

ANALYSIS OF METHODS TO HANDLE MEDICAL SENSOR DATA TOWARDS HEALTH DISORDER IDENTIFICATION

Introduction

Medical sensors are important for monitoring the health condition of a person. The huge amount of data generated by these sensors, has a great potential for early detection of disorders. This study focuses on analysing methods to handle medical sensor data, particularly the data from Photo-plethysmograph (PPG). PPG quantifies the volumetric change of the heart by measuring the light transmission or reflection on arteries.

Objective

- Analyse various pre-processing approaches to handle issues arising in medical datasets, especially the imbalance of data items.
- Compare various classifiers to identify an optimal classifier for medical sensor data.

Dataset

- This study has used PPG-BP dataset containing features namely, Systolic Blood Pressure, Diastolic Blood Pressure, Heart Rate, Age, Sex and Body Mass Index (BMI) which were collected from 219 persons.
- Systolic blood pressure, diastolic blood pressure and heart rate were calculated from the PPG recordings.
- The dataset contains data for normal persons and persons with four disorders namely; hypertension, diabetes, cerebral infarction, and cerebrovascular disease.



S. Meruja and E.Y.A. Charles Department of Computer Science, Faculty of Science, University of Jaffna merujajoseph@gmail.com, charles.ey@univ.jfn.ac.lk

Methodology Dataset using sci-kit learn library. Train - Test Preprocess the dataset Training **Testing Data** Data particular disorder. Model Evaluation Development classifiers. Performance measure

→ Analyse the data and identify appropriate methods to represent data suitable for Machine Learning.

 \rightarrow Identify appropriate methods to handle class imbalances in the data set. The dataset selected for this study is highly imbalanced. Dataset contains a higher number of records from normal persons. To avoid the effect of *overfitting* during learning, this study used: Undersampling technique: Taking a random number of data with target value 'Normal' according to the ratio of other classes.

Duplicating the dataset: Make a copy of dataset and scaled it and added to the dataset. Again it leads to overfitting.

Undersampling and duplicating: Taking a random number of data with target value 'Normal' according to the ratio of other classes and a copy of the dataset is scaled and added to the dataset.

- → Identify appropriate Machine learning methods suitable for the selected dataset. This study selected SVM, K-Nearest Neighbor and Naïve Bayes algorithms.
- → Construction of classifiers Relationship between attributes and target values. Four different models were constructed, each for one disorder.
- → Evaluate the constructed models and revise the whole process to identify methods and models for optimal performance.

Testing Result

Accuracy, Precision, F1 score and Recall are used as the performance metric to validate the models and confusion matrix to analyse the performance of the proposed approaches.

```
Accuracy = (TP + TN)/(TP + TN + FP + FN)
Precision = TP / (TP + FP)
              = TP / (TP + FN)
Recall
```

F1 Score = 2*[(Precision * Recall)/(Precision + Recall)]

- True Positive, TN – True Negative, FP – False Positive, FN – False Negative TP Table 1: The classification rate of the proposed method for each disease.

Classifier	Hypertension		Diabetes		Cerebrovascular		Cerebral Infarction	
	Overall	N-Fold	Overall	N-Fold	Overall	N-Fold	Overall	N-Fold
SVM-OvA	85.45%	86.79%	63.16%	59.89%	66.67%	88.57%	90.00%	87.50%
SVM-OvO	98.18%	96.82%	63.16%	59.89%	71.43%	74.29%	90.00%	80.00%
K-Nearest								
Neighbour	83.64%	83.22%	84.21%	57.81%	90.48%	62%	100.00%	83%
Naive Bayes	92.73%	90.91%	52.63%	65.89%	71.43%	85.71%	90.00%	92.50%

References

[1]. N.D.K.G. Dharmasiri and S. Vasanthapriyan, "Approach to Heart Diseases Diagnosis and Monitoring through Machine Learning and iOS Mobile Application," Proceedings of the International Conference on Advances in ICT for Emerging Regions ICTers, vol. 18, pp. 407–412, 2018. [2]. Liang, Y. et al, "A new, short-recorded photoplethysmogram dataset for blood pressure monitoring in China," Sci. Data 5:180020 doi: 10.1038/sdata.2018. [3]. Arrozag Ave, Hamdan Fauzan S. Rhandy Adhitya and Hasballah Zakaria, "Earlier detection of cardiovascular diseases with photoplethysmogram (PPG) sensor," Proceedings of the International Conference on Electrical Engineering and Informatics, pp. 676–681, 2015.



Experimental setup

Proposed methods and models were implemented and tested



The data set is divided into four parts, each for identifying a

As the standard practice, 70% of the dataset was used to train the models and 30% was used for testing.

Feature Selection : Recursive Feature Elimination - To reduce the number of parameters used to construct the

Classifiers: Linear SVM-OvA, Linear SVM-OvO, K-Nearest Neighbour, and Naive Bayes

