

Neonatal Sepsis Detection Using Decision Tree Ensemble Methods: Random Forest and XGBoost

Marwan Al-Bardaji and Nahir Danho

Abstract—Neonatal sepsis is a potentially fatal medical condition due to an infection and is attributed to about 200 000 annual deaths globally. With healthcare systems that are facing constant challenges, there exists a potential for introducing machine learning models as a diagnostic tool that can be automatized within existing workflows and would not entail more work for healthcare personnel. The Herlenius Research Team at Karolinska Institutet has collected neonatal sepsis data that has been used for the development of many machine learning models across several papers. However, none have tried to study decision tree ensemble methods. In this paper, random forest and XGBoost models will be developed and evaluated in order to assess their feasibility for clinical practice. The data contained 24 features of vital parameters that are easily collected through a patient monitoring system. The validation and evaluation procedure needed special consideration due to the data being grouped based on patients and being imbalanced. The proposed method has the potential to be generalized to other similar applications. Finally, using the measure receiver-operating-characteristic area-under-curve (ROC AUC), both models achieved around ROC AUC = 0.84. Such results suggest that the random forest and XGBoost models are potentially feasible for clinical practice. Another gained insight was that both models seemed to perform better with simpler models, suggesting that future work could create a more explainable model.

Sammanfattning—Neonatal sepsis är ett potentiellt dödligt medicinskt tillstånd till följd av en infektion och uppges globalt orsaka 200 000 dödsfall årligen. Med sjukvårdssystem som konstant utsätts för utmaningar existerar det en potential för maskininlärningsmodeller som diagnostiska verktyg automatiserade inom existerande arbetsflöden utan att innebära mer arbete för sjukvårdsanställda. Herlenius forskarteam på Karolinska Institutet har samlat ihop neonatal sepsis data som har använts för att utveckla många maskininlärningsmodeller över flera studier. Emellertid har ingen prövat att undersöka beslutsträds ensemble metoder. Syftet med denna studie är att utveckla och utvärdera random forest och XGBoost modeller för att bedöma deras möjligheter i klinisk praxis. Datan innehåller 24 attribut av vitalparameterar som enkelt samlas in genom patientövervakningssystem. Förfarandet för validering och utvärdering krävde särskild hänsyn med tanke på att datan var grupperad på patientnivå och var obalanserad. Den föreslagna metoden har potential att generaliseras till andra liknande tillämpningar. Slutligen, genom att använda receiver-operating-characteristic area-under-curve (ROC AUC) måttet kunde vi uppvisa att båda modellerna presterade med ett resultat på ROC AUC = 0.84. Sådana resultat föreslår att både random forest och XGBoost modellerna kan potentiellt användas i klinisk praxis. En annan insikt var att båda modellerna verkade prestera bättre med enklare modeller vilket föreslår att framtida arbete skulle kunna vara att skapa en mer förklarlig maskininlärningsmodell.

Index Terms—Machine Learning, Sepsis, Neonatal Sepsis, Random Forest, XGBoost, Imbalanced Data, Binary Classification, Cross-Validation, Hyperparameter Tuning.

Supervisors: Antoine Honoré

TRITA number: TRITA-EECS-EX-2022:175

I. INTRODUCTION

The goal of the healthcare system is the maintenance and improvement of the health of a population. An essential step in that work is the diagnosis and detection of diseases, where physicians typically work using a combination of their own experience and medical guidelines using years of research [1]. This project focuses on neonatal sepsis, a possibly fatal medical condition due to an infection [2]. Globally sepsis is attributed to about 200 000 annual deaths [3]. There exist several scoring systems for the detection and prognosis estimation of sepsis built on international consensus; however, none is perfect [2], [4]. Meanwhile, healthcare systems are constantly facing issues such as rising costs [5], staffing shortages [6], and aging populations with ever-increasing healthcare needs. This could lead to situations where guidelines cannot be followed perfectly, thus endangering patient safety [1]. Any method that could increase diagnostic performance without inferring increased effort by healthcare personnel would be desirable. Machine learning models have been on the rise since the 20th century, and adoption is continuously increasing with more powerful computers and the improved ability to collect large datasets. There exist many machine learning algorithms whose goal is to classify different data and could therefore have the potential as a diagnostic tool in healthcare that can be automatized within existing workflows [7].

A. Problem Formulation

This project will use neonatal sepsis data provided by the Herlenius Research team at Karolinska Institutet. Earlier work for the detection of sepsis using machine learning models has been conducted using the same data. Examples of previous models that have been used are Markov models, logistic regression, naïve Bayes, multi-layer perceptrons, Gaussian mixture models, and normalizing flows [8], [9]. Many of the models achieved promising results. Nevertheless, no work has studied the performance of models that use decision tree ensemble methods.

B. Project Goal and Scope

The goal of the project is to study the feasibility of using decision tree ensemble methods to detect neonatal sepsis. Specifically, the performance of random forest and XGBoost

models will be evaluated and compared. Moreover, insightful learnings from the development of the models will be collected and discussed. The project is limited to a proof of concept and does not create a model ready for real-life implementation.

II. BACKGROUND

A. Machine Learning

Machine Learning (ML) is a subset of artificial intelligence that deals with learning in the sense that the algorithm's performance on future tasks can be improved by making observations of the world [7]. Machine learning algorithms are generally divided into three categories (1) supervised learning, (2) unsupervised learning, and (3) reinforcement learning. Supervised learning is predictive, where the algorithm maps an input to an output. Hidden training data is provided for which the mapping is known beforehand, also known as labeled data. A trained model can later be used to predict the output of future inputs with unknown outputs. Unsupervised learning analyzes unstructured data and tries to find a pattern on its own; here, there exists only input and no output. On the other hand, reinforcement learning is an algorithm that reacts to the environment rewarding desired behaviors and punishing undesired ones [7]. There exists a myriad of algorithms under each category; moreover, there also exist algorithms that do not fall under one category, some of which are semi-supervised learning. This project will deal with a classification problem within supervised learning.

B. Medical Background

Sepsis is a possibly fatal medical condition which, in layman's terms, often somewhat incorrectly is called blood poisoning, thus not revealing the entire truth. According to international consensus [?], sepsis should be defined as "life-threatening organ dysfunction caused by a dysregulated host response to infection. Left untreated, sepsis can turn into sepsis shock. sepsis shock is defined as [?] "a subset of sepsis in which underlying circulatory and cellular metabolism abnormalities are profound enough to increase mortality substantially." Many survivors develop permanent neurologic impairment [10]. Sepsis does not have obvious symptoms, especially in the early course of the condition. To identify sepsis in clinical practice, international consensus recommends using the Sequential Organ Failure Assessment (SOFA) score or the quick SOFA (qSOFA) score [4], [11], [12]. These scoring systems use a combination of vital parameters, such as partial pressure of oxygen in the blood and blood pressure; blood tests, such as platelet count and bilirubin; and neurological status [11]. Moreover, after a sepsis suspicion has arisen, there exist international guidelines created by the "Surviving Sepsis Campaign" for the most appropriate tests and treatments to continue with [12]. Suspected sepsis can be confirmed through microbiologic blood cultures; the culture needs to be obtained before any antibiotic treatments. In principle, there exist two types of treatment (i) antimicrobial treatment, which usually starts with empiric broad-spectrum antibiotics, and (ii) organ-supportive such as fluid therapy and vasoactive medication [12]. An inherent limitation of using the SOFA scoring system

is that doctors and nurses conduct additional tasks such as taking blood tests that do not necessarily need to happen for all patients [11].

This project will study neonatal sepsis, meaning sepsis in the first weeks of an infant's life. In neonates, sepsis is difficult to diagnose clinically since they may be asymptomatic until organ dysfunction is prominent [13]. There exists an adapted version of the SOFA score called the nSOFA score; however, this score still requires invasive testing such as blood tests which both are costly and take time compared to collecting vital parameters. Moreover, the nSOFA scoring is not widely adopted [14]. A study by Fairchild suggests that there exists potential with using heart rate variability, heart rate characteristics, and other vital signs in the detection of Neonatal Sepsis. These measurements are non-invasive and can easily be collected through a monitoring machine [10]. Current challenges in the treatment of neonatal sepsis include late detection, overuse of antibiotics, and difficulties with invasive testing [10].

C. Classification Problems

A classification problem within supervised machine learning refers to a predictive modeling problem where the aim is to predict an output given an input. The following definition of a classification problem will assume a one-dimensional output, also called the label. Each data input will be provided as a feature vector \vec{x}_i with its corresponding output y_i . Since the input is in vector form, it may include multiple data points corresponding to several features.

