

1 **TITLE**

2 Applying data mining and Computer Vision Techniques to MRI to estimate
3 quality traits in Iberian hams.

4 **AUTHORS**

5 Trinidad Pérez-Palacios^{a,*}, Daniel Caballero^b, Andrés Caro^b, Pablo G.
6 Rodríguez^b, Teresa Antequera^a.

7 ^a Tecnología de los Alimentos, Facultad Veterinaria, Universidad de
8 Extremadura, Av. Universidad s/n, 10003 Cáceres, Spain

9 ^b Departamento de Ingeniería de Sistemas Informáticos y Telemáticos, Escuela
10 Politécnica, Universidad de Extremadura, Av. Universidad s/n, 10003 Cáceres,
11 Spain

12

13 * Corresponding author. Tel.: +34 927 257123; fax: +34 927 257110. e-mail
14 address: triny@unex.es (T. Pérez-Palacios).

15

16 **ABBREVIATIONS**

17 KDD: knowledge discovery in databases; R: raw hams; SA: end of salting; PS:
18 end of post-salting; D: end of drying; DC: dry-cured hams; P-C: physico-
19 chemical; MRI-CVT: Magnetic Resonance Imaging and Computer Vision
20 Techniques; B: *Biceps femoris* muscle; S: *Semimembranosus* muscle.

21 **ABSTRACT**

22 This study aims to forecast quality characteristics of Iberian hams by using
23 non-destructive methods of analysis and data mining. Magnetic Resonance
24 Imaging and Computer Vision Techniques were conducted on hams
25 throughout their processing. Physico-chemical parameters were also measured
26 in these products. Information from these analyses was integrated in a
27 database. First, deductive techniques of data mining were applied to these
28 data. Multiple linear regression allows for the estimation of information from
29 magnetic resonance imaging, Computer Vision Techniques and physico-
30 chemical analysis. This enables the completion of the initial database. Then,
31 predictive techniques of data mining were applied. Both, multiple linear
32 regression and isotonic regression achieved the prediction of weight, moisture
33 and lipid content of hams as a function of features obtained by Magnetic
34 Resonance Imaging and Computer Vision Techniques. Thus, data mining,
35 magnetic resonance imaging and Computer Vision Techniques could be used
36 to estimate the quality traits of Iberian hams. This allows for the improvement of
37 the process control without destroying any piece.

38 **KEY WORDS**

39 Data mining; MRI and Computer Vision Techniques; deduction; prediction;
40 quality parameters; Iberian ham.

41

42 1. INTRODUCTION

43 Quality attributes of dry-cured hams depend on characteristics of raw
44 material and processing conditions. Throughout the processing of hams,
45 changes on the physico-chemical (P-C) characteristics of the thighs take
46 place, also influencing the quality of the final product. Thus, not only
47 characteristics of thighs but also their modifications during the processing are
48 important parameters to control the technological process of dry-cured hams
49 (Pérez-Palacios et al., 2011a).

50 Temperature and relative humidity conditions during the processing lead
51 to ham dehydration and, hence, to weight loss. The ham industry estimates the
52 optimal ripening time by the percentage of weight loss, related to the amount
53 of water contained in the ham muscles (Martín et al., 1998). Raw material for
54 ham production should contain plenty of intramuscular fat, which is an
55 important characteristic, due to its positive influence on quality parameters on
56 the final product, such as marbling, juiciness, odour, and aroma (Ruiz et al.,
57 2002).

58 Usual methods for evaluation of the P-C characteristics (i.e. weight loss,
59 moisture, fat content) of dry-cured hams throughout the whole processing are
60 tedious and time-consuming, and sometimes involve the destruction of the
61 pieces. In this sense, the use of non-destructive techniques, such as computed
62 tomography (CT), near infra-red reflectance spectroscopy (NIRs) and Magnetic
63 Resonance Imaging (MRI), has been proposed for determining quality
64 parameters in this product. Studies on salt content by means of CT have been
65 carried out by several authors (Fulladosa et al., 2010; Haseth et al., 2012;
66 Picouet et al., 2013; Santos-Garcés et al., 2010; Vestergard et al., 2005). CT has
67 also been applied for predicting the water content throughout the process of
68 hams (Fulladosa et al., 2010; Santos-Garcés et al., 2010), and the weight and
69 lean content of the raw material (Picouet et al., 2010). In pig carcass, Furnols et
70 al. (2009) estimated the lean meat content by using CT. Collell et al. (2011) used
71 NIRs to predict moisture, water activity and NaCl content at the surface of dry-
72 cured ham during the process. Results obtained by Pérez-Juan et al. (2010)

73 showed the accuracy of NIRs to predict the fatty acid composition of ham
74 subcutaneous fat.

75 MRI is a non-destructive, non-invasive, non-intrusive, **non-ionizing** and
76 innocuous technique. Thus, as an alternative to P-C procedures, MRI has also
77 been proposed to study some characteristics in hams. Fantazinni et al. (2009)
78 used this technique to obtain information on moisture and salt distribution
79 during the processing of Parma hams. Recently, predictive models have been
80 proposed for estimating water activity, moisture, salt content and proteolysis
81 extent in S. Daniele hams on the basis of the MR signal intensity (Manzoco et al.,
82 2013).

83 The implementation of active contours in MRI can be used to recognize
84 the *Biceps femoris* and *Semimembranosus* muscles in Iberian hams, determine
85 the volume of the muscle and estimate ham weight and moisture (Antequera
86 et al., 2007; Caro et al., 2001). MRI and computational texture features allowed
87 for the classification of fresh and dry-cured Iberian hams as a function of pig
88 feeding background (Pérez-Palacios et al., 2010a, 2011b). Sensory traits in
89 Iberian dry-cured hams were predicted from computational texture
90 characteristics obtained from MRI of fresh hams (Pérez-Palacios et al., 2010b).

91 The calculation of intramuscular fat levels of Iberian ham has also been
92 attempted by using MRI applications (Ávila et al., 2005; Caro et al., 2003),
93 obtaining reasonable, but not very high, correlation coefficients (around 0.50 –
94 0.63), which shows the potential of this technique for determining intramuscular
95 fat level in Iberian hams. In these studies, database obtained from P-C, and
96 MRI and computational analysis are processed by applying usual statistical
97 tools such as Pearson's correlation coefficients or principal components analysis
98 (Pérez-Palacios et al., 2010b, 2011b). The integration of heterogeneous P-C
99 information with computer vision data, and the analysis of this new data set by
100 data management and database applications would be innovative and could
101 give accurate results, playing an increasing role in furthering food research
102 (Cortez et al., 2009; Holmes et al., 2012).

