

Dimensionality reduction in a database related with viticulture crops using wrapper techniques

R. Fernandez-Martinez, J. Fernandez-Ceniceros, A. Sanz-Garcia, R. Lostado-Lorza, F.J. Martinez-De-Pison-Ascacibar

Abstract— Recent advance in environmental monitoring technologies allows that every day, major amount of agricultural productions have a support to be controlled better. Coverall, thanks to manufacturing advances of new sensors, which allow realizing acquisition of physical variables with almost no limitation. This entails the existence of great amount of stored data, distributed in different variables that make really complicated work with them. In these circumstances, the problem arises at the time of building models when it works with a large number of variables. In order to solve it, feature selection methods are used to reduce this large number, improving building, training and validation models processes based on machine learning techniques. The methods used due to their satisfactory results, in the practical case of several viticulture crops, have been wrappers methods.

Keywords— feature selection, ripening grape berries, viticulture crop, wrapper methods.

I. INTRODUCTION

AGRICULTURAL monitoring systems gather crowds of physical variables during the growing and ripening season of crops, especially now that new technologies give the possibility to have sensors that can obtain data from almost any physical variable in order to be measured [13] [25]. This means that it is possible to perform a collection of information about what is happening at all development stages. And thanks to a forthcoming study of collected data, it is possible to extract previously unknown knowledge to help to improve the process and get a better crop management based in

Manuscript received January 31, 2011; Revised version received March 8, 2011. (Write the dates on which you submitted your paper for review as well as the revised version). This work was supported in part by La Rioja University through FPI fellowships and the Autonomous Government of La Rioja under Grant FOMENTA 2010/13.

R. Fernandez-Martinez is with EDMANS Group, Mechanical Engineering Department at La Rioja University, Logroño, La Rioja, 26004, Spain. (roberto.fernandez@unirioja.es).

J. Fernandez-Ceniceros is with EDMANS Group, Mechanical Engineering Department at La Rioja University, Logroño, La Rioja, 26004, Spain. (julio.fernandezc@unirioja.es).

A. Sanz-Garcia is with EDMANS Group, Mechanical Engineering Department at La Rioja University, Logroño, La Rioja, 26004, Spain. (andres.sanz@unirioja.es).

R. Lostado-Lorza is with EDMANS Group, Mechanical Engineering Department at La Rioja University, Logroño, La Rioja, 26004, Spain. (ruben.lostado@unirioja.es).

F. J. Martinez-De-Pison-Ascacibar is with EDMANS Group, Mechanical Engineering Department at La Rioja University, Logroño, La Rioja, 26004, Spain. (fjmartin@unirioja.es).

decision support systems [4].

With all this amount of data, classified by variables, it is performed a study trying to obtain information to assist farmers in order to improve the process [31]. These studies are conducted with data mining techniques that turn data into information that helps to control crop [5] [32] [35]. All these techniques are one of the research areas where more progress is being made, although it should be noted that many of data collected may not have useful information for this study purpose. And when the amount of data is extremely high is very complicated to detect which is the data that do not contain useful information and which does [36].

To solve this problem it is worked on algorithms that perform a selection of the most interesting variables for the selected goal [1] [3] [26] [27]. The gathered data can be studied with these algorithms and determine which variables have the most significant information and which provide little or nothing.

This kind of algorithms is used in the monitoring of a crop over several productive seasons, to get the variables that have some influence on the cultivation. A large number of variables are associated with weather conditions, as they have significant influence on most crops. Although not all variables that can be collected by weather stations, which are now on the market, are useful to show a correlation with the crop evolution. So they are selected variables that can actually have meaningful correlations between measures from weather stations and information collected by farmers about the crop.

Using learning algorithms is possible to work with several methods that can be used to relate all this information, but when the amount of information is high, worse conclusions are obtained in many times. Besides computational cost is increased according the amount of variables and data, that is involved in the process.

To improve the efficiency of these methods is carried out a previous step with data. It is made a feature selection that allows reducing the execution time and simplifying the calculation process when applying learning algorithms [24]. It is developed a reduction of variables and data gathered from weather stations, in order to obtain a better correlation of these variables with those obtained in the process of crop maturation.

II. PROBLEM FORMULATION

It is worked in the prediction of some physical and chemical variables (weight, sugar, acidity, ...) to learn how those features evolve during the maturation process in a vineyard and how they are influenced by different environmental parameters. The data used in this study were collected from some different study areas of La Rioja (Spain). Once collected, it is generated all the necessary variables to control this process, according to the recommendations of several authors [6] [14] [15] [29]. Concluding that the needed variables used to predict these features are 29. Furthermore, it is had a great amount of data of each of these variables, because the study is conducted with data collected over 7 years.

Before using feature selection and transformation processes a study of available variables was realized with the

The work done with all these data generates several problems that thanks to the feature selection algorithms can be resolved [23]:

- 1) The results quality is worse, because the initial data can have many variables, and some of them can confuse the algorithms with irrelevant information, that makes the algorithms work wrongly.
- 2) Irrelevant variables and redundant information, that may not provide an increase in the information quality available within of each class, make to increase the learning time used by the algorithms
- 3) The predictive accuracy is much worse since these data make the calculations generating some higher errors.

