



Simplifying the Interpretation of Steroid Metabolome Data by a Machine-Learning Approach

Tarik Kirkgoz¹, Semih Kilic², Zehra Yavas Abali¹, Ali Yaman³, Sare Betul Kaygusuz¹, Mehmet Eltan¹, Serap Turan¹, Goncagul Haklar³, Mahmut Samil Sagiroglu⁴, Abdullah Bereket¹, Tulay Guran¹

¹Marmara University, School of Medicine, Department of Pediatric Endocrinology and Diabetes, Istanbul, Turkey

²Politecnico di Milano, Department of Electronics, Information and Bioengineering, Milan, Italy

³Marmara University, School of Medicine, Department of Biochemistry, Istanbul, Turkey

⁴Genpute Computation Technologies, Istanbul, Turkey

P1-03

Background: Liquid chromatography-mass spectrometry (LC-MS) based panels of steroid hormones and their precursors offer a distinct pattern of steroid metabolome for various disorders of adrenal and gonadal steroidogenesis. However, it may not be easy to handle this high throughput data rapidly in clinical setting which requires expert opinion for correct interpretations. Analytical results of steroid panelling can be allied to automated review systems to simplify the complexity of data for disease-related interpretation.

Methods: We have implemented a machine-learning algorithm for a time-saving and experience-independent review and interpretation of analytical results. We have tested the performance of this algorithm using our archived data of quantitation of 16 steroid hormones and precursors by an LC-MS/MS based panel in 500 healthy controls and 427 treatment-naive children with a disorder of adrenal steroidogenesis. This cohort included classic CYP21A2 ($n=75$), non-classic CYP21A2 ($n=19$), CYP11B1 ($n=66$), mutation-positive HSD3B2 ($n=31$), mutation-negative HSD3B2 ($n=21$), CYP11B2 ($n=19$), CYP17A1 ($n=11$), POR ($n=7$) deficiencies and non-CAH PAI ($n=21$). Due to the relatively low numbers of some of the conditions in the patient cohort, the number of samples in one class has outnumbered the other one. This imbalance has been overcome by utilizing data sampling and boosting algorithms, specifically Random Oversampling Boosting (RUSBoost).

Results: Dataset of 415 patients fed to the algorithm with 10-fold cross validation to prevent overfitting. For discrimination of patients from the healthy controls; the sensitivity and specificity of the RUSBoost algorithm was 97.7% and 92.6%, respectively. The differentiation of each disorder could be achieved with overall accuracy of up to 95% independent of age and sex.