103 Data mining is an important part of a larger process known as KDD
104 ("Knowledge Discovery in Databases") (Fayyad et al., 1996). It is associated with

105 large data. The main goal of data mining consists in extracting hidden
106 information from a data set. This can be achieved by the automatic or semi-
107 automatic analysis of large amounts of data, which allows for the extraction of
108 interesting and previously unknown patterns (Hastie et al., 2001). These patterns
109 can be groups of data records (cluster analysis), unusual records (anomaly
110 detection) and dependencies among data (association rules). Thus, the
111 patterns can be seen as a summary of the input data, and can be used for
112 further analysis.

113 Interest in data mining has recently grown because of the rapidly
114 decreasing cost of large storage devices and increasing ease in data
115 collection over networks. Other factors include, the development of robust and
116 efficient algorithms to process this data, and the increase in computing power,
117 enabling the use of intensive computational methods for data analysis
118 (Mitchell, 1999).

119 To our knowledge, few studies apply data mining to food. Song et al.,
120 (2002) and Cortez et al. (2006) used this computing technique to predict quality
121 traits in beef and lamb, respectively. It has also been used to predict the
122 oxidation of menhaden fish oil (Klaypradith et al., 2010) or to model wine
123 preferences (Cortez et al., 2009). Holmes et al. (2012) applied data mining to
124 detect fruit and vegetables contaminated with pesticide and to identify these
125 products as a function of their home country.

126 For this study, data obtained from the MRI-CVT (volume) and P-C analysis
127 (moisture, lipid and weight) of a homogeneous Iberian ham batch were used
128 to construct a database. Several data mining techniques were applied to this
129 database in order to i) estimate values for the analysed parameters in a higher
130 number of samples and ii) predict moisture, lipid content and weight
131 throughout the processing of the Iberian ham.

132

133 **2. MATERIAL AND METHODS**

134 **2.1. Experimental design**

135 This study was carried out with 15 Iberian thighs, which were processed
136 following the traditional processing as described in Antequera et al., (2007).
137 Four stages were considered: raw hams (R), 0 days; end of post-salting (PS), 90
138 days; end of drying (D), 270 days; and dry-cured hams (DC), 660 days. This
139 experimental design is shown in Figure 1.

140 At each stage, 6 hams were scanned for obtaining MR images. After then,
141 three hams were destroyed at each stage for the P-C analysis, having 12, 9 and
142 6 hams at PS, D and DC stages, respectively. Ham weights were recorded at
143 these four stages and also at the end of the salting step (SA). This is how the P-C
144 data set is formed.

145 In this work MRI has been used as a non-invasive technique only to
146 acquire images of the hams without destroying them. Then, our own active
147 contour algorithms were applied to recognize the *Biceps femoris* and
148 *Semimembranosus* muscles, in order to compute their volumes, as described in
149 Antequera et al. (2007).

150 Numerical data is extracted by Data Mining from the data sets obtained
151 by our MRI-CVT and from the data sets obtained by P-C. Figure 2 describes the
152 whole process.

153 **2.2. MRI acquisition**

154 Magnetic resonance images were generated at the "Infanta Cristina"
155 University Hospital (Badajoz, Spain). A MRI scanner (Philips Gyroscan NT Intera
156 1.5 T) was used, with a quadrature whole-body coil. Sequences of T1 were
157 applied with the following parameters: 120 x 85 mm for field-of view (FOV), 20
158 ms for echo time (TE), 500 ms for repetition time (TR), 2 mm thick slices, 90° for flip
159 angle, i.e. a T1-weighted spin echo (SE), 0.23 x 0.20 mm per pixel resolution. Sixty
160 slices per ham piece were obtained. The MRI acquisition was done at 20° C
161 and it took 28 min for each ham. All the images were in DICOM format, with a
162 512 x 512 resolution, and 256 grey levels.

163 **2.3. Computer Vision Techniques**

164 After the images were acquired, our own computer vision algorithms were
165 applied to extract numerical data from these images. Then, data mining
166 techniques were tested over these data to obtain prediction equations.

167 The automated procedure was run as described in Figure 2. First, a
168 previous image pre-processing stage was carried out. Then, the *Biceps Femoris*
169 and *Semimembranosus* muscles (B and S, respectively) were recognized
170 distinctly by using Active Contours, applying a greedy algorithm method
171 (Antequera et al., 2007). The surface and volume for all the contours is
172 calculated by relying on classical methods in analytical geometry. **Volume is**
173 **expressed in voxel (volume per element), which is $0.23 \times 0.2 \times 2 \text{ mm}^3$.**

174 **2.4. Physico-chemical analysis**

175 At each stage of the processing, ham weight was recorded and the B
176 and S muscles of three hams were dissected, weighed and analyzed for
177 moisture (AOAC, 2000; reference 935.29) and lipid content (Pérez-Palacios et
178 al., 2008). Analyses were done in triplicate.

179 **2.5. Data mining**

180 The free software WEKA (Waikato Environment for Knowledge Analysis)
181 (<http://www.cs.waikato.ac.nz/ml/weka/>) was used for carrying out the data
182 mining analysis. The primary groups in data mining tasks are descriptive and
183 predictive techniques. The first ones include deductive techniques, which have
184 the ability to infer new values based on actual data. In predictive techniques,
185 future models can be predicted from current data by trend analysis (Witten &
186 Frank, 2005; Wu et al., 2008). Both, descriptive and predictive techniques were
187 applied in this study.

188 Multiple linear regression was used for the deductive tasks. The dependent
189 variable to be estimated was always unique and numerical and this method
190 enables the removal of collinear attributes. In addition, regression techniques
191 seem to be the most appropriate to forecast values, as it allows inferring
192 numerical data from the available numerical values. The M5 method of
193 attribute selection and a ridge value of 1×10^{-4} were applied. This method steps
194 through the attributes, and removes the one with the smallest standardised

195 coefficient until no improvement is observed in the estimate of the error given
196 by the Akaike information criterion (Hastie et al., 2001).

