

A Novel Recommender System for Stroke Risk Stratification

Nour Dekhil

National Engineering School of Sfax
nour.dekhil@enis.tn

Yasin Mamatjan, Safwat Hassan

Thompson Rivers University
ymamatjan@tru.ca, shassan@tru.ca

Mira Salih

Beth Israel Deaconess Medical Center
mirasalih28@gmail.com

Abstract—Stroke is one of the leading causes of neurological deficits and disability. Early prediction or detection of strokes is required to reduce mortality and morbidity rates. In this study, we propose a novel recommender approach to identify key risk factors impacting risk stratification for a particular patient. Our proposed approach achieves high accuracy in predicting patient’s risk level with an average AUC of 0.98 and an average F1 score of 0.91. Our approach adopted model interpretation techniques (i.e., LIME and SHAP) to help patients understand the risk factors associated with the predicted risk level. Finally, we compare LIME and SHAP results for prediction explanation. The results showed strong agreement between both models on the selection of top risk factors for a particular risk stratification. The obtained results show that our proposed approach can help patients and clinicians better understand patient specific risk factors associated with stroke risk severity that eventually prevents strokes by controlling these risk factors.

I. INTRODUCTION

Stroke is a major cause of mortality. Early stroke risk stratification helps to understand stroke risk to prevent stroke and its outcome. Additionally, identifying key stroke risk factors, such as cholesterol and blood pressure, and quantifying their effects may help provide appropriate care to patients. This study introduces a novel approach that uses patient data (e.g. age and blood pressure) to automatically predict the patient’s stroke risk level and rank patient specific risk factors that are most impacting the decision. First, our approach uses CatBoost Classifier [1] to predict the patient’s stroke risk level. Then, our approach uses Local Interpretable Model-Agnostic Explanations (LIME) [2] and SHapley Additive exPlanations (SHAP) [3] to quantify the importance of each risk factor.

Prior work uses machine learning techniques to predict stroke level and assist clinicians in their diagnosis [4], [5]. However, the lack of explainability of these black-box models makes it hard to leverage the obtained results. Therefore, our primary goal is to provide meaningful insights of why a particular decision was made. A study has been conducted by ElShawi et al. [6] to address this challenge in healthcare by introducing interpretable machine learning techniques such as LIME, Anchors, and SHAP. Chang et al. [7] applied these interpretability techniques on acute kidney injury prediction to understand individualized predictions. However, these techniques have not been applied to stroke risk stratification yet. This work aims to develop a novel approach to predict stroke risk level and identify patient specific risk factors.

II. METHODOLOGY

In this paper, we propose a smart health recommender system with a novel technique to identify and rank risk factors that are most important when predicting stroke severity level. Therefore, the objective of this research is twofold. First, we used CatBoost Classifier to predict the stroke severity class of an individual. Second, we quantified the impact of each risk factor on the prediction using SHAP and LIME. CatBoost, SHAP and LIME are briefly described in the next subsections.

A. Dataset

The stroke analysis dataset [8] is used in this study. The dataset contains 4,798 information about 4,798 patients and 12 risk factors and stroke risk level output. The output response is a stroke risk presented in three levels such as no risk, low risk, moderate risk, and high risk. The 12 risk factors used for the prediction of stroke risk are as follows: age, National Institutes of Health Stroke Scale (NIHSS), Modified Rankin Scale (mRS), systolic blood pressure, diastolic blood pressure, glucose level, paralysis, smoking status, Body Mass Index (BMI), and cholesterol.

B. Stroke Risk Level Classifier

We used CatBoost to build the classification model for four risk level classes. CatBoost is a machine learning algorithm that can handle categorical features by converting them into numbers without any pre-processing effort. In addition, CatBoost uses Bayesian estimators to reduce the risk of overfitting caused by traditional gradient boosting algorithms.

C. SHapley Additive exPlanations (SHAP)

SHAP is a model interpretation method that uses the concept of game theory to interpret a particular prediction by determining the contribution of each feature in terms of its Shapley value. SHAP helps in reverse-engineering any black-box model’s results, which makes it model agnostic [9]. This technique provides both local and global interpretability, but the focus of our study is on identifying what the prediction locally depends on. SHAP presents three main classes - KernelExplainer for linear models, DeepExplainer for deep neural networks, and TreeExplainer for tree-based models. In this study, the TreeExplainer was used to interpret the predictions of the CatBoost classifier.

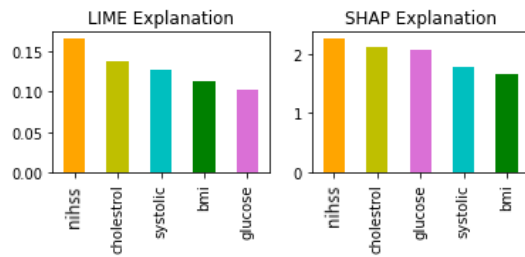


Fig. 1. Top 5 risk factors selected by LIME (on the left) and SHAP (on the right) using the color bars ranked in order for patient’s prediction result.

D. Local Interpretable Model-Agnostic Explanations

LIME is a method that aims to evaluate the impact of features on the prediction made for a particular instance. The main idea behind LIME is training a surrogate model to approximate the prediction of a complex model. It starts by perturbing the original data point and feeding it to an interpretable model such as linear regression model. Then, the perturbed data points are weighted according to their similarity to the instance to be explained. Finally, assigning a coefficient to each feature by linear regression.

