

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

A GP approach for precision farming

This is the author's manuscript

Original Citation:

Availability:

This version is available <http://hdl.handle.net/2318/1781793> since 2021-03-22T10:58:50Z

Publisher:

Institute of Electrical and Electronics Engineers Inc.

Published version:

DOI:10.1109/CEC48606.2020.9185637

Terms of use:

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

A GP Approach for Precision Farming

Francesca Abbona ^{1,4,*}, Leonardo Vanneschi ^{2,3}, Marco Bona ⁴ and Mario Giacobini ¹

¹ Department of Veterinary Sciences, University of Torino, Turin, Italy.

² NOVA Information Management School (NOVA IMS), Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisboa, Portugal.

³ LASIGE, Departamento de Informática, Faculdade de Ciências, Universidade de Lisboa, 1749-016 Lisboa, Portugal

⁴ ANABORAPI, Associazione Nazionale Allevatori Bovini Razza Piemontese, Carrù, Italy

* Corresponding author.

E-mail addresses:

francesca.abbona@unito.it (F. Abbona), lvanneschi@novaims.unl.pt (L. Vanneschi), marco.bona@anaborapi.it (M. Bona), mario.giacobini@unito.it (M. Giacobini).

1

2 **Abstract**—Livestock is increasingly treated not just as food containers, but as animals that can be susceptible to stress
3 and diseases, affecting, therefore, the production of offspring and the performance of the farm. The breeder needs a
4 simple and useful tool to make the best decisions for his farm, as well as being able to objectively check whether the
5 choices and investments made have improved or worsened its performance. The amount of data is huge but often
6 dispersive: it is therefore essential to provide the farmer with a clear and comprehensible solution, that represents an
7 additional investment. This research proposes a genetic programming approach to predict the yearly number of weaned
8 calves per cow of a farm, namely the measure of its performance. To investigate the efficiency of genetic programming
9 in such a problem, a dataset composed by observations on representative Piedmontese breedings was used. The results
10 show that the algorithm is appropriate, and can perform an implicit feature selection, highlighting important variables
11 and leading to simple and interpretable models.

12 **Keywords**—Genetic Programming, Precision Livestock Farming, Cattle Breeding, Piedmontese Bovines.

13 1. Introduction

14 In this article, the performance of the breeding farms of *Piemontese* bovines are investigated. The considered cattle
15 farms are located in Piedmont, a region in Northwestern Italy. The Piedmontese cattle derives its name from this region,
16 its cradle of origin, even if today it is spreading in several foreign countries. The bovines are usually bred in beef intensive
17 farms, which are therefore provided with the installation of stables to control the animals, grazing for fattening
18 purposes, the addition of different artificial fodder on feed and curative intents, and particular attention to the
19 reproduction of the livestock. The main information that represents the yield of a Piedmontese cattle farm is given by
20 the *count of calves per cow per year* [1, 2]. It is a quantity that is basically predicted considering the average calving
21 interval *intp*, expressed in days, and the average perinatal mortality of the farm *m*, referred to the previous 12 months:

$$22 \quad Y_a = \frac{365}{intp} \left(1 - \frac{m}{100} \right) \quad (1)$$

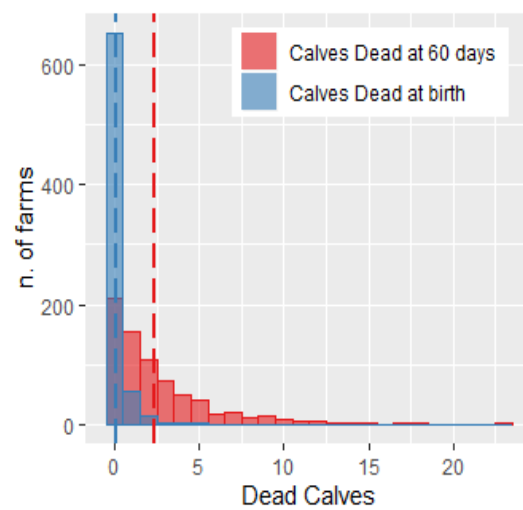
23
24 However, this expression does not take into account the period following the birth of the calf. Calf mortality is an
25 important cause of economic losses in Piedmontese cattle farms [2]. It represents for the farmer the loss of the
26 economic value of the calf and the reduction of the herd's genetic potential. Furthermore, the high mortality rate
27 reduces the number of young animals to be used to increase the size of the breeding. From the analysis of the recordings
28 in 2017, described later in the manuscript, the difference between the number of dead calves at birth and those that
29 did not survive through the weaning period is straightforward extremely significant. In Figure 1, the distributions of
30 dead calves are represented for the selected group of breedings. Most of the farms report no death at birth (mean
31 value: 0.13), whereas during weaning period records show up to 23 deaths per farm in the considered year (mean value
32 3.02), entailing that many of the newborns were not able to survive. This issue can lead to the need of buying animals
33 and, therefore, to additional costs. Perinatal mortality is related to birth and the first few hours after it: it is mainly due
34 to the delivery itself (difficulty of parturition of the cow and its health condition) and to the difficulty of birth and the
35 weight of the newborn.

