression and mood disorders increases SI in stroke patients, and that there is a significant positive correlation between depression and SI in stroke patients [9, 10]. Therefore, a clinical data prediction model is necessary to reduce SI in patients after a stroke.

Most of the developed stroke prediction models are reported in studies on diagnosis, sequela, mortality, and physical function, and cannot be conveniently used practically owing to the associated invasive measurements and analyses [11–14]. Additionally, while studies on predictive model development for stroke-related emotional disorders, such as post-stroke anxiety and PSD have been conducted [15, 16], the

predictors used in these models were assessed at one-month post-stroke, at which point full depressive symptoms may not be present. Additionally, procedures need to be devised for the comparison of different machine learning models to select the best among them.

This study presents a stroke patient SI prediction model independent of biochemical data that are not routinely collected and aims to differentiate various SIs. To date, there have been no similar studies, and most of the developed models require image data and invasive test data, which are difficult to collect. Therefore, this study aims to identify SI risk factors from demographic factors, medical history, cognitive function, activities of daily living (ADL) function, and psychological factors. To the best of our knowledge, this study is also the first to apply the best model selected after comparing the performance of three boosting models using data collected from a sample of subacute and chronic stroke patients in an attempt to create an SI prediction tool.

## Methods

### Setting, data description, and pre-processing

This study aimed to develop and validate an SI prediction model using clinical data and identified SI predictors. For this purpose, we used the data collected from a specialized hospital in Daegu Metropolitan City, Republic of Korea, to predict high or low levels of SI outcomes in patients with stroke. A total of 385 stroke patients were screened for eligibility between October 2012 and March 2014. The eligibility criteria were as follows: diagnoses confirmed based on the results of magnetic resonance imaging and computed tomography images evaluated by a physician; patients in the age range of 18–80 years; a diagnosis of ischemic and hemorrhage stroke type; and patients with an onset of subacute stroke between one and six months and chronic stroke over six months. The collected anonymized sample data included information on demographics, hospital admission, cognitive function, motor function, ADL, and emotion assessment results. The ethics committee of our Institutional Review Board reviewed this study. This is a retrospective study using anonymized data obtained with written consent from all patients. This study has been the ethics committee of Daegu University Institutional Review Board (IRB) approved this study (1040621-202111-HR-079) and all methods were performed in accordance with the relevant guidelines and regulations.

The features obtained from pre-processing were then divided into five domains based on the assessment for which they were collected. All the potential predictors, including sociodemographic factors, cognitive function, motor function, ADL, and emotional parameters, were extracted from the hospital's electronic medical records and experimental data. Assessments included the Scale for SI [17, 18], the Korean version of the Mini-Mental State Examination (MMSE-K) [19], the Manual Function Test (MFT) [20], the Korean version of the Modified Bathel Index (K-MBI) [21], Self-Efficacy Scale [22], the Rehabilitation Motivation Scale (RMS) [23], the Beck Anxiety Inventory (BAI) [24], the Beck Depression Inventory (BDI) [25]. The study data indicated that the assessment outcome had high reliability and validity.

Demographic features included sex, age, phase, type, affected side, dominant hand, socioeconomic level, marital status, hypertension, diabetes, family/past history, smoking and drinking, education, and transfer. Cognitive function was measured using the MMSE-K, motor function using the MFT, and ADL using the K-MBI. Finally, positive emotions were measured using the Self-Efficacy scale and the RMS, and negative emotions were measured using the BAI and BDI.

Variables for demographic features, cognitive function, motor function, and ADL, as well as numerical variables for emotion were included in the dataset. The target variable was the SSI Scale score. To transform the problem into a binary classification one and to compare our results directly with those obtained by existing methods, we discretized the SSI into two classes: high SI group ( $> 14$ ) and low SI group ( $\leq 14$ ) [17, 18]. This particular discretization is medically relevant because it helps to distinguish between stroke patients who will be able to live an independent life from those with a significant suicide risk.

The age variables were transformed into categorical variables. Two pre-processing methods were used to eliminate the outliers and missing values. After data cleaning, the resulting dataset contained 23 features, and the data of 304 patients who met the inclusion criteria were included in the datasets, which were then used for model training and validation (Fig. 1). All the stroke patients included in the study were screened, and anonymized data were used for a retrospective study comprising two groups: high SI group ( $n = 165$ ) and low SI group ( $n = 139$ ).

## Statistical analysis

The data were analyzed using the IBM Statistical Package for Social Sciences (SPSS) version 25.0. Frequency analysis and chi-square test were performed, and a normality test was performed to determine normality of the distributions. The age variable was converted into a categorical variable for anonymization; it was converted into a 10-year interval based on the original data. The Mann–Whitney U (two tailed) test was conducted to determine the difference in the SI between the two groups. Differences were considered statistically significant at  $p < 0.05$  (Fig. 1). Three models (Xgboost, CatBoost, LGBM) were compared and the one with the best performance, that is the CatBoost model, was selected.

