

SOAP: Efficient Feature Selection of Numeric Attributes

Roberto Ruiz, Jesús S. Aguilar-Ruiz, and José C. Riquelme

Department of Computer Science. University de Seville.
Avda. Reina Mercedes S/n. 41012 Sevilla, Spain.
{rruiz,aguilar,riquelme}@lsi.us.es

Abstract. The attribute selection techniques for supervised learning, used in the preprocessing phase to emphasize the most relevant attributes, allow making models of classification simpler and easy to understand. Depending on the method to apply: starting point, search organization, evaluation strategy, and the stopping criterion, there is an added cost to the classification algorithm that we are going to use, that normally will be compensated, in greater or smaller extent, by the attribute reduction in the classification model. The algorithm (SOAP: Selection of Attributes by Projection) has some interesting characteristics: lower computational cost ($O(mn \log n)$ m attributes and n examples in the data set) with respect to other typical algorithms due to the absence of distance and statistical calculations; with no need for transformation. The performance of SOAP is analysed in two ways: percentage of reduction and classification. SOAP has been compared to CFS [6] and ReliefF [11]. The results are generated by C4.5 and 1NN before and after the application of the algorithms.

1 Introduction

The data mining researchers, especially those dedicated to the study of algorithms that produce knowledge in some of the usual representations (decision lists, decision trees, association rules, etc.), usually make their tests on standard and accessible databases (most of them of small size). The purpose is to independently verify and validate the results of their algorithms. Nevertheless, these algorithms are modified to solve specific problems, for example real databases that contain much more information (number of examples) than standard databases used in training. To accomplish the final tests on these real databases with tens of attributes and thousands of examples is a task that takes a lot of time and memory size.

It is advisable to apply to the database preprocessing techniques to reduce the number of attributes or the number of examples in such a way as to decrease the computational time cost. These preprocessing techniques are fundamentally oriented to either of the next goals: feature selection (eliminating non-relevant attributes) and editing (reduction of the number of examples by eliminating some of them or calculating prototypes [1]). Our algorithm belongs to the first group.

In this paper we present a new method of attribute selection, called SOAP (Selection of Attributes by Projection), which has some important characteristics:

- Considerable reduction of the number of attributes.
- Lower computational time $O(mn \log n)$ than other algorithms.

- Absence of distance and statistical calculations: correlation, information gain, etc.
- Conservation of the error rates of the classification systems.

The hypothesis on which the heuristic is based is: "place the best attributes with the smallest number of label changes". The next section discusses related work. Section 3 describes the SOAP algorithm. Section 4 presents the results. Which deal with several databases from the UCI repository [4]. The last section summarises the findings.

2 Related Work

Several authors defined the feature selection by looking at it from various angles depending on the characteristic that we want to accentuate. In general, attribute selection algorithms perform a search through the space of feature subsets, and must address four basic issues affecting the nature of the search: 1) Starting point: forward and backward, according to whether it began with no features or with all features. 2) Search organization: exhaustive or heuristic search. 3) Evaluation strategy: wrapper or filter. 4) Stopping criterion: a feature selector must decide when to stop searching through the space of feature subsets. A predefined number of features are selected, a predefined number of iterations reached. Whether or not the addition or deletion of any feature produces a better subset, we also stop the search, if an optimal subset according to some evaluation function is obtained.

Algorithms that perform feature selection as a preprocessing step prior to learning can generally be placed into one of two broad categories: wrappers, Kohavi [9], which employs a statistical re-sampling technique (such as cross validation) using the actual target learning algorithm to estimate the accuracy of feature subsets. This approach has proved to be useful but is very slow to execute because the learning algorithm is called upon repeatedly. Another option called filter, operates independently of any learning algorithm. Undesirable features are filtered out of the data before induction begins. Filters use heuristics based on general the characteristics of the data to evaluate the merit of feature subsets. As a consequence, filter methods are generally much faster than wrapper methods, and, as such, are more practical for use on data of high dimensionality. FOCUS [3], LVF [18] use class consistency as an evaluation meter. One method for discretization called Chi2 [17]. Relief [8] works by randomly sampling an instance from the data, and then locating its nearest neighbour from the same and opposite class. Relief was originally defined for two-class problems and was later expanded as ReliefF [11] to handle noise and multi-class data sets, and RReliefF [16] handles regression problems. Other authors suggest Neuronal Networks for an attribute selector [19]. In addition, learning procedures can be used to select attributes, like ID3 [14], FRINGE [13] and C4.5 [15]. Methods based on the correlation like CFS [6], etc.

