The Smaller, the Better: Selecting Refined SVM Training Sets Using Adaptive Memetic Algorithm

Jakub Nalepa Institute of Informatics Silesian University of Technology Akademicka 16 44-100 Gliwice, Poland jakub.nalepa@polsl.pl

ABSTRACT

Support vector machine (SVM) is a supervised classifier which has been applied for solving a wide range of pattern recognition problems. However, training of SVMs may easily become a bottleneck, because of its time and memory requirements. Enduring this issue is a vital research topic, especially in the era of big data. In this abstract, we present our adaptive memetic algorithm for selection of refined (significantly smaller) SVM training sets. The algorithm—being a hybrid of an adaptive genetic algorithm and some refinement procedures—exploits the knowledge about the training set vectors extracted before the evolution, and attained dynamically during the search. The results obtained for several real-life, benchmark, and artificial datasets showed that our approach outperforms the other state-of-the-art techniques, and is able to extract very high-quality SVM training sets.

CCS Concepts

•Computing methodologies \rightarrow Support vector machines; Bio-inspired approaches; Artificial intelligence;

Keywords

SVM; training set selection; memetic algorithm; PCA; adaptation

1. INTRODUCTION

The process of the SVM training involves solving a constrained quadratic programming optimization problem of $O(t^3)$ time and $O(t^2)$ memory complexity, where t denotes the cardinality of the training set **T**. It entails selecting important **T** vectors, referred to as *support vectors* (SVs). These vectors define the position of the decision hyperplane, which is later used to classify the incoming (unseen) data. Additionally, the classification time linearly depends on the number of SVs (s). Hence, minimizing s (ideally without

GECCO'16 Companion July 20-24, 2016, Denver, CO, USA

© 2016 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-4323-7/16/07.

DOI: http://dx.doi.org/10.1145/2908961.2930950

Michal Kawulok Institute of Informatics Silesian University of Technology Akademicka 16 44-100 Gliwice, Poland michal.kawulok@polsl.pl

affecting the classification performance) is of a great importance too, especially in the case of real-time applications.

There are two main streams of research towards dealing with the time-consuming SVM training. The first one encompasses techniques which aim at enhancing the training process itself, whereas the other one includes algorithms which retrieve subsets of T (denoted as T''s) for feeding them into the training procedure. Since SVs constitute a subset of T, training SVMs with T''s containing only SVs may be expected to give the same model as if they are trained with the entire T. It is beneficial to extract significantly smaller refined sets—this will not only speed up the training, but can also improve the classifier performance (if the "misleading", e.g., mislabeled, vectors are removed from T), and decrease the classification time (if s is reduced).

In our very recent work published in the Neurocomputing journal [2], we proposed the adaptive memetic algorithm (abbreviated as PCA²MA—Principal Component Analysis Adaptive Memetic Algorithm) for selection of refined training sets for SVMs. In PCA²MA, the information about potentially important vectors from T is extracted before the evolution in the pre-processing step, and is coupled with the knowledge attained during the optimization to better guide the search towards high-quality refined sets. Thanks to the efficient, parameter-less adaptation techniques applied in PCA²MA, this algorithm dynamically updates its crucial parameters (the size of refined sets being the most important) on the fly, thus does not require undertaking the time-consuming tuning process. The comprehensive experiments performed on several real-life, benchmark, and artificial datasets revealed that PCA²MA significantly outperforms not only our previous evolutionary approaches (including the memetic algorithm presented at GECCO '14 [1]). but also other state-of-the-art training set selection techniques belonging to various groups of algorithms. The experimental study included visualizations of retrieved T''s along with the annotated SVs, and the statistical tests (two-tailed Wilcoxon tests) to prove the significance of the results.

2. THE ALGORITHM

In PCA²MA, the population of individuals (each representing a refined training set) evolves in time. This evolution is preceded by the pre-processing step, in which the input space is transformed using principal component analysis to extract most discriminating dimensions. Each dimension is then divided into intervals in such a way, that the intervals contain the same number of T vectors. The geometry

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

of these vector groups is further analyzed to form the set of *candidate vectors* (CVs), which will contain as heterogeneous T vectors as possible. The CVs are exploited at various PCA²MA steps, to (i) create the initial population, and to (ii) guide the evolutionary process.

During the evolution, the individuals undergo the genetic operations (selection, crossover, and mutation), as well as the memetic ones. The latter operations encompass the education, in which "weak" vectors from an individual are replaced with those residing in the support vectors pool. The pool gathers the vectors which were appended to at least one chromosome, and were selected as SVs at some point during the optimization. Hence, they can be considered as important vectors (they influenced the position of the SVM decision hyperplane). The pool is also utilized to create super individuals (chromosomes composed of the so-far-determined SVs only). Determining the fitness of an individual consists in training an SVM using the corresponding T' (containing t' vectors, where $t' \ll t$), and quantifying its performance for the set T—the area under the receiver operating characteristic curve (AUC) is used for this purpose.