## SI prediction model

We used an ML approach to develop the SI prediction models for stroke patients. The three boosting models (Xgboost, CatBoost, and Light GBM [gradient booting model]) apply an algorithm based on gradient boosted decision trees. Xgboost implements the gradient boosting algorithm, which combines numerous decision trees for elaborate classification, in a fast and generalized manner [26]. XGBoost also applies a sparsity-aware algorithm to find the best split faster than the other methods. Light GBM (LGBM) is an advanced implementation of gradient boosting. This algorithm differs from the other algorithms in the growth of the tree in-depth or by leaves. LGBM handles large amounts of data with the lowest memory requirements [27, 28]. Almost all the modern gradient-based methods work well with numerical attributes. If the dataset contains both numerical and categorical variables, then the categorical ones must be converted to numerical ones; this however leads to a potential decrease in the model's accuracy.

CatBoost is a gradient enhancement library whose main advantage lies in that it works well with categorical features [29]. One-hot encoding is used for processing categorical features, but this method incurs more computational complexity and memory owing to its high cardinalities. Therefore, an effective way to process categorical features is to use the CatBoost algorithm based on modified target statistics.

## Model performance evaluation

In the previous section, the variables that showed a significant difference between the two groups were selected through a traditional statistical analysis. The stroke SI model was tested using the K-fold validation dataset [30]. The overall model predictive performance was assessed using the area under the receiver operating characteristic (ROC) curve. The performance characteristics of the stroke SI model indicate sensitivity, specificity, accuracy, positive predictive value (PPV), and negative predictive value (NPV) results. A sub-assessment analysis was performed by selecting the model with the highest performance. For each assessment, a separate ROC curve was constructed according to the predictions obtained from the highest-performance model and the outcomes within each assessment. The importance and relationship of stroke SI variables were derived through Shapley additive explanation (SHAP) values. The red and blue dots indicate that the variables at that point had positive and negative effects on the SI occurrence, respectively. The data were analyzed using Python 3.7.12 (Python Software Foundation).

## Results

The demographic data of the stroke patients are presented in Table 1. The variables that showed a significant difference between the two groups based on the SI outcome criterion were age, onset, type, socioeconomic level, and education level ( $p < 0.05$ ). The high SI group had a higher frequency of older adults over 65 years of age than the low SI group. The onset group had a higher frequency when the stroke onset was less than 6 months, the socioeconomic level was poor, and the education level was low.

Table 1  
Demographic and clinical characteristics based on suicidal ideation

| Variables                  | Total (n = 304) |      | SI low (n = 165) |      | SI High (n = 139) |      | p    |
|----------------------------|-----------------|------|------------------|------|-------------------|------|------|
|                            | n               | %    | n                | %    | n                 | %    |      |
| Sex                        |                 |      |                  |      |                   |      | .090 |
| Male                       | 220             | 72.4 | 126              | 76.4 | 94                | 67.6 |      |
| Female                     | 84              | 27.6 | 39               | 23.6 | 45                | 32.4 |      |
| Age                        |                 |      |                  |      |                   |      | .016 |
| Under 45                   | 76              | 25.0 | 49               | 29.7 | 27                | 19.4 |      |
| 45 ~ 54                    | 131             | 43.1 | 67               | 40.6 | 64                | 46.0 |      |
| 55 ~ 64                    | 66              | 21.7 | 39               | 23.6 | 27                | 19.4 |      |
| 65 over                    | 31              | 10.2 | 10               | 6.1  | 21                | 15.1 |      |
| Onset                      |                 |      |                  |      |                   |      | .015 |
| Subacute                   | 174             | 57.2 | 84               | 50.9 | 90                | 64.7 |      |
| Chronic                    | 130             | 42.8 | 81               | 49.1 | 49                | 35.3 |      |
| Type                       |                 |      |                  |      |                   |      | .009 |
| Ischemic                   | 235             | 77.3 | 137              | 83.0 | 98                | 70.5 |      |
| Hemorrhage                 | 69              | 22.7 | 28               | 17.0 | 41                | 29.5 |      |
| Affected side              |                 |      |                  |      |                   |      | .084 |
| Right                      | 178             | 58.6 | 104              | 63.0 | 74                | 53.2 |      |
| Left                       | 126             | 41.4 | 61               | 37.0 | 65                | 46.8 |      |
| Dominant hand              |                 |      |                  |      |                   |      | .242 |
| Right                      | 291             | 95.7 | 160              | 97.0 | 131               | 94.2 |      |
| Left                       | 13              | 4.3  | 5                | 3.0  | 8                 | 5.8  |      |
| Socioeconomic level        |                 |      |                  |      |                   |      | .001 |
| High<br>(Health insurance) | 267             | 87.8 | 154              | 93.3 | 113               | 81.3 |      |

