# AN APPLICATION OF SOFT COMPUTING TECHNIQUES FOR A HOTEL DECISION SUPPORT SYSTEM

Francisco Javier Martínez-de-Pisón Ascacíbar

Rubén Escribano García

Julio Fernández Ceniceros

Roberto Fernández Martínez

*Grupo EDMANS (http://www.mineriadatos.com)*

*Universidad de La Rioja. España*

José Ignacio Pérez Moneo

*Angestur Consultores e Innovación S.L (HotelOptimizer) (http://hoteloptimizer.com/)*

## Abstract

In last years, due to the spreading of websites specialized on online hotel booking, many clients wait until the last moment to make the reservations. Thus, at the last moment, they can compare the room prices in different closely hotels, with almost the same services, in order to save money. As a consequence, for hotel managers, it is a challenge to maximize profits by continuously changing room prices to present them more attractive than the competition hotels. This paper shows the results of a research project which principal aim is the application of soft computing techniques to obtain models useful to predict future bookings for each calendar day, by means of the use of the hotel historical reservation data and the characteristics of each day (month, day of the week, season, festivities in closely regions or cities, the weather, etc.). Built models will be part of a hotel decision support software which is been developed for *HotelOptimizer* company with EDMANS group support.

*Keywords*: hotel decision support systems; soft-computing; hotel booking systems; data mining

## Resumen

En los últimos años, debido a la proliferación de portales Web especializados en reservas *On-Line* de habitaciones de hotel, muchos clientes esperan hasta el último momento para realizar sus reservas. Así, pueden ahorrarse dinero comparando los precios de última hora de hoteles cercanos que ofrecen prácticamente los mismos servicios. Como consecuencia de esto, para los gerentes de hotel es un reto poder maximizar los beneficios mediante el ajuste continuo de los precios de las habitaciones de forma que éstos sean más atractivos para los clientes frente a los de hoteles de la competencia. En este artículo, se muestran los resultados de un proyecto de investigación donde se han aplicado técnicas de *soft-computing* para desarrollar modelos capaces de predecir las futuras reservas de cada día del año para un hotel cualquiera a partir de la base de datos de históricos de reservas del mismo y las características de cada día (mes, día de la semana, temporada del año, fiestas en comunidades o ciudades cercanas, climatología, etc.). Los modelos desarrollados formarán parte de un software de apoyo a la toma de decisiones para la gestión de hoteles que está desarrollando la empresa *HotelOptimizer* en colaboración con el Grupo EDMANS.

*Palabras clave:* sistemas de apoyo a la toma de decisiones en hoteles; computación flexible; reservas de hotel; minería de datos

## 1. Introduction

Nowadays, due to the increase of Websites to online hotel booking, hotel managers need to change continuously the room prices in order to be competitive with other competition hotels. In this case, the aim consists to predict the precise room price to obtain maximum profit in each moment.

*HotelOptimizer* (http://www.hoteloptimizer.com) is a new company that is developing new decision support tools to help to the hotel managers. For example, they present a free Web-tool which permits to analyze a 14 days prediction of hotel prices in different Spanish regions (figure 1.).

**Figure 1: Web-tool where is shown a prediction of hotel prices for the next 14 days in Spain**
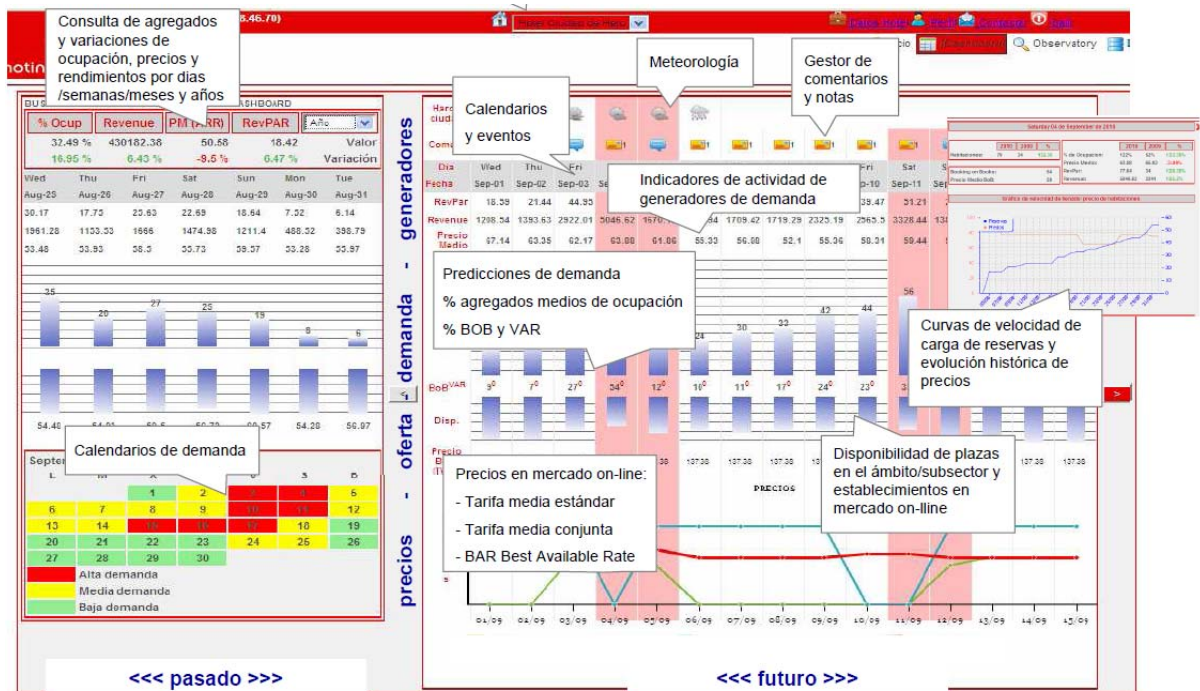


Also, they have a new commercial product, named "Hop Yield & Revenue Management Assistant (Hop YRMA)" (figure 2), which permits to:

1. Obtain business positioning in the market and management analysis occupation money: developments, trends, schedules and forecasts of demand for employment, reservations and price elasticity, price analysis and strategies of competitors, anticipating to the fluctuations in the availability of places in the area, planning of prices, availability, conditions and restrictions and updating of distribution channels.

