

## **Suggestion of statistical validation on feature importance of machine learning**

Youngro Lee<sup>1</sup>

Jongmo Seo<sup>2</sup>

<sup>1</sup> Electrical and Computer Engineering, Seoul National University

<sup>2</sup> Electrical and Computer Engineering, Seoul National University / Seoul National University Hospital

Feature importance method is used as both primary and subsidiary tool in machine learning analysis for medical datasets in these days. Feature importance methods can be used for selecting biomarkers or markers indicating target disease. If number of markers are not clear, it can give field experts hint for the mechanism of disease. If it corresponds to the field knowledge, it can be the way for validating the model adding to the performance. However, listing up features with its feature importance rank does not provide clear utility for this aspect. Unlike statistical approach, there is no clear definition which value of feature importance can be defined as statistically significant features. Although statistical significance cannot be calculated inside of model algorithm, the value can be estimated by comparing between feature importance values. In this paper, we suggested simple ways to evaluate the statistical significance of feature importance, and select the number of biomarkers. We showed the example by applying the methods for heart failure dataset, which is public open dataset.