III. PROBLEM SOLUTION

To solve the problem deriving from working with so many variables there are two methods that allow reducing this quantity. Feature transformation and feature selection (Fig. 3):

- 1) Feature transformation: According to this method, new variables are created from a transformation of the initial group. These transformations are made from linear or nonlinear combinations in order to reduce the dimensionality of the dataset and lose as little information as possible [22].
- 2) Feature selection. However, according to this other method, the variables are not transformed but the most significant of all are chosen to make up the database [7] [11] [17]. Feature selection consists of selecting what type of characteristics or traits are the best suited to describe the variables that it is wanted to predict. In order to do this, it is located the features that affect the problem in a more crucial way.

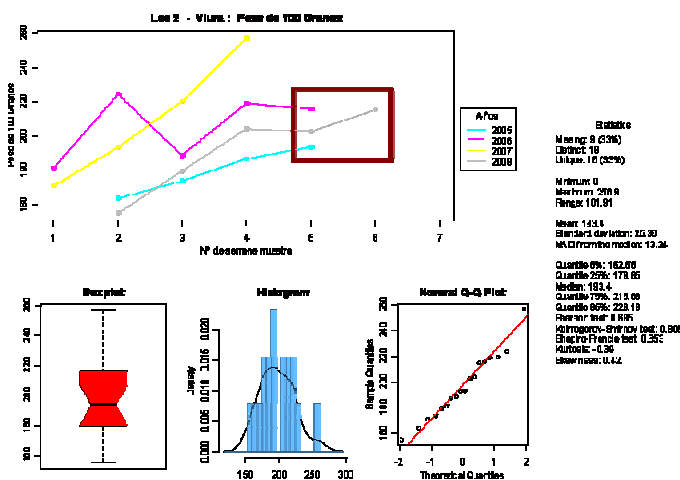


Fig. 1 Example of analysis done before using feature selection and transformation processes

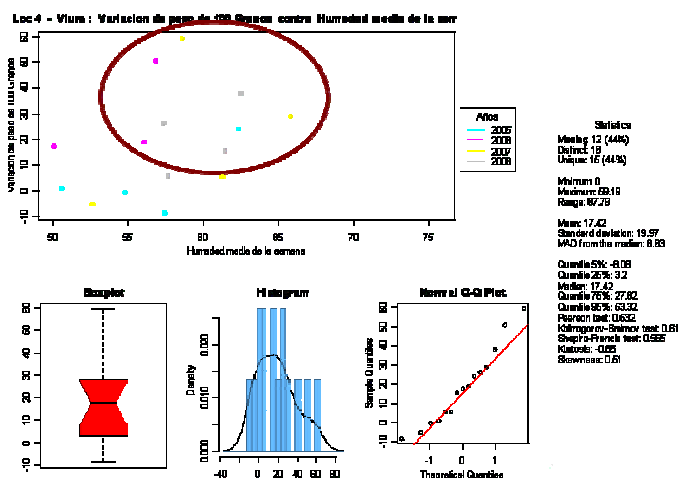


Fig. 2 Example of analysis performed before feature transformation and selection to test the influence of some variables with others.

information stored (Fig. 1 and 2).

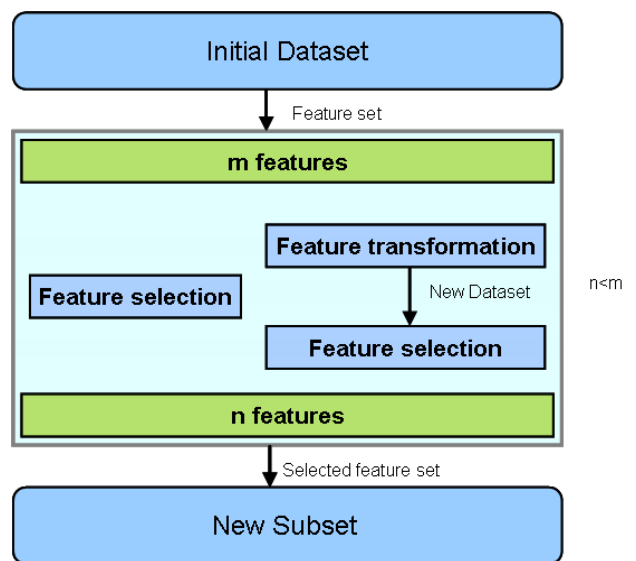


Fig. 3 Feature transformation and feature selection

Of these two methods to reduce variables, it is considered in this case, that the second is more useful for researchers because it does not transform the variables and shows more clearly how each one affects to the final result. The advantages generated by this method are:

- 1) Improved compression of the generated models, because the models are composed of fewer variables.
- 2) Reduction of computational cost when the final model is generated. Especially in more complex models.
- 3) Improved accuracy of final model, because reduction of non-significant variables improves its accuracy.