To train the machine learning model, a dataset \mathcal{D} with n examples and m features is provided.

$$\mathcal{D} = \{(\vec{x}_i, y_i)\} \quad (|\mathcal{D}| = n, \vec{x}_i \in \mathbb{R}, y_i \in \mathbb{R}) \quad (1)$$

(Binary classification is the case when $y_i \in \mathbb{B} = \{0, 1\}$). A prediction function ϕ gives the prediction.

$$\hat{y}_i = \phi(\vec{x}_i) \quad (2)$$

The behavior of ϕ depends on internal model parameters that are unique for each machine learning implementation. The procedure of finding the internal parameters is what training a model means [7]. The internal parameters of a model are fit by minimizing the objective function \mathcal{L} which typically is of the following form.

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \Omega(\theta) \quad (3)$$

Here l is a differentiable convex loss function that measures the difference between y_i and \hat{y}_i . The second term Ω is a function of the internal parameters which penalize complex models. There exist several loss and penalty functions depending on the implementation that shares the concept of distance and model complexity, respectively. [7]

If an unseen data point \vec{x}^{unseen} is provided without a label then there exists no way to know whether the prediction $\hat{y}^{\text{unseen}} = \phi(\vec{x}^{\text{unseen}})$ is correct. However, if a theoretical label y^{real} is allowed to exist then the goal of any classification algorithm is that $\hat{y} = y^{\text{real}}$. To imitate this procedure, the

original dataset \mathcal{D} is commonly divided into a test set and a training set.

The training set is used for fitting the internal parameters of the selected model, while the test set is used to evaluate the performance of the model. Sometimes the original dataset is divided into yet another set called the validation set to optimize external parameters, also called hyperparameters, on unseen data before evaluating the results on the test set. External parameters affect a model's behavior and do not change depending on the training data [7].

D. Decision Trees

A decision tree is a representation of a function that maps a feature vector to a single output value which is called the decision [7]. In essence, it is similar to a flowchart of questions that are commonly used within the healthcare system [15]. The decision tree starts with a root node. Each node may, in turn, split into several other nodes. A node that does not split into other nodes is called a leaf and contains the final decision. At each node, a question regarding the question is asked. In implementations using numerical data the question are comparisons of the types $<$, \leq , $=$, \geq , $>$.

The prediction functions ϕ works by using an algorithm that finds the feature and the question to ask that provide the highest "importance." Importance is measured using information gain, which is defined in terms of entropy. These quantities are fundamental in information theory. To train the decision tree, all the feature vectors in the training data set to go through the entire decision tree, and the number of samples and their class are calculated at the leaf nodes [7], which means that each node sorts the incoming data into smaller sets.

Entropy is a measure of uncertainty of a random variable; the more information, the less entropy. In general, the entropy of a random variable V with values v_k having the probability $P(v_k)$ is defined as [7].

$$\text{Entropy} = H(V) = - \sum_k P(v_k) \log_2 P(v_k) \quad (4)$$

Entropy is measured in bits and corresponds to the expected number of 50/50 guesses it would require to narrow down to a specific value v_k . For example, a fair coin flip has an entropy of $H(\text{Fair coinflip}) = -2 \cdot (0.5 \log_2 0.5) = 1$ and a fair six sided die has an entropy of $H(\text{Fair six-sided die}) = -6 \cdot (1/6 \log 1/6) \approx 2.6$ [7].

The information gain at a node is calculated as the expected reduction in entropy by sorting the incoming data. The information gain on attribute A and data S is defined as [7].

$$\text{Information gain} = H(S) - \sum_{i \in v} \frac{|S_i|}{|S|} H(S_v) \quad (5)$$

where S is the incoming data to the node, v is a set of mutually exclusive questions, and S_i is the sorted version of S after asking a question i . S is considered to be a random variable with values and probabilities according to the distribution in the data [7].

The maximization of the information is thus the objective function of the decision tree. However, in order to limit the

complexity of the model and reduce bias, it is possible to penalize complex models with a regularization parameter Ω [7].

E. Random Forests

Random forests are an ensemble of decision trees that uses bootstrap aggregating to reduce variance in a noisy dataset by training multiple different trees. A majority vote of all trees then decides the output [16].

A random forest classifier R is defined as

$$R = \text{Majority vote of } \{h(\vec{x}, \Theta_k), k = 1, \dots\} \quad (6)$$

Where h is a decision tree, and Θ_k is an independent random vector that the decision tree k uses in the construction of the tree [16].

The primary source of randomness in a random forest is feature subsampling, where a random number of features are selected for each tree. This reduces the bias by increasing the probability that the trees are uncorrelated. The objective function and penalization of complex trees are analogous to the decision tree [16].

F. XGBoost

XGBoost, an implementation of gradient boosted decision trees, has shown state-of-the-art results in many machine learning challenges. XGBoost is an ensemble method of regression trees that uses boosting, which aims to improve performance by creating a strong classifier from many weak classifiers. Regression trees use the same concepts as a decision tree but have a continuous target variable. Moreover, XGBoost is designed with system performance in mind and is easily scalable [17].

1. Objective Function

For a given dataset with n examples and m features

$$\mathcal{D} = \{(\vec{x}_i, y_i)\} \quad (|\mathcal{D}| = n, \vec{x}_i \in \mathbb{R}, y_i \in \mathbb{R}) \quad (7)$$

A tree ensemble method uses K additive functions to predict the output according to

$$\hat{y} = \phi(\vec{x}_i) = \sum_{k=1}^K f_k(\vec{x}_i), \quad f_k \in \mathcal{F} \quad (8)$$

where $\mathcal{F} = \{f(\vec{x}) = w_{q(\vec{x})}\} \quad (q : \mathbb{R}^m \rightarrow T, w \in \mathbb{R}^T)$ is the space of regression trees. T is the number of leaves in each tree and q represents the structure of each tree that maps an input to the leaves. Each tree structure q and leaf weights w correspond to a specific f_k .

To learn the weights w used in the model, the following regularized objective is minimized [17].

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k) \quad (9)$$

$$\text{where } \Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2 \quad (10)$$

Here l is a differentiable convex loss function, \hat{y}_i is the prediction and y_i is the target. The second term penalizes

a complex model regarding the size of the weights and the number of leaves [17].

2. Gradient Tree Boosting

Since the equation

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k) \quad (11)$$

contains functions as parameters it cannot be optimized using traditional optimization methods in Euclidean space and has to train in an additive manner. The prediction of instance i at iteration t is defined as $\hat{y}_i^{(t)}$. Then $\hat{y}_i^{(t)}$ is calculated by adding $f_t(\vec{x}_i)$ to $\hat{y}_i^{(t-1)}$ [17]. Therefore the objective to minimize turns into

$$\mathcal{L}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(\vec{x}_i)) + \Omega(f_t) \quad (12)$$

The greedy approach is to add the f_t that improves the model the most. By using a second-order approximation of the objective, it turns into

$$\mathcal{L}^{(t)} \simeq \sum_{i=1}^n [l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(\vec{x}_i) + \frac{1}{2} h_i f_t^2(\vec{x}_i)] + \Omega(f_t) \quad (13)$$

$$\text{where } g_i = \frac{\partial l(\hat{y}_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} \text{ and } h_i = \frac{\partial^2 l(\hat{y}_i, \hat{y}_i^{(t-1)})}{\partial (\hat{y}_i^{(t-1)})^2} \quad (14)$$

A simplified objective function $\tilde{\mathcal{L}}^{(t)}$ at step t is obtained by removing the constant terms

$$\tilde{\mathcal{L}}^{(t)} = \sum_{i=1}^n [g_i f_t(\vec{x}_i) + \frac{1}{2} h_i f_t^2(\vec{x}_i)] + \Omega(f_t) \quad (15)$$

Define $I_j = \{i | q(\vec{x}_i) = j\}$ as the instance set of leaf j meaning all the instances that correspond to the leaf. The equation can be rewritten by expanding Ω

$$\begin{aligned} \tilde{\mathcal{L}}^{(t)} &= \sum_{i=1}^n [g_i f_t(\vec{x}_i) + \frac{1}{2} h_i f_t^2(\vec{x}_i)] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\ &= \sum_{j=1}^T \left[\left(\sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma T \end{aligned} \quad (16)$$

For a fixed structure $q(\vec{x})$ the optimal weight w_j^* of leaf j is calculated by

$$w_j^* = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda} \quad (17)$$

The corresponding optimal value is given by

$$\tilde{\mathcal{L}}^{(t)} = - \frac{1}{2} \sum_{j=1}^T \frac{(\sum_{i \in I_j} g_i)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma T \quad (18)$$

Equation 15 can be used as a scoring function to measure the quality of the tree structure q . Usually, it is impossible to enumerate all the possible tree structures q . A greedy algorithm is developed by starting from a single leaf and iteratively adding. Assume that I_R and I_L are instance

sets of the right and left nodes after a split. Letting $I = I_R \cup I_L$ then the loss reduction after the split is given by [17]

$$\mathcal{L}_{\text{split}} = \frac{1}{2} \left[\frac{(\sum_{i \in I_R} g_i)^2}{\sum_{i \in I_R} h_i + \lambda} + \frac{(\sum_{i \in I_L} g_i)^2}{\sum_{i \in I_L} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} - \gamma \right] \quad (19)$$

Equation 16 is the formula that is used in practice in XGBoost [17].