197 Again, multiple linear regression was used for the experiments of
198 prediction. This technique obtains a linear regression equation, which can be
199 used to predict future values (Hastie et al., 2001). The M5 method of attribute
200 selection and a ridge value of 1×10^{-4} were also applied.

201 Isotonic regression was also tested for prediction. When the values of the
202 database are highly correlated, the use of non-linear regression is
203 recommended. In these cases, the isotonic regression is considered as a good
204 option. Isotonic regression provides a set of values from the information stored
205 on a database. It is based on estimating ordered values for an independent
206 variable (i.e. weight) as a function of one of the input parameters (attributes of
207 the database). Thus, the ham weight is predicted as a function of the volume or
208 the maturation stage. Only the input parameters providing better adjustment
209 results (for example, the stage) will be selected. Finally, an interpolation
210 function is established (polynomial trend line) to compare the provided set
211 data with original values in the database, obtaining the prediction equation
212 (Borge, 1985; Barlow, 1972).

213 **2.6. Databases**

214 An initial database was built with data obtained throughout the ham
215 processing: i) stage of the ham processing, ii) P-C analysis (ham, B and S
216 weight; moisture and lipid content of the B and S), and iii) MRI-CVT (ham, B and
217 S volume) (see Figure 3).

218 As previously explained (Figure 1), this study was carried out with 15 Iberian
219 hams and three of them were discarded at each stage. Thus, the number of
220 pieces at R, SA, PS, D and DC stages were 15, 12, 12, 9 and 6, respectively. The
221 initial database contained 54 records, with each record treated as a data set
222 obtained from a ham. Although this database might be regarded as small, it
223 should be noted that each Iberian ham presents considerable costs, about 30
224 Euros per kilo plus lab work.

225 Since the 15 hams were not analysed at all the ripening stages, this initial
226 database presents incomplete records (Figure 3A). After applying data mining
227 techniques (multiple linear regression), the values for all analysed parameters
228 were estimated. The records thus completed made up the whole database, as
229 can be observed in Figure 3B.

230 **2.7. Statistical design**

231 Differences throughout the processing of Iberian hams with parameters
232 determined by P-C analysis and MRI-CVT were analysed by one-way analysis of
233 variance (ANOVA). When significant differences ($p < 0.05$) were found, the
234 Tukey's test was conducted. Analyses were done by using the SPSS package
235 (v.18.0).

236

237 **3. RESULTS AND DISCUSSION**

238 **3.1. Physico-chemical and MRI-Computer Vision Techniques.**

239 Table 1 shows results on ham weight, moisture content, lipid content, and
240 weight of B and S muscles in Iberian hams throughout the processing. Weight
241 and lipid content in B are known to be greater in comparison to the S muscle
242 (Pérez-Palacios et al., 2008b, 2010c), which is corroborated in this study. As
243 expected, ham and muscle weight and moisture decreased during the
244 processing due to water loss (Martín et al., 1998; Pérez-Palacios et al., 2011b).
245 The significant increase in the percentage of lipid content in the B muscle
246 during the processing can be also related to the water loss, since the
247 percentage of dry matter (as fat) increased as the water content decreased.

248 Moisture loss during the processing occurs more in the S muscle than in B,
249 above all at the last stages of the processing. This fact agrees with previous
250 results in Andrés et al. (2005). This phenomenon is related to muscle location in
251 the ham (the B muscle is an internal muscle, while S is external), since water loss
252 takes place from the inner to the outer part. Thus, water loss is facilitated in
253 external muscles, such as S.

254 Volume of ham, B and S muscles at R, PS, D and DC stages achieved by
255 MRI-CVT are shown in Table 2. These three objects of study decreased during
256 processing, which coincides with the changes found in ham and muscle
257 weight. The accuracy of volume estimation for the muscles is very high, as can
258 be examined in Antequera et al. (2007). There was a high correlation ($R^2 =$
259 0.992) between the data obtained by physical measurement and sizes
260 measured on MRI by computer vision methods.

261 **3.2. Data mining for deduction.**

262 As previously explained, a database with 54 records was built (Figure 3A).
263 A record is the data set of a ham, which includes i) the stage of the processing,
264 ii) data from P-C analysis (ham weight, B and S muscles weight, moisture and
265 lipid content), and iii) data from MRI-CVT (ham, B and S muscles volume). Most
266 of the records in the database were incomplete. By applying multiple linear
267 regression, the unknown information of the records in the database is
268 estimated. Hence, a database of 54 full records is computed (Figure 3B). This
269 process could be seen as a type of data reconstruction: data that did not exist
270 is reconstructed by using various algorithms with some degree of confidence.

271 Correlation index R^2 is used to prove the correctness and precision of the
272 estimated values by using multiple linear regression. Table 3 shows the
273 correlation coefficients between real and predicted data for the features
274 analysed: ham weight; B and S muscles weight, moisture and lipid content;
275 ham, B and S muscle volume. As can be seen, high correlations ($R^2 > 0.900$)
276 have been obtained for all traits, except for lipid content of the S muscle ($R^2 =$
277 0.665). This lower correlation could be related to the high variability of fat
278 content in Iberian ham. Particularly noteworthy is the high correlation obtained
279 for moisture in the two muscles (> 0.990).

280 Table 4 displays the value range of the predicted features, which can be
281 compared to the real values shown in Table 1, for the P-C characteristics, and
282 Table 2, for data obtained by MRI-CVT, in order to corroborate the good
283 correlation between real and predicted data. For example, at the R stage, the
284 average moisture of the B muscle was 68.69 % (BM value at Raw in Table 1) and
285 the values predicted for this characteristic range between 64.97 and 71.95 %

286 (BM value at Raw in Table 4); at the D stage, the real value for ham volume was
287 64.50 voxel, and its predicted values were 58.94-67.21 voxel.

288 To the best of our knowledge, deductive methods from data mining
289 techniques have not been applied at all in food science. This fact is really
290 important since this approach yields a large number of data from a small and
291 incomplete database. In the case of Iberian ham production, the application
292 of deductive methods of data mining would be an interesting tool due to the
293 high cost of this product.

294 **3.3. Data mining for prediction.**

295 The prediction of ham quality parameters (weight, moisture content, and
296 lipid content in the B and S muscles) was also tested. Predictive techniques from
297 data mining were applied to information retrieved from MRI-CVT (BV, SV and
298 HV) procedures. Two methods in data mining were used, multiple linear
299 regression and isotonic regression.