Abbreviation: SI, suicidal ideation

| Variables             | Total ( <i>n</i> = 304) |      | SI low ( <i>n</i> = 165) |      | SI High ( <i>n</i> = 139) |      | p    |
|-----------------------|-------------------------|------|--------------------------|------|---------------------------|------|------|
|                       | n                       | %    | n                        | %    | n                         | %    |      |
| Low<br>(Medical care) | 37                      | 12.2 | 11                       | 6.7  | 26                        | 18.7 |      |
| Hypertension          |                         |      |                          |      |                           |      | .870 |
| Yes                   | 180                     | 59.2 | 97                       | 58.8 | 83                        | 59.7 |      |
| No                    | 124                     | 40.8 | 68                       | 41.2 | 56                        | 40.3 |      |
| Diabetes              |                         |      |                          |      |                           |      | .295 |
| Yes                   | 90                      | 29.6 | 53                       | 32.1 | 37                        | 26.6 |      |
| No                    | 214                     | 70.4 | 112                      | 67.9 | 102                       | 73.4 |      |
| Marital status        |                         |      |                          |      |                           |      | .056 |
| Married               | 222                     | 73.0 | 124                      | 75.2 | 98                        | 70.5 |      |
| Unmarried             | 59                      | 19.4 | 34                       | 20.6 | 25                        | 18.0 |      |
| Divorced/Widowed      | 23                      | 7.6  | 7                        | 4.2  | 16                        | 11.5 |      |
| Family history        |                         |      |                          |      |                           |      | .092 |
| Yes                   | 170                     | 55.9 | 85                       | 51.5 | 85                        | 61.2 |      |
| No                    | 134                     | 44.1 | 80                       | 48.5 | 54                        | 38.8 |      |
| Past history          |                         |      |                          |      |                           |      | .421 |
| Yes                   | 119                     | 39.1 | 68                       | 41.2 | 51                        | 36.7 |      |
| No                    | 185                     | 60.9 | 97                       | 58.8 | 88                        | 63.3 |      |
| Smoking               |                         |      |                          |      |                           |      | .408 |
| Yes                   | 174                     | 57.2 | 98                       | 59.4 | 76                        | 54.7 |      |
| No                    | 130                     | 42.8 | 67                       | 40.6 | 63                        | 45.3 |      |
| Drinking              |                         |      |                          |      |                           |      | .791 |
| Yes                   | 262                     | 86.2 | 143                      | 86.7 | 119                       | 85.6 |      |
| No                    | 42                      | 13.8 | 22                       | 13.3 | 20                        | 14.4 |      |
| Education             |                         |      |                          |      |                           |      | .017 |
| Uneducated            | 12                      | 3.9% | 1                        | 0.6  | 11                        | 7.9  |      |

Abbreviation: SI, suicidal ideation

| Variables                           | Total ( <i>n</i> = 304) |       | SI low ( <i>n</i> = 165) |      | SI High ( <i>n</i> = 139) |      | p    |
|-------------------------------------|-------------------------|-------|--------------------------|------|---------------------------|------|------|
|                                     | n                       | %     | n                        | %    | n                         | %    |      |
| Elementary                          | 7                       | 2.3%  | 4                        | 2.4  | 3                         | 2.2  |      |
| Middle School                       | 18                      | 5.9%  | 8                        | 4.8  | 10                        | 7.2  |      |
| High School                         | 216                     | 71.1% | 125                      | 75.8 | 91                        | 65.5 |      |
| University                          | 51                      | 16.8% | 27                       | 16.4 | 139                       | 17.3 |      |
| Transfer                            |                         |       |                          |      |                           |      | .108 |
| Wheelchair                          | 91                      | 29.9  | 43                       | 26.1 | 48                        | 34.5 |      |
| Ambulation                          | 213                     | 70.1  | 122                      | 73.9 | 91                        | 65.5 |      |
| Abbreviation: SI, suicidal ideation |                         |       |                          |      |                           |      |      |

The results presented in Table 2 indicate a significant difference in ADL and emotions in both the groups ( $p < 0.05$ ). In particular, there was a significant difference between the two groups in the emotional domain ( $p < 0.001$ ). Cognition and motor functions, on the other hand, did not differ between the two groups.

Table 2

Comparison of cognitive functions, motor functions, ADL, emotion functions between both groups

