Contents lists available at ScienceDirect

# Artificial Intelligence in Medicine

# Building interpretable predictive models for pediatric hospital readmission using Tree-Lasso logistic regression

Milos Jovanovic [a], Sandro Radovanovic [a], Milan Vukicevic [a,*], Sven Van Poucke [b], Boris Delibasic [a]

[a] *University of Belgrade, Faculty of Organizational Sciences, Jove Ilica 154, 11010 Vozdovac, Belgrade, Serbia*
[b] *Department of Anesthesiology, Critical Care, Emergency Medicine and Pain Therapy, Ziekenhuis Oost-Limburg, Schiepse Bos 6, B-3600 Genk, Belgium*

## ARTICLE INFO

## ABSTRACT

*Objectives:* Quantification and early identification of unplanned readmission risk have the potential to improve the quality of care during hospitalization and after discharge. However, high dimensionality, sparsity, and class imbalance of electronic health data and the complexity of risk quantification, challenge the development of accurate predictive models. Predictive models require a certain level of interpretability in order to be applicable in real settings and create actionable insights. This paper aims to develop accurate and interpretable predictive models for readmission in a general pediatric patient population, by integrating a data-driven model (sparse logistic regression) and domain knowledge based on the international classification of diseases 9th—revision clinical modification (ICD-9-CM) hierarchy of diseases. Additionally, we propose a way to quantify the interpretability of a model and inspect the stability of alternative solutions.

*Materials and methods:* The analysis was conducted on >66,000 pediatric hospital discharge records from California, State Inpatient Databases, Healthcare Cost and Utilization Project between 2009 and 2011. We incorporated domain knowledge based on the ICD-9-CM hierarchy in a data driven, Tree-Lasso regularized logistic regression model, providing the framework for model interpretation. This approach was compared with traditional Lasso logistic regression resulting in models that are easier to interpret by fewer high-level diagnoses, with comparable prediction accuracy.

*Results:* The results revealed that the use of a Tree-Lasso model was as competitive in terms of accuracy (measured by area under the receiver operating characteristic curve—AUC) as the traditional Lasso logistic regression, but integration with the ICD-9-CM hierarchy of diseases provided more interpretable models in terms of high-level diagnoses. Additionally, interpretations of models are in accordance with existing medical understanding of pediatric readmission. Best performing models have similar performances reaching AUC values 0.783 and 0.779 for traditional Lasso and Tree-Lasso, respectfully. However, information loss of Lasso models is 0.35 bits higher compared to Tree-Lasso model.

*Conclusions:* We propose a method for building predictive models applicable for the detection of readmission risk based on Electronic Health records. Integration of domain knowledge (in the form of ICD-9-CM taxonomy) and a data-driven, sparse predictive algorithm (Tree-Lasso Logistic Regression) resulted in an increase of interpretability of the resulting model. The models are interpreted for the readmission prediction problem in general pediatric population in California, as well as several important subpopulations, and the interpretations of models comply with existing medical understanding of pediatric readmission. Finally, quantitative assessment of the interpretability of the models is given, that is beyond simple counts of selected low-level features.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

An increased availability of pediatric health data facilitates the investigations and efforts for improved quality of care [1] and enables improvement in several areas, including the optimization of treatments, reduction of adverse events and readmission rates and earlier identification of populations in need.

\* Corresponding author.
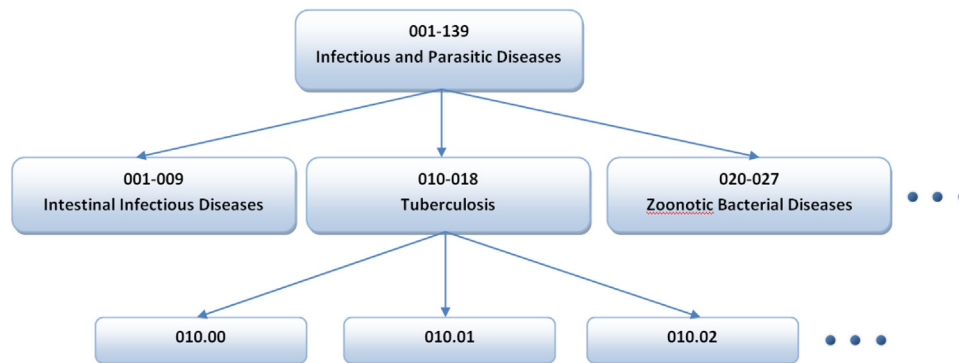*E-mail address:* milan.vukicevic@fon.bg.ac.rs (M. Vukicevic).

**Fig. 1.** Excerpt of ICD-9-CM hierarchy.

Hospital readmission (admission to a hospital within 30 days of discharge) is disruptive to both patients and healthcare providers and is often associated with higher costs and penalties. Modern care standards require effective discharge planning including the transfer of information at discharge, patient and parent education, and coordination of care after discharge. In pediatrics, the analysis of hospital readmission continues to be challenging based on the multitude of influencing factors (e.g. seasonal variations) and is considered a critical metric of the quality and cost of healthcare [2,3]. Based on a recent report [4], readmission rate within 30 days is 19.6%. Additionally, 34.0% of the pediatric patients return to the hospital within 90 days and 56.1% within one year following discharge. According to the Institute for Healthcare Improvement [4], approximately 76% of the 5 million U.S. hospital readmissions are more or less preventable [5]. Therefore, predictive modeling of the readmission risk deserves more attention to be elaborated.

Predictive algorithms can be used to identify the risk of possible readmission, and make an early warning system for the patient at the risk. Moreover, patterns in patient data as recognized by predictive algorithms could provide additional insights into the factors influencing readmission.

However, learning algorithms often fail to capture dependencies between high dimensional health related data and readmission. This could be directly related to the high-dimensionality of the data [1]. Another challenge is the high class-imbalance, meaning that majority of the patients are not readmitted within 30 days resulting in a predictive model biased in the direction of predicting a negative output (not readmitted). Another problem is that state-of-the-art predictive algorithms (that often provide highly accurate models) usually do not provide interpretable models (i.e. neural networks or support vector machines). Such models could potentially be used as early warning systems, but do not provide any actionable insights and reasons that lead to potential readmission. Finally, the notion of interpretability is subjective and in predictive modeling its analysis is most often limited to more interpretable models (i.e. Decision Trees).

This research aims to develop an accurate and interpretable predictive model for hospital readmission based on the integration of data from electronic health records (EHRs) and domain knowledge represented by the international classification of diseases 9th revision—clinical modification (ICD-9-CM) code hierarchy [6]. In this paper, we applied a Tree-Lasso logistic regression [7] in order to integrate ICD-9-CM hierarchy and logistic regression. This kind of integration of domain knowledge is interesting for research [8–10] since it improves the generalizability of predictive models and reduces the number of features, resulting in more interpretable models. Additionally, we proposed a method for quantification of interpretability of sparse predictive models.

The main contributions of this research are: (1) integration of domain knowledge (in the form of ICD-9-CM taxonomy) and learning algorithms, to improve interpretability of models; (2) a quantitative assessment of the interpretability of predictive models, beyond simple counts of features; (3) evaluation and interpretation of predictive models for readmission prediction for State Inpatient Databases (SID) pediatric patient data in California; (4) the resulting models are interpreted for general pediatric population, as well as several important subpopulations while interpretation of models comply with existing medical knowledge.

Experiments are conducted on SID California data consisting of 67,000 pediatric patients admitted between 2009 and 2011. Each admission is described by maximum 15 diagnoses leading to over 15,000 binary features.

## 2. Background

Model interpretability is fairly abstract and subjective. However, a model can be defined as interpretable if the behavior of that model can be explained verbally and that the model can be used for reasoning. For medical applications, model interpretability is without any doubt a very important property, because of the underlying complexity of the phenomena that are analyzed, and the potential impact of wrong decisions. In this context, medical doctors as decision makers require a fundamental understanding of predictive models if they are implemented as clinical decision support [11].

Datasets with thousands of features are common in medicine. The interpretability of a model is related to (i) the number of features and (ii) the information provided by the features. The number of features is intuitively evident as interpretability measure. The higher the dimensionality, the more complex it becomes for human beings to analyze the relative impact of features and patterns with the potentially important in decisions. Therefore, using a reduced set of features might lead to more interpretable models. On the other hand, the contextual information provided by the features is important regardless to dimensionality. If a model is based on a limited number of features but the model is considered a black box by the human interpreter, then the model is not interpretable. This is the reason predictive modeling in the medical domain usually relies on traditional algorithms like Logistic regression (parameters of logistic regression can be interpreted as the logarithm of odds ratio) or Decision Trees (that can be interpreted as hierarchical rule set). In this research, we utilized Tree-Lasso logistic regression [7] which harnesses both elements of model interpretability. First, it forces model parameters to be zero resulting in a smaller number of features (sparse models). Second, it selects features based on group similarity (in this case to groups of diagnoses provided by ICD-9-CM codes hierarchy) facilitating the applicability of a model in a medical environment. Additionally, we proposed a method for quantification and comparison of interpretability of sparse predictive models.