3. Split Finding Algorithms Using equations 14 and 15, two algorithms for split finding can be written according to the following pseudocode [17].

Algorithm 1 Exact Greedy Algorithm for Split Finding

Input: I , instance set of current node

Input: d feature dimension

$gain \leftarrow 0$

$G \leftarrow \sum_{i \in I} g_i$, $H \leftarrow \sum_{i \in I} h_i$

for $k=1$ **to** m **do**

$G_L \leftarrow 0$, $H_L \leftarrow 0$

for j **om** sorted(I , by \vec{x}_{jk}) **do**

$G_L \leftarrow G_L + g_j$ $H_L \leftarrow H_L + h_j$

$G_R \leftarrow G - G_L$, $H_R \leftarrow H - H_L$

end for

end for

Output: Split with maximum score

Algorithm 2 Approximate Algorithm for Split Finding

for $k=1$ **to** m **do**

Propose $S_k = \{s_{k1}, s_{k2}, s_{k3}, \dots, s_{kl}\}$ by percentiles on feature k

Proposal can be done per tree (global), or per split(local).

end for

for $k=1$ **to** m **do**

$G_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \vec{x}_{jk} > s_{k,v-1}\}} g_j$

$H_{kv} \leftarrow \sum_{j \in \{j | s_{k,v} \geq \vec{x}_{jk} > s_{k,v-1}\}} h_j$

end for

Follow same step as in previous section to find maximum score only among proposed splits.

III. METHODS

A. Selection of Machine Learning Models

A binary classifier was to be created that could distinguish sepsis using patient data. The random forest and XGBoost models were selected to study their potential for this endeavor using the implementations of the sci-kit learn and XGBoost libraries in python, respectively [18], [19].

B. Study population

Data was acquired from the Herlenius Research team at Karolinska Institutet that are currently conducting research regarding neonatal healthcare. The population consisted of

very low birth weight infants (< 1500 g) hospitalized in the Neonatal Intensive Care Unit (NICU) at Karolinska University Hospital, Stockholm, Sweden [9].

C. Data Description

Time-series data were collected from all patients was collected from a high-frequency Phillips IntelliVue MX800 Patient Monitor [9]. The data initially contained monitor data sampled at 1 Hz. The times-data was split into windows of 70 minutes where 19 different features were extracted. If a patient received a sepsis diagnosis at a specific time instance, all the windows 24 hours before the time of the diagnosis was categorized as sepsis-like. Sepsis-like windows were set to the value 1, and non-sepsis-like time windows were set to the value 0, and these values corresponded to the target y_i in our binary classification problem. Moreover, five more features were also provided that included parameters regularly updated in the medical history. In total, the data consisted of 24 features. Table I provides a list of the target label and all the features that were provided in the data.

Data is provided in a tabular manner where each data row corresponds to a specific time window. The data initially contained 134668 data rows from 118 different patients. A total of 10 patients experienced sepsis during their hospitalization together, totaling 556 rows included sepsis-like characteristics, thus having a target_y = 1.

It can be noticed that many features seem to contain large negative values down to -99999 ; also there appears to exist missing data for the `feats_cirk_vikt` feature. A more detailed description of the data, including the median and percentiles, can be found in the appendix of the project.

D. Preprocessing

The data in its unmodified state is not universally usable due to several reasons. The following issues have been identified and need to be dealt with. When discussing the following naming convention is used in this project: the feature vector corresponds to a row, and one specific value in the feature vector is called a data point.

- **Erroneous data points** - Some data points have been identified where the value of some features are exactly -99999 or -9999 and are far outliers relative to the rest of the data. The reason for this data is related to errors during the patient monitoring data collection.
- **Categorical data** - The unmodified data includes the `feats_group_uid` feature, which is categorical. Many machine learning algorithms do not work with categorical data [7].
- **Large variations in the data** - There exist large variations in the data which may affect the weights of the internal parameters in distance-based models [7].
- **Missing data** - One of the features has missing data that needs to be dealt with. The unmodified data contains 6595 missing data points in only the `feats_cirk_vikt` feature. If all the rows containing a missing data point were removed, it would result in an additional 158280 data points lost. This is equal to almost 5% of all rows.

With a background in solving the aforementioned issues, the following data preprocessing steps are taken. These steps aim to prepare the data for the random forest and XGBoost algorithms but also prepare for other eventual algorithms that would want to be tested in the future, including distance-based machine learning algorithms.

1. **Label Encoding** - The 'group_uid' for the entire dataset is encoded using a Label Encoder, creating a 1-to-1 mapping from the String domain to the integer domain. This is to facilitate future slicing and grouping of the dataset.
2. **Removal of Erroneous data points** - Some of the data points from features that initially were collected from the patient monitor was exactly equal to -99999 and -9999 . Since the probability of achieving an exact integer in any continuous distribution is zero, there exists no risk of real data being removed.
3. **Scaling** - This part is not necessary for random forests and XGBoost. However, regarding distance-based algorithms, the internal parameters are affected by the values of the datapoint. Having a large variance in the sizes of the data points makes it more difficult to fit appropriate internal parameters [7]. All the data except 'target_y' and 'group_uid' is scaled using a standard scaling that removes each data point's mean within a feature and scales to unit variance.
4. **Missing Data** - The remaining missing data points are assigned values based on the neighboring data points. More specifically, a KNN imputer is used, which stands for K nearest neighbor imputer. Initial testing resulted in the selection of $K = 5$ due to not significantly changing the mean and the standard deviation of the `feats_cirk_vikt` feature. In other words, the algorithm looks at the five nearest neighbors. The specifics of how the algorithm works are explored in further detail in the appendix to the project.

Table I contains a description of the data after removing erroneous data points. A more detailed description describes the data after each step also, including the median and percentiles, can be found in the appendix to the project.

E. Cross-Validation for Highly Imbalanced Grouped Time Series Data

Data is needed both to train the internal parameters and to test the performance of the model data. Learning the internal parameters of a prediction function and testing on the same data results in a methodological mistake. The model could simply repeat the labels that it has seen without learning any patterns [7]. A trivial method to carry out such a split would be to randomly select a specific percentage in the training set and the rest to be in the testing set. Typical values are to select about 70% to 80% in the training set [20]. However, when trying different hyperparameters of the models, there may exist situations where part model finds patterns in the training set that does not exist in the testing set. This situation is called overfitting, where the patterns found do not display the entire truth of the data. Moreover, indirect knowledge about