300 To validate the predicted results, the coefficient correlation R^2 of the two
301 explored data mining methods was computed (Table 5). For weight, moisture
302 content in B and S muscles and lipid content in B, high correlation coefficients
303 (0.87- 0.99) were obtained. Very few differences were found between
304 correlation coefficients achieved by multiple linear regression and isotonic
305 regression methods. The computational cost of both techniques is similar, and
306 yet, isotonic regression is not automatic and needs a subsequent interpolating
307 step by using a spreadsheet. Thus, the use of multiple linear regression for
308 deducing these P-C parameters seems to be more comfortable.

309 In the case of lipid content in the S muscle, no good correlations were
310 obtained when applying multiple linear regression, but accurate results were
311 achieved ($R^2 = 0.817$) with isotonic regression. As previously explained, this could
312 be related to the high variability of fat content in Iberian ham. In fact, the use
313 of isotonic regression is indicated when having non-linear dependent data
314 (Barlow et al., 1972).

315 Figure 4 presents the adjustment between real and predicted values of
316 lipid content in the S muscle by the two deductive techniques applied in this

317 study. Isotonic regression shows higher accuracy in comparison to multiple
318 linear regression for predicting the lipid content of S.

319 Table 6 shows prediction equations for weight, moisture and lipid content
320 in the B and S muscles by multiple linear regression and isotonic regression. Thus,
321 by using data obtained non-destructively by MRI-CVT (HV, BV and SV) weight,
322 moisture and lipid content can be now reliable estimated. These
323 determinations have always been carried out in Iberian hams, but the
324 traditional methods are time-consuming and require the destruction of the
325 sample. Therefore, our equations could be considered as a useful tool.

326 **CONCLUSIONS**

327 **To the best of our knowledge this work has been the first** to apply data
328 mining to Iberian ham information obtained from P-C analysis, weight, moisture
329 and lipid content, and MRI-CVT techniques, volume.

330 The application of deductive techniques from data mining, multiple linear
331 regression, to information from MRI-CVT and P-C analysis allows for the
332 accurate estimation of more records of the analysed traits: weight, moisture
333 content, lipid content, and volume in Iberian hams.

334 Multiple linear regression and isotonic regression are accurate methods of
335 data mining for predicting weight, moisture and lipid content in Iberian ham as
336 a function of features obtained from MRI-CVT techniques.

337 Data mining and MRI-CVT have been used as a pioneering approach to
338 study the features of hams. These tools can be useful for calculating P-C
339 parameters related to ham quality and for improving the control of the
340 processing without destroying meat pieces.

341

342 **ACKNOWLEDGMENTS**

343 The authors wish to acknowledge the funding received for this research
344 from both the Junta de Extremadura (Regional Government Board – Research
345 Projects 3PR05B027 and PDT08A021; Consejería de Economía, Comercio e

346 Innovación and FEDER– economic support for research groups: GRU09148 and
347 GRU09025) and from the Spanish Government (National Research Plan) and
348 the European Union (FEDER funds) by means of the grant reference TIN2008-
349 03063. We also wish to thank the “Hermanos Roa” company from Villar del Rey
350 (Badajoz), as well as the “Infanta Cristina” University Hospital Radiology Service,
351 specially to the Dr. Ramón Palacios, for their contribution and support.

352

353 **CONFLICTS OF INTERESTS**

354 Authors state that there are no known conflicts of interest associated with this
355 publication.

356

357 **REFERENCES**

- 358 Andres, A.I., Ventanas, S., Ventanas, J., Cava, R., & Ruiz, J. (2005). Physicochemical
359 changes throughout the ripening of dry cured hams with different salt content
360 and processing conditions. *European Food Research and Technology*, 221, 30–35.
- 361 Antequera, T., Caro, A., Rodriguez, P.G., & Pérez-Palacios, T. (2007). Monitoring the
362 ripening process of Iberian Ham by computer vision on magnetic resonance
363 imaging. *Meat Science*, 76, 561-567.
- 364 Association of Official Analytical Chemist (AOAC) (2000). Official Methods of Analysis of
365 the Association of Official Analytical Chemists. 17th ed. Gaithersburg, Maryland.
- 366 Ávila, M., Durán, M.L., Caro, A., Antequera, T., & Gallardo, R. (2005). Thresholding
367 methods on MRI to evaluate intramuscular fat level on Iberian ham. Lectures
368 Notes in Computer Science (LNCS 3523). *Pattern Recognition and Image Analysis*,
369 697–704.
- 370 Barlow, R.E., Bartholomew, D., Bremner, J.M. & Brunk, H.D (1972). *Statistical inference*
371 *under order restriction: The theory and application of Isotonic Regression*. Wiley,
372 New York.
- 373 Borge, L. (1985). Estimacion y contrastes de hipótesis en el modelo lineal general con
374 restricciones de desigualdad. Doctoral thesis. University of Valladolid, Spain.
- 375 Caro, A., Durán, M.L., Rodríguez, P., Antequera, T., & Palacios, R. (2003). Mathematical
376 morphology on MRI for the determination of Iberian ham fat content. Lecture
377 Notes in Computer Science (LNCS 2905). *Progress in Pattern Recognition, Speech*
378 *and Image Analysis*, 359-366.

379 Caro, A., Rodríguez, P.G., Cernadas, E., Durán, M.L., & Villa, D. (2001). Applying active
380 contours to muscle recognition in Iberian ham MRI. In *IASTED International*
381 *Conference Signal Processing, Pattern Recognition and Applications*, Rhodes,
382 Greece.

383 Collell, C., Gou, P., Arnau, J., & Comaposada, J. (2011). Non-destructive estimation of
384 moisture, water activity and NaCl at ham surface during resting and drying using
385 NIR spectroscopy. *Food Chemistry*, 129, 601–607.

386 Cortez, P., Cedeira, A., Almeida, F., Matos, T., & Reis, J. (2009). Modeling wine
387 preferences by data mining from physicochemical properties. *Decision Support*
388 *System*, 47, 547-553.

389 Cortez, P., Portelinha, S., Rodrigues, S., Cadavez, V., & Teixeira, A. (2006). Lamb Meat
390 Quality Assessment by support vector machines. *Neural Processing Letters*, 24, 41-
391 51.