A. Feature selection methods: filters and wrappers

There are two main groups of feature selection methods. One called filter or indirect approach, and another one called wrapper or direct approach.

B. Filter

These methods make use of heuristic algorithms to determine the optimal subset of features. Where the final solution is not determined directly but it is got making attempts. To obtain it, a group of possible solutions are generated during the process according to a given pattern. These solutions are tested using the criteria that characterize the solution. And of all the solutions generated, invalid solutions are not taken into account [18] [20].

The main advantage of this algorithm is its faster speed in the calculation and the computational cost savings.

C. Wrapper

In this method, each data subset is evaluated by means of learning algorithms. This gives greater accuracy, but also carries a higher computational cost [19].

Decision trees are often used to examine the attributes that are not used or used in fewer rules.

D. Used methods

On many occasions, filters methods are most commonly used for its speed and acceptable results. There are even a greater number of techniques based on these methods. But in this study it is studied wrapper methods of subset selection type, as it is considered they are the ones that produce more accurate results. Also it is considered acceptable that they have a higher duration and a higher computational cost in his generation.

Wrappers methods used to carry out this study have been implemented within the WEKA tool and are:

- 1) ClassifierSubsetEval [34]: Evaluates subsets of training data attributes or a set of independent test, using classifiers. Uses a classifier to estimate the 'merit' of a set of attributes.
- 2) WrapperSubsetEval [19]: Evaluates subsets of attributes using classifiers. Internal cross validation is used to estimate the accuracy of the learning system in each set. And later assign a merit to each of given groups of

variables. It is also developed several trees to evaluate each data set to get merit estimation with less error

Several classifiers are used with different learning strategies, for each of the selected methods:

- 1) Linear Regression (LR) [33]: There are a number of independent variables, which, when taken together, produce a result of a dependent variable. The regression model is then used to predict the result of an unknown dependent variable, given the values of the independent variables. In this case, it uses the Akaike criterion for model selection, and is able to deal with weighted instances.
- 2) Multilayer Perceptron (MLP) [12] [8]: A classifier and predictor that uses backpropagation to classify instances. All nodes in this network are sigmoid, except when the class is numeric. In the latter case, the output nodes become unthresholded linear units. Training is performed with networks that have between 1 and 30 neurons in the hidden layer. In this case the network was developed with five neurons due to it was the one that gets best results.

Table 1. Initial variables group used to put into practice feature selection methods.

Vineyard variables	
Variety	Var
Vineyard age (year)	Age
Altitude (m)	Altit
Environmental variables related to the amount of rainfall	
Total rainfall over the preceding week (mm.)	RFW
Total rainfall over the preceding two weeks (mm)	RF2W
Total rainfall over the preceding three weeks (mm)	RF3W
Total rainfall since the beginning of the year (mm)	RFY
Total rainfall since bud break (mm)	RFBB
Total rainfall during the penultimate week (mm)	RFW2
Total rainfall during the penultimate and antepenultimate week (mm)	RF2W2
Total rainfall between bud break and flowering (mm)	RFBBF
Total rainfall between flowering and setting (mm)	RFFS
Total rainfall between setting and véraison (mm)	RFSV
Total rainfall between véraison and harvest (mm)	RFVH
Environmental variables related to wind, humidity and weight	
Prevailing wind direction over the preceding week (N,S,E,W)	Dir
Average relative humidity over the preceding week (%)	Hum
Minimum relative humidity over the preceding week (%)	HumMin
Maximum relative humidity over the preceding week (%)	HumMax
Average wind speed in Km/h over the preceding week (Km/h)	Speed
Maximum wind speed in Km/h over the preceding week (Km/h)	SpeedMax
Weight of 100 berries	W100B
Environmental variables related to temperature	
Average temperature over the preceding week (°C)	Temp
Minimum temperature over the preceding week (°C)	TempMin
Maximum temperature over the preceding week (°C)	TempMax
Aggregate of average daily temperatures since the beginning of the year (°C)	STemp
Days with maximum temperatures above 40° C	D40
Days with average temperatures above 18° C during maturation	DM18
Days with maximum temperatures above 30° C during maturation	DM30
Average differences between maximum and minimum daily temperature during maturation (°C)	DDN

Table 2. Results of applying the method ClassifierSubSetEval.