2. Demand forecast: Hop YRMA has a demand forecasting system through complex formulas and identifying patterns using data mining techniques.

3. See details of the competition: Learn the strategies of competitors knowing their prices for the coming days and weeks with a single click (figure 3).

4. See demand schedules: Hop YRMA has a comprehensive database of events and issues that could to affect to the hotel.

5. Can be integrated with any Property Management Systems (PMS).

6. Have a Channel Manager connected to the hotel PMS and distribution channels allowing automatically link availability in all distributors and, even able to process reservations through the system connected with the hotel's own website.

One of the most important parts of this product is the capacity to forecast the future hotel demand. The basic idea was to use soft computing techniques to develop models useful to predict future room bookings for each calendar day, by means of the use of the hotel historical reservation data and the characteristics of each day (month, day of the week, season, festivities in closely regions or cities, weather, etc.).

**Figure 2: Hop Yield & Revenue Management Assistant (Hop YRMA)**
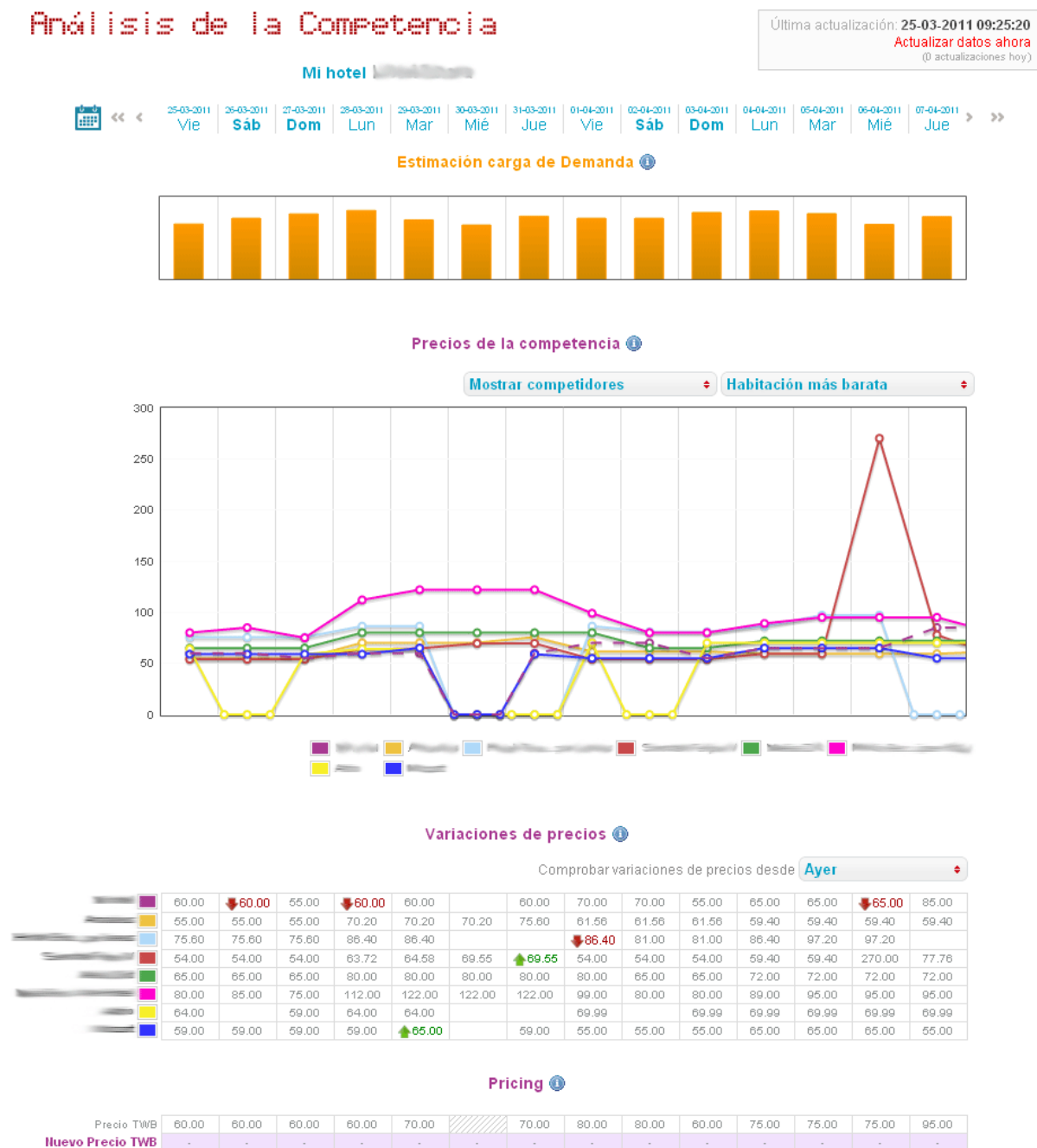


In this paper the results of the research project, named "CIERZO", are presented in which its principal objective was to build useful models to create demand calendars.