392 Fantazinni, P., Gombia, M., Schembri, P., Simoncini, N., & Virgili, R. (2009). Use of
393 Magnetic Resonance Imaging for monitoring Parma dry-cured ham processing.
394 *Meat Science*, 82, 219-227.

395 Fayyad, U., Pietetsky-Shapiro, G., & Smyth, P. (1996). From data mining to knowledge
396 discovery in databases. *AI Magazine*, 17, 37-54.

397 Fulladosa, E., Santos-Garcés, E., Picouet, P., & Gou, P. (2010). Prediction of salt and
398 water content in dry-cured hams by computed tomography. *Journal of Food*
399 *Engineering*, 96, 80–85.

400 Furnols, M.F., Teran, M.F., & Gispert, M. (2009). Estimation of lean meat content in pig
401 carcasses using X-ray Computed Tomography and PLS regression. *Chemometrics*
402 *and Intelligent Laboratory Systems*, 98, 31–37.

403 Haseth, T. T., Sørheim, O., Høy, M., & Egelanddal, B. (2012). Use of computed
404 tomography to study raw ham properties and predict salt content and
405 distribution during dry-cured ham production. *Meat Science*, 90, 858–864.

406 Hastie, T., Tibshirani, R., Friedman, J. (2001). *The Elements of Statistical Learning: Data*
407 *Mining, Inference and Prediction*. Springer-Verlag, New York.

408 Holmes, G., Fletcher, D., Hernández, J., Ramírez, M.J., & Ferri, C. (2007). *Introducción a la*
409 *minería de datos*. Prentice-Hall.

410 Reutermann, P. (2012). An application of data mining to fruit and vegetable sample
411 identification using gas chromatography-mass spectrometry. In *Proceedings of*
412 *International Congress of Enviromental Modelling and Software Managing*
413 *Resources of a Limited Planet*, Leipzig, Germany.

414 Klaypradith, W., Kerdpiboon, S., & Singh, R.K. (2010). Application of artificial neural
415 networks to predict the oxidation of menhaden fish oil obtained from Fourier

416 transform infrared spectroscopy method. *Food and Bioprocess Technology*, 4,
417 475-480.

418 Manzoco, L., Anese, M., Marzona, S., Innocente, N., Lagazio, C., & Nicoli, M. C. (2013).
419 Monitoring dry-curing of S. Daniele ham by magnetic resonance imaging. *Food*
420 *Chemistry*, 141, 2246-2252.

421 Martin, L., Córdoba, J.J., Antequera, T., Timon, M.L., & Ventanas, J. (1998). Effects of salt
422 and temperature on proteolysis during ripening of Iberian ham. *Meat Science*, 49,
423 145-53.

424 Mitchell, T.M. (1999). Machine learning and data mining. *Communications of the ACM*,
425 42, 30-36.

426 Pérez-Juan, M., Afseth, N.K., González, J., Díaz, I., Gispert, M., Furnols, M.F., Oliver, M.A.,
427 & Realini, C.E. (2010). Prediction of fatty acid composition using a NIRS fibre optics
428 probe at two different locations of ham subcutaneous fat. *Food Research*
429 *International*, 43, 1416-1422.

430 Pérez-Palacios, T., Ruiz, J., Martín, D., Barat, J. M., & Antequera, T. (2011a). Pre-cure
431 freezing effect on physicochemical, texture and sensory characteristics of Iberian
432 ham. *Food Science and Technology International*, 17, 127-133.

433 Pérez-Palacios, T., Antequera, T., Durán, M. L., Caro, A., Rodríguez, P. G., & Palacios, R.
434 (2011b). MRI-based analysis of feeding background effect on fresh Iberian ham.
435 *Food Research International*, 43, 248-254.

436 Pérez-Palacios, T., Antequera, T., Durán, M.L., Caro, A., Rodríguez, P.G., & Ruiz, J.
437 (2010a). MRI-based analysis, lipid composition and sensory traits for studying
438 Iberian dry-cured hams from pigs fed with different diets. *Food Chemistry*, 126,
439 1366-1372.

440 Pérez-Palacios, T., Antequera, T., Molano, R., Rodríguez, P.G., & Palacios, R. (2010b).
441 Sensory traits prediction in dry-cured hams from fresh product via MRI and lipid
442 composition. *Journal of Food Engineering*, 101, 152-157.

443 Pérez-Palacios, T., Ruiz, J., Dewettinck, K., Trung Le., T., & Antequera, T. (2010c).
444 Individual phospholipid classes from Iberian pig meat as affected by diet. *Journal*
445 *of Agricultural and Food Chemistry*, 58, 1755-1760.

446 Pérez-Palacios, T., Ruiz, R., Martin, D., Muriel, E., & Antequera, T. (2008a). Comparison of
447 different methods for total lipid quantification. *Food Chemistry*, 110, 1025-1029.

448 Pérez-Palacios, T., Ruiz, J., & Antequera, T. (2008b). Perfil de ácidos grasos de la grasa
449 subcutánea e intramuscular de credos ibéricos cebados en montanera y con
450 pienso "alto oleico". *Eurocarne*, 1-10.

451 Picouet, P. A., Gou, P., Fulladosa, E., Santos-Garcés, E., & Arnau, J. (2013). Estimation of
452 NaCl diffusivity by computed tomography in the Semimembranosus muscle

453 during salting of fresh and frozen/thawed hams. *LWT - Food Science and*
454 *Technology*, 51, 275-280.

455 Ruiz, J., García, C., Muriel, E., Andrés, A.I., & Ventanas, J. (2002). Influence of sensory
456 characteristics on the acceptability of dry-cured ham. *Meat Science*, 61, 347-354.

457 Santos-Garcés, E., Gou, P., Garcia-Gil, N., Arnau, J., & Fulladosa, E. (2010). Non-
458 destructive analysis of aw, salt and water in dry-cured hams during drying process
459 by means of computed tomography. *Journal of Food Engineering*, 101, 187-192.

460 Song, Y.H., Kim, S.J., & Lee, S.K. (2002). Evaluation of ultrasound for prediction of carcass
461 meat yield and meat quality in Korean native cattle. *Asian Journal Animal*
462 *Science*, 15, 591-595.

463 Vestergaard, C., Erbou, S. G., Thauland, T., Adler-Nissen, J., & Berg, B. (2005). Salt
464 distribution in dry-cured ham measured by computed tomography and image
465 analysis. *Meat Science*, 69, 9-15.