Variables	ClassifierSubSetEval														
	Best First					Genetic search					Linear Forward selection				
	DS	LR	M5P	MLP	REP	DS	LR	M5P	MLP	REP	DS	LR	M5P	MLP	REP
Var	90	100	40	80		90	100	70	100		90	100	40	80	
Age		100	100	70	80	20	100	80	80	100		100	100	70	80
Altit	10	100	10	80	100	40	100	40	100	100	10	100	10	80	100
Dir		30	80	40	60	20	10	90	60	30		30	80	40	60
Hum		100	100	50	100	10	100	100	80	90		100	100	50	100
HumMin		100	100	80	100		100	90	70	90		100	100	80	100
HumMax		100	100	80	90		100	100	80	90		100	100	80	90
Speed			100	50	80		10	70	80	80			100	50	80
SpeedMax		100	100	70	70	20	100	100	70	50		100	100	70	70
W100B			20	80		20		30	100				20	80	
RFW		60	80	80	30	20	30	50	60	60		60	80	80	30
RF2W		30	90	60	30	40	50	60	20	20		30	90	60	30
RF3W		60	70	80	50	10	30	60	70	60		60	70	70	50
RFY		40	80	50		30	50	100	60	60		40	80	60	
RFBB		100	100	70	100	20	100	90	70	60		100	100	70	100
RFW2		20	60	40	30	30	70	60	40	30		20	60	40	30
RF2W2		10	20	30	60	10	20	50	30	60		10	20	30	60
RFBBF	60	100	60	80	20	90	80	80	80	60	60	100	60	80	20
RFFS			70	80	40	20	50	50	70	70			70	80	40
RFSV		100	90	80	50	10	100	90	80	50		100	90	80	50
RFVH	40	80	30	30	80	40	70	50	60	90	40	80	30	30	80
Temp		100	20	60	70	40	100	50	60	30		100	20	60	70
TempMin		100	80	60	40	50	100	100	60	40		100	80	60	40
TempMax		100	60	60	30	50	100	30	70	50		100	60	60	30
STemp	100	20	100	60	80	100		100	70	60	100	20	100	60	80
D40		100	60	30	50		100	60	20	30		100	60	30	50
DM18		100	90	90	40		100	100	60	60		100	90	90	40
DM30		100	100	60	30		100	100	50	40		100	100	60	30
DDN		100	20	60	20		100	40	70	40		100	20	60	20

Table 3. Results of applying the method WrapperSubSetEval.

Variables	WrapperSubSetEval														
	Best First					Genetic search					Linear Forward selection				
	DS	LR	M5P	MLP	REP	DS	LR	M5P	MLP	REP	DS	LR	M5P	MLP	REP
Var	30	100	20	100		30	100	40	100		30	100	20	100	
Age		50	40	100		30	50	70	90			50	40	100	
Altit	70	100		100		70	100	30	100		70	100		100	
Dir			80	70	70	10		60	60	60			80	70	70
Hum		100	100	70	50		100	100	70	50		100	100	70	50
HumMin		50	100	100	70		50	100	100	70		50	100	100	70
HumMax		100	100	90	30	10	100	100	90	60		100	100	90	30
Speed		50	90	70	80		50	60	80	70		50	90	70	80
SpeedMax		100	100	100	50	10	100	100	90	40		100	100	100	50
W100B		50		100			60	30	80			60		100	
RFW		60	60	60	40	10	50	50	60	20		60	60	60	40
RF2W		40	80	80	60	20	40	70	60	40		40	80	80	60
RF3W		10	30	80	60	20	10	70	80	50		10	30	80	60
RFY	20	90	90	90	80	30	60	100	100	70	20	90	90	90	80
RFBB		100	100	90	100	20	100	100	100	80		100	100	90	100
RFW2			70	30	30		10	40	30	40			70	30	30
RF2W2		10	40	40	20	10	40	10	90	30		10	40	40	20
RFBBF	30	100	100	90	90	30	90	70	100	90	30	100	100	90	90
RFFS		20	80	90	60	20		100	100	90		20	80	90	60
RFSV		100	100	90	100	10	100	90	100	80		100	100	90	100
RFVH	50	10	30	60	80	40	50	10	70	60	50	10	30	60	80
Temp		100	60	40	20	20	100	60	70	60		100	60	40	20
TempMin		100	70	90	50	20	100	90	90	90		100	70	90	50
TempMax		100	90	50	40	40	100	70	50	70		100	90	50	40
STemp	100	10	100	100	100	100	10	100	100	100	100	10	100	100	100
D40		80	60	40	40		80	50	60	50		80	60	40	40
DM18		100	90	100	70		100	70	80	70		100	90	100	70
DM30		100	100	70	60		100	100	90	100		100	100	70	60
DDN		100	30	70	10		100	40	40	40		100	30	70	10

- 3) Decision stump (DS) [2]: Class for building and using a decision stump. Usually used in conjunction with a boosting algorithm. Does regression (based on mean-squared error) or classification (based on entropy). Missing is treated as a separate value.
- 4) M5P (M5P) [28]: Implementation of base routines for generating M5Model trees. A decision list for regression problems is generated using separate-and-conquer. It builds a model tree in each iteration using M5 algorithm and makes the 'best' leaf into a rule. Quinlan's M5P can learn such piece-wise linear models. M5P also generates a decision tree that indicates when to use which linear model.
- 5) RepTree (REP) [34]: Fast decision tree learner. Builds a decision/regression tree using information gain/variance and prunes it using reduced-error pruning (with backfitting). Only sorts values for numeric attributes once. Missing values are dealt with by splitting the corresponding instances into pieces (i.e. as in C4.5).

E. Search methods

Like feature selection methods, it is unclear which search method is the most appropriate for each case.

And so that several studies, comparing different search algorithms are made by several authors [16] [21], although in this study is only used the following methods.