36 However, the complete development of the calf occurs in the 60 days following birth. During the weaning period, the
37 physiological development process of the animal reaches completion and it is straightforward that the gestational phase
38 alone is not exhaustive: it is therefore crucial to consider neonatal mortality, outlining the calf's ability to survive. Thus,
39 in addition to the genetic factors previously mentioned, it is inevitable to consider congenital calf's defects, such as
40 arthrogryposis or macroglossia. Together with environmental and food conditions, they affect the quality of life of the

41 newborn, denoting an important source of stress that can compromise the immune response, the growth rate, the
42 disease resistance, and the well-being of the animals. It is therefore necessary to incorporate in expression (1) the
43 factors that encapsulate the effect of the weaning period of the newborn. This issue leads to the reformulation of the
44 problem into the following question:

45 *“How many weaned calves per cow are produced per year?”*

46 It is straightforward that a proper model should be formulated, with the encapsulation of other parameters, among
47 those available in the dataset.



48
49 Fig.1 Distribution of the number of dead calves at birth and during the weaning period in 2017. Mean values are
50 represented with the dashed line at the two different time reference. The data derive from the dataset described in
51 Section III. All the breedings (725) show extremely different values between the dead calves at birth (in blue) and (in
52 red) at 60 days after it (Kruskal-Wallis test: p-value $\ll 0.001$).

53
54 This study aims hence at investigating the production performances of Piemontese calves and its optimization for
55 fattening purposes but also for the calf's reproductive career. In particular, the intention is to extend the horizon by
56 investigating which administrative and production variables available in the dataset may influence the production of
57 calves. Studies conducted so far within the association are based on traditional statistical identification approaches.
58 Actual modelling involves only two variables, without exploiting the huge number of parameters in the dataset [1, 2, 3].
59 Without making a priori bio-, epi-, or eco- logical assumptions about data or the relationship between the response and
60 the independent variables, even if still uncommon for this class of problems, Machine Learning (ML) techniques may
61 provide interesting feature selection characteristics, representing a flexible and robust alternative in predictors

62 identification. Specifically, the potential of Genetic Programming (GP) [4, 5] is investigated to create and to analyze
63 predictive models for the number of weaned calves in Piedmontese cattle breedings, which could improve the analysis
64 of Piedmontese breeding performance. Inside the ML arena, we chose to use GP, because this technique has a set of
65 interesting characteristics, that distinguish it from many other methods. First of all, it assumes no hypothesis about the
66 shape of the final model, which is very important for the problem under investigation, considering that no a priori
67 knowledge is given. Secondly, using some precaution, GP can be able to generate readable and interpretable models,
68 which is crucial for our application. Finally, GP is able to perform an automatic feature selection, thus relieving us from
69 any pre-processing task. These models are compared with the predictive model that is currently adopted by the National
70 Association of Piedmontese Cattle Breeders ANABORAPI to monitor the progress of each farm.

71

72 The paper is organized as follows: in Section II, the background is described. Then, the dataset is analyzed and some
73 basic assumptions on the model are made in Section III. GP models and their performance are illustrated in Section IV,
74 where hypothesis tests and results are examined. Finally, discussions are presented, and further developments are
75 highlighted in Section V.

76 **2. Background**

77 The 'Piedmontese' is an Italian bovine breed native of Piedmont and represents a characteristic element of the territory.
78 It is the major bred breed among beef cattle in the region, showing both organoleptic and zootechnical remarkable
79 qualities. If, on one side, it results in greater tenderness of the meat, on the other hand, it is a breed with exceptional
80 character skills, such as meekness, maternal attitude, resistance to diseases, little stress, and great adaptation to
81 pasture. It, therefore, allows easy management and, not less important, the use and development of the local area [1,
82 3]. The association ANABORAPI is responsible for promoting the breed through the study of the productive, reproductive
83 and management processes of the Piedmontese breeding [6]. The activity is carried out with the management of the
84 Herd Book of the Race, a complex database that preserves the pedigrees of all the registered animals and a series of
85 additional information, such as validation of breed characters, reproductive career, morphological studies, and genetic
86 values. Nowadays, these activities must deal with new needs, increasingly connected to the sustainability of breeding
87 and well-being of animals, in the perspective of monitoring every animal. The contribution of each individual is the
88 concept behind Precision Livestock Farming (PLF). It is the solution to avoid imprecise or non-objective farmers
89 evaluations and to facilitate management methods, to obtain hence the best profit both for the individual and the
90 community. ANABORAPI offers to its members a wide section of statistics, which provide a detailed analysis of various

91 parameters of technical and economic efficiency of the farm and can contribute to identifying the breeding strengths
92 and critical points for possible improvements or developments. In particular, the average situation of breeding due to
93 the main fertility parameters is monitored, summarized by the average number of calves per cow produced in the last
94 year, net of mortality and calving interval (the period between two deliveries of a cow). This is then translated into a
95 brief economic summary, which compares the gross revenue with the mortality losses, providing the farmer with an
96 economic indicator of breeding performance.