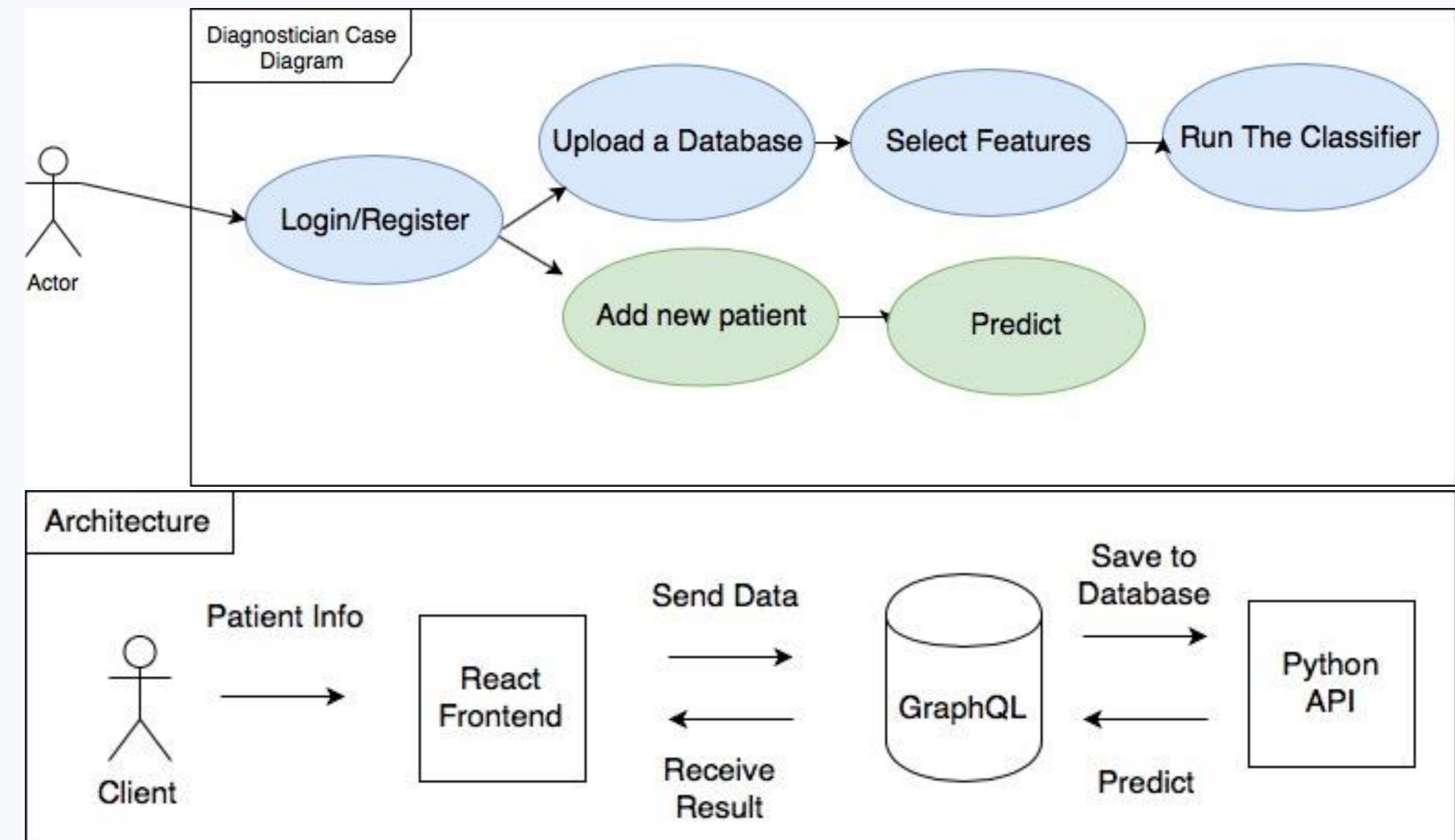


Figure 1) Schematization of programme

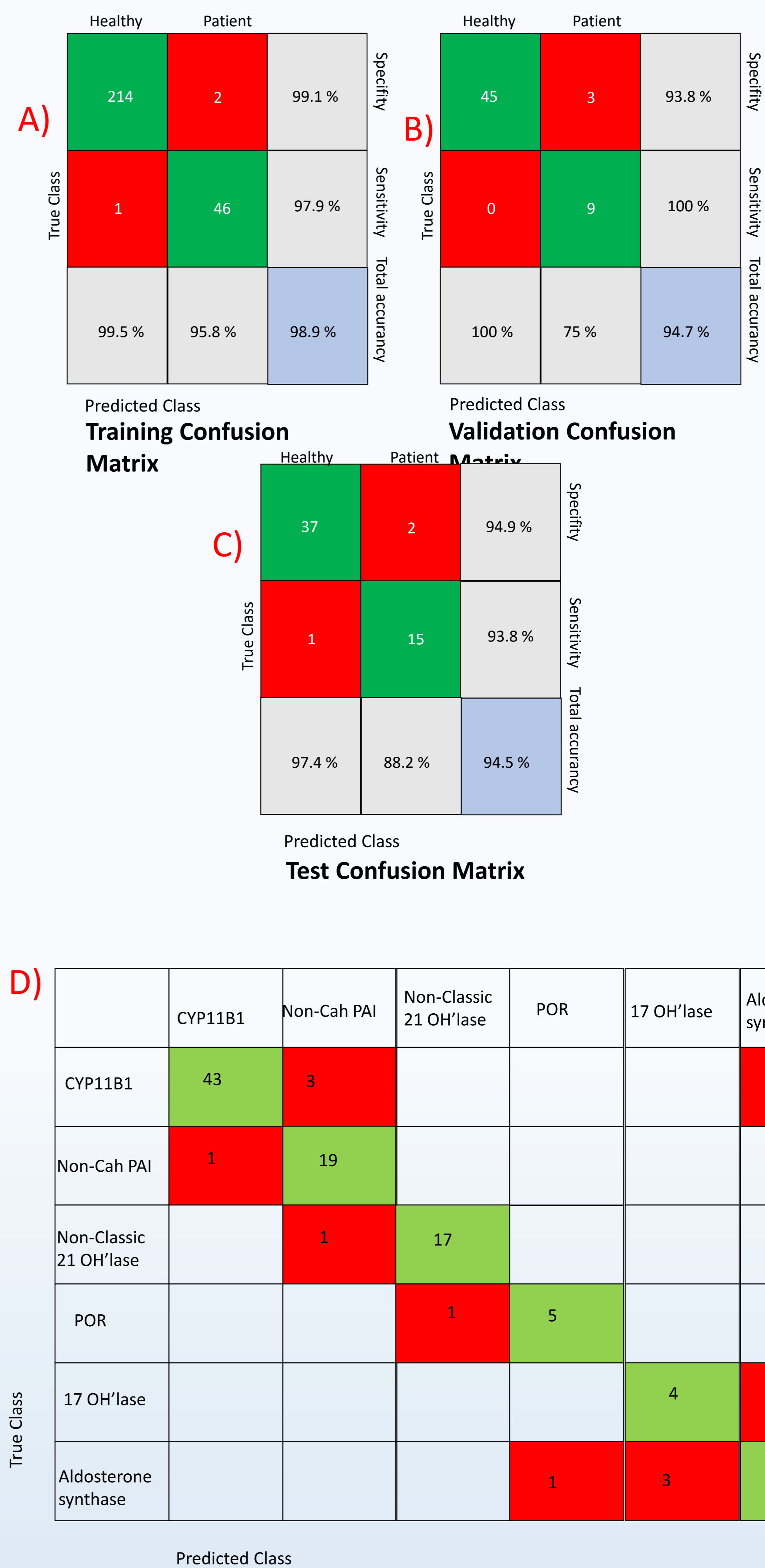


Figure 2) A,B,C) Confusion matrixes and results
D) Subgrouping patients data results

Conclusion :

✓ Application of RUSBoost machine learning algorithm enables a rapid and standardized review of complicated plasma steroid panelling data, which can widely be used by clinicians to make correct diagnosis for disorders of steroidogenesis.