3 SOAP: Selection of Attributes by Projection

3.1 Description

To describe the algorithm we will use the well-known data set IRIS, because of the easy interpretation of their two-dimensional projections.

Three projections of IRIS have been made in two-dimensional graphs. In Fig. 1 it is possible to observe that if the projection of the examples is made on the abscissas or ordinate axis we can not obtain intervals where any class is a majority, only can be seen the intervals [4.3,4.8] of Sepallength for the Setosa class or [7.1,8.0] for Virginica. In Fig. 2 for the Sepalwidth parameter in the ordinate axis clear intervals are not appraised either. Nevertheless, for the Petalwidth attribute is possible to appreciate some intervals where the class is unique: [0,0.6] for Setosa, [1.0,1.3] for Versicolor and [1.8,2.5] for Virginica. Finally in Fig. 3, it is possible to appreciate the class divisions, which are almost clear in both attributes. This is because when projecting the examples on each attribute the number of label changes is minimum. For example, it is possible to verify that for Petalength the first label change takes place for value 3 (setosa to Versicolor), the second in 4.5 (Versicolor to Virginica), there are other changes later in 4.8, 4.9, 5.0 and the last one is in 5.1.

SOAP is based on this principle: to count the label changes, produced when crossing the projections of each example in each dimension. If the attributes are in ascending order according to the number of label changes, we will have a list that defines the priority of selection, from greater to smaller importance. SOAP presumes to eliminate the basic redundancy between attributes, that is to say, the attributes with interdependence have been eliminated. Finally, to choose the more advisable number of features, we define a reduction factor, RF, in order to take the subset from attributes formed by the first of the aforementioned list.

Before formally exposing the algorithm, we will explain with more details the main idea. We considered the situation depicted in Fig. 2: the projection of the examples on the abscissas axis produces a ordered sequence of intervals (some of them can be a single point) which have assigned a single label or a set of them: {[0,0.6] Se, [1.0,1.3] Ve, [1.4,1.4] Ve-Vi, [1.5,1.5] Ve-Vi, [1.6,1.6] Ve-Vi, [1.7,1.7] Ve-Vi, [1.8,1.8] Ve-Vi, [1.9,2.5] Vi}. If we apply the same idea with the projection on the ordinate axis, we calculate the partitions of the ordered sequences: {Ve, R, R, Ve, R, R, R, R, R, R, R, R, R, Se, R, Se, R, Se}, where R is a combination of two or three labels. We can observe that we obtain almost one subsequence of the same value with different classes for each value from the ordered projection. That is to say, projections on the ordinate axis provide much less information than on the abscissas axis.

In the intervals with multiple labels we will consider the worst case, that being the maximum number of label changes possible for a same value.

The number of label changes obtained by the algorithm in the projection of each dimension is: Petalwidth 16, Petalength 19, Sepallenth 87 and Sepalwidth 120. In this way, we can achieve a ranking with the best attributes from the point of view of the classification. This result agrees with what is common knowledge in data mining, which states that the width and length of petals are more important than those related to sepals.