97

98 A large amount of data is now collected through the use of sensors, ear tags, collars, images and video recordings in
99 many fields, and livestock sector is not different [7-9]. It is increasingly common in farms to monitor each animal: as
100 already mentioned, the PLF approach aims for greater accuracy on the quantity and quality of information, to achieve
101 the economic and environmental sustainability of farms. The breeder must generally deal with animals' problems like
102 their health conditions and social behavior, that affect the quality of the product, the life of the animal, and the
103 performance of the farm. Indeed, the PLF approach provides the offset of incurred costs, as these issues are identified
104 in advance, allowing decisions to be made in time [10, 11]. The creation of prediction models on a specific result in the
105 zootechnical field is increasingly addressed with the use of ML techniques [10-20]. These approaches are suitable for
106 the management of large data sets and are used to predict livestock issues such as the time of disease events, risk
107 factors for health conditions, and failure to complete a production cycle. Studies have been conducted, based on the
108 application of ML techniques, to model the individual intake of cow feed [12], optimizing health and fertility, to predict
109 the rumen fermentation pattern from milk fatty acids [13], which influence the quantity and composition of the milk
110 produced but also the sensorial and technological characteristics of the meat. The use of ML techniques is also often
111 exploited to identify potential disease predictors, e.g. Bovine Viral Diarrhoea Virus (BVDV), Infectious Bovine
112 Rhinotracheitis (IBR), Bovine Tuberculosis (TB), lameness, and mastitis [14-16, 21], to classify grazing and social behavior
113 [17-19], and to predict carcass conformation [20], an important component of price negotiations between beef
114 producers and market operators. These works are mostly carried out on dairy cattle, which are more critical to manage
115 from a health point of view. Dairy animals generally have a shorter average life compared to the lifespan of beef bovines
116 and are usually affected by diseases and metabolic problems. In the beef cattle sector, and in particular in the
117 Piedmontese cattle breed, animals are more resistant and exposed to fewer stress factors. This is an explanation why
118 meat farms show moderate use of devices. However, individual information is already recorded and loaded by the
119 technicians during the checks, and therefore the management of big data is necessary.

120 **3. The Dataset**

121 The content of the dataset elaborated by the ANABORAPI system covers a total of over 4000 active farms, keeping
122 historical records for all of them. The elaboration processed by the ANABORAPI system to evaluate Y_a (see Equation (1))
123 goes back 365 days, starting from the last check, to process the average summaries. A first restriction is therefore the
124 isolation of the data of a whole year (in our case 2017) and to consider the target we want to infer for the following year
125 (2018). Since the performance of the farm mainly focuses on fertility, the data concerning multiparae cows were
126 considered to elaborate the number of deliveries and the calving intervals. In the same way, data on bulls used for
127 artificial insemination were maintained (i.e. selection indices, that represent namely estimations of the additive genetic
128 effect of a subject). Information referred to inbreeding levels between animals were not incorporated into the study,
129 since they required more investigations. However, they will be included in future developments, for a more accurate
130 inspection on the consanguinity of unborn calves. Finally, restrictions on farms were imposed to obtain a solid
131 representative subset: filters on breeding located in Piedmont with at least 30 cows and percentage of artificial
132 insemination between 90% and 100% were applied. This last condition means that a part of the considered farms
133 actually own bulls and carry out natural impregnations. Thereby, two main groups of farms result from the selection: a
134 smaller one, containing 330 farms, and a larger one, consisting of 395 breedings, resulting in a total of 725 breedings.
135 The difference between the two sets results in a major use of the breeding bull: this means that instead of recording
136 the date on which the insemination took place, breedings belonging to the second group use to set a period of several
137 days, followed by the diagnosis of the pregnancy. As both datasets are representative for the Piedmontese breeding
138 reality, where the second dataset features a more diffused situation and the first one the most accurate one, we used
139 both groups in the study, as propaedeutic to the objective. Since the aim is the building of predictive models via a ML
140 technique, we therefore decided to designate the first set of farms (size 330) as a learning set, as the algorithm can
141 learn on precise recordings, while the second set (size 395) was designated as a test set. Each record of the final datasets
142 stands for a single farm and variables {1 – 19} refer to year 2017, whereas Y is the actual number for weaned calves
143 recorded in 2018 (Table I). All variables can only assume positive values.

144

145 **4. Application of GP**

146 *GP main traits.*

147 The GP technique is a tree-based algorithm, in which the initial population evolves through mechanisms of selection,
148 mutation, and recombination of individuals (i.e. mutation and crossover), as in a biological evolutionary process.

149 Subtrees at each generation are recombined and recursively evaluated. The best candidates are eligible for the new
150 generation and they are on average fitter than previously generated individuals, i.e. show a smaller error. The error is
151 measured with a fitness function, that is an objective function used to evaluate the distance from the experimental
152 target. with a regression problem we have chosen as fitness function the Root Mean Square Error (RMSE) between the
153 expected and predicted numbers as the measure for the fitness of the models: lower values represent better solutions
154 (i.e. expressions fitting well correspond to low error levels).

