

# An Ensemble Machine Learning Approach for Stroke Prediction

P. Premisha\*

*Dept. of Information Communication Technology  
University of Vavuniya, Sri Lanka  
premishaprem@vau.ac.lk*

Mauran Kanagarathnam  
*Department of Physical Science  
University of Vavuniya, Sri Lanka  
maurank@vau.ac.lk*

Senthan Prasanth

*Department of Physical Sciences and Technology  
Sabaragamuwa University of Sri Lanka  
sprasanth@appsc.sab.ac.lk*

Kuhaneswaran Banujan  
*Department of Computing and Information Systems  
Sabaragamuwa University of Sri Lanka  
bhakuha@appsc.sab.ac.lk*

**Abstract** - Nowadays, one out of four people above 25 will suffer from a stroke. Especially this year, with the highest count of around 13.7 million people discovered with stroke for the first time. Out of 13.7 million, 5.5 million were fatalities. This was stated in a recent WHO study. It is estimated that if no action is taken, the number of fatalities will rise to 6.7 million yearly. The pandemic situation of COVID-19 will play a significant cause in the expanded death rate of stroke. Even for adults and patients with minor risk factors affected by stroke rather than in previous years. This study predicts the impact level of stroke with the development of an ensemble model by combining the various classifiers performed well in isolation. Predicting the stroke status in patients would help the physicians determine the prognosis and assist them in providing the targeted therapy in a limited time. During this study, an ensemble model was built by considering the base, bagging, and boosting classifiers: Support Vector Machine, Naïve Bayes, Decision Tree, Logistic Regression, Artificial Neural Network, Random Forest, XGBoost, LightGBM, and CatBoost. The dataset consists of 5110 patient details, along with 12 attributes that were analyzed in this research. The final ensemble model was developed by carrying out the methodology in two phases. During the first and second phases, the classifiers mentioned above were trained without hyperparameter tuning and with hyperparameter tuning and tested against the fundamental evaluation matrices. During each phase, the classifier that produces the highest classification accuracy is discovered from the base, bagging, and boosting categories. From the results obtained, the final ensemble model was constructed using the Max Voting approach, which yielded an accuracy of 95.76%.

**Keywords** - bagging and boosting, ensemble, Machine Learning, medical informatics, stroke