TABLE I
FEATURES THAT WERE PROVIDED AND A DESCRIPTION OF THE DATA AFTER ERRONEOUS DATA POINTS HAVE BEEN REMOVED

Feature Name	Description	Count	Mean	Std	Min	Max
feats_btb_mean	The mean value of the beat to beat interval over a time window.	95849	0.000203	0.00104	-0.0284	0.0239
feats_rf_mean	The mean value of the respiratory frequency over a time window.	95849	50.2	11.5	8.01	97.0
feats_spo2_mean	The mean value of the oxygen saturation over a time window.	95849	93.5	3.47	35.6	100.0
feats_btb_std	The standard deviation value of the beat to beat interval over a time window.	95849	0.0201	0.0144	0.000556	0.241
feats_rf_std	The mean value of the respiratory frequency over a time window.	95849	13.5	4.30	0.0	38.1
feats_spo2_std	The mean value of the oxygen saturation over a time window.	95849	4.31	2.73	0.0	27.8
feats_btb_maximum	The maximum value of the beat to beat interval over a time window.	95849	0.146	0.149	0.0016	1.64
feats_rf_maximum	The maximum value of the respiratory frequency over a time window.	95849	89.7	16.52	15.3	163
feats_spo2_maximum	The maximum value of the oxygen saturation over a time window.	95849	99.7	0.911	79.9	100.0
feats_btb_minimum	The minimum value of the beat to beat interval over a time window.	95849	-0.0430	0.0265	-0.808	-0.00140
feats_rf_minimum	The minimum value of the respiratory frequency over a time window.	95849	20.1	6.85	0.394	55.3
feats_spo2_minimum	The minimum value of the oxygen saturation over a time window.	95849	74.5	14.7	0.0333	100.0
feats_btb_skew	The sample skewness of the beat to beat interval over a time window.	95849	2.18	2.93	-5.24	14.8
feats_rf_skew	The sample skewness of the respiratory frequency over a time window.	95849	0.151	0.626	-17.8	12.7
feats_spo2_skew	The sample skewness of the oxygen saturation over a time window.	95849	-1.34	1.19	-23.2	2.51
feats_btb_kurtosis	The kurtosis of the beat to beat interval over a time window. A measure of the "tailedness" of the sampling points.	95849	20.4	31.8	-1.33	240
feats_rf_kurtosis	The kurtosis of the respiratory frequency over a time window. A measure of the "tailedness" of the sampling points.	95849	0.0564	4.33	-3.0	341
feats_spo2_kurtosis	The kurtosis of the oxygen saturation over a time window. A measure of the "tailedness" of the sampling points.	95849	4.01	10.3	-3.0	592
feats_btb_sampAs	The sample asymmetry of the beat to beat interval over a time window. Assessing the asymmetry of the sampling points [10].	95849	3.64	5.72	0.0652	276
feats_btb_sampEn	The sample entropy of the beat to beat interval over a time window. Assessing the complexity and thus possible deterioration [10].	95849	0.410	0.150	0.0124	1.01
feats_cirk_vikt	The weight of the patient during the time window. (Not sampled through the monitor)	91550	1.28	0.470	0.497	3.72
feats_bw	The birth weight of the patient. (Not sampled through the monitor)	95849	839	258	400	150
feats_sex	Categorical data with values of biological sex at birth. Either 1 or 2 for males and females. (Not sampled through the monitor)	95849	1.56	0.496	1.0	2.0
feats_pnage_days	Patient age in days since birth. Negative days are possible due to inconsistencies with registration of exact birth time. (Not sampled through the monitor)	95849	31.3	22.2	0.0250	13
feats_group_uid	An anonymized personalized string for each patient. (Not sampled through the monitor) ⁴					
target_y	1 for sepsis-like behavior over a time window, 0 for non-sepsis like.					

the testing set through iterative evaluation of the model may "leak" into the choice of hyperparameters, and in such an instance, the testing set does not show an unbiased picture of the model's performance. A possible option to overcome this problem is to partition the original dataset into three parts where a validation set is added. This data is not touched until the final valuation of the model after the hyperparameters have been found. Nevertheless, this partitioning drastically reduced the number of samples available for training the model [21]. Since the available data was heavily imbalanced, with a small proportion of the data having positive labels, a partition of the dataset into three parts would result in partitions with only about 100 to 200 sepsis-like rows.

A solution to reduce overfitting without using a validation set is to use cross-validation. The principle idea is that the model is trained several times for a specific hyperparameter choice with different choices of test sets each time. Every instance of model training is called a split. If the model would overfit one of the splits, then an eventual loss of generalization due to overfitting would affect the performance of the other splits. There exist several types of methods of cross-validation, and the choice of the cross validator depends on the data. The data used in this report had the following characteristics that need to be taken into account. (i) The data is heavily imbalanced, with less than one percent of the data points having a positive label. If a split is selected entirely at random, there exists a probability that no positive labels are included in either the training or the testing set of that split. (ii) The data included time-series data grouped on a patient basis. It is salient for any cross-validation to take into not having data rows from the same patient in both the training set and test set of the split. This would result in a methodological error since the model possibly would train on past data and test on future data from the same patient finding patient-specific patterns instead of general patterns.

The following cross-validation iterators were considered from which one was selected to be used.

- k -fold - This cross-validation iterator randomly divides the data into k groups of samples where all the samples. In total there are k splits where one of the groups is selected as the testing set for each split. In total the entire dataset will have been used in testing across all the splits. However, this cross-validation iterator does not solve either issue (i) or (ii) [21].

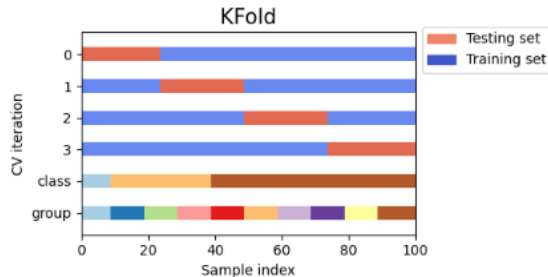


Fig. 1. How the data can be split in a GroupKFold. Source: [21]

- Group k -fold - This cross-validation iterator works in the same way as k -fold; however it ensures that that the same group is not represented in both the training and testing set. The division of groups is made possible by providing an id with each data row. In the data used in this report, the id corresponds therefore to the `feats_group_uid`, meaning each patient. Therefore this cross-validation iterator takes into account issue (ii) but not (i). Due to an imbalance of the data the sizes of the training and testing sets in each split are not necessarily the same [21].

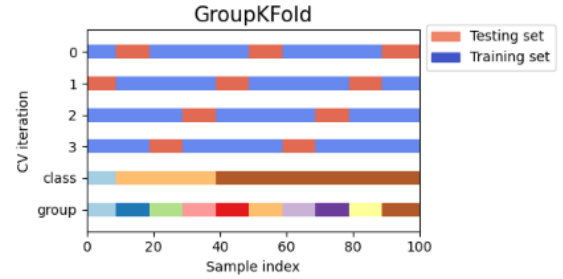


Fig. 2. How the data can be split in a GroupKFold. Source: [21]

- Stratified k -fold - This cross-validation iterator is a variation of the k -fold where each set in the splits contains approximately the same percentage of samples of each target class as the complete set. This means that the iterator tries to preserve the ratios of the classes in both the training and the test set. This cross-validation iterator solves the issue (i) but not (ii). Due to imbalances of the data, the ratios are not necessarily the same but given that it is possible to divide the data a split will not result in a set missing any positive labels [21].

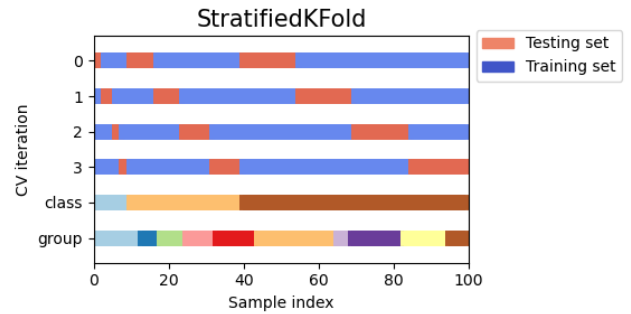


Fig. 3. How the data can be split in a StratifiedKFold. Source: [21]

- Stratified group k -fold - This cross-validation iterator combines both the group k -fold and stratified k -fold cross-validation iterators. Thus, this iterator preserves both the ratios of classes in each split and keeps each group within either the test or training split within a single split. This cross-validation iterator solves both issues (i) and (ii).

Finally, the best cross-validation iterator that caters to the data that does not provide methodological errors is the Stratified

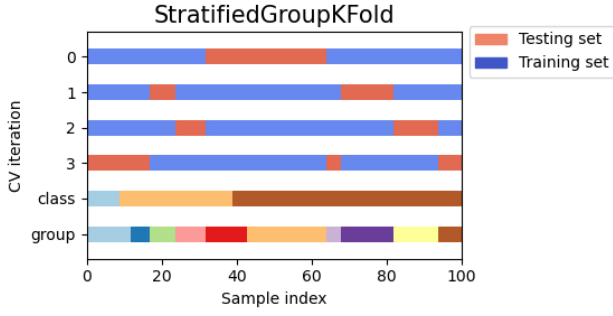


Fig. 4. How the data can be split in a StratifiedGroupKFold. Source: [21]

Group k -fold cross-validation iterator. Henceforth, when referring to a cross-validation iterator it is assumed that stratified group k -fold is used.

F. Model Evaluation for Highly Imbalanced Data

The four possible outcomes when the model makes a prediction are the following:

- 1) True Positive (TP): True positives occur when the model correctly predicts that a patient has sepsis.
- 2) True Negative (TN): True negatives occur when the model correctly predicts that a patient does not have sepsis.
- 3) False Positives (FP): False positives occur when the model incorrectly predicts that a patient has sepsis.
- 4) False Negatives (FN): False Negatives occur when the model incorrectly predicts that a patient does not have sepsis.

With these four outcomes, they can be used to determine the overall quality of the model. Typically, the following metrics are used [22].