| Domain  | SI low ( <i>n</i> = 165) |         | SI High ( <i>n</i> = 139) |         | p    |
|---|--------------------------|---------|---------------------------|---------|------|
|   | Mean (SD)                | IQR     | Mean (SD)                 | IQR     |      |
| Cognition   |                          |         |                           |         |      |
| MMSE  | 23.99 (1.54)             | 23–25   | 23.91 (1.70)              | 23–25   | .784 |
| Motor   |                          |         |                           |         |      |
| MFT   | 22.28 (2.04)             | 21–24   | 22.03 (1.78)              | 21–23   | .207 |
| ADL   |                          |         |                           |         |      |
| MBI   | 72.45 (5.99)             | 69–77   | 71.00 (6.59)              | 66–77   | .023 |
| Emotion   |                          |         |                           |         |      |
| Self-efficacy   | 189.18 (30.61)           | 181–216 | 154.24 (30.89)            | 131–156 | .001 |
| Rehabilitation<br>motivation  | 90.77 (14.61)            | 81–100  | 73.14 (14.37)             | 63–82   | .001 |
| BAI   | 15.53 (4.04)             | 12–17   | 21.83 (5.22)              | 18–23   | .001 |
| BDI   | 14.35 (4.47)             | 13–18   | 20.08 (4.62)              | 19–24   | .001 |
| Abbreviations: SI, suicidal ideation; SD, standard deviation; IQR, interquartile range; MMSE, mini-mental state examination; MFT, manual function test; MBI, modified bathel index; BAI, beck anxiety inventory; BDI, beck depression inventory |                          |         |                           |         |      |

Among the three models, the area under the AUC value was higher for the CatBoost model than the other two models, and most values (sensitivity, NPV) outperformed the XGBoost and LGBM scores. The specificity, PPV, accuracy values were the highest in Xgboost (Supplementary information 2). As shown in Table 3, emotional features such as depression, anxiety, self-efficacy, and rehabilitation motivation showed the best results in the CatBoost model. Sensitivity was the Rehabilitation motivation, specificity was depression and anxiety, NPV was self-efficacy, PPV was the Rehabilitation motivation, and depression showed the highest accuracy. Additionally, as for the cut-off points, BDI showed a mild depressive state, and K-MBI showed a cut-off point of moderate dependence, whereas BAI showed a normal level cut-off point.

Table 3  
Result of the CatBoost model based on emotion and ADL data

| Variables  | Sensitivity | Specificity | PPV  | NPV  | Accuracy | Cut off value |
|--|-------------|-------------|------|------|----------|---------------|
| MBI  | .317        | .861        | .656 | .599 | .612     | 68            |
| Self-efficacy  | .799        | .824        | .792 | .829 | .812     | 177           |
| Rehabilitation<br>motivation   | .899        | .636        | .675 | .882 | .757     | 86            |
| BAI  | .791        | .794        | .763 | .818 | .793     | 18            |
| BDI  | .813        | .794        | .768 | .834 | .803     | 19            |
| Abbreviations: ADL, activity daily living; PPV, positive predict value; NPV, negative predict value; MBI, modified bathel index; BAI, beck anxiety inventory; BDI, beck depression inventory |             |             |      |      |          |               |

Supplementary information 1 shows the ROC curves of the CatBoost classifier for the five functional assessments. The AUC values were ordered as per the order presented in Table 3: first, negative emotion evaluation, such as evaluation of depression and anxiety; second, positive emotion evaluation; and third, ADL assessment. Furthermore, the AUC value, which includes the demographic variables that indicated a significant difference between the two groups, as well as the exercise and emotion evaluation, showed the highest result.

Regarding SHAP, depression was found to be the most important predictor for SI in stroke patients, followed by emotional variables such as self-efficacy, anxiety, and rehabilitation motivation. In the SHAP summary plot result (Figure 2), it was seen that the higher the negative emotions such as depression and anxiety, the higher the SI. Conversely, the lower the positive emotions such as self-efficacy and rehabilitation motivation, the higher the SI.

Using the SHAP dependence plot, the results of the interaction relationship between anxiety, rehabilitation motivation, self-efficacy, and ADL that exhibited significant differences were derived based on depression, which demonstrated the greatest importance for SI in stroke patients. Negative emotions, such as anxiety and depression, showed a positive relationship, and positive emotions, such as rehabilitation motivation and self-efficacy, exhibited an inverse relationship with SI. There was no evident association between depression and ADL function (Figure 3).

## Discussion

In this study, using stroke patients' data from a rehabilitation hospital, we developed and validated a model for SI prediction in stroke patients within a post-onset period. Using the statistically significant predictors that a stroke patient can report in a direct interview and survey, performance was compared for the three boosting models.

Using the chi-square test for the demographic variables used in this study, statistically significant differences were observed between the two groups divided on the basis of age, onset, stroke type, and economic and education level. Among them, the high SI group had a high proportion of participants aged 65 years, an onset of less than six months, hemorrhagic stroke, and low economic, and education levels. This suggested that risk factors for SI in stroke patients increased in various pathologies due to rapid changes that take place associated with old age, loss and maladaptation immediately after onset [9], hemorrhagic stroke, severe pain, poor prognosis [31], low socioeconomic level, and low educational level. This can be seen as a low-income group [32, 33]. Additionally, there was a significant difference for widowed or divorced patients, which showed an approximate result (Table 1). This finding was consistent with a previous study that indicated a large difference depending on whether the support of the family or spouse was present [34].