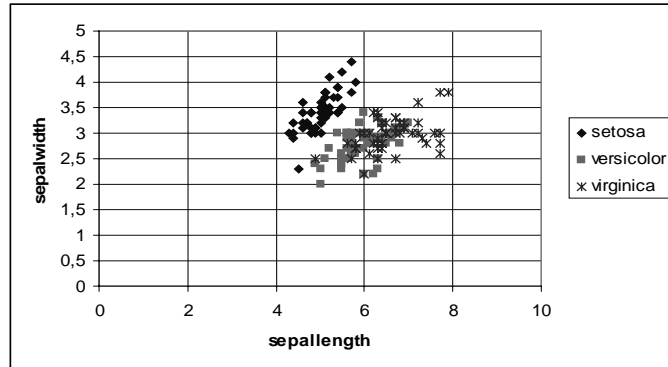


Fig. 1. Two-dimensional representation

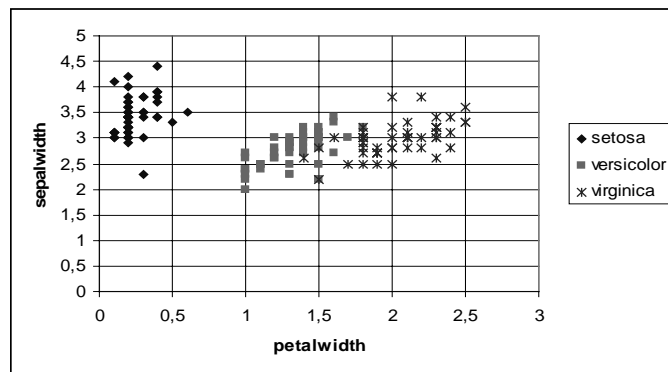


Fig. 2. Two-dimensional representation

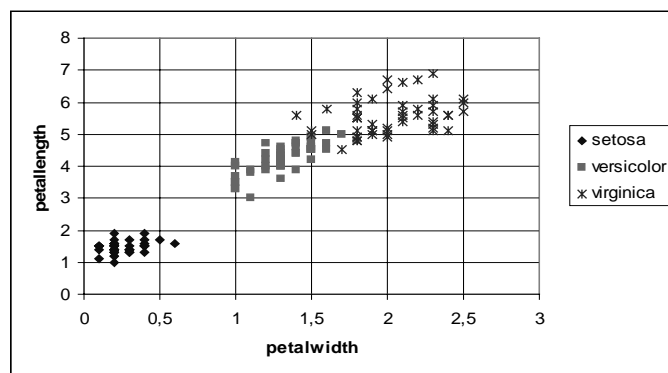


Fig. 3. Two-dimensional representation

3.2 Definitions

Definition 1: Let the attribute A_i be a continuous variable that takes values in $I_i=[\min_i, \max_i]$. Then, A is the attributes space defined as $A=I_1 \times I_2 \times \dots \times I_m$, where m is the number of attributes.

Definition 2: An example $e \in E$ is a tuple formed by the Cartesian product of the value sets of each attribute and the set C of labels. We define the operations att and lab to access the attribute and its label (or class): $att: E \times N \rightarrow A$ and $lab: E \rightarrow C$, where N is the set of natural numbers.

Definition 3: Let the universe U be a sequence of example from E . We will say that a database with n examples, each of them with m attributes and one class, forms a particular universe. Then $U=\langle u[1], \dots, u[n] \rangle$ and as the database is a séquence, the access to an example is achieved by means of its position. Likewise, the access to j -th attribute of the i -th example is made by $att(u[i], j)$, and for identifying its label $lab(u[i])$.

Definition 4: An ordered projected sequence is a sequence formed by the projection of the universe onto the i -th attribute. This sequence is sorted out in ascending order.

Definition 5: A partition in subsequences is the set of subsequences formed from the ordered projected sequence of an attribute in such a way as to maintain the projection order. All the examples belonging to a subsequence have the same class and every two consecutive subsequences are disjointed with respect to the class. Henceforth, a subsequence will be called a partition.

Definition 6: A subsequence of the same value is the sequence composed of the examples with identical value from the i -th attribute within the ordered projected sequence.

3.3 Algorithm

The algorithm is very simple and fast, see Fig. 4. It operates with continuous variables as well as with databases which have two classes or multiple classes. In the ascending-order-task for each attribute, the QuickSort algorithm is used [7]. This algorithm is $O(n \log n)$, on average. Once ordered by an attribute, we can count the label changes throughout the ordered projected sequence. NumberChanges in Fig. 5, considers whether we deal with different values from an attribute, or with a subsequence of the same value. In the first case, it compares the present label with that of the following value. Whereas in the second case, where the subsequence is of the same value, it counts as many label changes as are possible (function ChangesSameValue).

The k first attribute which NCE (number of label changes) under NCE_{lim} will be selected. NCE_{lim} is calculated applying the follow equation:

$$NCE_{lim} = NCE_{min} + (NCE_{max} - NCE_{min}) * RF \quad (1)$$

RF: reduction factor.

```

Input: E training (n examples, m attributes)
Output: E reduced (n examples, k attributes (k<=m))
  For each attribute  $a_i$  with  $i$  in  $\{1..m\}$ 
     $E_i \leftarrow \text{QuickSort}(E_i, a_i)$ 
     $\text{NCE}_i \leftarrow \text{NumberChanges}(E_i, a_i)$ 
  NCE Attribute Ranking
  Select the k first

```

Fig. 4. SOAP algorithm

```

Input: E training (n examples, m attributes)
Output: number of label changes
  For each example  $e_j \in E$  with  $j$  in  $\{1..n\}$ 
    If  $\text{att}(u[j], i) \in \text{Subsequence same value}$ 
       $\text{labelChanges} += \text{ChangesSameValue}()$ 
    Else
      If  $\text{lab}(u[j]) \neq \text{lab}(u[j+1])$ 
         $\text{labelChanges}++$ 