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}} \quad (20)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (21)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (22)$$

$$\text{Specificity} = \frac{\text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}} \quad (23)$$

$$\text{True Positive Rate} = \text{TPR} = \frac{\text{TP}}{\text{P}} \quad (24)$$

$$\text{False Positive Rate} = \text{FPR} = \frac{\text{FP}}{\text{P}} \quad (25)$$

However, due to the data being highly imbalanced, the accuracy would be high just by creating a model that predicts all inputs as non-sepsis-like. Given the distribution of classes in the provided data, such a model would have an accuracy of about 99.5% even if $\text{TP} = 0$. Moreover, these metrics are threshold-dependent, which is not of importance in the process of discriminating between sepsis-like and non-sepsis-like time windows. A better approach is thus to look at the receiver operating characteristic, which is received by plotting the true positive rate against the false positive rate. In other words,

the ROC curve shows the true positive rate for the classifier given an accepted value of the false positive rate. A random classifier follows a linear slope in which the true positive rate and the true negative rate are always equal. The area under the ROC, also called the ROC AUC, can be calculated to be used as a scoring metric. A ROC AUC of 0.5 indicates that there is no discrimination between classes, and 1 indicates perfect discrimination [22]. Henceforth, the ROC AUC will be used as the scoring metric in all machine learning models.

G. Hyperparameter Tuning

1) *Workflow*: An important part of training a machine learning model is to find the parameters that generate the best performance according to an evaluation method, which in this instance is the best ROC AUC. The workflow of finding the best hyperparameters consists of first defining a hyperparameter space meaning all potential values of hyperparameters that aims to be tested. Combinations from the hyperparameter space are tested together with a cross-validation iterator. A brute force solution would consist of testing all the possible combinations of parameters from the parameters space to select the best performing combination. After having found the retrained model, final cross-validation over the entire dataset is tested to find a final evaluation of the performance of the model.

The specific hyperparameters that are available to work within the parameter space are the machine learning model and the specific implementation in the library.

2) *Random Forest Hyperparameters*: The following hyperparameters for the random forest model are tuned according to their implementation in the sci-kit-learn library [18].

- *n estimators* - The number of trees in the forest. The default value is 100. Typical values to consider are integers ranging from 10 to 300 [18].
- *maximum depth* - The maximum depth of the tree. The default value is unlimited. Typical values to consider are integers ranging from 1 to 30 [18].
- *maximum features* - The number of features to consider when looking for the best. The default value is the same as the number of features which correspond to 24 for the provided dataset. Typical values to consider are integers ranging from 1 to the number of features = 24 [18].
- *minimum samples split* - The minimum number of samples required to split an internal node. The default value is 2. Typical values to consider are integers ranging from 2 to 30 [18].
- *minimum samples leaf* - The minimum number of samples required to be at a leaf node. The default value is 1. Typical values to consider are integers ranging from 1 to 30 [18].

3) *XGBoost Hyperparameters*: The following hyperparameters for the random forest model are tuned according to their implementation in the XGBoost library [19].

- *n estimators* - The number of trees [19]. Increasing this will make the model more complex and more likely to overfit. The default value is 100. Typical values to

consider in the hyperparameter space are integers ranging from 10 to 500 [19].

- **maximum depth** - Maximum depth of a tree. Increasing this value will make the model more complex and more likely to overfit. The default value is 6. Typical values to consider in the hyperparameter space are integers ranging from 1 to 10 [19].
- **learning rate** - Step size shrinkage used in the update to prevent overfitting. After each boosting step, weights of new features can be obtained directly, and the learning rate shrinks the feature weights to make the boosting process more conservative. The default value is 0.3. Typical values to consider in the hyperparameter space are values ranging from 0.01 to 1 in 0.01 step [19].
- **minimum child weight** - Minimum sum of instance weight (hessian) needed in a child. If the tree partition step results in a leaf node with the sum of instance weight less than minimum child weight, then the building process will give up further partitioning. The default value is 1. Typical values to consider in the hyperparameter space are values ranging from 0.5 to 1.5 in 0.01 steps [19].
- **gamma** - Minimum loss reduction required to make a further partition on a leaf node of the tree. The larger gamma is, the more conservative the algorithm is. The default value is 0. Typical values to consider in the hyperparameter space are values ranging from 0 to 0.5 in 0.01 steps [19].
- **column sample by tree** - The subsample ratio of columns when constructing each tree. Subsampling occurs once for every tree constructed. The default value is 1. Typical values to consider in the hyperparameter space are values ranging from 0.5 to 1 in 0.01 steps [19].
- **subsample** - subsample ratio of the training instances. Setting it to 0.5 means that XGBoost would randomly sample half of the training data before growing trees and this will prevent overfitting. Subsampling will occur once in every boosting iteration. The default value is 1. Typical values to consider in the hyperparameter space are values ranging from 0.5 to 1 in 0.01 steps [19].
- **lambda** - L2 regularization term on weights. Increasing this value will make the model more conservative, punishing complex models. The default value is 1. Typical values to consider in the hyperparameter space are values ranging from 0.5 to 1.5 in 0.01 steps [19].
- **alpha** - L1 regularization term on weights. Increasing this value will make the model more conservative, punishing complex models. The default value is 0. Typical values to consider in the hyperparameter space are values ranging from 0 to 0.5 in 0.01 steps [19].
- **scale positive weight** - Control the balance of positive and negative weights, useful for unbalanced classes. A typical value to consider: $\text{sum}(\text{negative instances}) / \text{sum}(\text{positive instances})$ [19].

4) *Hyperparameter spaces notation*: The following notation is used to define the entire hyperparameter space H of q

different parameters:

$$H = P_1 \times \cdots \times P_q = \bigotimes_i P_i \quad (26)$$

Each hyperparameter space P_i corresponds to hyperparameter i and consists of the set of possible hyperparameter values that can be chosen. The following notation for the sets is used:

$$[a, b, s] = \text{All real values from } a \text{ to } b \text{ (including } b) \text{ with step sizes of } s \text{ starting from } a \quad (27)$$

The size or cardinality of the hyperparameter space is given by $|H|$.

5) *The cardinality of initial hyperparameter spaces*: Using the notation in equations (26) and (27) the cardinality of the hyperparameter space that uses the entire range of typical values would be: $|H_{\text{Typical values Random Forest}}| \sim 10^8$ and $|H_{\text{Typical values XGBoost}}| \sim 10^{14}$. Assuming a generous computation time of 0.01 seconds for each cross-validation, an entire exhaustive search of the entire hyperparameter space would take about a month and 20 thousand years, respectively. An exhaustive search is therefore not an option.

6) *Hyperparameter Search*: To reduce the number of parameter combinations that are searched for from the initial total hyperparameter spaces $H_{\text{Typical values Random Forest}}$ and $H_{\text{Typical values XGBoost}}$ a combination of the following methods is used.

- **Cross validated grid search** - This method exhaustively goes through all the hyperparameter combinations in H and runs an iteration of cross-validation to calculate an average scoring value using an evaluator function which is ROC AUC [23].
- **Cross validated random halving grid search** - This method starts with selecting one "resource" parameter that has the property that it correlates with the increased resource usage of the model as it increases. For both the random forest and XGBoost a suitable hyperparameter is the n estimators. The minimum value and maximum values of the n estimators follow the limits of the set $P_{n \text{ estimators}}$. From $H \setminus P_{n \text{ estimators}}$ random values are sampled. During each iteration of the search, the cross-validated score of each chosen combination is calculated and the best third combinations are chosen to continue to the next iteration. This continues until there only exists one combination left. The number of candidates in the initial iteration is chosen so that the last iteration uses as many resources as possible within the limits of the maximum value for the "resource" parameter [24].

The chosen approach to finding the hyperparameter spaces in this project is a combination of, (1) a cross-validated random halving search to find a starting position using $H_{\text{Typical values Random Forest}}$ and $H_{\text{Typical values XGBoost}}$, (2) one or more narrow cross-validated grid searches with H s that use the previous best parameters $\pm \approx 10\%$ of range in $H_{\text{typical values}}$ for each P_i . The hyperparameter space is extended if the best parameter found is on the boundary of the H . In this step, only two to four hyperparameters are searched for in each iteration of the grid search using the previous best value when continuing with the next search.

This limitation is needed to reduce the cardinality of the H s. When finding the hyperparameters for the random forest all the grid searches are divided into two groups in the following order $[n \text{ estimators, maximum depth, maximum features}]$, $[\text{minimum samples split, minimum samples leaf}]$. The reasoning is that the first group mainly deals with the complexity of the model and the second group deals with node splits [18].

When finding the hyperparameters for the XGBoost model the grid searches are divided into three groups in the following order $[n \text{ estimators, maximum depth, learning rate, minimum child weight}]$, $[\text{gamma, column sample by tree, subsample}]$, $[\text{lambda, alpha}]$. The reasoning is that the first group mainly deals with general tree parameters, the second group deals with node construction, and the third group deals with general regularization parameters that penalize complex models. The regularization parameters are also left out during the initial cross-validated random halving search.

IV. RESULTS

A. Random Forest

1) *Before hyperparameter Tuning:* Using the default parameters of the random forest model before any tuning resulted in a mean ROC AUC of 0.67 with a standard error of 0.04 in the final cross-validation iteration. The ROC curve is presented in Figure5.

Receiver operating characteristic Random Forest before tuning hyperparameter

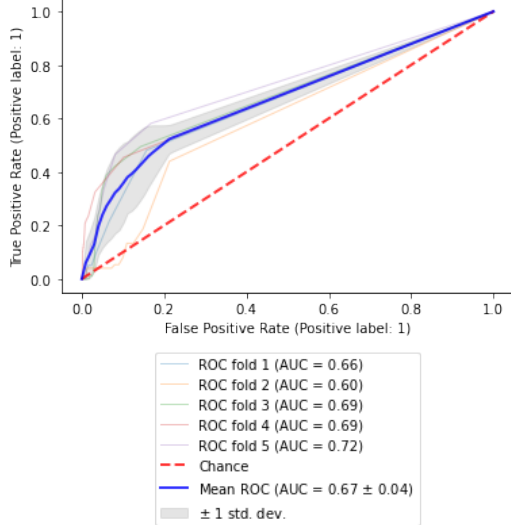


Fig. 5. Results of the Random Forests model before tuning for hyperparameters

2) *Final Hyperparameters:* The following hyperparameters were found as the best hyperparameters after following the tuning procedure described in the methods section. See the

appendix for a detailed description of all steps.

$$\begin{cases} n \text{ estimators} = 190 \\ \text{maximum depth} = 1 \\ \text{maximum features} = 2 \\ \text{minimum samples split} = 3 \\ \text{minimum samples leaf} = 23 \end{cases} \quad (28)$$

The final performance of the model received a mean ROC AUC of 0.842 with a standard error of 0.10 in the cross-validation iteration. The ROC curve is presented in Figure6.

Receiver operating characteristic Random Forest with Tuned Hyperparameters

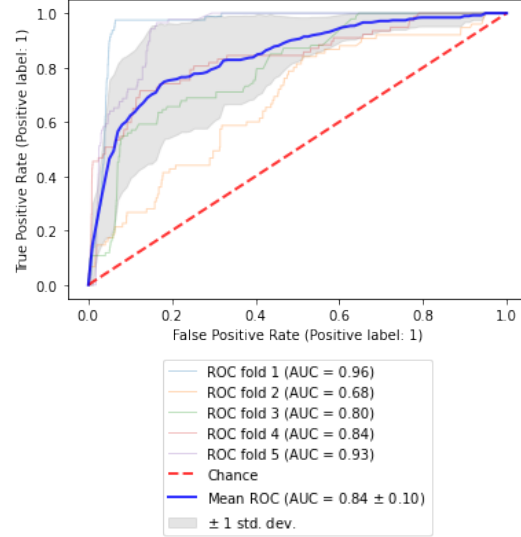


Fig. 6. Results of the Random Forests model after tuning for hyperparameters

B. Results XGBoost

1) *Before hyperparameter Tuning:* Using the default parameters of the XGBoost model before any tuning, except for the scale positive weight hyperparameter that, according to the typical convention, was set to $\text{sum}(\text{negative instances}) / \text{sum}(\text{positive instances})$, which for the dataset was equal to 232.78, resulted in a mean ROC AUC of 0.75 with a standard error of 0.14 in the final cross-validation iteration. The ROC curve is presented in Figure5.

2) *Final hyperparameters:* The following hyperparameters were found as the best hyperparameters after following the tuning procedure described in the methods section. See the

Receiver operating characteristic XGBoost before tuning hyperparameters

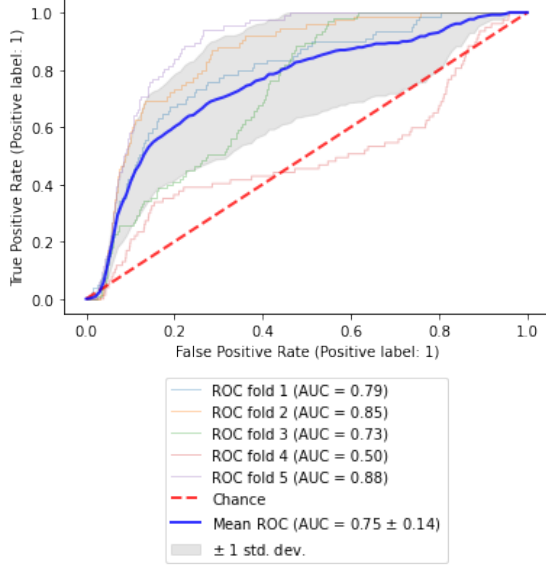


Fig. 7. Results of the XGBoost model before tuning for hyperparameters

Receiver operating characteristic XGBoost with tuned hyperparameters

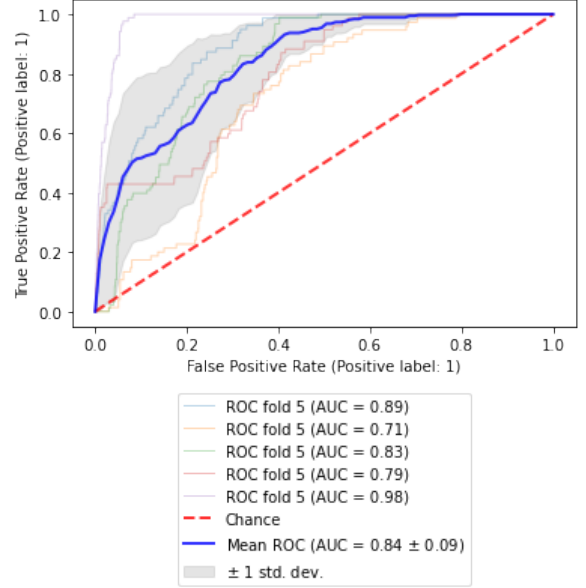


Fig. 8. Results of the XGBoost model after tuning for hyperparameters

appendix for a detailed description of all steps.

$$\left\{ \begin{array}{l} n \text{ estimators} = 140 \\ \text{maximum depth} = 1 \\ \text{learning rate} = 0.34 \\ \text{minimum child weight} = 0 \\ \text{gamma} = 0 \\ \text{column sample by tree} = 0.92 \\ \text{subsample} = 0.6 \\ \text{lambda} = 0.86 \\ \text{alpha} = 0.17 \\ \text{scale positive weight} = 232.78 \end{array} \right. \quad (29)$$

The final performance of the model received a mean ROC AUC of 0.840 with a standard error of 0.10 in the cross-validation iteration. The ROC curve is presented in Figure 8.

V. DISCUSSION

A. Model Performance

The final random forest model demonstrated a mean ROC AUC $\overline{\text{AUC}} = 0.84$ with a standard deviation $s = 0.1$ and the final XGBoost model demonstrated a mean ROC AUC $\overline{\text{AUC}} = 0.84$ with a standard deviation $s = 0.09$. Any value for the ROC AUC that is greater than 0.5 indicates that the model has the ability to distinguish between sepsis-like and non-sepsis-like time windows [22]. It is natural to perform a hypothesis test to evaluate whether $\overline{\text{AUC}}$ differs significantly from 0.5. The test statistic given by $(\overline{\text{AUC}} - 0.5)/d$ is approximately normally distributed [25]. The standard error $d = s/\sqrt{n}$ where $n = 5$ is equal to the number of folds in the cross-validation iterator. With the null and alternative hypotheses defined as $H_0 : \overline{\text{AUC}} = 0.5$ versus $H_1 : \overline{\text{AUC}} \neq 0.5$ yields a test statistics with p -values of $7 \cdot 10^{-12}$ and $9 \cdot 10^{-14}$ for the random forest model and XGBoost model respectively. Thus,

significantly rejecting the null hypothesis, showing that both models can distinguish sepsis-like windows. The grade of a

TABLE II
GRADES OF DIAGNOSTIC TESTS ACCORDING TO THE ROC AUC [25]

ROC AUC	Grade
0.9 - 1	Outstanding
0.8 - 0.9	Excellent
0.7 - 0.8	Acceptable
0.6 - 0.7	Poor
≤ 0.6	Unacceptable

diagnostic test predictor in medicine can be defined according to the categories in table II [25]. Testing whether a model at least has a specific grade g can be done by formulating the null hypotheses $H_0 : \overline{\text{AUC}} < l_g$ against the alternative hypothesis $H_1 : \overline{\text{AUC}} \geq l_g$ where l_g is the lower limit in the ROC AUC for the grade g . The p values for models that at least achieve different grades are given in Table III It is shown