III. RESULTS

A. Model Evaluation

The dataset was split into independent training and test sets, containing 3,838 and 960 patients respectively. To build the model, we set the number of iterations to 1,000 and the learning rate to 0.3. The classification model achieved an average AUC (Area under the ROC Curve) of 0.98 using one versus all approach. A Cohen’s Kappa score of 0.74 and a weighted average F1 score of 0.91 were also obtained.

B. Risk Factors Ranking

We provided insights of stroke risk stratification (CatBoost prediction) by adopting LIME and SHAP methods. We ranked the feature importance based on their respective SHAP values and feature coefficient (LIME). The top five risk factors found by each approach were presented in two bar graphs. Longer bars indicate a greater feature impact on making a particular prediction. Fig. 1 shows the top five risk factors contributing to the prediction made for a particular patient. This patient was a 72 year-old man with the following clinical attributes: NIHSS 42, mRS 6, Systolic 185, Diastolic 131, Glucose 293, Paralysis 3, Smoker, BMI 43, Cholesterol 243, and TOS 3. The predicted stroke risk level by CatBoost classifier was class 3 (High-Risk). SHAP and LIME agreed that the factors detected by CatBoost that increased the stroke risk level for this patient were NIHSS, Cholesterol, Systolic, BMI, and Glucose.

C. Models Consistency Analysis

Model consistency analysis was performed in Fig. 2. We calculated the number of similar risk factors among the top five ones generated by each model explainer (i.e., LIME and SHAP). For all risk levels, the explainers found at least two

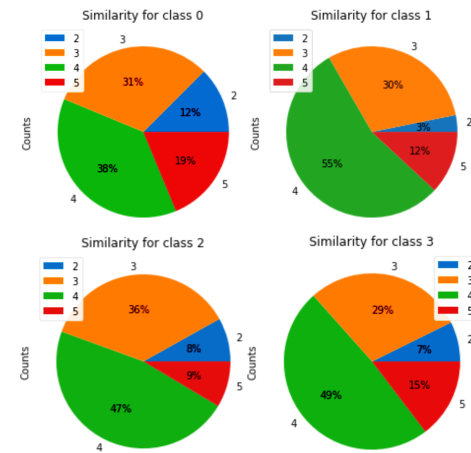


Fig. 2. Percentage similarity between SHAP and LIME in terms of top 2 to 5 overlapping risk factors for class 0 to 3 in order.

similar risk factors. In addition, LIME and SHAP showed high similarity of 88%, 97%, 92%, and 93% for class 0 to 3 respectively.

IV. CONCLUSION

In this paper, we proposed a novel recommender system for stroke risk stratification that can be used by both clinicians and patients to identify patient specific risk factors, which eventually prevent stroke by controlling these risk factors. We introduced SHAP and LIME techniques to rank risk factors after quantifying their effect on stroke risk stratification. A high level of similarity was achieved after applying both models to 960 patients from the test set. In addition, the CatBoost classifier made a reliable and accurate prediction of stroke risk score that can be applied in the clinic. Therefore, we are developing a recommender system based online tool to assist physicians and patients to detect strokes early.

REFERENCES

- [1] A. V. Drogush, V. Ershov, and A. Gulin, “CatBoost: gradient boosting with categorical features support,” *CoRR*, vol. abs/1810.11363, 2018.
- [2] M. T. Ribeiro, S. Singh and C. Guestrin, “Model-agnostic interpretability of machine learning,” *ArXiv*, 2016.
- [3] Scott M. Lundberg and Su-In Lee, “A Unified Approach to Interpreting Model Predictions,” *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 2017.
- [4] T. Tazin, M. N. Alam, N. N. Dola, M. S. Bari, S. Bourouis, and M. Monirujjaman Khan, “Stroke disease detection and prediction using robust learning approaches,” *Journal of Healthcare Engineering*, vol. 2021, pp. 1–12, 2021.
- [5] H. K V, H. P, G. Gupta, V. P, and P. K B, “Stroke prediction using machine learning algorithms,” *International Journal of Innovative Research in Engineering & Management*, vol. 8, no. 4, 2021.
- [6] R. ElShawi, Y. Sherif, M. Al-Mallah, and S. Sakr, “Interpretability inHealthcare: A Comparative Study of Local Machine Learning Interpretability Techniques,” *Comput. Intell.* 2020. doi:10.1111/coin.12410
- [7] C. Hu, Q. Tan, Q. Zhang, Y. Li, F. Wang, X. Zou, Z. Peng, “Application of interpretable machine learning for early prediction of prognosis in acute kidney injury,” *Computational and Structural Biotechnology Journal*, Volume 20, 2022, Pages 2861-2870
- [8] V. Bandi, D. Midhunchakkaravarthy, D. Bhattacharyya, “Stroke Analysis”, *Mendeley Data*, V1, 2020. doi: 10.17632/jpb5tds9f6.1
- [9] M.T. Ribeiro, S. Singh, C. Guestrin, “Model-agnostic interpretability of machine learning,” arXiv:1606.05386, 2016.