466 Witten, I.H., & Frank, E. (2005). *Data Mining: Practical Machine Learning Tools and*
467 *Techniques with Java Implementations*. Morgan Kaufmann, San Francisco.

468 Wu, X., Kumar, V., Ross Quinlan, J., Ghosh, J., Ghosh, J., Yang, Q., Motoda, H.,
469 McLachlan, G.J., Ng, A., Liu, B., Yu, P.S., Zhou, Z-H., Steinbach, M., Hand, D.J., &
470 Steinberg, D. (2008). Top 10 algorithms in data mining. *Knowledge and*
471 *Information Systems*, 14, 1-37.

472

FIGURE CAPTIONS

Caption Figure 1. Sampling throughout the Iberian ham processing for the physico-chemical analysis (P-C) and the MRI acquisition.

Caption Figure 2. Acquisition of Iberian ham data (from physico-chemical analysis and MRI and computer vision techniques) used to estimate quality parameters by applying data mining.

Caption Figure 3. Initial database with incomplete records (A) and with all records filled after applying data mining (B). The data set of each record is composed by i) processing stages (Stage) (raw hams = 1; salting = 1.5; post-salting = 2; drying = 3; dry-cured ham = 4), ii) physico-chemical parameters (ham weight = HW, *Biceps femoris* muscle weight, moisture and lipid content = BW, BM and BL, respectively, *Semimembranosus* muscle weight, moisture and lipid content = SW, SM and SL, respectively), and iii) MRI and computer vision techniques (ham, *Biceps femoris* and *Semimembranosus* volume = HV, BV, and SV, respectively).

Suspension dots indicate that the database is greater and it has been cut.

N = ham identifier.

HW, BW and SW are expressed in grams; BM, BL, SM and SL are expressed in g/100g sample; HV, BV and SV are expressed in voxel.

Caption Figure 4. Adjustment between real (\blacklozenge) and predicted values of the lipid content of *Semimembranosus* muscle by using multiple linear regression (---) and isotonic regression (—) as a function of the *Semimembranosus* volume (expressed in voxel).

Figure 1.

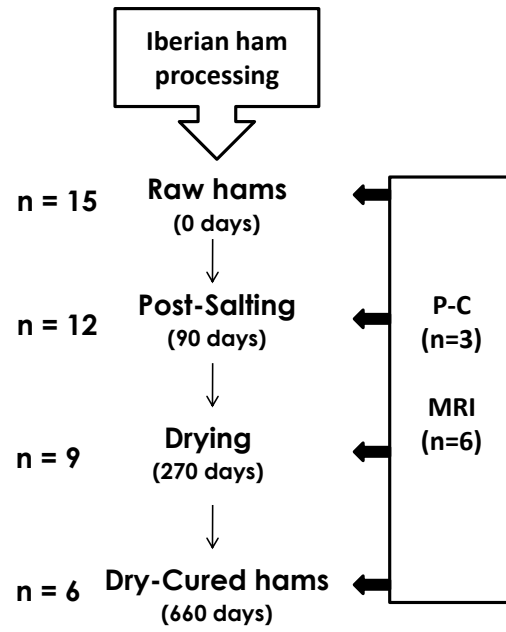


Figure 2.

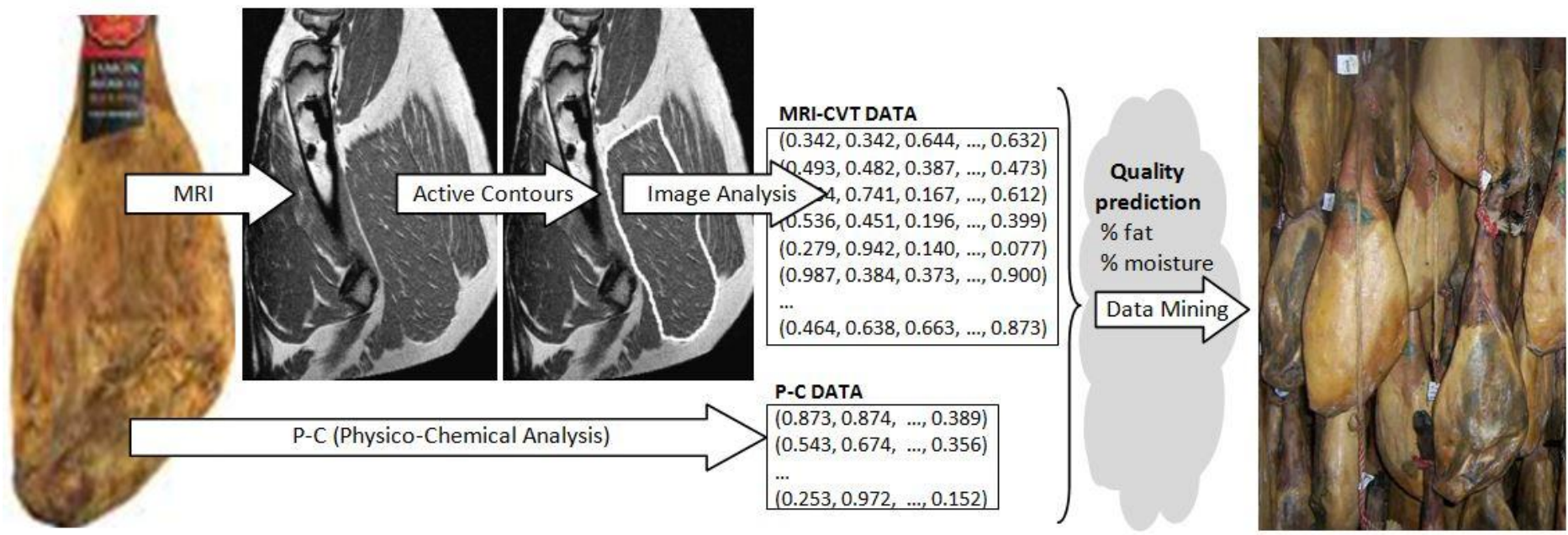


Figure 3

Figure 3.