155

156 *Dataset partitioning.*

157 For our experimental study, we used the GPLab Toolbox of MATLAB [4]. As mentioned in the previous section, the first
158 group of farms (size 330) was used as a learning set, while the second one (size 395) as a test set. We considered the
159 possibility of dividing the datasets through a k -fold cross validation approach. However, the reduced set of data does
160 not allow us to find a suitable k value: for instance, if we chose a k smaller than 10, we would obtain a small number of
161 subsets, leading to a small number of runs (i.e. less than 10). On the contrary, with a k greater than 10, we would have
162 a restrained number of records within the test sets (i.e. less than 39 test farms for each run). We used the splitting of
163 the learning set into 30 different subsets, with constant training-validation partitioning (75%-25%). Each division was
164 carried out with a random choice of records at each run with uniform distribution and without repetition, keeping
165 separate training and validation. In other words, among the total 330 learning records, 83 records were chosen to form
166 the validation set, and the remaining 247 were labeled as training ones, reiterating the process with different sets for
167 all the 30 runs. For each run, the individuals obtained on the training set were evaluated on the validation set, in order
168 to select the best ones (i.e. models with the lowest error among the validation set). Finally, the generalization ability of
169 the latters was checked, by analyzing the respective error achieved on the test set.

170

171 *Terminals, nodes and operators*

172 The GP individuals were generated using a tree-based representation, where the trees were built using a set of terminal
173 symbols T and a set of primitive functional symbols F . The set T was composed by the previously described variables,
174 plus a set of random constants between 0 and 1 generated during the initialization process. The set F was equal to $\{plus;$
175 $minus; times; mydivide\}$, where *plus*, *minus* and *times* indicate the usual operators of binary addition, subtraction and
176 multiplication, respectively, while *mydivide* represents the protected division, that returns the numerator when the
177 denominator is equal to zero. In order to limit overfitting and maintain the models as simple as possible, besides

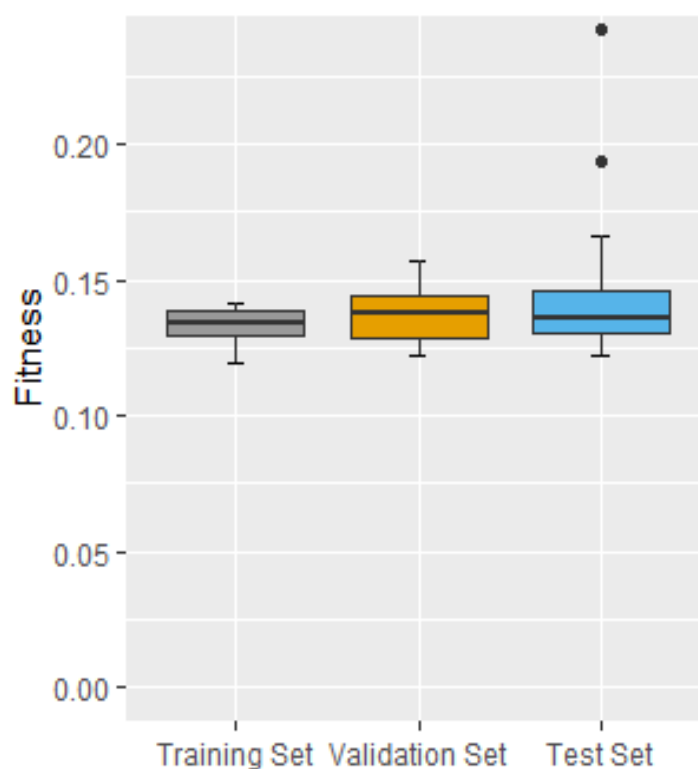
178 crossover and mutation, operators such as *shrinkmutation* and *swapmutation* (predefined in GPLab) were used. These
179 two operators respectively exchange a subtree with a terminal node and permute binary non-commutative functions'
180 elements. Table II reports the employed experimental setting.

181

182 *Performance evaluation.*

183 The performance of the simulations is reported in Figure 2, where the fitness among the 30 runs on the training, the
184 validation and the test sets are presented. The Lilliefors test, performed with significance level $\alpha=0.05$, showed that a
185 normal distribution can be assumed only on the training set. Hence, we applied a Kruskal-Wallis test ($\alpha=0.05$), under
186 the alternative hypothesis that, at the end of the runs, the RMSEs do not have equal medians. Results entailed that
187 there is no significant difference between the three distributions: given a p-value $p=0.17$, the null hypothesis was not
188 rejected, that is the median values of the errors committed on the three sets are not different. The median value
189 obtained on the test set allows us to affirm that the obtained models are able to generalize well, on unseen data.

190



191

192 Fig.2 Performance of the best 30 selected models, respectively, on the training, validation and test sets. There is no
193 significant difference between the results (Kruskal-Wallis test: $p = 0.17$, with $\alpha=0.05$), i.e. the median values of the errors
194 committed on the three phases are not different.

195

196 Table I. Final set of variables used in the studied dataset. The last line (variable Y) represents the dependent variable,
 197 target of the predictive models generated by GP.