TABLE III
P VALUES FOR THE MODELS TO AT LEAST HAVE A SPECIFIC GRADE

At least ...	XGBoost p -value	Random forest p -value
Outstanding	0.91	0.88
Excellent	0.19	0.20
Acceptable	$1.1 \cdot 10^{-2}$	$2.2 \cdot 10^{-2}$
Poor	$7.1 \cdot 10^{-8}$	$6.0 \cdot 10^{-7}$

that both the random forest and the XGBoost models with significance have at least an acceptable grade as a diagnostic test. Although measures have been taken to reduce bias and overfitting, it is still not guaranteed that the model performs at the calculated level. Only patients from Sweden have been used in the data set, which could have introduced selection bias. Introducing a validation set would have provided more information on overfitting. A common denominator regarding

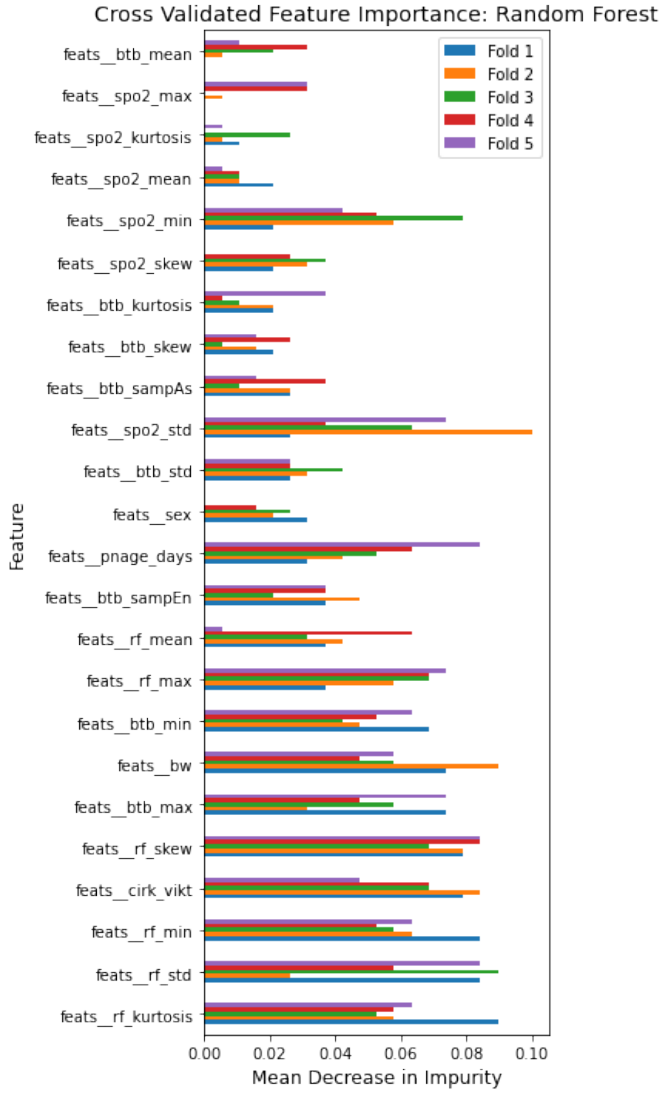


Fig. 9. Feature importance of the Random Forest model

the pitfalls of the results is that they could be partially eluded if the data set had been larger.

1) *Comparison with the SOFA score:* The SOFA score has a highly distinctive ability to predict sepsis with a ROC AUC of 0.89 [26]. However, machine learning models that are close to that level would still be of interest, especially if they do not require invasive methods and increased tasks for healthcare personnel, as in the case of the SOFA score [4].

B. Feature Importance

Feature importances for the final random forest and XGBoost models were calculated by finding the mean decrease in impurity when a feature is removed. The (Gini) impurity is a metric similar to entropy from information theory. A large decrease when a feature is removed signifies greater importance. Calculations were performed on a cross-validation iteration to get a sense of whether the importances were random or not. The importance of each parameter is found in Figures 9 and 10.

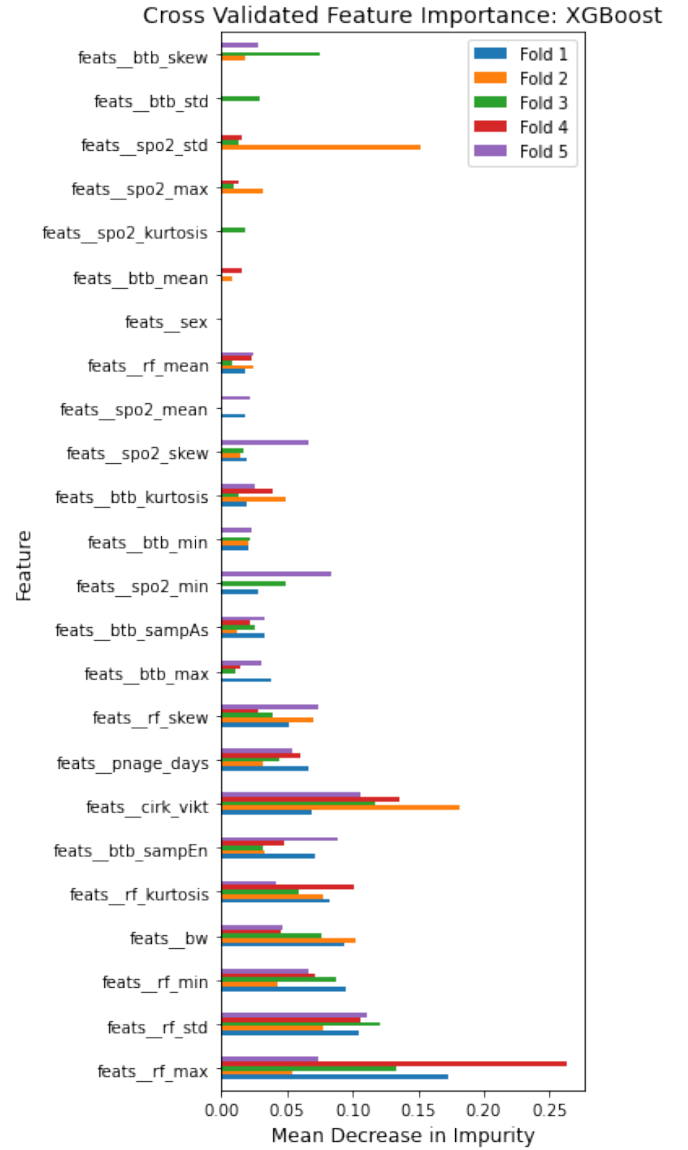


Fig. 10. Feature importance of the XGBoost model

1) *Insights:* For both models, the feature importances for the different folds fluctuated, indicating a component of randomness within the models. However, it is possible to find some general patterns that apply to both models.

- Features of high importance - All features regarding the respiratory frequency and body weight.
- Features of low importance - All features regarding oxygen saturation and sex.
- Features with varying importance - The beat-to-beat interval features achieved varying importance with the maximum and sample entropy features scoring on the higher side within the group.

The results go somewhat against the SOFA score that takes into account oxygen saturation, which was shown to be less important, but none of the respiratory frequencies, bodyweight measurements, or beat-to-beat intervals were shown to be important [11].

C. Model Complexity

Both the random forest and XGBoost models ended up with hyperparameters toward the less complex side, the most significant hyperparameter being a maximum depth of only 1. This signifies that only one comparison of one feature is carried out in every tree. It is up for discussion whether a single node even should be called a tree. An analysis of the model performance dependence on hyperparameter values was performed using data from the hyperparameter tuning procedure. Some of the discovered relationships are as follows.

- Keeping all other hyperparameters the same, the performance of the XGBoost model was greatly reduced with depths of 5 having ROC AUCs of 0.65 to 0.7, while depths of 1 and 2 performed between 0.8 and 0.85.
- Keeping all other hyperparameters the same, the performance of the random forest model was greatly reduced for maximum features greater than 20 with a ROC AUC of 0.65 to 0.7, while depths of maximum features less than 5 performed closer to 0.8.
- For both models, performance plateaued for n estimators larger than 150.

See the appendix for more relationships and graphs of the dependence of the ROC AUC score on the hyperparameters for both the training and the testing sets.

1) *Implications for explainability:* Since both models tended towards the less complex side, a relevant question is whether the models chosen are too complex. On the one hand, it can be argued that there exists room for adding more features and finding more complex patterns. On the other hand, one could see the possible implications for explainability. It is highly debated whether artificial intelligence in healthcare should be explainable, which in essence means that a human should understand why the model makes its decisions [27]. In the current state of the world, it is easier to overcome legal and ethical hurdles with explainable models [27]. Together with the fact that simpler models tend to be more explainable, this implies that having a simple model that is the best performing model brings hope for easy adoption [28]. Since both the random forest and XGBoost are built on the concept of a decision tree that asks simple questions at each node, it is theoretically possible to reverse engineer the decision tree and come up with an easier set of questions that could replace the SOFA scoring system.

