



Modeling medical decision making by support vector machines, explaining by rules of evolutionary algorithms with feature selection

Ruxandra Stoean*, Catalin Stoean

Department of Computer Science, Faculty of Exact Sciences, University of Craiova, A.I. Cuza Str., No. 13, 200585 Craiova, Romania

ARTICLE INFO

Keywords:

Support vector machines
Evolutionary algorithms
Rule extraction
Cooperative coevolution
Feature selection
Breast cancer diagnosis

ABSTRACT

Machine learning support for medical decision making is truly helpful only when it meets two conditions: high prediction accuracy and a good explanation of how the diagnosis was reached. Support vector machines (SVMs) successfully achieve the first target due to a kernel-based engine; evolutionary algorithms (EAs) can greatly accomplish the second owing to their adaptable nature. In this context, the current paper puts forward a two-step hybridized methodology, where learning is accurately performed by the SVMs and a comprehensible emulation of the resulting decision model is generated by EAs in the form of propositional rules, while referring only those indicators that highly influence the class separation. An individual highlighting of the medical attributes that trigger a specific diagnosis for a current patient record is additionally obtained; this feature thus increases the confidence of the physician in the resulting automated diagnosis. Without loss of generality, we aim to model three breast cancer instances, for reasons of both high incidence of the disease and the large application of state of the art artificial intelligence methods for this medical task. As such, the prediction of a benign/malignant condition as well as the recurrence/nonrecurrence of a cancer event are studied on the Wisconsin corresponding data sets from the UCI Machine Learning Repository. The proposed hybridization reached its goals. Rule prototypes evolve against a SVM consistent training data, while diversity among the different classes is implicitly preserved. Feature selection eventually leads to a resulting rule set where only the significant medical indicators together with the discriminating threshold values are referred, while individual relevance of attributes can be additionally obtained for each patient. The gain is thus dual: the EA benefits from a noise-free SVM preprocessed data and the resulting SVM model is able to output rules in a comprehensible, concise format for the physician.

© 2012 Elsevier Ltd. All rights reserved.

1. Introduction

In the context of the highly increasing amount of machine learning methods to assist medical diagnosis, a good prediction accuracy alone is no longer sufficient for some given approach to be considered of truly reliable decision support to the physician. While this is certainly a prerequisite, supplementary information on how a verdict had been reached, based on the medical indicators at hand, is necessary for the learning model to be fully trusted as a second opinion.

Such knowledge can be expressed in the form of rules that discover attribute values connected to each possible outcome. The expert can then look for those indicators and values that distinguish the possible diagnoses for a case.

Within this context, kernel-based techniques like neural networks (NNs) and support vector machines (SVMs) outperform

rule-centered classifiers in prediction accuracy, while only the learning engine of the latter is transparent. The compromise is to combine such two methodologies to achieve both reasonable accuracy and an expressive logic of the decision process.

SVMs are as such one of the best classifiers when it comes to an increased accuracy of prediction. Their weak point yet lies in their opaque learning machine. Existing implementations nonetheless offer the possibility to also extract the coefficients of the decision hyperplane. However, this merely provides a formula where the importance of each attribute is weighted. Moreover, indicators are included in totality in such an expression and no particular guidelines of the decision making process can be guessed.

In this respect, the paper puts forward a novel interplay between SVMs and evolutionary algorithms (EAs) to achieve comprehensible evolutionary rules that explain a SVM model of medical diagnosis. SVMs originality among kernel machines is due to their geometrical learning fashion and the particular support vectors that delineate the decision surface (Vapnik, 1995). Rule definition and manipulation can be very flexible through EA evolution (Eiben & Smith, 2003), while the multiple distinct prototypes to comprise

* Corresponding author. Tel./fax: +40 251 413728.

E-mail addresses: ruxandra.stoean@inf.ucv.ro (R. Stoean), catalin.stoean@inf.ucv.ro (C. Stoean).

the final set can be easily created and maintained through a multimodal treatment, like cooperative coevolution (CC) (Potter & De Jong, 2000). What is more, of all the combinations between opaque predictors and rule generators, hybridizations between SVMs and EAs have been the least often explored.

Nevertheless, even with an account of the attribute threshold values that discriminate between the outcomes, it is usual that the medical exams undertaken for an accurate assessment of the illness of a patient generate a large number of attributes. This poses problems both to the classifiers (the curse of dimensionality) and to the human interpretation of the resulting rules. Therefore, the EA within the proposed hybridized approach is endowed with an additional mechanism for feature selection. The SVM part still has to deal with the entirety of indicators, but this theoretically constitutes no problem, as their working is supposed to be independent of the number of features involved in a decision problem (Joachims, 1998). For the EA classifier, however, this both eases rule generation and induces a dynamic dual evolution between the chosen attributes and their determined thresholds. What is though more important from the practical point of view is that the physician can now concentrate only on those indicators (with given threshold values) that are significant for the respective output, without losing the focus on the redundant or misleading ones.