```

Fig. 5. NumberChanges algorithm

4 Experiments

In order to compare the effectiveness of SOAP as a feature selector for common machine learning algorithms, experiments were performed using twelve standard data sets from the UCI collection [4]. The data sets and their characteristics are summarized in Table 3. The percentage of correct classification with C4.5 and 1NN, averaged over ten ten-fold cross-validation runs, were calculated for each algorithm-data set combination before and after feature selection by SOAP (RF 0.35 and 0.25), CFS and ReliefF (threshold 0.05). For each train-test split, the dimensionality was reduced by each feature selector before being passed to the learning algorithms. The same fold were used for each feature selector-learning scheme combination.

To perform the experiment with CFS and ReliefF we used the Weka¹ (Waikato Environment for Knowledge Analysis) implementation.

Table 1 shows the results for attribute selection with C4.5 and compares the size (number of nodes) of the trees produced by each attribute selection scheme against the size of the trees produced by C4.5 with no attribute selection. Smaller trees are preferred as they are easier to interpret, but accuracy is generally degraded. The table shows how often each method performs significantly better (denoted by \circ) or worse (denoted by \bullet) than when performing no feature selection (column 2 and 3). Throughout we speak of results being significantly different if the difference is statistically at the 5% level according to a paired two-sided t test. Each pair of points

¹ <http://www.cs.waikato.ac.nz/~ml>

consisting of the estimates obtained in one of the ten, ten-fold cross-validation runs, for before and after feature selection. For SOAP, feature selection degrades performance on four datasets, improves on one and it is equal on seven. The results are similar to ReliefF and a little worse than those provided by CFS. From this table it can be seen that SOAP produces the smallest trees, it improves C4.5's performances on nine data sets and degrades it on one.

Tables 1 and 2 show the average for two execution of SOAP (RF 0.35 and 0.25, equation 1).

Table 1. Result of attribute selection with C4.5. Accuracy and size of trees. ○,● Statistically significant improvement or degradation (p=0.05).

Data Set	Original		SOAP		CFS		RLF	
	Ac.	Size	Ac.	Size	Ac.	Size	Ac.	Size
balance-scale	78,18	81,08	57,94 ●	6,28 ○	78,18	81,08	78,29	81,54
breast-w	95,01	24,96	94,84	21,62 ○	95,02	24,68	95,02	24,68
diabetes	74,64	42,06	74,34	8,56 ○	74,36	14,68 ○	65,10 ●	1,00 ○
glass	68,18	46,34	66,78	46,10	69,35	40,90 ○	68,97	30,32 ○
glass2	78,71	24,00	78,90	16,32 ○	79,82	14,06 ○	53,50 ●	1,70 ○
heart-stat	78,11	34,58	79,56	28,20 ○	80,63 ○	23,84 ○	82,33 ○	14,78 ○
ionosphere	89,83	26,36	90,06	22,52 ○	90,26	23,38 ○	89,91	22,72 ○
iris	94,27	8,18	94,40	8,12	94,13	7,98	94,40	8,16
segment	96,94	80,98	90,94 ●	110,68 ●	96,35 ●	73,92 ○	96,93	80,66
sonar	74,28	27,98	70,72 ●	13,18 ○	74,38	28,18	70,19 ●	9,74 ○
vehicle	71,83	139,34	52,84 ●	22,26 ○	66,42 ●	106,60 ○	66,22 ●	137,42
waveform	75,36	592,92	77,47 ○	485,26 ○	77,18 ○	513,78 ○	75,51	217,72 ○
Average (35)	81	94	77	66	81	79	78	53
Average (25)			77	59				

Table 2. Result of attribute selection with 1NN. Average number of features selected, the percentage of the original features retained and the accuracy. ○,● Statistically significant improvement or degradation (p=0.05).