## 2. Developing Hotel Demand Calendars Using Soft Computing Techniques
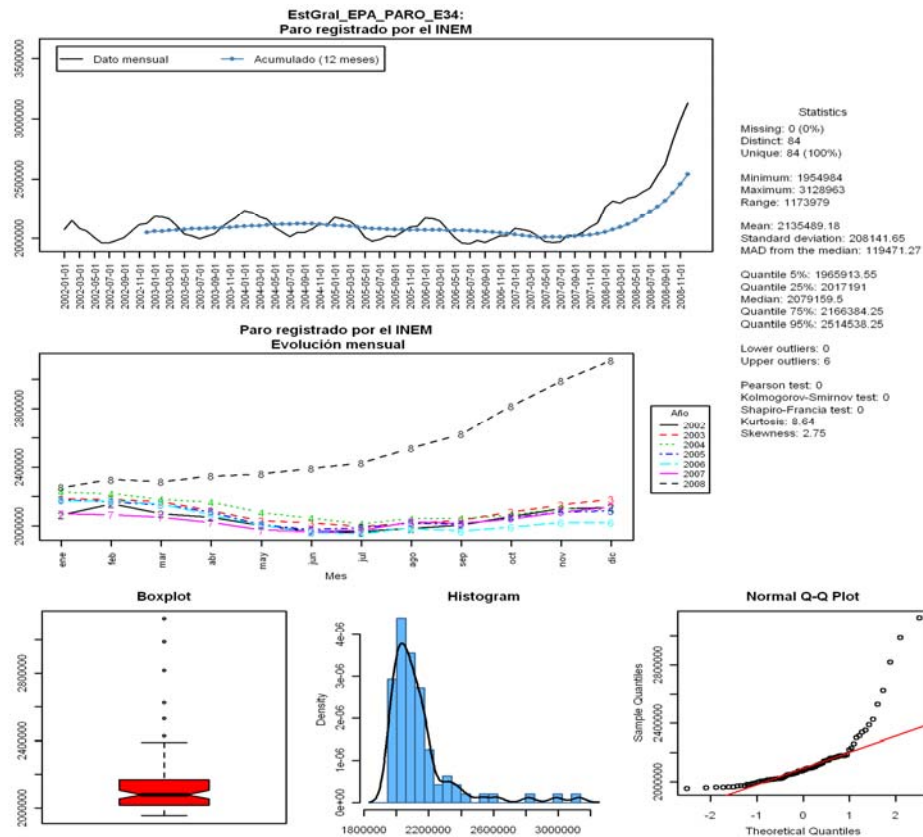
### 2.1 Attributes Selection

The first step was to analyze visually many variables extracted from different sources: Macroeconomic and Social variables from the "*Instituto Nacional de Estadística (INE)*" databases, Meteorological databases from Websites, local and global festivities and other type of data included manually from different sites (figure 4).

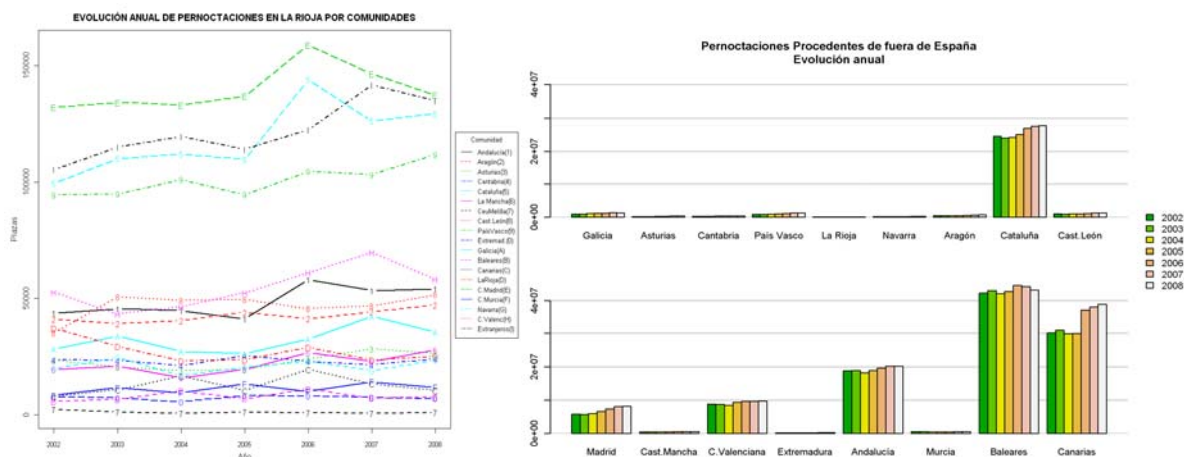**Figure 3: Example of analysis of competitors**



Also, evolution of hotel room prices and others hotels indicators were studied in different Spanish regions (figure 5). Finally, 119 of the most interesting attributes were selected.