Based on the study results, a statistically significant difference between the two groups in the variables of ADL and emotional function was noted. In previous studies, cognitive dysfunction was found to be associated with suicide [33, 35], which was not observed in the results of the current study. The cognitive function evaluation tool used in this study, the MMSE, is simple and efficient; however, we believe that it may have been affected by low sensitivity, as it is a screening tool for mild cognitive impairment [36]. In the case of MFT, lower extremity functions, such as gait function [37, 38], that affect depression in stroke patients were not included, and so, there was no significant difference between the two groups. In contrast, depression can be viewed as the biggest risk factor for SI according to previous studies' results [39], and has previously showed a strong correlation with ADL, anxiety, self-efficacy, and Rehabilitation motivation [39, 40]. Therefore, it is thought that there was a significant difference between the two groups in ADL and emotional variables.

Only statistically significant demographic and functional domain variables were applied to the three boosting models to derive their respective performances [41, 42]. After comparing the performance of the three models, it was found that LGBM had the most inferior performance, whereas Xgboost showed the best performance in terms of specificity, PPV, and accuracy. Further, CatBoost showed the best performance in terms of sensitivity, NPV, and AUC (Supplementary information 2). While XGBoost and LightGBM offer several advantages, it must be noted that 16 out of the 23 variables of the stroke data used in this study were categorical. When a large number of categorical features are present in the dataset, then CatBoost may offer a more efficient performance [43, 44]. In addition, LGBM is disadvantageous in that its application to small datasets (i.e., fewer than 10,000 cases) leads to leaf-wise growth, which, in turn, causes significant overfitting, whereas XGBoost cannot handle categorical features on its own [45, 46]. Additionally, the classification performance improved when more features were added to the classifiers (Supplementary information1). The predicted results can be used to take the necessary precautions and improve the function of stroke patients. Further, the AUC of the best classifiers was approximately 0.900. This value can be said to be sufficient for the reliable prediction of patients' functional outcomes [47].

Figure 2 (a) shows the absolute influence of each variable of CatBoost through SHAP on the model. Notably, it is crucial for physicians to understand the effect of various factors on the SI of stroke patients. The variable that showed the greatest influence on stroke occurrence in patient SI was “depression,” followed by “self-efficacy,” “anxiety,” “rehabilitation motivation,” and so forth. The emotion function level had a significant influence on the occurrence of SI in stroke patients. Figure 2 (b) is a SHAP summary showing the degree of influence of each variable on stroke patient SI prediction [42]. Thus, higher levels of “depression” and “anxiety” meant that the probability of SI occurrence increased [48]. Therefore, the higher the “self-efficacy” and “rehabilitation motivation,” the lower the probability of SI occurrence, thereby exhibiting an inverse relationship with each other. Figure 3 is a SHAP partial dependent plot showing the correlation between depression, the most influential SI predictor in stroke patients, and other important predictive factors. Positive emotions, such as rehabilitation motivation, and self-efficacy, are observed to have a negative correlation (Figure 3 b, c). The results thus obtained were identical to those reported in previous studies on depression, anxiety, rehabilitation motivation, and self-efficacy in stroke patients; negative and positive emotions were found to be the main factors affecting the SI of stroke patients; further, it was found that the two had opposite effects on each other [49, 50].

The stroke SI prediction model developed in this study can therefore be used to classify stroke patients into low- and high-risk SI groups based on routinely collected medical data and self-report questions. Furthermore, improved characterization of low and high risk for stroke-related SI can be achieved by analyzing the importance and correlation of the model’s prediction features. The implementation of a stroke SI prediction model in public health systems may facilitate early stroke SI detection and intervention programs, thereby reducing suicidal ideation. Additionally, it should be noted that a prediction model is only a tool to support the clinician and therefore cannot be used to replace personal judgment.

## Limitations

This study has some limitations. First, prospective clinical trials are needed to demonstrate a clear clinical benefit of the addition of a stroke SI prediction model to the clinical intervention system. Clearer information about risk predictors can be provided by collecting additional data. Second, the study results cannot be generalized for all stroke features, such as biochemical indices and lesion location, which are also considered risk factors. Future studies should combine these to reveal the interactions of pathophysiological risk factors [15]. In a follow-up study, the model may benefit from the inclusion of as yet unavailable contributing predictors, such as invasive test data like quantitative brain structural and functional imaging data of stroke patients.

## Conclusion

We constructed a comprehensive risk prediction model for SI in stroke patients based on clinical and psychological features. The model indicated that psychological factors were important for identifying SI risk in subacute and chronic stroke patients and contributed to post-stroke rehabilitation and mental health. Furthermore, the prediction model ultimately works as a decision tool to help clinicians identify

the SI risk early, which will allow the optimization of stroke patients' suicide prevention strategies in personalized medicine.