VI. CONCLUSION

In conclusion, the use of random forests and XGBoost for the detection of neonatal sepsis is feasible, offering performance within the range of the SOFA scoring system. At the same time, requiring less invasive methods and being automated in the sense that healthcare personnel is not burdened with more tasks. However, both models tended to perform better with reduced complexity. This suggests that it could be possible to gain benefit from adding more features, thus discovering more complex patterns in the data. Or else generate even simple models that focus on the explainability of the data to ease the adoption of the algorithm. One possible method would be to reverse engineer the questions asked at

specific nodes in the decision tree ensemble models.

The project has also found appropriate methods for validation and evaluation that could be generalized to all types of highly unbalanced grouped tabular data that can generalize to other similar applications.

A. Future Work

A possible extension of this project is to calibrate the final models. Currently, due to the usage of non-threshold-dependent scoring systems, the model does not find the exact probability of input being sepsis-like. Generating the probabilities from the model predictions would require calibration with prior probabilities [7]. Another extension is to reverse engineer the models to find simple questions that enhance explainability. Other future work includes acquiring a larger dataset with more sample points and features and testing other diseases or conditions.

ETHICS STATEMENT

All patient data were anonymized by removing personal identification numbers and not providing the exact date and time of the collected data. Therefore, precluding any attempt to identify specific patients. Accessing the data required multi-factor authentication and the generation of one-time passwords for each login session.

ACKNOWLEDGMENT

The authors would like to express gratitude to Antoine Honoré for his support and assistance throughout the course of the project. Furthermore, expressing a special appreciation for the Herlenius Research Team for their work in collecting, working with, and providing the dataset for this project.

SOURCE CODE

The source code for the project and all plots can be found in the following GitHub repository: <https://gits-15.sys.kth.se/mmmab/Neonatal-Sepsis-Detection-Using-Decision-Tree-Ensemble-Methods-Random-Forest-and-XGBoost>.

REFERENCES

- [1] A.-M. Audet, S. Greenfield, and M. Field, "Medical practice guidelines: current activities and future directions," *Annals of Internal Medicine*, vol. 113, no. 9, pp. 709–714, 1990.
- [2] M. Singer, C. S. Deutschman, C. W. Seymour, M. Shankar-Hari, D. Annane, M. Bauer, R. Bellomo, G. R. Bernard, J.-D. Chiche, C. M. Coopersmith, R. S. Hotchkiss, M. M. Levy, J. C. Marshall, G. S. Martin, S. M. Opal, G. D. Rubenfeld, T. van der Poll, J.-L. Vincent, and D. C. Angus, "The Third International Consensus Definitions for Sepsis and Septic Shock (Sepsis-3)," *JAMA*, vol. 315, no. 8, pp. 801–810, 02 2016. [Online]. Available: <https://doi.org/10.1001/jama.2016.0287>
- [3] C. Fleischmann, F. Reichert, A. Cassini, R. Horner, T. Harder, R. Markwart, M. Tröndle, Y. Savova, N. Kissoon, P. Schlattmann *et al.*, "Global incidence and mortality of neonatal sepsis: a systematic review and meta-analysis," *Archives of Disease in Childhood*, vol. 106, no. 8, pp. 745–752, 2021.
- [4] E. P. Raith, A. A. Udy, M. Bailey, S. McGloughlin, C. MacIsaac, R. Bellomo, D. V. Pilcher, for the Australian, N. Z. I. C. S. A. C. for Outcomes, and R. E. (CORE), "Prognostic Accuracy of the SOFA Score, SIRS Criteria, and qSOFA Score for In-Hospital Mortality Among Adults With Suspected Infection Admitted to the Intensive Care Unit," *JAMA*, vol. 317, no. 3, pp. 290–300, 01 2017. [Online]. Available: <https://doi.org/10.1001/jama.2016.20328>

- [5] B. Cummings, "Rising healthcare costs are a rising concern," *Journal of Financial Planning*, vol. 35, no. 2, pp. 19–19, 2022.
- [6] M. Marć, A. Bartosiewicz, J. Burzyńska, Z. Chmiel, and P. Januszewicz, "A nursing shortage—a prospect of global and local policies," *International nursing review*, vol. 66, no. 1, pp. 9–16, 2019.
- [7] R. N. Stuart Peter, *Artificial Intelligence A Modern Approach*, 3rd ed. Harlow, England: Pearson Education Limited, 2016.
- [8] A. Honore, D. Liu, D. Forsberg, K. Coste, E. Herlenius, S. Chatterjee, and M. Skoglund, "Hidden markov models for sepsis detection in preterm infants," in *Proceedings IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020* :, ser. International Conference on Acoustics Speech and Signal Processing ICASSP. Institute of Electrical and Electronics Engineers (IEEE), 2020, pp. 1130–1134, qC 20210324.
- [9] A. Honoré, D. Forsberg, K. Jost, K. Adolphson, A. Stålhammar, E. Herlenius, and S. Chatterjee, "Classification and feature extraction for neonatal sepsis detection," 2022.
- [10] K. D. Fairchild and T. M. O'Shea, "Heart rate characteristics: physiologic markers for detection of late-onset neonatal sepsis," *Clinics in perinatology*, vol. 37, no. 3, pp. 581–598, 2010.
- [11] A. Oscarson, C. Bjurman, J. E. Wallér, and M. Werner, "Sepsis hos vuxna – tidig upptäckt och initial behandling," *Läkartidningen*, Mar 2017.
- [12] A. Rhodes, L. E. Evans, W. Alhazzani, M. M. Levy, M. Antonelli, R. Ferrer, A. Kumar, J. E. Sevransky, C. L. Sprung, M. E. Nunnally *et al.*, "Surviving sepsis campaign: international guidelines for management of sepsis and septic shock: 2016," *Intensive care medicine*, vol. 43, no. 3, pp. 304–377, 2017.
- [13] P.-Y. Iroh Tam and C. M. Bendel, "Diagnostics for neonatal sepsis: current approaches and future directions," *Pediatric research*, vol. 82, no. 4, pp. 574–583, 2017.
- [14] J. L. Wynn and R. A. Polin, "A neonatal sequential organ failure assessment score predicts mortality to late-onset sepsis in preterm very low birth weight infants," *Pediatric research*, vol. 88, no. 1, pp. 85–90, 2020.
- [15] M. J. Aspinall, "Use of a decision tree to improve accuracy of diagnosis," *Nursing Research*, vol. 28, no. 3, pp. 182–185, 1979.
- [16] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [17] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," ser. KDD '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 785–794. [Online]. Available: <https://doi.org/10.1145/2939672.2939785>
- [18] SKlearn, "Sklearn.ensemble.randomforestclassifier," 2022. [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
- [19] XGBoost, "Xgboost parameters," 2021. [Online]. Available: <https://xgboost.readthedocs.io/en/stable/parameter.html>
- [20] J. Brownlee, "Train-test split for evaluating machine learning algorithms," Aug 2020. [Online]. Available: <https://machinelearningmastery.com/train-test-split-for-evaluating-machine-learning-algorithms/>
- [21] Scikit, "3.1. cross-validation: Evaluating estimator performance," 2022. [Online]. Available: https://scikit-learn.org/stable/modules/cross_validation.html
- [22] X.-H. Zhou, D. K. McClish, and N. A. Obuchowski, *Statistical methods in diagnostic medicine*. John Wiley & Sons, 2009.
- [23] SKlearn, "Sklearn.model_selection.gridsearchcv," 2022. [Online]. Available: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
- [24] Scikit, "Sklearn.model_selection.halvingrandomsearchcv," 2022. [Online]. Available: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.HalvingRandomSearchCV.html
- [25] J. N. Mandrekar, "Receiver operating characteristic curve in diagnostic test assessment," *Journal of Thoracic Oncology*, vol. 5, no. 9, pp. 1315–1316, 2010.
- [26] A. K. Toker, S. Kose, and M. Turken, "Comparison of sofa score, sir's, qsofa, and qsofa+ 1 criteria in the diagnosis and prognosis of sepsis," *The Eurasian Journal of Medicine*, vol. 53, no. 1, p. 40, 2021.
- [27] J. Amann, A. Blasimme, E. Vayena, D. Frey, and V. I. Madai, "Explainability for artificial intelligence in healthcare: a multidisciplinary perspective," *BMC Medical Informatics and Decision Making*, vol. 20, no. 1, pp. 1–9, 2020.
- [28] P. Linardatos, V. Papastefanopoulos, and S. Kotsiantis, "Explainable ai: A review of machine learning interpretability methods," *Entropy*, vol. 23, no. 1, p. 18, 2020.