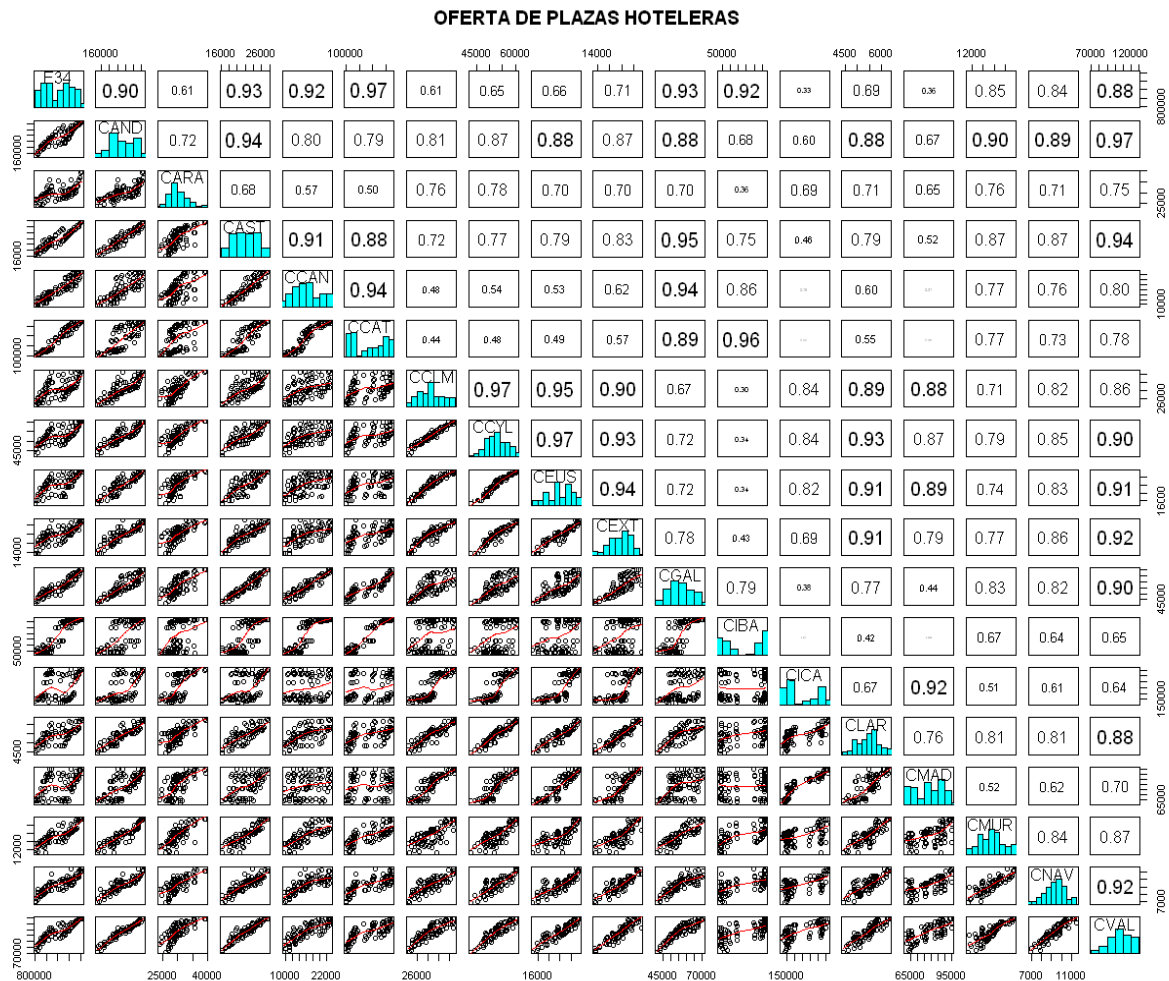**Figure 4: Example of unemployment evolution in Spain**



**Note: It is interesting to observe the strong increment from the end of 2007 when government was saying that it had not been crisis.**

**Figure 5: Evolution of overnights in La Rioja of people from different Spanish Regions (left)**

**and Evolution of overnights in each Spanish region of people from different countries**



After this, advanced scatter plots (figure 6) were used to identify correlation between Macroeconomic indexes, meteorological data, local and global festivities, etc. in front of hotel overnights attributes and from different Spanish regions. At the end of this step, the list of the most important variables was reduced significantly.

**Figure 6: Example of scatter plot of room vacancy in different Spanish Regions**

OFERTA DE PLAZAS HOTELERAS

## 2.2 Calendar Forecast for one specific Hotel with Data Mining Techniques

After the selection of the most important variables (Economical, Social, Festive, Meteorological and Hotel Indicators), the next aim was to search a methodology for develop models useful in order to create calendar forecast for a specific hotel.

The basic idea was to use the historical reservation data from each hotel to create specific models to predict the future bookings for each future calendar day. The design of the models was structured in 22 input attributes that define the characteristics of each day (month, day of the week, season, festivities in closely regions or cities, festivities in Spain, the weather, etc.) and one output variable with the number of bookings for the day. Due to confidentiality is not possible to show the name and type of each selected attribute.

First, models were created using data of the years 2007 to 2009. Data from before 2007 were not considered because of the different Spanish economic situation.

In order to find models that generate a low prediction error, a battery of Data Mining algorithms were used:

- M5P algorithm (M5P): Implements base routines for generating M5Model trees. A decision list for regression problems is generated using separate-and-conquer. In each iteration, it builds a model tree using M5 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 (Quinlan, 1992).

- Multilayer Perceptron (MLP): 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 (Haykin, 1999). Training is performed with networks that have between 1 and 40 neurons in the hidden layer.

- Linear Regression (LINREG): A class for using linear regression for prediction. It uses the Akaike criterion for variable selection and is able to deal with weighted instances (Wilkinson and Rogers, 1973).

- Simple Linear Regression (SIMPLR): Uses only the best attribute to obtain the model. It is useful for comparing with other algorithms.

- LeastMedSq (LMSQ): Implements a least median squared linear regression to make predictions. Least squared regression functions are generated from random sub-samples of the data. The least squared regression that has the lowest median squared error is chosen as the final model (Portnoy and Koenke, 1997).

- IBk (IBk): A version of the k-nearest neighbour algorithm. K is the number of neighbours to be used. It also permits the use of distance weighting. As it is a lazy algorithm, there is no training time (Aha and Kibler, 1991).

To obtain the best precision, ten models of each type of algorithm configuration were trained with 70% of the data from the training database and the remaining data (30%) were used to validate each model. By generating 10 models of each algorithm configuration, the influence of local minima was reduced and much more realistic errors were obtained.

The purpose of this work was to determine the algorithm configuration that provide the best rooms booking prediction or, in other words, the algorithm configuration that yields the lowest Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) for other different days not used for model construction. These errors were:

$$RMSE = \sqrt{\frac{1}{n}\sum_{k=1}^{n}(y(k)-\hat{y}(k))^2} \qquad\qquad (1)$$

and

$$MAE = \frac{1}{n}\sum_{k=1}^{n}\left|y(k)-\hat{y}(k)\right| \qquad\qquad (2)$$

where $y$ and $\hat{y}$ were, respectively, the measured and predicted rooms booking and $n$ was the number of examples in the database used to validate the models.

The results of the training and validation process are shown in table 1. This table provides a summary of the validation errors arranged by the RMSE corresponding to ten trained models for the algorithm configurations. This table presents the mean (MEAN) and standard deviation (SD) of RMSE and MAE validation of ten models of each type of algorithm configuration.

As it is observed, the best algorithm was a Multilayer Perceptron with 20 neurons in the hidden layer but the errors were not acceptable (RMSEMEAN over 40%).

**Table 1: Results of the Training and Validation Process**

| Algorithm | RMSEMEAN | RMSESD | MAEMEAN | MAESD | TIME |
|---|---|---|---|---|---|
| MLP (Neurons=20) | 0.404 | 0.004 | 0.311 | 0.008 | 8628.589 |
| MLP (Neurons=15) | 0.405 | 0.003 | 0.313 | 0.006 | 1804.589 |
| MLP (Neurons=10) | 0.405 | 0.002 | 0.316 | 0.010 | 1951.41 |
| IBk (K=10) | 0.405 | 0.004 | 0.312 | 0.003 | 0.02 |
| LMSQ | 0.405 | 0.003 | 0.329 | 0.002 | 22.165 |
| MLP (Neurons=5) | 0.405 | 0.004 | 0.314 | 0.010 | 899.565 |
| MLP (Neurons=40) | 0.406 | 0.004 | 0.311 | 0.012 | 13739.393 |
| MLP (Neurons=30) | 0.406 | 0.003 | 0.312 | 0.011 | 8228.685 |
| MLP (Neurons=3) | 0.406 | 0.003 | 0.322 | 0.011 | 715.824 |
| MLP (Neurons=7) | 0.407 | 0.002 | 0.312 | 0.008 | 1009.002 |
| M5P | 0.407 | 0.005 | 0.316 | 0.004 | 4.122 |
| IBk (K=6) | 0.409 | 0.004 | 0.308 | 0.002 | 0.017 |
| IBk (K=5) | 0.412 | 0.003 | 0.308 | 0.002 | 0.018 |
| LINREG | 0.425 | 0.004 | 0.352 | 0.009 | 2709.669 |
| IBk (K=3) | 0.430 | 0.004 | 0.308 | 0.003 | 0.017 |
| SIMPLR | 0.443 | 0.004 | 0.307 | 0.003 | 0.02 |

## 2.3 Designing the Final Algorithm using Soft Computing Techniques

From the conclusions raised in the previous works, a new model was contemplated using soft computing techniques.

Soft computing consists in the use of computational techniques and intelligent systems in order to solve inexact and complex problems (Sedano et al., 2010). The involved computational techniques are stochastic and therefore suited to investigate real-world problems (Banerjee et al., 2010; Corchado and Herrero, 2011).

Evolutionary algorithms (EA), which comprise a fundamental component of soft computing, are said to be inspired by nature. One of the most representative techniques of soft computing is the genetic algorithms (GA): a large number of systematic methods used to solve optimization problems applying the principles of biological evolution, namely survival of the 'fittest', sexual reproduction and mutation (Mitchell, 1998). From this conception of GA, it has been possible to solve real-world problems which were impossible to tackle previously with the classic techniques.

The new algorithm was created with the following phases: selection of the most important variables, identification of similar days and average values obtained with robust estimation of reserves. In this case, we sought to optimize the parameters of the proposed model using genetic algorithms, with the idea of getting a model as accurate as possible.

The proposed model was based on an instance-based learning algorithm where the day's Euclidean distance depended on: if there are parties or in nearby cities, If there are parties or not in nearby regions, the month of the year, the week of the month, the day of the week, the weather, etc.

The final algorithm was formed for these steps:

- Each attribute is normalized between 0 and 1 and multiplied with a weighting coefficient.

- Euclidean distance matrix is created with the computation of the Euclidean distance for each day with the others.

- K-nearest days are chosen for each date from the Euclidean distance matrix.

In order to obtain robust values, the mean room booking is calculated using the mean of the K-values that falls into the range limited by PERCENTIL and 1-PERCENTIL.

## 2.4 Searching the Best Algorithm's Parameters with Genetic Algorithms

Weighting coefficients, K, PERCENTIL and the most important attributes used to calculate the day's Euclidean distance were optimized with genetic algorithms.

First of all, 100 individual algorithm's parameters were initialised with random values.

For each individual the training data was randomly selected from the 70% of the database. It was used for calculates day's Euclidean distance matrix. Then, for each day, the most K-closed days were selected using this matrix. The other 30% of the data (test data) was used to calculate the correlation between the predictions of the room reservations and the actuals. In each case, correlation was computed ten times with different test and training data. The final objective function, to be minimised, was one minus the mean of the ten correlations.

Those of the many individuals in generation 0 which have the lowest objective function were selected and used as the basis for obtaining the next generation (generation 1) by means of crosses and mutations.

The new generation was made up as follows:

- 25% comprises the best individuals from the previous generation (parents of the new generation);

- 60% comprises individuals obtained by crossovers from selected parents. The crossover process involves changing various digits in the chromosomes of the variables to be modified. These chromosomes were made up of the digits for the variables with the decimal points removed, joined together in a single set.

- The remaining 15% was obtained by mutation, through the creation at random of chromosomes within the ranges established. The aim was to find new solutions in areas not previously explored.

This process was repeated over several generations until the objective function of the best individual was observed to remain the same or not to drop significantly from one generation to the next. If the error was suitable, the best individual was selected as the final solution.

## 3. Results

After 300 generations of individuals and 50 days of calculating, the maximum correlation achieved for the best model was 1-0.27142 = 0.72858. In figure 7 is possible to observe the evolution of the 100 individuals of the generations 1, 100, 200 and the final generation 300.

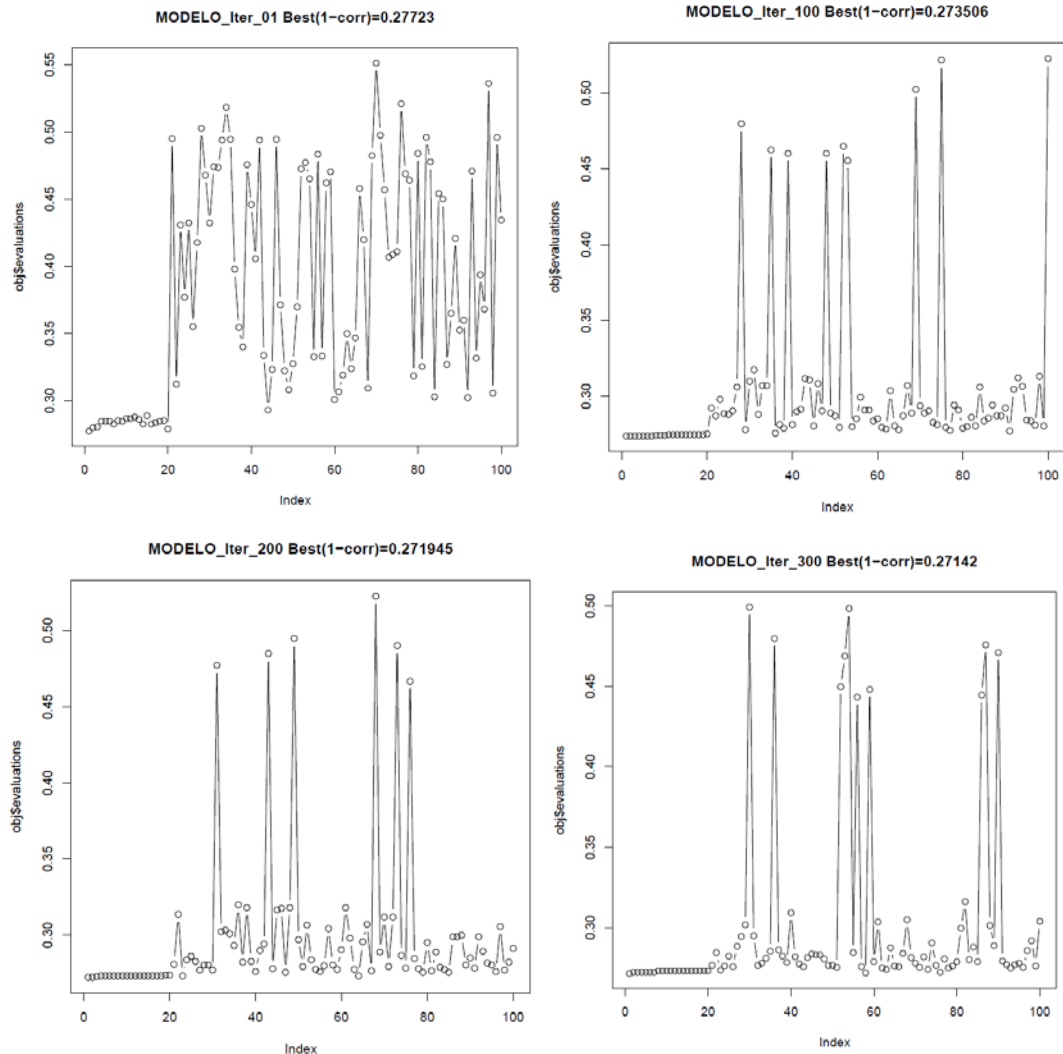**Figure 7: Generations: 1, 100, 200 and 300 (Final)**



Figure 8 presents the predictions of the best model for the first seven months of the year 2010 contrasting to the real room bookings. Calendar shows the prediction for each day on this way:

- The number on the left is the month day.

- The number in brackets shows the prediction of reservations for that day as estimated by the model using the database for the years 2007 to 2009.

- The number between brackets indicates the absolute error or absolute difference of the value that predicts the model and the actual bookings that were made for that day.

- In colors is showed the error respect the maximum value. Box is blue if the error is below 20%, green if error is between 20% and 40%, and so on like is presented in the box on the right.

As can be seen, most of the days of the calendar had obtained less than 20% error. Fourteen (14) of them had an error between 20% and 40% (green), and five (5) yellow errors (range between 40% and 60%).