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

	Reference Year	Variable	Description
1	2017	<i>COWS</i>	Consistency for cow (n. of cows in the farm)
2	2017	<i>C_AGE</i>	Mean age of cows, expressed in days.
3	2017	<i>INT_P</i>	Mean value of calving interval, i.e. the average number of days that elapse between a parturition and the following one.
4	2017	<i>C_PAR</i>	Mean number of parturitions of cows.
5	2017	<i>N_PAR</i>	Number of occurred deliveries.
6	2017	<i>C_EASE</i>	Number of easy parturitions, that did not require human intervention and that did not cause stress to the cow nor the calf.
7	2017	<i>C_GRAVID</i>	Number of pregnant cows.
8	2017	<i>C_INS</i>	Number of inseminated cows.
9	2017	<i>BIRTHW_M</i>	Mean birth weight of male calves.
10	2017	<i>BIRTHW_F</i>	Mean birth weight of female calves.
11	2017	<i>IND_PAR</i>	Mean Genetic selection index referred to facility of parturition of the cows.
12	2017	<i>TFA_BIRTH</i>	Mean Genetic selection index referred to facility of birth of the bulls, which semen has been used on artificial inseminations.
13	2017	<i>TFA_PAR</i>	Mean Genetic selection index referred to facility of parturition with which the bulls, which semen has been used on artificial inseminations, have been born.
14	2017	<i>N_ELIM</i>	Number of calves dead within 60 days after birth.
15	2017	<i>N_TOT</i>	Total number of newborns.
16	2017	<i>N_BALIVE</i>	Total number of calves born alive.
17	2017	<i>N_CORRECT</i>	Percentage of calves born without birth defects, such as Macroglossia or Arthrogryphosys.
18	2017	<i>ABORT</i>	Percentage of abortions.
19	2017	<i>MORT</i>	Mean neonatal mortality.
20	2018 range= [0.26;1.24]	Y	Number of calves per cow per year. It is obtained on data from 2018 with the following: $Y = \frac{N_{BALIVE} - N_{ELIM}}{COWS}$

Table II. Parameters used in our experimental study

Parameter	Description
Maximum number of generations	20
Population size	500
Selection Method	Lexicographic Parsimony Pressure
Elitism	Keepbest
Initialization Method	Ramped half and half
Tournament Size	2
Subtree Crossover Rate	0.8
Subtree Mutation Rate	0.1
Subtree Srinkmutation Rate	0.05
Subtree Swapmutation Rate	0.05

Key role variables result in non-null median frequencies among the best solutions on all the runs, whereas negligible ones correspond to null estimations: values greater than zero suggest that the corresponding variables were used in over 50% of the final solutions, namely the number of cows (*COWS*), the number of occurred deliveries in the farm during the year (*N_PAR*), and the number of calves that were born alive (*N_BALIVE*). The information was confirmed also by the equivalent percentage, reported in the second column of Table III.

Finally, we investigated the interpretability of the expressions, considering the number of variables involved in each one of the best final models and the corresponding fitness. In order to compare the performance of the GP models, we examined the number of parameters encapsulated in each one, paying attention to the corresponding fitness obtained on the test set (Table IV). Observing Table IV, we can identify a general trend: models that use less variables tend to have a worse fitness (i.e. a larger error) on the test set than those that use more variables. Among the 19 variables in the dataset, the obtained models encapsulate from a minimum of 3 to a maximum of 10 variables. An intermediate situation is represented by models involving 4 of these parameters, since in this case the error is small and, as shown later, the expression is interpretable. We selected two models in order to make comparisons, the one showing the best fitness among all the evolved expressions (*GP3* in Figure 3) and the one with the best fitness among the models that use 4 variables (*GP8* in Figure 3). The choice of the second model was entailed, as shown below, as a consequence of its interpretability. For both Models *GP3* and *GP8*, the distance values between predictions based on 2017 and target values Y_i recorded in 2018 are represented through boxplots, that is:

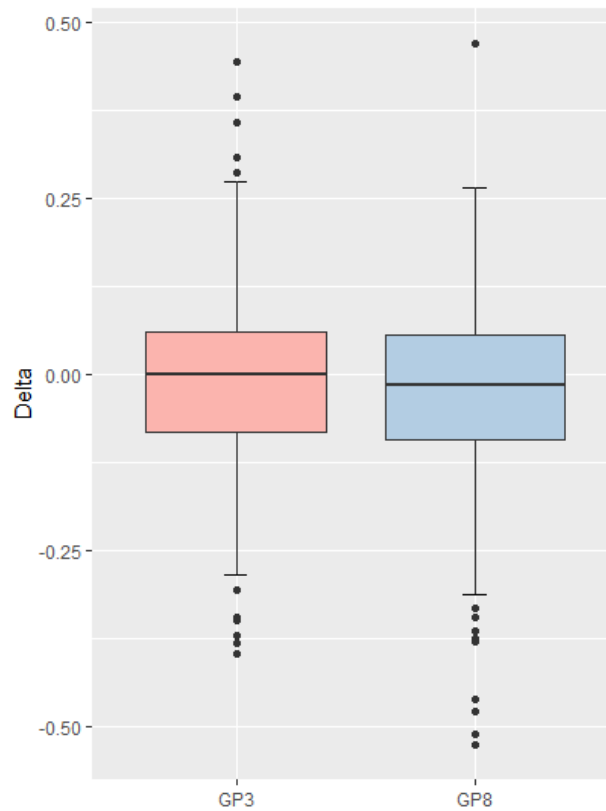
$$\Delta_{model,i} = Y_{model,i} - Y_i$$

for each record $i = 1, \dots, 395$ in the test set. Predictions obtained with the two models *GP3* and *GP8* are not significantly different (Kruskal Wallis: p-value = 0.2372).

253 Table III. Median frequencies and percentage of use of each variable among the best 30 individuals found by GP.

254
255
256
257
258
259
260
261
262
263
264
265

Variable	Median	% of use on 30 runs
<i>X₁ - COWS</i>	1	73%
<i>X₂ - C_AGE</i>	0	27%
<i>X₃ - INTP</i>	0	43%
<i>X₄ - C_PAR</i>	0	27%
<i>X₅ - N_PAR</i>	1	53%
<i>X₆ - C_EASE</i>	0	40%
<i>X₇ - C_GRAVID</i>	0	23%
<i>X₈ - C_INS</i>	0	17%
<i>X₉ - BIRTHW_M</i>	0	13%
<i>X₁₀ - BIRTHW_F</i>	0	10%
<i>X₁₁ - IND_PAR</i>	0	37%
<i>X₁₂ - TFA_BIRTH</i>	0	13%
<i>X₁₃ - TFA_PAR</i>	0	23%
<i>X₁₄ - N_ELIM</i>	0	37%
<i>X₁₅ - N_TOT</i>	0	43%
<i>X₁₆ - N_BALIVE</i>	0.5	50%
<i>X₁₇ - N_CORRECT</i>	0	37%
<i>X₁₈ - ABORT</i>	0	23%
<i>X₁₉ - MORT</i>	0	13%



266

267

268 Fig.3 Comparisons between GP models on the test set. Distributions of the differences between predicted and real
 269 values are plotted. Both GP predicted values are not significantly different (Kruskal-Wallis: p - value = 0.2372). GP3
 270 shows a median value equal to -0.0005928782, smaller than the median value obtained with GP8 (-0.0146762341).

271 Table IV. Fitness on the test set, number of involved variables and corresponding percentage are reported for each
 272 model evolved by GP in each one of the 30 performed runs.

<i>Prediction model</i>	<i>Fitness on Test</i>	<i>N. of variables</i>	<i>% of variables</i>
<i>model 1</i>	0.1379	5	26%
<i>model 2</i>	0.1418	3	16%
<i>model 3</i>	0.1218	9	47%
<i>model 4</i>	0.1354	8	42%
<i>model 5</i>	0.1660	3	16%
<i>model 6</i>	0.1290	8	42%
<i>model 7</i>	0.1370	4	21%
<i>model 8</i>	0.1321	4	21%
<i>model 9</i>	0.1258	8	42%
<i>model 10</i>	0.1357	3	16%
<i>model 11</i>	0.2422	9	47%
<i>model 12</i>	0.1461	3	16%
<i>model 13</i>	0.1286	7	37%
<i>model 14</i>	0.1548	4	21%
<i>model 15</i>	0.1320	9	47%
<i>model 16</i>	0.1261	7	37%
<i>model 17</i>	0.1285	8	42%
<i>model 18</i>	0.1371	9	47%
<i>model 19</i>	0.1610	3	16%
<i>model 20</i>	0.1571	4	21%
<i>model 21</i>	0.1355	9	47%
<i>model 22</i>	0.1450	3	16%
<i>model 23</i>	0.1291	7	37%
<i>model 24</i>	0.1426	4	21%
<i>model 25</i>	0.1935	5	26%
<i>model 26</i>	0.1330	10	53%
<i>model 27</i>	0.1305	6	32%
<i>model 28</i>	0.1543	3	16%
<i>model 29</i>	0.1308	7	37%
<i>model 30</i>	0.1361	9	47%

273

274 We therefore concluded that the two models, whose expression is provided in Equations (3) and (5), perform likewise,
 275 incorporating different variables with respect to Y_a (see Equation (1)). Parameters such as *MORT* and *N_ELIM* used in
 276 Equation 1 were encapsulated also in GP expressions, i.e. mortality at 60 days (*GP8*) and number of calves born alive
 277 (*GP3* and *GP8*). Regarding *GP3*, the expression in infix notation to obtain the predictions is:

278

279
$$Y_{GP3} = \frac{X_{11}}{X_{17} + \frac{X_3}{X_{16}} + \frac{X_3}{X_6 \cdot \frac{2 \cdot X_{18} + X_{16}}{X_9 + X_1}}}, \quad (3)$$

280

281 where

282
283
284
285
286
287

X_1 -COWS,
X_3 -INTP,
X_6 -C_EASE,
X_9 -BIRTHW_M,
X_{11} -IND_PAR,
X_{16} -N_BALIVE,
X_{17} -N_CORRECT,
X_{18} -ABORT,
X_{19} -MORT.

288 In model *GP3*, the denominators of *mydivide* operator do not meet existence conditions, that is they can assume null
289 values (e.g. perinatal mortality X_{19} is null for some records). It is not possible to assert that the *mydivide* operator is
290 actually a division and the previous expression (3) cannot be further simplified. Contrarily to *GP3*, the model for *GP8* is
291 comprehensible:

292
$$Y_{GP8} = \frac{X_5}{\frac{(X_5 \cdot X_{14} + X_{16})}{X_1} + X_1} \quad (4)$$

293 Since we previously set the constraint in the dataset on farms with more than 30 cows, and the other variables can even
294 assume only positive values, the denominators of *mydivide* that appear in the latter model are also positive (in Model
295 4, the mentioned values cannot reach null levels, since the number of cows is added to a quantity, greater than zero).
296 Existence conditions are in this case always verified and therefore the function *mydivide* is a division, leading to a
297 simplified version:

298

299
$$Y_{GP8} = \frac{X_1 \cdot X_5}{X_1^2 + X_5 \cdot X_{14} + X_{16}}, \quad (5)$$

300 where

X_1 - COWS
X_5 - N_PAR
X_{14} - N_ELIM
X_{16} - N_BALIVE.