- 1) Best first: Searches the space of attribute subsets by greedy hillclimbing augmented with a backtracking facility. Setting the number of consecutive non-improving nodes allowed controls the level of backtracking done. Best first may start with the empty set of attributes and search forward, or start with the full set of attributes and search backward, or start at any point and search in both directions (by considering all possible single attribute additions and deletions at a given point) [30].
- 2) Genetic search: Search using a simple genetic algorithm [9]. Genetic Algorithms are efficient and robust search methods that are being employed in a plethora of applications with extremely large search spaces. The directed search mechanism employed in Genetic Algorithms performs a simultaneous and balanced, exploration of new regions in the search space and exploitation of already discovered regions.
- 3) Linear forward selection: Extension of Best First. Takes a restricted number of k attributes into account. Fixed-set selects a fixed number k of attributes, whereas k is increased in each step when fixed-width is selected. The search uses either the initial ordering to select the top k attributes, or performs a ranking (with the same evaluator the search uses later on). The search direction can be forward, or floating forward selection (with optional backward search steps) [10].

F. Results

The variables selected initially are shown in Table 1. And from a prescribed number of variables, this study shows that

as a preliminary step to calibrate a model, there are several techniques for reducing the number of variables.

In Tables 2 and 3 are shown, for different search methods and wrappers methods used, which are the variables selected. To see how feature selection affects to the problem, a study is made of how much influence has a model calibration with all variables or just using the variables selected by the feature selection algorithms. The finally selected variables were the ones which have shown best figures with all used methods.

In this case, the variables that result more suitable according to the algorithms are: Var, Age, Altit, Hum, HumMin, HumMax, SpeedMax, RFY, RFBB, RFBBF, RFSV, TempMin, STemp, DM18 and DM30. The accumulated percentages of all the variables, from which these features are selected, are shown in Table 4. A test is done to develop a regression model of sugar concentration that owns a grain of grape of a vineyard during its maturation, in order to verify that the work carried out with the feature selection methods is useful. The regression model chosen is a neuronal network formed by five neurons. This model is calibrated and tested, as much with all the initial data, like just by the variables selected by the methods used in this experiment. Not only it contributes to diminish the time used in the model generation but it improves the forecast results of this regression (Table 5 and Fig. 4).

Table 4. Accumulated percentages of used methods.

Vineyard variables	
Var	1750
Age	1700
Altit	1820
Environmental variables related to the amount of rainfall	
RFW	1350
RF2W	1360
RF3W	1330
RFY	1750
RFBB	2260
RFW2	910
RF2W2	810
RFBBF	2230
RFFS	1450
RFSV	2130
RFVH	1520
Environmental variables related to wind, humidity and weight	
Dir	1260
Hum	2040
HumMin	2070
HumMax	2110
Speed	1540
SpeedMax	2060
W100B	830
Environmental variables related to temperature	
Temp	1530
TempMin	1920
TempMax	1690
STemp	2280
D40	1370
DM18	2000
DM30	1920
DDN	1290

The studied error indices are the following:

1) Correlation coefficient (CORR)

$$CORR = \frac{\sqrt{\sum_{k=1}^n (p_k - m_k)^2}}{\sqrt{\sum_{k=1}^n (m_k - \bar{m})^2}} \quad (1)$$

2) Root mean squared error (RMSE)

$$RMSE = \sqrt{\frac{1}{n} \sum_{k=1}^n (m_k - p_k)^2} \quad (2)$$

3) Mean absolute error (MAE)

$$MAE = \frac{1}{n} \sum_{k=1}^n |m_k - p_k| \quad (3)$$

4) Root relative squared error (%) (RRSE)

$$RRSE = \frac{\sqrt{\sum_{k=1}^n (p_k - m_k)^2}}{\sqrt{\sum_{k=1}^n (m_k - \bar{m})^2}} \quad (4)$$

5) Relative absolute error (%) (RAE)

$$RAE = \frac{\sum_{k=1}^n |(p_k - m_k)|}{\sum_{k=1}^n |(m_k - \bar{m})|} \quad (5)$$

where m and P are, respectively, the measured and predicted outputs, n is the number of points of the database

used to validate the models, $\bar{m} = \frac{1}{n} \sum_{k=1}^n m_k$ and $\bar{p} = \frac{1}{n} \sum_{k=1}^n p_k$.

Also it is possible to observe like comparing the results obtained in the correlation between the real data and the calculated, it is obtained better results using only the selected variables (Fig. 5 and 6).