**Figure 8: Calendar with the room reservation predictions for the first seven months for 2010**
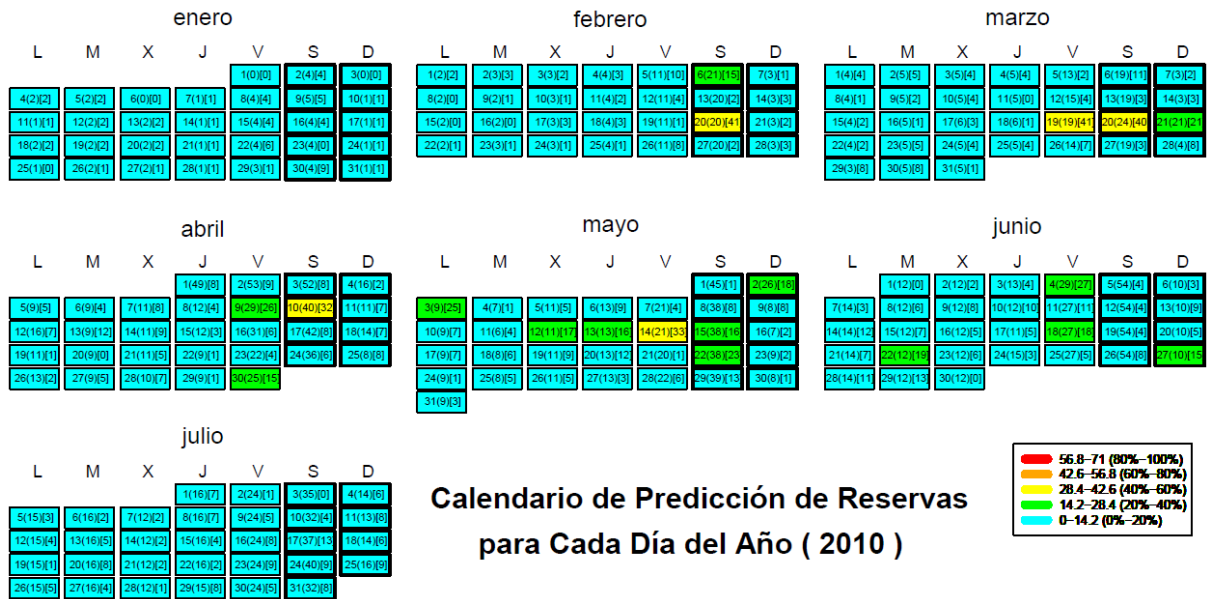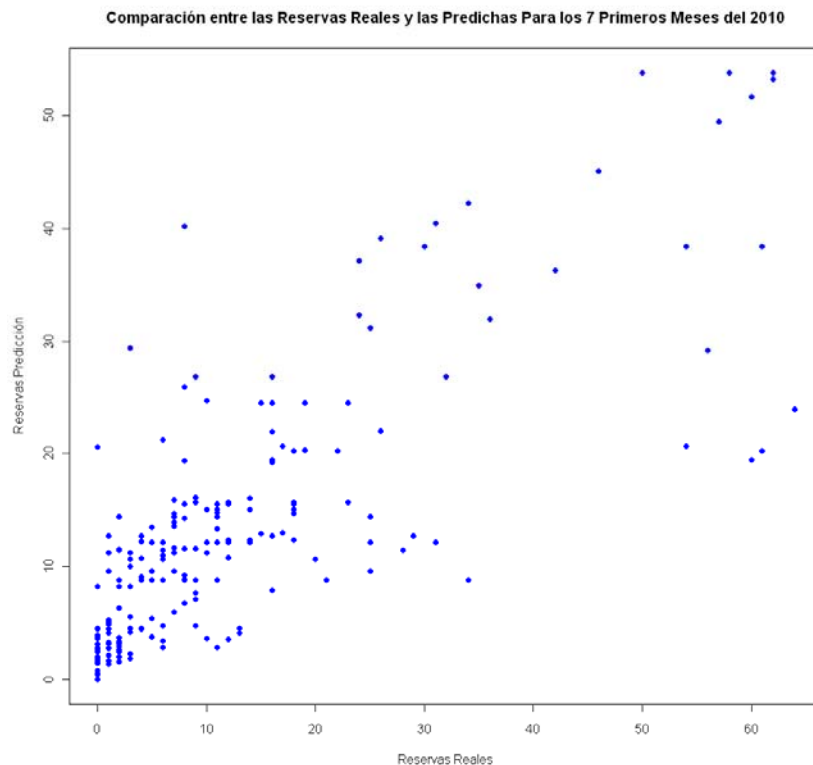


Figure 9 shows an XY plot which represents the actual reserves of these months compared with that predicted values. In this case, the points tend to approach the main diagonal of a correlation graph obtained corroborating close to 0.80 (80%) that can be considered as excellent.