304 Model (5) can further be rewritten as

305
$$Y_{GP8} = \left(\left(\frac{N_PAR}{COWS} \right)^{-1} + \frac{N_ELIM}{COWS} + \left(\frac{N_BALIVE}{COWS} \cdot \frac{1}{N_PAR} \right) \right)^{-1} \cdot (6)$$

306 The first term can be expressed as the invers of the number of mean value of the yearly deliveries occurred in the farm,
307 since the number of all parturitions is divided by the total number of cows ($\overline{N_PAR}$). Likewise, the second and third
308 terms contain, respectively, the yearly number of calves per cow that did not survive during the weaning period
309 ($\overline{N_ELIM}$) and the yearly number per cow of calves born alive ($\overline{N_BALIVE}$), that is:

$$Y_{GP8} = \left(\frac{1}{\overline{N_PAR}} + \overline{N_ELIM} + \frac{\overline{N_BALIVE}}{\overline{N_PAR}} \right)^{-1}. \quad (7)$$

311 Stated otherwise, by renaming the terms and performing basic operations, we obtained the following:

$$312 \quad 1 = n_j v_{1,j} + n_j v_{2,j} + n_j v_{3,j}, \quad (8)$$

313 for $j=1, \dots, 725$, since we considered the complete dataset with all the selected farms (see Section III), and where:

314

315

316

317

318

$$\begin{aligned} n_j &= Y_{GP8,j}, \\ v_{1,j} &= (\overline{N_PAR}_j)^{-1}, \\ v_{2,j} &= \overline{N_ELIM}_j, \\ v_{3,j} &= \frac{\overline{N_BALIVE}_j}{\overline{N_PAR}_j}. \end{aligned}$$

319 It is straightforward that Equation (8) can be formulated as the sum of rescaled variables

320

$$321 \quad 1 = \tilde{v}_{1,j} + \tilde{v}_{2,j} + \tilde{v}_{3,j}, \quad (9)$$

322 where $\tilde{v}_{i,j} = n_j v_{i,j}$ for $i=\{1,2,3\}$. Thereby, it was possible to measure the contribution of each term in the sum expressed

323 in Equation (9). The distributions of each $\tilde{v}_{i,j}$ was statistically analyzed and the three boxplots were displayed (Figure

324 4). Extremely significant difference is verified between all variables (Wilcoxon test with Bonferroni correction: $\alpha=0.017$,

325 $p < 0.001$). Moreover, we inspected how far the mean value of each variable is from the unit. We compared, one by

326 one, the three distributions via a single sample Wilcoxon test. We set alternative hypothesis that the distribution shows

327 a mean value $\mu \neq 1$, with $\alpha=0.05$. Once again, we found an extremely statistical difference between the mean value of

328 $\tilde{v}_{i,j}$ from the value 1. Similarly, we compared the distributions with respect to 0. The results of the test were analogous

329 to the previous ones: with extremely significant p-values ($p < 0.001$), we could deduce that $\tilde{v}_{2,j}$ and $\tilde{v}_{3,j}$ remain relevant

330 parameters, even assuming values close to zero, providing hence a minimal contribute in Equation (8). In other words,

331 we could assert that all the variables in Equation (9) are influent: in particular, $\tilde{v}_{1,j}$ is the most important one, since its

332 mean value was $\mu_1=0.951$, whereas $\tilde{v}_{2,j}$ and $\tilde{v}_{3,j}$ respectively showed $\mu_2=0.032$ and $\mu_3=0.021$. Model (5) can be

333 simplified, to the point of being expressed as the sum of three parameters. We verified that these three parameters are

334 the average number of parts occurred during the year in the herd ($(\overline{N_PAR})^{-1}$), the number of calves per cow that have

335 not passed the weaning phase ($\overline{N_ELIM}$) and finally the number of calves per cow live births compared to the total

336 number of parts of the herd during the year ($\overline{N_BALIVE}/\overline{N_PAR}$). From a zoological point of view, these are actually the

337 main parameters that intuitively can give an idea of the economic performance of the farm. All of them play a significant

