

## **Sepsis prediction:**

**Approach 1:** This approach included two steps to analyse sepsis data: features selection, features engineering and then ensembling.

The features engineering schema separated all the covariates in the data into two clusters: first cluster of covariates with the minimum number of missing values with threshold of 10%.

For the features with low missing values:

- Aggregation using Sliding windows of 5 and 11 hours frames while applying different methods:
  - Min, max, mean, median, and variance.
  - Quantile of 95%, 99%, 5%, and 1%
- Other features to capture Long and Short Term Dependencies like Shannon Entropy energy, mean for the first differences and the length of stay for a specific patient.

Morteza et al. H while the second cluster grouped the remaining covariates. The total number of features is around 410 features.

In order to reduce the number of covariates and reduce the bias and variance, a second layer of features selection was applied for two main objectives: select best performing features and five best hyperparameters.

The described procedure above was applied on only 10% of the data. One interesting point about the features selection in this approach is that it was performed using BoostARoot algorithms based on Xgboost. The procedure is that for each feature we create another shadow feature and use only both of them on the training. Then, the feature importance will be presented. If shadow features reported as more important than the original features, this latter will be eliminated.

Once it's decided about which are the best performing parameters and best performing covariates, the training starts using the remaining 90% of the training data. In fact, in this approach we construct five randomly disjoint sets and apply an undersampling technique to balance class in each set separately. At the end, 5 Xgboost models are trained using each set with 5-fold cross-validation. The final output is calculated using the geometric mean of the five outputs of the trained model.

As results, this model training structure has achieved an AUC score of 0.833 and an accuracy of 0.8440.