## List Of Abbreviations

Post-stroke depression (PSD)

Suicidal ideation (SI)

Activities of daily living (ADL)

Mini-Mental State Examination (MMSE-K)

Manual Function Test (MFT)

Rehabilitation Motivation Scale (RMS)

Modified Bathel Index (K-MBI)

Beck Anxiety Inventory (BAI)

Beck Depression Inventory (BDI)

Light GBM (LGBM)

Positive predictive value (PPV)

Negative predictive value (NPV)

Receiver operating characteristic (ROC)

Shapley additive explanation (SHAP)

## Declarations

### Ethics approval and consent to participate

The ethics committee of Daegu University Institutional Review Board (IRB) approved this study (1040621-202111-HR-079). This is a retrospective study using anonymized data obtained with written informed consent from all patients. This study has been independently reviewed and approved by an IRB.

### Consent for publication

Not applicable

### Availability of data and materials

Due to privacy/ethical restriction, data are available from the corresponding author on reasonable request.

### **Competing interests**

The authors declare that they have no competing interests.

### **Funding**

This research was supported by a grant of the Korea Health Technology R&D Project through the Korea Health Industry Development Institute (KHIDI), funded by the Ministry of Health & Welfare, Republic of Korea [grant number: HI21C1074]. The funding source had no role in the design of the study; collection, analysis and interpretation of data; in the writing of the report; and in the decision to submit the article for publication.

### **Authors' contributions**

S.I.S., S.B.L contributed to design the study, and was a major contributor in writing the manuscript. S.I.S., H.T.H collected and pre-processing the stroke patient data regarding the rehabilitation. S.I.S., S.B.L, C.L analyzed and interpreted the stroke patient data regarding the rehabilitation. All authors read and approved the final manuscript.

### **Acknowledgements**

We would like to express our deepest gratitude to Prof. Eui kyu Chie and Prof. Hwan Kim for their support. Furthermore, we would like to thank Editage ([www.editage.com](http://www.editage.com)) for English language editing.

## **References**

1. Umphred DA, Lazaro RT. Neurological rehabilitation. Elsevier Health Sciences; 2012
2. De Ryck A, Brouns R, Geurden M, Elseviers M, De Deyn PP, Engelborghs S. Risk factors for poststroke depression: identification of inconsistencies based on a systematic review. *J Geriatr Psychiatry Neurol.* 2014;27:147–58. <https://doi.org/10.1177/0891988714527514>
3. Pohjasvaara T, Vataja R, Leppävuori A, Kaste M, Erkinjuntti T. Suicidal ideas in stroke patients 3 and 15 months after stroke. *Cerebrovasc Dis.* 2001;12:21–6. <https://doi.org/10.1159/000047676>
4. Fuller-Thomson E, Tulipano MJ, Song M. The association between depression, suicidal ideation, and stroke in a population-based sample. *Int J Stroke.* 2012;7:188–94. <https://doi.org/10.1111/j.1747-4949.2011.00702.x>
5. Park SM. The impact of the COVID-19 pandemic on mental health among population. *kjhep.* 2020;37:83–91. <https://doi.org/10.14367/kjhep.2020.37.5.83>
6. Poudel K, Subedi P. Impact of COVID-19 pandemic on socioeconomic and mental health aspects in Nepal. *Int J Soc Psychiatry.* 2020;66:748–55. <https://doi.org/10.1177/0020764020942247>