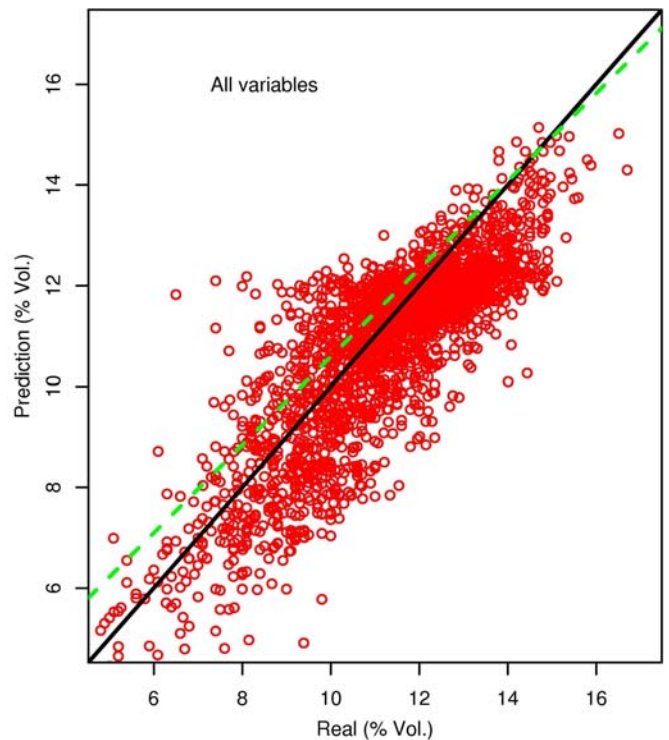


Fig. 5 Existing correlation between the real data and the predicted ones for all the variables.

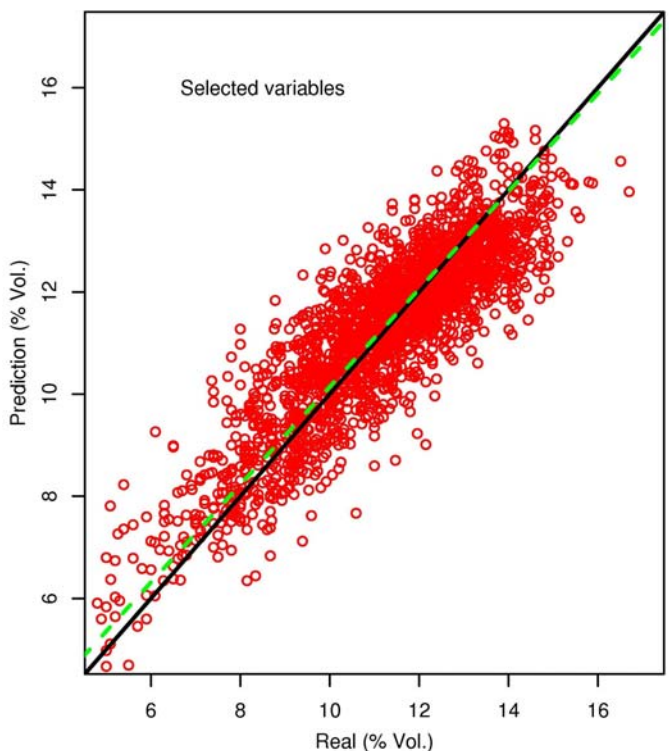


Fig. 6 Existing correlation between the real data and the predicted ones for only the selected variables.

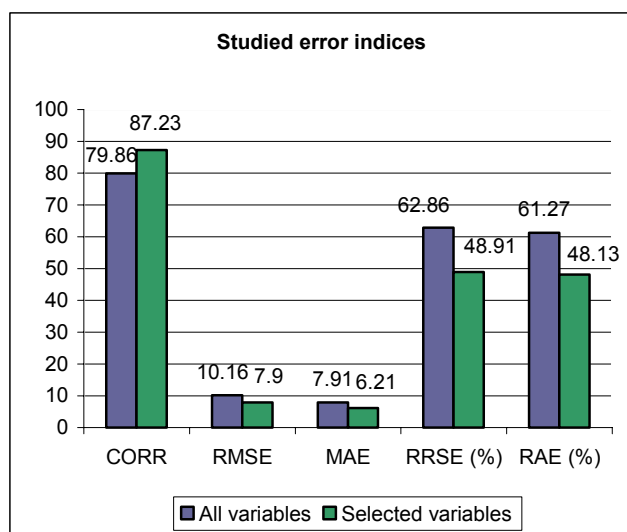


Fig. 4 Comparison of error rates obtained for the methods used effectiveness.

Table 5. Accumulated percentages of used methods.

	CORR	RMSE	MAE	RRSE (%)	RAE (%)
All variables	0.7986	0.1016	0.0791	62.8631	61.2768
Selected variables	0.8723	0.079	0.0621	48.9106	48.139

IV. CONCLUSIONS

It is observed that the use of feature selection algorithms is effective since when reducing the input variables, the understanding of the generated models improves. And since the model is compound of less number of variables, the problem is defined by the most significant variables in a clearer way. And as it is not used feature transformation methods these variables are more understandable.

Also a reduction of computational cost is verified when generating the final model. Mainly in the most complex models, since he is not the same to calibrate models with many variables that do not contribute significant information, that to calibrate only with the most significant. In the studied case the time of operation is reduced in a 34%.

In addition the precision of the final model has improved. The reduction of no significant variables improves the precision in the five studied indices to verify the error.

It is also verified, that all the used algorithms to reduce features do not get the same conclusions, although realising an analysis of all the methods jointly allow us to obtain the best solution for the problem resolution.