**Figure 9: Predictions vs. actual room bookings of the best model**

With the fitted model is possible to define a booking calendar like in the figure 10. This case is a Forecast Calendar for 2010 using data from years 2007 to 2009.
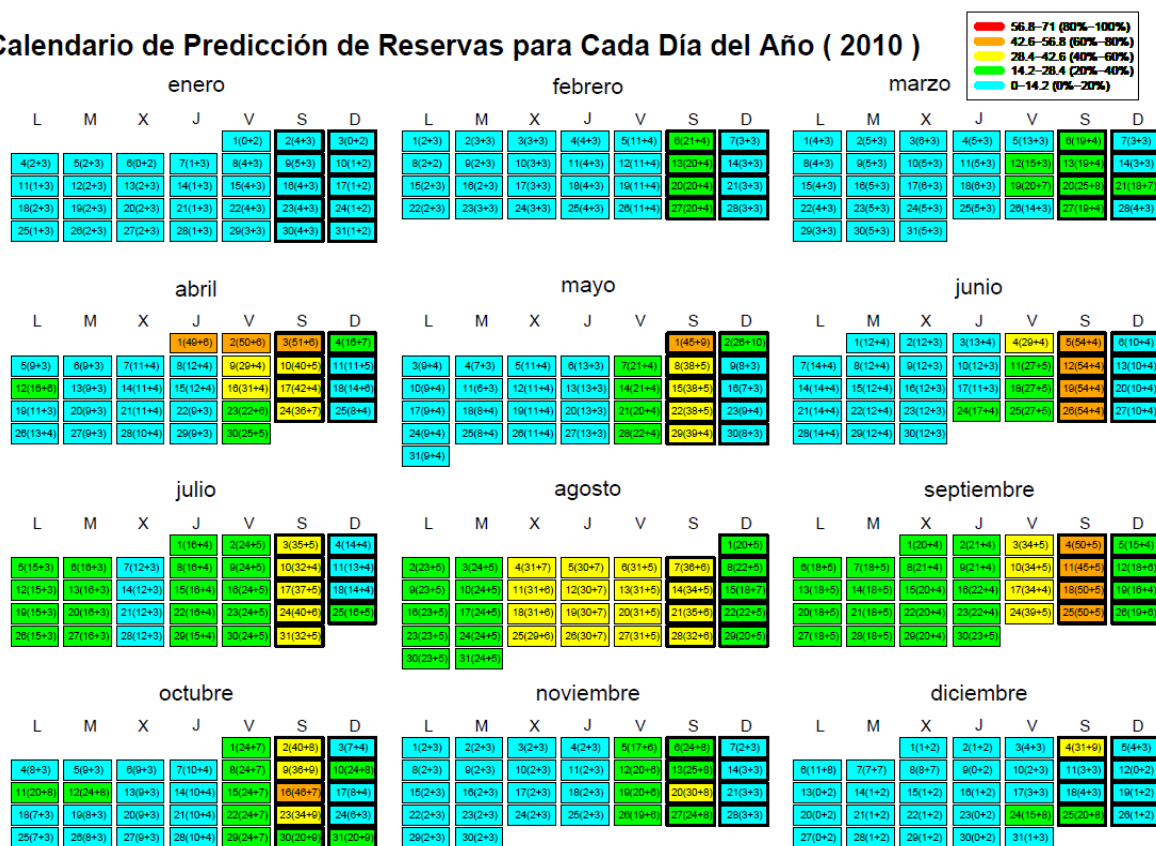
In brackets is shown the reservations prediction for each day (left) + the interval confidence at 95% of level (right). That is, a value 20+4 could be interpreted as: room bookings prediction for that day is 20±4, or whatever is the same, 95% of the days used to calculate the average value 20 have fallen into the interval [20-4, 20+4] = [16.24].

In colors, is shown the class of the day depending on the final reservations value predicted (cyan=very low, green=low, yellow=medium, orange=high, red=very high).

In this case, the maximum interval did not exceed the value of 10, but that does not mean they could be higher in the prediction errors due to abnormal days.

**Figure 10: Final forecast calendar to 2010**



## 4. Conclusions

Considering the high difficult to predict future room reservations, it could be considered than the results with the last model are quite good. Correlation with new data, not used to build the model, was 0.80. In contrast, the use of classical and advanced data mining algorithms, directly with the database, did not give satisfactory results.

The final model has been a heuristic algorithm based on the combination of three actions: selection of the best attributes, identification of similar days and average values obtained with robust estimation of reserves. Algorithm parameters have been optimized with genetic algorithms.

The optimization process was computed until the 500 generations but, in the last 200 generations, the best individual didn't improve (since the 300th generation). That comes to clearly indicate that the problem, for this type of hotel and these type of data, it had its own limitations because there was a high randomness in the rooms bookings. This came to mean that the models could be improved slightly but did not get very high precision due to the inherent properties of the data and the high uncertainty of the process.

Obviously, the time necessary for calculates 500 generations was very high (almost four months) but is possible to observe that in a few number of generations, the models' correlations were stable in the three principal decimal digits. In this manner, the methodology can be applied in order to obtain, in a few days, models that could be efficient enough to be used in hotel decision support systems.

## 5. Acknowledges

## 6. References

Aha, D. & Kibler, D. (1991). Instance based learning algorithms. *Machine Learning*, 6, 37-66.

Banerjee, T., Das, S., Roychoudhury, J. & Abraham, A. (2010). Implementation of a new hybrid methodology for fault signal classification using short -time fourier transform and support vector machines. In: Soft Computing Models in Industrial and Environmental Applications, 5th International Workshop (SOCO 2010), Advances in Intelligent and Soft Computing, vol 73, 219–225

Corchado, E. & Herrero. A., (2011). Neural visualization of network traffic data for intrusion detection. *Applied Soft Computing*, 11, (2), 2042-2056.

Haykin, S. (1999). *Neural networks, a comprehensive foundation (2nd ed.).* EE.UU, New Jersey: Prentice Hall.

Koenker, R. & Portnoy, S. (1997). The Gaussian Hare and the Laplacian Tortoise: Computability of squared-error vs. absolute-error estimators, with discussion. *Statistical Science,* 12, 279-300.

Mitchell, M. (1998). *An introduction to genetic algorithms.* Cambridge: The MIT Press.

Quinlan, J. R. (1992). Learning with Continuous Classes. 5th Australian Joint Conference on Artificial Intelligence, Singapore, 343-348.

Sedano, J., Curiel, L., Corchado, E., De la Cal, E. & Villar, J. (2010). A soft computing method for detecting lifetime building thermal insulation failures. *Integrated Computer-Aided Engineering*, 17, (2), 103–115.

Wilkinson, G. N. & Rogers, C. E. (1973). Symbolic description of factorial models for analysis of variance. *Applied Statistics*, 22, 392-399.

**Correspondence:**

Dr. Francisco Javier Martínez de Pisón Ascacíbar
Grupo EDMANS.  URL: http://www.mineriadatos.com
Área de Proyectos de Ingeniería. Departamento de Ingeniería Mecánica
Edificio Departamental. ETSII de Logroño. C/ Luís de Ulloa, 20, 26004 Logroño (España).
Phone: +34 941 299 232
Fax: + 34 941 299 794
E-mail: fjmartin@unririoja.es.