Data Set	Original			SOAP			CFS			RLF		
	Atts	Ac.		Atts	%	Ac.	Atts	%	Ac.	Atts	%	Ac.
balance-scale	4	86,56	1,39	35	57,98 ●	4,00	100	86,56	4,00	100	86,56	
breast-w	9	95,25	6,00	67	94,16 ●	8,97	100	95,24	8,05	89	95,35	
diabetes	8	70,35	2,99	37	70,16	3,11	39	70,07	0,00	0	34,90 ●	
glass	9	70,28	3,94	44	73,04 ○	6,30	70	74,25 ○	3,39	38	63,83 ●	
glass2	9	77,79	4,72	52	80,37	3,95	44	83,07 ○	0,32	4	54,29 ●	
heart-stat	13	75,59	7,11	55	77,74	6,26	48	78,37 ○	6,27	48	78,89 ○	
ionosphere	34	86,78	31,55	93	87,07	12,30	36	89,72 ○	30,88	91	87,49	
iris	4	95,27	2,00	50	96,33	1,93	48	95,60	4,00	100	95,27	
segment	19	97,13	7,00	37	91,29 ●	5,66	30	97,00	15,04	79	97,19	
sonar	60	84,47	5,42	9	70,63 ●	17,84	30	83,56	3,89	6	68,61 ●	
vehicle	18	69,48	1,09	6	46,50 ●	7,45	41	62,86 ●	5,81	32	61,28 ●	
waveform	40	73,59	12,99	32	79,33 ○	14,85	37	79,13 ○	5,77	14	73,09	
Average (35)	19	82	7	43	77	8	52	83	7	50	75	
Average (25)			6	35	75							

Table 2 shows the average number of features selected, the percentage of the original features retained and the accuracy of INN. SOAP is a specially selective algorithm compared with CFS and RLF. If SOAP and CFS are compared, only in one dataset (ionosphere) is the number of characteristics significantly greater than those selected by CFS. In five data sets there are no significant differences, and in six, the number of features is significantly smaller than CFS. Compare to RLF, only in glass2 and diabetes, SOAP obtains more parameters in the reduction process (threshold 0.05 is not sufficient). It can be seen (by looking at the fifth column) that SOAP retained 43% (35%) of the attributes on average. Figure 6 shows the average number of feature selected by SOAP, CFS and ReliefF as well as the number present in the full data set.

Table 3. Discrete class data sets with numeric attributes. Time in milliseconds.

Data Set	Original			SOAP t-ms	CFS t-ms	RLF t-ms
	Instances	Atts	Classes			
1 balance-scale	625	4	3	10	17455	561
2 breast-cancer	699	9	2	10	40	1322
3 diabetes	768	8	2	10	30	1422
4 glass	214	9	7	0	20	160
5 glass2	163	9	2	0	10	80
6 heart-statlog	270	13	2	10	10	281
7 ionosphere	351	34	2	10	120	1202
8 iris	150	4	3	0	10	40
9 segment	2310	19	7	40	521	29362
10 sonar	208	60	2	10	100	771
11 vehicle	846	18	4	10	70	3956
12 waveform	5000	40	3	210	2434	282366
Sum				320	20820	321523

It is interesting to compare the speed of the attribute selection techniques. We measured the time taken in milliseconds² to select the final subset of attributes. SOAP is an algorithm with a very short computation time. The results shown in Table 3 confirm the expectations. SOAP takes 320 milliseconds in reducing 12 datasets whereas CFS takes more than 20 seconds and RLF almost 6 minutes. In general, SOAP is faster than the other methods and it is independent of the classes number, a factor that excessively affects CFS, as it is possible to observe in the set “segment” with seven classes. Also it is possible to be observed that ReliefF is affected very negatively by the number of instances in the dataset, it can be seen in “segment” and “waveform”. Eventhough these two datasets were eliminated, SOAP is more than 200 times faster than CFS, and more than 100 times than ReliefF.

5 Conclusions

In this paper we present a deterministic attribute selection algorithm. It is a very efficient and simple method used in the preprocessing phase. A considerable reduction of the number of attributes is produced in comparison to other techniques. It does not

² This is a rough measure. Obtaining true cpu time from within a Java program is quite difficult.

need distance nor statistical calculations, which could be very costly in time (correlation, gain of information, etc.). The computational cost is lower than other methods $O(m \cdot n \cdot \log n)$.

In later works, we will focus our research on the selection of the subset of attributes once they have been obtained. Finally we will try to adapt SOAP to databases with discrete attributes where redundant features have not been eliminated.