Finally, a determination of the relevant indicators of the reached diagnosis for a particular patient is obtained taking into account the similarity between the gathered rules and the individual values. This trait answers the practical motivation of building such an automated diagnosis system: that of truly making the physician understand and trust the judgement of the fellow computer when establishing a (sometimes vital) medical diagnosis.

Without losing the generality of the model, the particular medical problem tackled by the present methodology is the diagnosis of breast cancer for reasons of both increased incidence and numerous related machine learning studies on the task. Hence, the potential of an accurate, comprehensible and short explanatory rule set is investigated on three distinct data instances related to this type of cancer.

The paper is structured as follows. The description of the breast cancer instances is made in Section 2. The proposed machine learning approach is detailed in Section 3: the SVM model and output are outlined in Section 3.2, the CC framework and the consequential rule extraction EA engine are presented in Section 3.3, the inclusion of the simultaneous feature selection is explained in Section 3.4 and the methodology of explaining an individual classification, fact of high importance for the medical expert, is included in Section 3.5. The experimental task, setup, obtained results and conclusions are discussed in Section 4, while some final remarks are comprised in Section 5.

2. Materials

Breast cancer ranks second in the list of common cancer types, as estimated by the American National Cancer Institute, with a forecasted 226,870 (female) – 2,190 (male) new cases and 39,510 (female) – 410 (male) deaths in the United States in 2012. It is therefore needless to say that a reliable computational framework to assist breast cancer diagnosis and recurrence would aid in the early detection and treatment of this disease as soon as the results of the specific medical exams are available. Such a methodology would rapidly match the relationship between the present indicators and the possible outcomes (based on its prior learning of many different other cases) and would provide its reached decision together with an explanation of the underlying reasoning. The physician would then also be aware of a second, trustworthy (machine)

opinion and interpretation of the result and could then take his/her final decision on the diagnosis.

The breast cancer topic has been extensively studied by the machine learning community by appointing various state of the art techniques for many available data sets reflecting the disease. Among the most recent published studies, one can see SVMs (Zeng & Liu, 2010; Li, Liu, & Hu, 2011; Akay, 2009), NNs (Gorunescu, Gorunescu, El-Darzi, & Gorunescu, 2010) and fuzzy approaches (Ghazavi & Liao, 2008) are present as some of the prominent players in the field.

In this paper, we selected three data sets that are available at the UCI Machine Learning Repository (Frank & Asuncion, 2010), namely breast cancer Wisconsin original, diagnostic and prognostic cases. The potential of the proposed approach will thus be more thoroughly examined, as the combination of medical indicators is different with each problem as well as one of them treats the equally important task of the recognition of a recurrent cancer event.

The original data has 9 attributes, detailed below with their corresponding domains:

- Clump thickness: 1–10.
- Uniformity of cell size: 1–10.
- Uniformity of cell shape: 1–10.
- Marginal adhesion: 1–10.
- Single epithelial cell size: 1–10.
- Bare nuclei: 1–10.
- Bland chromatin: 1–10.
- Normal nucleoli: 1–10.
- Mitoses: 1–10.

The two possible outcomes are benign and malignant with 458 and 241 cases each.

The features of the diagnostic collection describe characteristics of the cell nuclei present in a digitized image of a fine needle aspirate (FNA) of a breast mass (Frank & Asuncion, 2010). Every cell nucleus is defined by ten traits and for every trait the mean, the standard error and the worst (mean of the three largest values) are computed, resulting in a total of 30 features for each image:

- Radius (mean of distances from center to points on the perimeter): 10.95–27.22.
- Texture (standard deviation of gray-scale values): 10.38–39.28.
- Perimeter: 71.90–182.10.
- Area: 361.60–2250.
- Smoothness (local variation in radius lengths): 0.075–0.145.
- Compactness ($\text{perimeter}^2/\text{area} - 1.0$): 0.046–0.311.
- Concavity (severity of concave portions of the contour): 0.024–0.427.
- Concave points (number of concave portions of the contour): 0.020–0.201.
- Symmetry: 0.131–0.304.
- Fractal dimension (coastline approximation – 1): 0.050–0.097.

The class distribution for this data set is 357 benign and 212 malignant.

The prognostic problem has two outcomes (nonrecurrent with 151 samples and recurrent with 47) and has the same 30 attributes measured for breast images in the diagnostic situation, plus three more, i.e.:

- Time (recurrence time if class is recurrent, disease-free time if nonrecurrent): 1–125.
- Tumor size – diameter of the excised tumor in centimeters: 0.400–10.00.

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات