Chapter 46 Comprehensible Explanation of Predictive Models

Marko Robnik-Šikonja

University of Ljubljana, Slovenia

ABSTRACT

The most successful prediction models (e.g., SVM, neural networks, or boosting) unfortunately do not provide explanations of their predictions. In many important applications of machine learning, the comprehension of the decision process is of utmost importance and dominates the classification accuracy (e.g., in business and medicine). This chapter introduces general explanation methods that are independent of the prediction model and can be used with all classification models that output probabilities. It explains how the methods work and graphically explains models' decisions for new unlabeled cases. The approach is put in the context of applications from medicine, business, and macro-economy.

INTRODUCTION

In many areas where machine learning methods are applied the practitioners and users of produced prediction models are interested in comprehensible explanation of their predictions. Unfortunately, the best performing predictive models do not offer an intrinsic introspection into their decision processes or provide explanations of their prediction. This is true for Support Vector Machines (SVM), Artificial Neural Networks (ANN), and all ensemble methods (for example, boosting, random forests, bagging, stacking and multiple adaptive regression splines). Approaches that do offer an intrinsic introspection such as decision trees or decision rules do not perform so well or are not applicable in many cases (Meyer et al., 2003). The areas where models' transparency is of crucial importance are for example most of business and marketing applications where the executives are just as interested in the comprehension of the decision process, explanation of the existing and new customers' needs and expectations in a given business case, as in the classification accuracy of the prediction model. The same is true for many areas of business intelligence, finance, marketing, insurance, medicine, science, policy making, and strategic planning where knowledge discovery dominates prediction accuracy.

DOI: 10.4018/978-1-5225-7362-3.ch046

To alleviate this problem two types of solutions have been proposed. The first type is based on internal working of each particular learning algorithm and exploits its learning process to gain insight into the presumptions, biases and reasoning leading to final decisions. A well-known example of such an approach are random forests for which several visualizations exist mostly exploiting the fact that during bootstrap sampling some of the instances are not selected for learning and can serve as an internal validation set. With the help of this set important features can be identified and similarity between objects can be measured. The second type of explanation approaches are general and can be applied to any predictive model. Examples of this approach are methods EXPLAIN (Robnik-Šikonja & Kononenko, 2008) and IME (Štrumbelj et al., 2009). These two methods are based on efficient implementation of input perturbations. They can explain models' decision process for each individual predicted instance as well as the model as a whole. As both methods are efficient, offer comprehensible explanations, can be visualized, and are readily available in R package ExplainPrediction (Robnik-Šikonja, 2015) they are the focus of this article. Other general explanation methods are discussed in the background Section.

The objective of the article is to explain how EXPLAIN and IME explanation methods work and to show their practical utility in several real world scenarios. The first aim is achieved through explanation of their working principle and graphical explanation of models' decisions on a well-known data set. Two types of explanations are demonstrated, predictions of new unlabeled cases and the functioning of the model as a whole. This allows inspection, comparison, and visualization of otherwise opaque models. The practical utility of the methodology is demonstrated with short description of several applications: in medicine (Štrumbelj et al. 2010), macro economy (Pregeljc et al., 2012) and business consultancy (Bohanec et al., 2015).

BACKGROUND

In a typical data science problem setting, users are concerned with both prediction accuracy and the interpretability of the prediction model. Complex models have potentially higher accuracy but are more difficult to interpret. This can be alleviated either by sacrificing some prediction accuracy for a more transparent model or by using an explanation method that improves the interpretability of the model. Explaining predictions is straightforward for symbolic models such as decision trees, decision rules, and inductive logic programming, where the models give an overall transparent knowledge in a symbolic form. Therefore, to obtain the explanations of predictions, one simply has to read the rules in the corresponding model. Whether such an explanation is comprehensive in the case of large trees and rule sets is questionable.

For non-symbolic models there are no such straightforward explanations. A lot of effort has been invested into increasing the interpretability of complex models. A taxonomy of explanation methods and a review of neural network explanation approaches is given by Jacobsson (2005). For Support Vector Machines an interesting approaches is proposed by Hamel (2006). Many approaches exploit the essential property of additive classifiers to provide more comprehensible explanations and visualizations, e.g., (Jakulin et al., 2005) and (Poulin et al. 2006).

Visualization of decision boundaries is an important aspect of model transparency. Barbosa et al. (2016) present a technique to visualize how the kernel embeds data into a high-dimensional feature space. With their Kelp method they visualize how kernel choice affects neighborhood structure and SVM decision boundaries. Schultz et al. (2015) propose a general framework for visualization of classifiers via

10 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/comprehensible-explanation-of-predictivemodels/212144

Related Content

Micropolitan Areas Creating Leadership in the New Economy: Developing Micropolitan Areas to Develop a New Economy

Kristin Joyce Tardif (2020). International Journal of Responsible Leadership and Ethical Decision-Making (pp. 1-18).

www.irma-international.org/article/micropolitan-areas-creating-leadership-in-the-new-economy/273056

Academic Community Manager: Manager of the Academic Community

Ariana Daniela Del Pino Espinoza (2017). Remote Work and Collaboration: Breakthroughs in Research and Practice (pp. 520-535).

www.irma-international.org/chapter/academic-community-manager/180119

Intuitive Knowledge Generation in Post-Bureaucratic Organizations

Marta Sinclair (2017). *Evolution of the Post-Bureaucratic Organization (pp. 383-400).* www.irma-international.org/chapter/intuitive-knowledge-generation-in-post-bureaucratic-organizations/174855

Create Value in Family Business SMEs in Colombia

Helder Barahona Urbano (2020). Handbook of Research on Increasing the Competitiveness of SMEs (pp. 305-328).

www.irma-international.org/chapter/create-value-in-family-business-smes-in-colombia/246467

A Study on the Wide-Ranging Ethical Implications of Big Data Technology in a Digital Society: How Likely Are Data Accidents During COVID-19?

Izabella V. Lokshinaand Cees J. M. Lanting (2021). *Journal of Business Ecosystems (pp. 32-57).* www.irma-international.org/article/a-study-on-the-wide-ranging-ethical-implications-of-big-data-technology-in-a-digitalsociety/270479