As final conclusion, it is determined that the application of these algorithms is useful with this kind of data and is advisable for future works.

REFERENCES

- [1] A.L. Blum, P. Langley, "Selection of relevant features and examples in machine learning", *Artificial Intelligence*, Vol. 97, 1997, pp. 245-271.
- [2] L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone, *Classification and Regression Trees*, Chapman & Hall, New York, 1984.
- [3] B. Bulut, N. Allahverdi, H. Kahramanli, S. Yalpir, "A Residential Real-Estate Valuation Model with Reduced Attributes". *International Journal Of Mathematical Models And Methods In Applied Sciences*, Issue 3, Volume 5, 2011, 586- 593.
- [4] R.A. Cardenas Tamayo, M.G. Lugo Ibarra, J.A. Garcia Macias, "Better crop management with decision support systems based on wireless sensor networks", *Electrical Engineering Computing Science and Automatic Control (CCE)*, 2010, pp. 412-417.
- [5] A. Ceglar, Z. Crepinsek, L. Kajfez-Bogataj, T. Pogacar, "The simulation of phenological development in dynamic crop model: The Bayesian comparison of different methods", *Agricultural and Forest Meteorology*, Vol. 151, 2011, pp. 101-115.
- [6] B.G. Coombe, "Research on Development and Ripening of the Grape Berry", *Am. J. Enol. Vitic.*, Vol. 43, 1992, pp. 101-110.
- [7] M. Dash, H. Liu, "Feature Selection for Classification", *Intelligent Data Analysis*, Vol. 1, 1997, pp. 131-156.
- [8] R. Furferi, L. Governi, Y. Volpe, "Neural Network based Classification of Car Seat Fabrics", *International Journal Of Mathematical Models And Methods In Applied Sciences*, Issue 3, Volume 5, 2011
- [9] D.E. Goldberg, *Genetic algorithms in search, optimization, and machine learning*, Addison-Wesley Pub. Co., Reading, Mass, 1989.
- [10] M. Gutlein, E. Frank, M. Hall, A. Karwath, "Large-scale attribute selection using wrappers", *Computational Intelligence and Data Mining, CIDM '09*, 2009, pp. 332-339.
- [11] I. Guyon, A. Elisseeff, "An Introduction to Variable and Feature Selection", *Journal of Machine Learning Research*, Vol. 3, 2003, pp. 1157-1182.
- [12] S. Haykin, *Neural Networks: A Comprehensive Foundation*, 2nd edn, Prentice Hall PTR: Upper Saddle River, NJ, USA, 1999.
- [13] J. Hwang, C. Shin, H. Yoe, "Study on an Agricultural Environment Monitoring Server System using Wireless Sensor Networks", *Sensors*, Vol. 10, 2010, pp. 11189-11211.
- [14] D.I. Jackson, P.B. Lombard, "Environmental and Management Practices Affecting Grape Composition and Wine Quality - A Review", *Am. J. Enol. Vitic.*, Vol. 44, 1993, pp. 409-430.
- [15] R.S. Jackson, *Wine Science Principles and Applications*, Third Edition, Elsevier Inc, 2008.
- [16] A. Jain, D. Zongker, "Feature selection: evaluation, application, and small sample performance", *Pattern Analysis and Machine Intelligence*, Vol. 19, 1997, pp. 153-158.
- [17] G.H. John, R. Kohavi, K. Pfleger, "Irrelevant feature and the subset selection problem". *Proceedings of the Eleventh International Conference on Machine Learning*, 1994, pp. 121-129.
- [18] K. Kira, L.A. Rendell, "The feature selection problem: traditional methods and a new algorithm". *Proceedings of the tenth national conference on Artificial intelligence*, 1992, pp. 129-134.
- [19] R. Kohavi, G.H. John, "Wrappers for feature subset selection", *Artif. Intell.*, Vol. 97, 1997, pp. 273-324.
- [20] D. Koller, M. Sahami, "Toward Optimal Feature Selection", *Proceedings of ICML-96, 13th International Conference on Machine Learning*, 1996, pp. 284-292.
- [21] M. Kudo, J. Sklansky, "Comparison of algorithms that select features for pattern classifiers", *Pattern Recognition*, Vol. 33, 2000, pp. 25-41.
- [22] H. Liu, H. Motoda, *Feature Selection for Knowledge Discovery and Data Mining*, Kluwer Academic Publishers, 1998.
- [23] H. Liu, H. Motoda, L. Yu, *Feature Extraction, Selection, and Construction, The Handbook of Data Mining*. Lawrence Erlbaum Associates, Inc. Publishers, 2003, pp. 409 - 423.
- [24] F. Mazzetto, A. Calcante, A. Mena, A. Vercesi, "Integration of optical and analogue sensors for monitoring canopy health and vigour in precision viticulture", *Precision Agriculture*, Vol. 11, 2010, pp. 636-649.
- [25] V. Olej, P. Hajek, "Modelling municipal rating by unsupervised methods," *WSEAS Transactions on Systems*, vol. 5, no. 7, pp. 1679 - 1786, 2006.
- [26] Y. Peng, Z. Wu, J. Jiang, "A novel feature selection approach for biomedical data classification", *Journal of Biomedical Informatics*, Vol. 43, 2010, pp. 15-23.