A

N	Stage	HW	BW	BM	BL	SW	SM	SL	HV	BV	SV
3	1	10,800									
4	1	11,000									
8	1	11,000									
12	1	11,200							80,542	19,614	24,639
13	1	11,200									
17	1	10,800							78,728	20,838	23,785
19	1	10,600									
22	1	11,200									
23	1	11,000							82,764	20,426	24,571
24	1	10,600	1,235	69.16	9.18	725	72.57	4.14			
27	1	10,800							84,275	19,025	23,073
28	1	11,000							82,220	20,114	31,994
32	1	11,200							80,624	20,856	25,150
34	1	11,000	1,435	64.97	9.52	755	70.36	4.24			
37	1	11,000	1,475	71.95	7.61	760	72.68	3.46			
3	1.5	10,600									
4	1.5	10,800									
8	1.5	10,800									
12	1.5	11,000									
13	1.5	10,900									
17	1.5	10,500									
19	1.5	10,300									
22	1.5	10,900									
23	1.5	10,800									
27	1.5	10,600									
28	1.5	10,800									
32	1.5	11,000									

⋮

B

N	Stage	HW	BW	BM	BL	SW	SM	SL	HV	BV	SV
3	1	10,800	1,391	69.37	8.49	728	69.92	6.18	68,276	20,244	23,632
4	1	11,000	1,388	69.17	8.32	724	73.19	6.27	71,383	20,441	23,176
8	1	11,000	1,386	69.15	8.32	715	72.80	6.27	73,473	20,315	23,704
12	1	11,200	1,392	69.14	7.92	712	73.96	6.40	80,542	19,614	24,639
13	1	11,200	1,393	69.39	8.01	733	72.49	6.16	79,040	20,354	25,425
17	1	10,800	1,378	69.27	8.66	724	70.56	5.95	78,728	20,838	23,785
19	1	10,600	1,323	69.29	8.99	720	69.82	6.12	76,628	20,089	23,866
22	1	11,200	1,402	68.99	8.10	721	73.89	5.99	83,822	20,123	26,180
23	1	11,000	1,392	69.36	8.22	731	71.20	6.14	82,764	20,426	24,571
24	1	10,600	1,235	69.16	9.18	725	72.57	4.14	74,918	20,384	22,895
27	1	10,800	1,385	68.96	8.12	695	70.63	6.29	84,275	19,025	23,073
28	1	11,000	1,380	69.08	9.06	732	71.32	5.64	82,220	20,114	31,994
32	1	11,200	1,382	69.01	8.20	749	72.68	6.19	80,624	20,856	25,150
34	1	11,000	1,435	64.97	9.52	755	70.36	4.24	84,922	19,949	26,507
37	1	11,000	1,475	71.95	7.61	760	72.68	3.46	92,725	20,027	27,599
3	1.5	10,600	1,296	65.68	8.78	664	67.33	6.58	69,061	18,083	21,587
4	1.5	10,800	1,330	64.24	9.44	670	68.05	6.57	72,089	18,680	21,736
8	1.5	10,800	1,315	63.98	9.44	662	67.91	6.54	74,128	18,643	22,272
12	1.5	11,000	1,262	65.87	9.01	702	70.82	6.41	76,602	18,497	23,505
13	1.5	10,900	1,442	66.11	8.96	678	69.55	6.66	72,824	18,214	23,983
17	1.5	10,500	1,246	65.20	8.99	656	65.24	6.48	75,279	18,044	22,087
19	1.5	10,300	1,279	64.95	9.72	654	65.09	6.27	73,231	18,449	22,648
22	1.5	10,900	1,290	65.05	8.98	661	68.61	6.57	83,029	18,491	24,184
23	1.5	10,800	1,347	65.23	9.79	689	67.58	6.91	79,061	18,512	24,112
27	1.5	10,600	1,271	65.45	9.39	682	67.11	6.28	80,382	18,479	24,588
28	1.5	10,800	1,285	65.18	9.05	674	66.86	6.31	85,031	18,628	25,043
32	1.5	11,000	1,292	65.04	8.91	657	68.90	6.56	89,310	18,470	25,793

⋮

Figure 4.

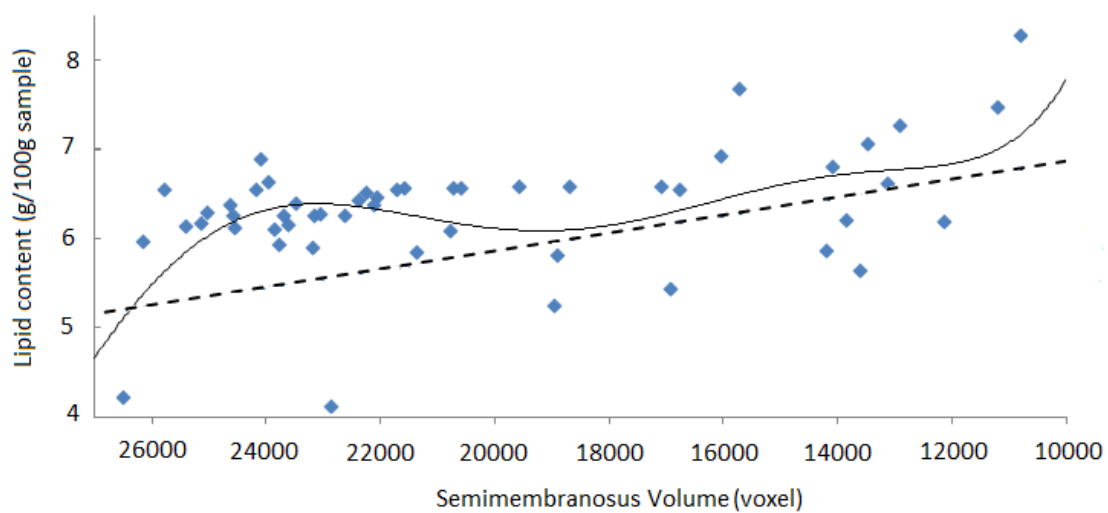


Table 1

Table 1. Results on physico-chemical analysis (ham weight = HW, *Biceps femoris* muscle weight, moisture and lipid content = BW, BM and BL, respectively, *Semimembranosus* muscle weight, moisture and lipid content = SW, SM and SL, respectively) at the different stages of the Iberian ham processing*.