7. Mash HBH, Ursano RJ, Kessler RC, Naifeh JA, Fullerton CS, Aliaga PA, et al. Predictors of suicide attempt within 30 days after first medically documented suicidal ideation in US Army soldiers. *Am J Psychiatry*. 2021;178:1050–9. <https://doi.org/10.1176/appi.ajp.2021.20111570>
8. Faber RA. Suicide in neurological disorders. *Neuroepidemiology*. 2003;22:103–5. <https://doi.org/10.1159/000068751>
9. Pompili M, Venturini P, Campi S, Seretti ME, Montebovi F, Lamis DA, et al. Do stroke patients have an increased risk of developing suicidal ideation or dying by suicide? An overview of the current literature. *CNS Neurosci Ther*. 2012;18:711–21. <https://doi.org/10.1111/j.1755-5949.2012.00364.x>
10. Shin KM, Cho SM, Hong CH, Park KS, Shin YM, Lim KY, et al. Suicide among the elderly and associated factors in South Korea. *Aging Ment Health*. 2013;17:109–14. <https://doi.org/10.1080/13607863.2012.702732>
11. Brugnara G, Neuberger U, Mahmutoglu MA, Foltyn M, Herweh C, Nagel S, et al. Multimodal predictive modeling of endovascular treatment outcome for acute ischemic stroke using machine-learning. *Stroke*. 2020;51:3541–51. <https://doi.org/10.1161/STROKEAHA.120.030287>
12. Heo J, Yoon JG, Park H, Kim YD, Nam HS, Heo JH. Machine learning–based model for prediction of outcomes in acute stroke. *Stroke*. 2019;50:1263–5. <https://doi.org/10.1161/STROKEAHA.118.024293>
13. Scrutinio D, Ricciardi C, Donisi L, Losavio E, Battista P, Guida P, et al. Machine learning to predict mortality after rehabilitation among patients with severe stroke. *Sci Rep*. 2020;10:20127. <https://doi.org/10.1038/s41598-020-77243-3>
14. Tozlu C, Edwards D, Boes A, Labar D, Tsagaris KZ, Silverstein J, et al. Machine learning methods predict individual upper-limb motor impairment following therapy in chronic stroke. *Neurorehabil Neural Repair*. 2020;34:428–39. <https://doi.org/10.1177/1545968320909796>
15. Liu R, Yue Y, Jiang H, Lu J, Wu A, Geng D, et al. A risk prediction model for post-stroke depression in Chinese stroke survivors based on clinical and socio-psychological features. *Oncotarget*. 2017;8:62891–9. <https://doi.org/10.18632/oncotarget.16907>
16. Wang J, Zhao D, Lin M, Huang X, Shang X. Post-stroke anxiety analysis via machine learning methods. *Front Aging Neurosci*. 2021;13:657937. <https://doi.org/10.3389/fnagi.2021.657937>
17. Folstein MF, Folstein SE, McHugh PR. “Mini-mental state”: a practical method for grading the cognitive state of patients for the clinician. *J Psychiatr Res*. 1975;12:189–98. [https://doi.org/10.1016/0022-3956\(75\)90026-6](https://doi.org/10.1016/0022-3956(75)90026-6)
18. Miyamoto S, Kondo T, Suzukamo Y, Michimata A, Izumi S-I. Reliability and validity of the Manual Function Test in patients with stroke. *Am J Phys Med Rehabil*. 2009;88:247–55. <https://doi.org/10.1097/PHM.0b013e3181951133>
19. Shah S, Vanclay F, Cooper B. Improving the sensitivity of the Barthel index for stroke rehabilitation. *J Clin Epidemiol*. 1989;42:703–9. [https://doi.org/10.1016/0895-4356\(89\)90065-6](https://doi.org/10.1016/0895-4356(89)90065-6)
20. Sherer M, Maddux JE, Mercandante B, Prentice-Dunn S, Jacobs B, Rogers RW. The self-efficacy scale: construction and validation. *Psychol Rep*. 1982;51:663–71.

<https://doi.org/10.2466/pr0.1982.51.2.663>

21. Kim H, Hwang Y, Yu J, Jung J, Woo H, Jung H. The correlation between depression, motivation for rehabilitation, activities of daily living, and quality of life in stroke patients. *The J Korean Soc Occup Ther.* 2009;17:41–53
22. Beck AT, Steer R. Beck Anxiety Inventory (BAI). Überblick Reliabilitäts Validitätsbefunde Klin Außerklinischen Selbst Fremdbeurteilungsverfahren. 1988;7
23. Beck AT, Steer RA, Carbin MG. Psychometric properties of the Beck Depression Inventory: twenty-five years of evaluation. *Clin Psychol Rev.* 1988;8:77–100. [https://doi.org/10.1016/0272-7358\(88\)90050-5](https://doi.org/10.1016/0272-7358(88)90050-5)
24. Beck AT, Kovacs M, Weissman A. Assessment of suicidal intention: the Scale for Suicide Ideation. *J Consult Clin Psychol.* 1979;47:343–52. <https://doi.org/10.1037/0022-006x.47.2.343>
25. Shin MS, Park KB, Oh KJ, Kim ZS. A study of suicidal ideation among high school students: the structural relation among depression, hopelessness, and suicidal ideation. *Korean J Clin Psychol.* 1990;9:1–19
26. Chen T, Xgboost GC. A scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*; 2016
27. Habib A-ZSB, Tasnim T, Billah MM. A study on coronary disease prediction using boosting-based ensemble machine learning approaches. Paper presented at the 2019 2nd International Conference on Innovation in Engineering and Technology (ICIET); 2019
28. Saber M, Boulmaiz T, Guermoui M, Abdrado KI, Kantoush SA, Sumi T, et al. Examining LightGBM and CatBoost models for wadi flash flood susceptibility prediction. *Geocarto Int.* 2021;1–26. <https://doi.org/10.1080/10106049.2021.1974959>
29. Dorogush AV, Ershov V, Gulin A. CatBoost: gradient boosting with categorical features support [Internet]. arXiv [Preprint]. Available from: arXiv:1810.11363
30. Chekroud AM, Bondar J, Delgadillo J, Doherty G, Wasil A, Fokkema M, et al. The promise of machine learning in predicting treatment outcomes in psychiatry. *World Psychiatry.* 2021;20:154–70. <https://doi.org/10.1002/wps.20882>
31. Jørgensen HS, Nakayama H, Raaschou HO, Vive-Larsen J, Støier M, Olsen TS. Outcome and time course of recovery in stroke. Part I: Outcome. The Copenhagen Stroke Study. *Arch Phys Med Rehabil.* 1995;76:399–405. [https://doi.org/10.1016/s0003-9993\(95\)80567-2](https://doi.org/10.1016/s0003-9993(95)80567-2)
32. Kim J-Y, Lee D-H, Hwang J-W, Lee K-U. Factors influencing suicidal ideation among lower-income group participating self-sufficiency Program in Gangwon Province, Korea. *J Korea Contents Assoc.* 2016;16:91–101. <https://doi.org/10.5392/JKCA.2016.16.12.091>
33. Park E. Suicide ideation and the related factors among Korean adults by gender. *J Agr Med Community Health.* 2014;39:161–75. <https://doi.org/10.5393/JAMCH.2014.39.3.161>
34. Morris PL, Robinson RG, Raphael B, Bishop D. The relationship between the perception of social support and post-stroke depression in hospitalized patients. *Psychiatry.* 1991;54:306–16. <https://doi.org/10.1080/00332747.1991.11024559>