Since the appropriate size of the refined sets (thus the size of the individuals) is not known *a priori*, it is dynamically updated during the evolution (along with the selection scheme). The process of updating t' is parameter-less, and is based on investigating the current (and desired) population characteristics only. This allowed for responding to the current search state and status very efficiently.

3. THE RESULTS

The experiments presented in [2] involved investigating the PCA^2MA performance for real-life, benchmark, and artificial datasets¹, and were divided into three groups: (i) the sensitivity analysis, the comparison with (ii) our previous evolutionary algorithms, and (iii) state-of-the-art techniques.

The results of the sensitivity analysis (in which various PCA^2MA variants were examined) revealed that the analysis of the T geometry (in the pre-processing step) is pivotal, and helps extract important training set vectors. Utilizing them allowed for improving the convergence capabilities of PCA^2MA . The new adaptation was shown much more efficient than our previous technique [1]. The PCA^2MA variant with both enhancements (pre-processing and the new adaptation) notably outperformed the other algorithm variants—the retrieved refined sets were of a higher quality.

The comparison with other evolutionary techniques proved that PCA^2MA is the best evolutionary algorithm for selection of T''s. SVMs trained with refined sets delivered by PCA^2MA gave the largest AUC values for the validation set. At the same time, the number of SVs was kept very small without affecting the classifier performance. It indicates that PCA^2MA selects most important T vectors, and refrains from adding other vectors unnecessarily. Importantly, the adaptive memetic algorithm was able to effectively balance the exploitation of smaller refined sets with the exploration of the larger ones. This balance is crucial in the case of challenging datasets, in which exploiting small T''s may not boost the classification performance due to a too small number of SVs. Finally, PCA^2MA appeared quite stable and led to very similar results in numerous independent algorithm runs for the investigated datasets.

The performance of PCA²MA was also compared with the selected state-of-the-art techniques, belonging to various groups of algorithms (randomized approaches and the one analyzing the geometrical properties of T). Additionally, SVMs were trained using the entire T's too. In most cases, the average AUC values obtained using SVMs trained with refined sets (the AUC values were calculated for the validation sets) retrieved using PCA²MA were larger than the best AUC values elaborated using other techniques. The sets obtained using PCA²MA contained significantly less samples compared with other approaches. It allowed for decreasing the number of SVs without lowering the classification score. We analyzed the convergence time of the proposed algorithm (i.e., the time after which the fitness of the best individual in the population could not be further improved) and showed, that it could be safely terminated faster. It is not possible in "one-pass" algorithms, which investigate the entire set Tin order to extract its valuable subset. The results proved the stability of PCA²MA. It is worth mentioning that the SVM training appeared impossible for massively-large reallife dataset (when the entire T was fed into the training procedure) due to its memory requirements.

Finally, the refined sets (along with the selected SVs) were visualized for all artificial sets. The investigation revealed that the refined sets obtained using PCA²MA do not necessarily follow geometrical patterns. Albeit the pre-processing examines the geometry of T, the vectors which were not annotated as CVs in this process can still be accessed and appended to individuals during the evolution. In many cases, it turned out that utilizing these vectors (which appeared not important at the first glance) allowed for boosting the SVM classification performance. This would not be possible for other methods, in which the rejected T vectors cannot be later re-analyzed and added to T''s.

4. CONCLUSIONS

In our recent paper [2], we proposed an adaptive memetic algorithm for selection of refined training sets for SVMs. This approach enhanced the evolutionary optimization of the refined sets using additional procedures exploiting the knowledge concerning the T vectors extracted before the execution (in the pre-processing step), and attained dynamically during the search. The crucial PCA²MA parameters were dynamically controlled on the fly to adapt to the search progress using the parameter-less adaptation schemes. The extensive experimental study revealed that PCA²MA is superb compared with other evolutionary techniques, and several algorithms from the literature. It allowed for (i) retrieving very high-quality T''s, (ii) decreasing s, and (iii) improving the process of selecting T''s. Finally, we showed that smaller training sets can be notably better than the larger ones, if they exclude mislabeled and noisy vectors from T.

5. REFERENCES

- J. Nalepa and M. Kawulok. A memetic algorithm to select training data for support vector machines. In *Proceedings of GECCO* '14, pages 573–580. ACM, 2014.
- J. Nalepa and M. Kawulok. Adaptive memetic algorithm enhanced with data geometry analysis to select training data for SVMs. *Neurocomputing*, 185:113 – 132, 2016.

¹The real-life and artificial datasets can be found at: http://sun.aei.polsl.pl/~jnalepa/SVM