

## Accuracy Optimization in Speech Pathology Diagnosis With Data Preprocessing Techniques

Joana Fernandes, Diamantino Rui Freitas and João P. Teixeira

Research Centre In Digitalization and Intelligent Robotics (CeDRI), Instituto Politécnico de Bragança, Faculdade de Engenharia da Universidade do Porto (FEUP)

Using acoustic analysis to classify and identify speech disorders non-invasively can reduce waiting times for patients and specialists while also increasing the accuracy of diagnoses. In order to identify models to use in a vocal disease diagnosis system, we want to know which models have higher success rates in distinguishing between healthy and pathological sounds. For this purpose, 708 diseased people spread throughout 19 pathologies, and 194 control people were used. There are nine sound files per subject, three vowels in three tones, for each subject. From each sound file, 13 parameters were extracted (jitta, jitter, Rap, PPQ5, ShdB, Shim, APQ3, APQ5, F0, HNR, autocorrelation, Shan-non entropy and logarithmic entropy). For the classification of healthy/pathological individuals, a variety of classifiers based on Machine Learning models were used, including decision trees, discriminant analyses, logistic regression classifiers, naive Bayes classifiers, support vector machines, classifiers of closely related variables, ensemble classifiers and artificial neural network classifiers. For each patient, 118 parameters were used initially. The first analysis aimed to find the best classifier, thus obtaining an accuracy of 81.3% for the Ensemble Sub-space Discriminant classifier. The second and third analyses aimed to improve ground accuracy using preprocessing methodologies. Therefore, in the second analysis, the PCA technique was used, with an accuracy of 80.2%. The third analysis combined several outlier treatment models with several data normalization models and, in general, accuracy improved, obtaining the best accuracy (82.9%) with the combination of the Greeps model for outliers treatment and the range model for the normalization of data procedure.

**Keywords:** Outliers · Normalization · Speech Pathologies · Machine Learning · Vocal Acoustic Analysis