35. Choi R, Moon H-J, Hwang B-D. The influence of chronic disease on the stress cognition, depression experience and suicide thoughts of the elderly. *The Korean. Health Serv Manag.* 2010;4:73–84
36. Kang Y, Na DL, Hahn S. A validity study on the Korean Mini-Mental State Examination (K-MMSE) in dementia patients. *J Korean Neurol Assoc.* 1997;15:300–8
37. Carod-Artal FJ, Egido JA. Quality of life after stroke: the importance of a good recovery. *Cerebrovasc Dis.* 2009;27;Suppl 1:204–14. <https://doi.org/10.1159/000200461>
38. Kim C, Koo K. The effects of physical activities of disabled men with stroke on depression and suicidal ideation. *KAHPERD.* 2017;56:657–64. <https://doi.org/10.23949/kjpe.2017.05.56.3.49>
39. Yu S-J, Kim H-S, Kim K-S, Baik H-G. The effects of community-based self-help management program by strengthening self-efficacy of post stroke elderly patients. *The Korean J Rehabil Nurs.* 2001;4:187–97
40. Diekstra RF. The epidemiology of suicide and parasuicide. *Acta Psychiatr Scand Suppl.* 1993;371:9–20. <https://doi.org/10.1111/j.1600-0447.1993.tb05368.x>
41. Choi J, Yang H, Oh H. Store sales prediction using gradient boosting model;25. *J Korea Institute Inf Commun Eng.* p. 171–7; 2021
42. Oh H-R, Son A-L, Lee Z. Occupational accident prediction modeling and analysis using SHAP. *dcs.* 2021;22:1115–23. <https://doi.org/10.9728/dcs.2021.22.7.1115>
43. Prokhorenkova L, Gusev G, Vorobev A, Dorogush AV, Gulin A. CatBoost: unbiased boosting with categorical features [Internet]. arXiv [Preprint]. 2017. Available from: arXiv:1706.09516
44. Chu Y, Knell G, Brayton RP, Burkhart SO, Jiang X, Shams S. Machine learning to predict sports-related concussion recovery using clinical data. *Ann Phys Rehabil Med.* 2022;65:101626. <https://doi.org/10.1016/j.rehab.2021.101626>
45. Ge X, Sun J, Lu B, Chen Q, Xun W, Jin Y. Classification of oolong tea varieties based on hyperspectral imaging technology and BOSS-LightGBM model. *J Food Process Eng.* 2019;42:e13289. <https://doi.org/10.1111/jfpe.13289>
46. Swalin A. CatBoost vs. Light GBM vs. XGBoost. *Towards Data Sci.* 2018;11
47. Muller MP, Tomlinson G, Marrie TJ, Tang P, McGeer A, Low DE, et al. Can routine laboratory tests discriminate between severe acute respiratory syndrome and other causes of community-acquired pneumonia? *Clin Infect Dis.* 2005;40:1079–86. <https://doi.org/10.1086/428577>
48. Forkmann T, Brähler E, Gauggel S, Glaesmer H. Prevalence of suicidal ideation and related risk factors in the German general population. *J Nerv Ment Dis.* 2012;200:401–5. <https://doi.org/10.1097/NMD.0b013e31825322cf>
49. Almhdawi KA, Alazrai A, Kanaan S, Shyyab AA, Oteir AO, Mansour ZM, et al. Post-stroke depression, anxiety, and stress symptoms and their associated factors: a cross-sectional study. *Neuropsychol Rehabil.* 2021;31:1091–104. <https://doi.org/10.1080/09602011.2020.1760893>
50. Robinson RG, Jorge RE. Post-stroke depression: a review. *Am J Psychiatry.* 2016;173:221–31. <https://doi.org/10.1176/appi.ajp.2015.15030363>

# Figures

## Figure 1

Stroke suicidal ideation prediction model.

## Figure 2

Feature importance based on SHAP values: (a) Mean absolute SHAP values (b) Summary

## Figure 3

Partial dependence plot by SHAP value. Relationship between (a) self-efficacy and depression (b) rehabilitation motivation and depression (c) anxiety and depression

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementaryinformation1.docx](#)
- [Supplementaryinformation2.docx](#)