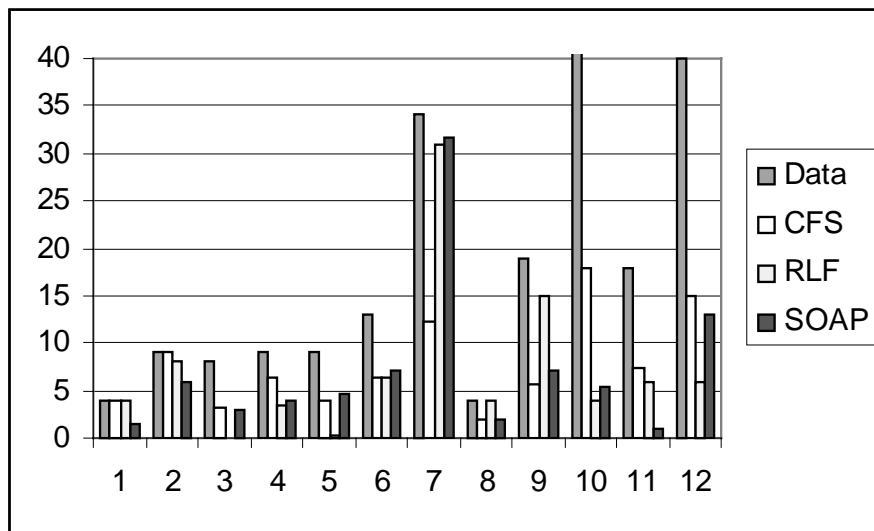


Fig. 6. Average number of feature selected.

Acknowledgments. This work has been supported by the Spanish Research Agency CICYT under grant TIC2001-1143-C03-02.

References

- [1] Aguilar-Ruiz, Jesús S., Riquelme, José C. and Toro, Miguel. Data Set Editing by Ordered Projection. *Intelligent Data Analysis Journal*. Vol. 5, nº5, pp. 1-13, IOS Press (2001).
- [2] Almuallim, H. and Dietterich, T.G. Learning with many irrelevant features. In *Proceedings of the Ninth National Conference on Artificial Intelligence*. pp. 547-552. AAAI Press (1991).
- [3] Almuallim, H. and Dietterich, T.G. Learning boolean concepts in the presence of many irrelevant features. *Artificial Intelligence*, 69(1-2):279-305 (1994).
- [4] Blake, C. and Merz, E. K. *UCI Repository of machine learning databases* (1998).
- [5] Brassard, G., and Bratley, P. *Fundamentals of algorithms*. Prentice Hall, New Jersey (1996).
- [6] Hall M.A. *Correlation-based feature selection for machine learning*. PhD thesis, Department of Computer Science, University of Waikato, Hamilton, New Zealand (1998).
- [7] Hoare, C. A. R. QuickSort. *Computer Journal*, 5(1):10-15 (1962).

- [8] Kira, K. and Rendell, L. A practical approach to feature selection. In Proceedings of the Ninth International Conference on Machine Learning. pp. 249-256, Morgan Kaufmann (1992).
- [9] Kohavi, R. and John, G. H. Wrappers for feature subset selection. *Artificial Intelligence*, 97, 273-324 (1997).
- [10] Koller, D. and Sahami, M. Toward optimal feature selection. In Proceedings of the Thirteenth International Conference on Machine Learning. pp. 284-292, Morgan Kaufmann (1996).
- [11] Kononenko, I. Estimating attributes: Analysis and extensions of relief. In Proceedings of the Seventh European Conference on Machine Learning. pp. 171-182, Springer-Verlag (1994).
- [12] Modrzejewski, M. Feature selection using rough sets theory. In Proceedings of the European Conference on Machine Learning. pp. 213-226, Springer (1993).
- [13] Pagallo, G. and Haussler, D. Boolean feature discovery in empirical learning. *Machine Learning*, 5, 71-99 (1990).
- [14] Quinlan, J. Induction of decision trees. *Machine Learning*, 1(1), 81-106 (1986).
- [15] Quinlan, J. C4.5: Programs for machine learning. Morgan Kaufmann (1993).
- [16] Robnik-Šikonja, M. and Kononenko, I. An adaptation of relief for attribute estimation in regression. In Proceedings of the Fourteenth International Conference on Machine Learning. pp. 296-304, Morgan Kaufmann (1997).
- [17] Setiono, R., and Liu, H. Chi2: Feature selection and discretization of numeric attributes. In Proceedings of the Seventh IEEE International Conference on Tools with Artificial Intelligence (1995).
- [18] Setiono, R., and Liu, H. A probabilistic approach to feature selection—a filter solution. In Proceedings of International Conference on Machine Learning, 319-327 (1996).
- [19] Setiono, R., and Liu, H. Neural network feature selectors. *IEEE Trans. On Neural Networks*, 8(3), 654-662 (1997).