	HW	BW	BM	BL	SW	SM	SL
	(g)	(g)	(g/100g sample)	(g/100g sample)	(g)	(g/100g sample)	(g/100g sample)
Raw hams	10960 ± 203a	1382 ± 129a	68.69 ± 3.17a	8.77 ± 0.98b	747 ± 19a	71.87 ± 1.34a	3.95 ± 0.45
Salting	10750 ± 211a	NA	NA	NA	NA	NA	NA
Postsalting	9683 ± 301b	1130 ± 30b	60.66 ± 1.16b	11.57 ± 1.79b	527 ± 15b	60.37 ± 1.65b	5.53 ± 0.32
Drying	8489 ± 401c	1030 ± 83b	54.43 ± 0.66c	10.63 ± 1.11b	552 ± 53b	34.24 ± 5.04c	6.66 ± 1.05
Dry-cured hams	7700 ± 110d	713 ± 21c	42.92 ± 2.49d	16.94 ± 1.44a	327 ± 70c	25.71 ± 2.76d	6.27 ± 1.99
p	<0.001	<0.001	<0.001	0.002	<0.001	<0.001	0.2

*Values are expressed as means ± standard deviation.

NA = not analysed.

In the same column, means with different letters differ significantly between stages.

Table 2. Results on MRI and Computer Vision Techniques (ham, *Biceps femoris* and *Semimembranosus* volume = HV, BV, and SV, respectively) at the different stages of the Iberian ham processing *.

	HV	BV	SV
	(voxel)	(voxel)	(voxel)
Raw hams	81520 ± 1950a	25530 ± 7200a	25530 ± 3240a
Salting	NA	NA	NA
Postsalting	75250 ± 2050b	21640 ± 8900b	21640 ± 1070b
Drying	64500 ± 2170c	15146 ± 1230c	15140 ± 1730c
Dry-cured hams	56990 ± 5630d	12130 ± 1270d	12130 ± 1710c
p	<0.001	<0.001	<0.001

*Values are expressed as means ± standard deviation.

NA = not analysed.

Table 3. Correlation coefficient (R^2) between real and predicted data obtained by data mining for the features analysed by physico-chemical analysis and MRI and Computer Vision Techniques.

	R^2
BW	0.975
SW	0.916
BM	0.994
SM	0.993
BL	0.908
SL	0.665
HV	0.975
BV	0.999
SV	0.993

See abbreviations in Figure 3.

Table 4. Minimum and maximum values for the features predicted by using data mining.

	BW	BM	BL	SW	SM	SL	HV	BV	SV
	(g)	(g/100g)	(g/100g)	(g)	(g/100g)	(g/100g)	(voxel)	(voxel)	(voxel)
Raw hams	1235-1475	64.97-71.95	7.61-9.52	695-760	69.82-73.96	3.46-6.40	68270-92720	19020-20850	22980-31990
Salting	1245-1442	63.98-66.11	8.78-9.79	654-702	65.09-70.82	6.27-6.91	69060-89310	18040-18680	21580-25790
Post-salting	1100-1187	59.79-62.31	59.79-62.31	510-608	53.80-62.38	5.26-6.60	57710-77050	15720-18220	16950-23190
Drying	880-1110	52.27-54.63	52.27-54.63	422-585	42.39-48.03	5.66-7.70	58940-67210	12120-15980	12160-17090
Dry-cured hams	690-756	40.66-45.88	40.66-45.88	256-423	29.11-40.21	3.88-8.29	48580-62330	8120-11340	9860-14200

See abbreviations in Figure 3.

Table 5. Correlation coefficient (R^2) for each physico-chemical characteristic predicted by applying data mining (multiple linear regression (MLR) and isotonic regression (IR)) on data achieved by MRI and Computer Vision Techniques (BV, SV and HV).

	BW	BM	BL	SW	SM	SL
MLR	0.954	0.966	0.871	0.937	0.969	0.035
IR	0.995	0.975	0.986	0.989	0.987	0.817

See abbreviations in Figure 3.

Table 6. Prediction equations of Iberian ham quality traits achieved by applying multiple linear regression (MLR) and isotonic regression (IR) on data achieved by MRI and Computer Vision Techniques (BV, SV and HV).

MLR	IR
$BW = 0.0445 * BV + 0.0131 * SV + 154.6591$	$BW = -4 \times 10^{-23} * HV^6 + 2 \times 10^{-17} * HV^5 - 3 \times 10^{-12} * HV^4 + 2 \times 10^{-7} * HV^3 - 0.0109 * HV^2 + 280.58 * HV - 3 \times 10^6$
$BM = 0.0021 * BV + 0.0002 * SV + 21.6042$	$BM = 2 \times 10^{-24} * SV^6 - 7 \times 10^{-19} * SV^5 + 5 \times 10^{-14} * SV^4 - 2 \times 10^{-9} * SV^3 + 3 \times 10^{-5} * SV^2 - 0.2673 * SV + 925.87$
$BL = -0.0007 * BV + 21.7736$	$BL = -2 \times 10^{-25} * HV^6 + 7 \times 10^{-20} * HV^5 - 1 \times 10^{-14} * HV^4 + 1 \times 10^{-9} * HV^3 - 8 \times 10^{-5} * HV^2 + 2.4017 * HV - 29387$
$SW = 0.0263 * BV + 0.0063 * SV + 28.4885$	$SW = 2 \times 10^{-21} * BV^6 - 3 \times 10^{-16} * BV^5 + 1 \times 10^{-11} * BV^4 - 3 \times 10^{-7} * BV^3 + 0.0044 * BV^2 - 29.538 * BV + 80337$
$SM = 0.0029 * BV + 0.0007 * SV - 4.2683$	$SM = 8 \times 10^{-23} * SV^6 - 1 \times 10^{-17} * SV^5 + 5 \times 10^{-13} * SV^4 - 1 \times 10^{-8} * SV^3 + 0.0002 * SV^2 - 1.2033 * SV + 3543.8$
$SL = -0.0001 * SV + 7.8575$	$SL = 2 \times 10^{-23} * SV^6 - 2 \times 10^{-18} * SV^5 + 1 \times 10^{-13} * SV^4 - 2 \times 10^{-8} * SV^3 + 3 \times 10^{-5} * SV^2 - 0.222 * SV + 636.56$

See abbreviations in Figure 3.

HIGHLIGHTS

- Data mining and MRI-CVT have been firstly used to study quality features of hams.
- Data mining tasks are appropriate to deduce and predict quality traits of hams.
- Physical-chemical and computer vision data are inferred by applying deductive tasks.
- Quality traits can be control by using predictive techniques and computer vision data.