- [27] S. Piramuthu, "Evaluating feature selection methods for learning in data mining applications", *European Journal of Operational Research*, Vol. 156, 2004, pp. 483-494.
- [28] J.R. Quinlan, "Learning with continuous classes", *Proceedings of the 5th Australian joint Conference on Artificial Intelligence*. Singapore, 1995, pp. 343-348.
- [29] P. Ribéreau-Gayon, D. Dubourdieu, B. Donèche, A. Lonvaud, *Chapter 10. The Grape and its Maturation, Handbook of Enology, Volume 1, 2nd Edition, The Microbiology of Wine and Vinifications*, John Wiley & Sons L, 2005.
- [30] L. Rios, L. Chaimowicz, "A Survey and Classification of A* Based Best-First Heuristic Search Algorithms", *Advances in Artificial Intelligence*, SBIA 2010, Springer Berlin / Heidelberg, 2011, pp. 253-262.
- [31] T.R. Sinclair, N.G. Seligman, "Crop Modeling: From Infancy to Maturity", *Agron. J.*, Vol 88, 1996, pp. 698-704.
- [32] C.O. Stöckle, M. Donatelli, R. Nelson, "CropSyst, a cropping systems simulation model", *European Journal of Agronomy*, Vol. 18, 2003, pp. 289-307.
- [33] G.N. Wilkinson, C.E. Rogers, "Symbolic description of factorial models for analysis of variance". *Journal of Applied Statistics*, Vol. 22, 1973, pp. 392-399.
- [34] I.H. Witten, E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques*, Second Edition edn, Morgan Kaufmann, 2005.
- [35] X. Zhang, M.A. Friedl, C.B. Schaaf, et al, "Monitoring vegetation phenology using MODIS", *Remote Sensing of Environment*, Vol. 84, 2003, pp. 471-475.
- [36] N. K. Zivadinovic, K. Dumicic, A. C. Casni, "Multivariate Analysis of Structural Economic Indicators for Croatia and EU 27," in *Proc. 2nd WSEAS Int. Conf. Multivariate Analysis and its Application in Science and Engineering*, Istanbul, 2009, pp. 134-139.

Roberto Fernández Martínez, MSc - graduated in Industrial Engineering from Universidad de La Rioja in 2006. He is working in his PhD degree at the same university in the area of industrial optimization through data mining techniques. He is researcher of the Universidad de La Rioja in the Mechanical Department. His main research interests are data mining; soft computing; pattern recognition; artificial intelligence; industrial optimization and feature selection. He has been researcher of several data mining Projects. Currently is finishing his thesis in the area of data mining and soft computing.

Julio Fernández Ceniceros received the M.Sc. degree in Industrial Engineering in 2009 and the Master Degree in Project Management from the University of La Rioja (Spain) in 2010. He is working at the University of La Rioja as a fellowship and his current research is applied numerical simulations (Finite Element Method) and Data Mining techniques in steel and concrete structures. His interests include the behavior of bolted connections, FEM simulations, failures modes, plasticity and damage.

Andrés Sanz García, received his B.S degree in Industrial Engineering in 1999 and M.S degree in Industrial Engineering from Universidad de La Rioja in 2002. Currently he is engaged on his PhD in the area of industrial optimization through soft computing techniques. He is registered Professional Engineer until 2011 in La Rioja, Spain. Recently, he is a professor and researcher of the Universidad de La Rioja in the Mechanical Department with EDMANS research group. He has published several papers in conferences. His research area includes but not limited to the use of finite element analysis to solve engineering problems, materials science, fracture mechanics, and failure analysis of engineering materials.

Rubén Lostado Lorza received the M.Sc. degree in engineering in 2003 and the Ph.D. degree in engineering from the University of La Rioja, Logroño, Spain in 2010. From 2003 to 2005 he was working on numerical simulation (Finite element method) at the department of Wind Energy in the company M-torres industrial designs. From 2005 to 2007 he was working on numerical simulation (Finite element method) in the Automotive Technological Center of Navarra (CITEAN). He is currently working at the University of La Rioja as professor of engineering projects and environmental technology. His current research is the combined modeling of industrial processes and products using the Finite Element Method and Data Mining

Francisco Javier Martínez de Pisón Ascacibar, PhD - graduated in Industrial Engineering from Universidad de La Rioja in 1999. In 2003 he completed a PhD degree at the same university in the area of industrial optimization through data mining techniques. He is teacher and researcher of the Universidad de La Rioja in the Mechanical Department. His main research interests are data mining; soft computing; pattern recognition; artificial intelligence and industrial optimization. He has published many papers in journals and is author of several books and industrial patents. He has been the main research of several national projects and has participated in other national and European data mining Projects. Currently supervises several MSc and PhD students in the area of data mining and soft computing.