

Instance Pruning with the SBL-PM-M-EKP Training Data Reducer

Karol Grudziński

E-mail: karol.grudzinski@wp.pl

Department of Physics, Kazimierz Wielki University
Plac Weysenhoffa 11, 85-072 Bydgoszcz, Poland
and

Institute of Applied Informatics
University of Economy
Garbary 2, 85-229 Bydgoszcz, Poland
Affiliation: Kazimierz Wielki University

Abstract. A prototype (case-based) way of data explanation is a powerful method for data analysis and understanding. Interesting instance vectors (prototypes) are usually generated by a training set pruning with various partial memory learners. This approach is the alternative to the rule induction techniques for knowledge discovery and understanding. In this paper a completely new system, SBL-PM-M-EKP is introduced. The study of suitability of SBL-PM-M-EKP for training data compression has been studied on ? datasets. As an underlying classifier we have chosen the well known IB1 system from the WEKA package. We compare the generalization ability of our system to the performance of IB1 trained on the entire training data. The results indicate that with only one prototype per class which was generated by SBL-PM-EKP system, on ? datasets we have obtained statistically indistinguishable results from those coming from IB1 or even the generalization ability has been improved by our system over the IB1 one in several cases.

1 Introduction.

Data mining is commonly used in many domains. A case-based way of data explanation is very popular among researchers. Such an approach to knowledge discovery and understanding is particularly often employed in medicine, where a medical doctor makes a diagnosis by referring to other similar cases in a database of patients.

Interesting instance vectors (prototypes) are usually generated by a training set pruning with various partial memory learners. The term ‘Partial Memory Learning’ (PML) is most often reserved for on-line learning systems that select and store a portion of the past learning examples. In this paper we stick to this naming convention (i.e. PML), but this methodology is called also ‘instance selection’, ‘training data compression, reduction or pruning’. The idea behind this machine learning paradigm is that only a small fraction of a usually much

larger, original training set is used for a final classification of unseen samples. [1–8].

The acronym SBL-PM-M-EKP is short for **S**imilarity-**B**ased-**L**earner-**P**artial-**M**emory-**M**inimization-**E**xactly-**k**-**P**rototypes. We however want to stress here that our new system is completely different from our earlier model, SBL-PM-M.

References

1. Maloof, M. A., Michalski, R. S.: AQ-PM: A System for Partial Memory Learning. Intelligent Information Systems, Ustroń, Poland (1999) 70-79
2. Maloof, M., Michalski, R. S.: Selecting Examples for Partial Memory Learning. Machine Learning, **41**, (2000) 27-52
3. Maloof, M., Michalski, R. S.: Incremental Learning with Partial Instance Memory. Proceedings of the Thirteenth International Symposium on Methodologies for Intelligent Systems, Lyon, France, (2002), In Foundations of Intelligent Systems, Lecture Notes in Artificial Intelligence, (2366), Berlin:Springer-Verlag, 16-27
4. Wilson, D. R., Martinez, T. R.: Instance Pruning Techniques. In Fisher, D.: Machine Learning: Proceedings of the Fourteenth International Conference. Morgan Kaufmann Publishers, San Francisco, CA, (1997), 404-417.
5. Wilson, D. R., Martinez, T. R.: Reduction Techniques for Instance-Based Learning Algorithms. Machine Learning, **38**, (2000), 257-286
6. Grochowski, M.: Selecting Reference Vectors in Selected Methods for Classification. MSc. Thesis, Nicolaus Copernicus University, Department of Applied Informatics, Toruń, Poland, (2003) (In Polish)
7. Jankowski N., Grochowski, M.: Comparison of Instances Selection Algorithms I: Algorithms Survey. Artificial Intelligence and Soft Computing ICAISC 2004, in Lecture Notes in Artificial Intelligence (LNAI 3070), (Springer), 598-603.
8. Grochowski, M., Jankowski N.: Comparison of Instance Selection Algorithms II: Results and Comments. Artificial Intelligence and Soft Computing ICAISC 2004, in Lecture Notes in Artificial Intelligence (LNAI 3070), (Springer), 580-585.
9. Grudziński, K.: SBL-PM-M: A System for Partial Memory Learning. Artificial Intelligence and Soft Computing ICAISC 2004, in Lecture Notes in Artificial Intelligence (LNAI 3070), (Springer), 586-591.
10. Nelder, J. A., Mead, R.: A simplex method for function minimization. Computer Journal **7** (1965), 308-313
11. Grudziński, K., Duch, W.: SBL-PM: A Simple Algorithm for Selection of Reference Instances for Similarity-Based Methods. Intelligent Information Systems, Bystra, Poland (2000), in Advances in Soft Computing, Physica-Verlag (Springer), 99-108
12. Ortega J., Koppel M., Argamon S.: Arbitrating Among Competing Classifiers Using Learned Referees. Knowledge and Information Systems 3 (2001), 470-490.
13. Bauer E., Kohavi R.: An empirical comparison of voting classification algorithms: bagging, boosting and variants. Machine Learning 36 (1999), 105-142
14. Duch, W.: Similarity Based Methods: a general framework for classification, approximation and association. Control and Cybernetics 29 **4**, (2000), 1-30
15. SBL, Similarity Based Learner, Software Developed by Karol Grudziński. Nicolaus Copernicus University: 1997-2002, Academy of Bydgoszcz: 2002-2006, University of Economy: 2005-2006.
16. Mertz, C. J., Murphy, P. M.: UCI repository of machine learning databases. <http://www.ics.uci.edu/pub/machine-learning-data-bases>.

17. Dudani S. A.: The distance-weighted k -nearest-neighbor-rule. IEEE Transactions on Systems, Man and Cybernetics 6 (4) 1975, 325-327.

Tables

Table 1. Results for the 10-fold CV Test on the Selected Datasets

Dataset	IB1%	std. dev %	SBL-PM-M-EKP:IB1 %	std. dev %
Breast-Cancer	68.6	7.5	73.6	6.4
Breast-Wisc.	95.7	2.4	95.5	2.5
Credit-Rating	81.6	4.6	79.6	6.1
Heart-Cleveland	76.1	6.8	80.3	7.2
Heart-Hungarian	78.3	7.5	82.5	6.5
Heart-Statlog	76.1	6.8	80.3	7.2
Hepatitis	81.4	8.6	80.7	9.2
Pima-Diabetes	70.6	4.7	70.4	5.9