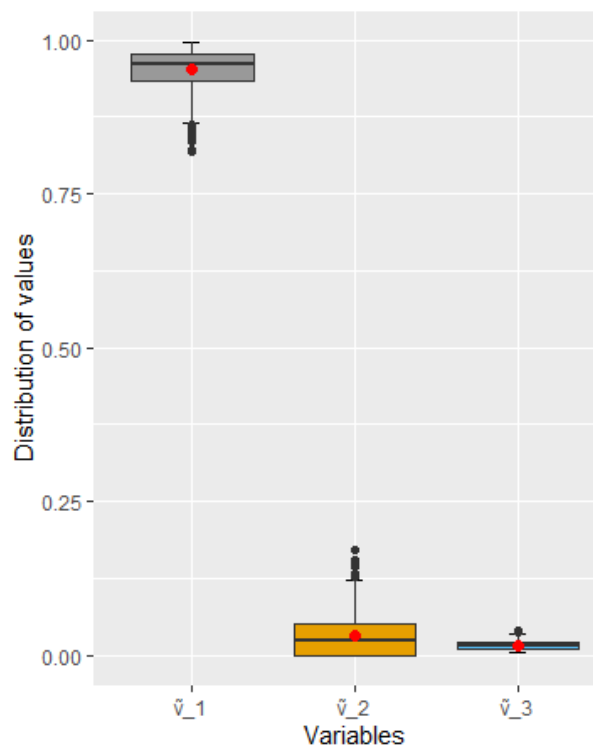
338 role with respect to the response variable: more importance is given to the parameter $(\overline{N_PAR})^{-1}$, associable to the

339 inverse of the mean calving interval (days between two deliveries) of the farm, whereas $\overline{N_ELIM}$ and $\overline{N_BALIVE}/N_PAR$
340 give a smaller contribute.

341

342 Summing up, the most frequent variables in the models, obtained with GP, are the number of cows in the farm (*COWS*),
343 the number of deliveries occurred in the breeding (*N_PAR*) and the number of calves born alive (*N_BALIVE*). The calving
344 interval (*INTP*) and the number of dead calves at 60 days (*N_ELIM*) are slightly less frequent. Perinatal mortality is not
345 so recurring, meaning that it could play a minor role in the prediction. The most frequent variables encapsulated in the
346 expression (3) are *COWS* and *N_BALIVE*, followed by *INTP* and *MORT*. Then there are 5 less frequent additional
347 parameters that could, therefore, be relevant in the refinement of the prediction. The median error of predictions
348 obtained with model *GP3* is slightly smaller than the one obtained with model *GP8*. The latter however processes less
349 variables, exploiting exactly the three most frequent ones, listed in Table III.

350



351

352

353 Fig.4 Boxplots of the distributions of the variables in Equation (9). Wilcoxon test with Bonferroni correction at
354 $\alpha=0.017$ reported $p<<0.001$. Hence, the variables are significantly different. The single sample Wilcoxon test, with
355 $\alpha=0.05$, showed for each distribution a mean value $\mu\neq 1(p<<0.001)$. Therefore, all the variables are extremely significant
356 in Equation (9). Mean values are respectively $\mu_1=0.951$, $\mu_2=0.032$ and $\mu_3=0.021$ (red dots).

357 **5. Conclusions and Future Works**

358 In this study, we investigated the performance of medium to large farms located in Piedmont of Piedmontese cattle,
359 starting from the model implemented in the systems of the National Association of Piedmontese Bovines (ANABORAPI)
360 [1-3]. The currently used model (reported in Equation (1)) predicts the number of calves per cow per year. However, it
361 is not completely suitable to represent the performance of the farms. In fact, during the weaning period, many calves
362 do not survive, entailing great losses to the economic revenues of the breedings. The reasons for those deaths are
363 various and difficult to identify objectively. It is hence necessary to take into account crucial parameters, that encompass
364 the calf's weaning in the output, as, for example, the number of calves born alive and those dead after the weaning
365 period, within 365 days. Although biologically acceptable, these hypotheses could not be sufficiently informative or they
366 could be informative enough, but not exhaustively combined in the formulation of a model. It is therefore difficult to
367 build a model with only zootechnical speculations and an automatic learning method has been applied, which can meet
368 the requirements. In addition, it is necessary to research and propose a simple model, which can be easily interpreted
369 by the breeder. The expression to target should be a simplification and an added value to the management of the farm.
370 The breeder should be able to easily read the information, in order to identify the critical points and strengths in
371 production.

372
373 Given its ability to perform an automatic feature selection, a Genetic Programming approach (GP) was used applied [4,
374 5] to build predictive models, trained, validated and tested on data recorded in 2017 and 2018. Accurate models were
375 achieved, and this means that GP can learn from a smaller dataset composed by representative farms and predict good
376 results on the selected test set. Moreover, the algorithm was able to select and process important variables, without
377 previous assumptions on the zoological aspect. The variety of expressions obtained by GP is composed of well-
378 performing models that involve more parameters, resulting in a more complex expression, hardly reducible to a simpler
379 one. However, other predictive models were also achieved that encapsulate fewer variables. Although these
380 expressions have a slightly larger error, their formula can be extremely simple and possibly easier to interpret from the
381 zoological point of view.

382
383 It is therefore worth investigating further the application of GP to a larger dataset. In this first study, we focused on data
384 directly referred to parturitions and artificial insemination, in order to process sound and solid data. The dataset was
385 filtered and resized, and 19 variables were kept among 210: many were duplicate fields, aggregates of several variables,

386 and even incomplete ones, because introduced lately in the database of ANABORAPI. Parameters such as those on
387 heifers, i.e. bovines that did not give birth yet, were not considered, since we focused on data directly referred to cows,
388 i.e. bovines that gave birth at least once. In breeding farms, heifers are mostly intended to the production of calves and
389 are going to contribute to the restock of the herd. The behavior of GP and its features selection ability among these
390 variables will be investigated, as well as among parameters on the bulls used for natural insemination. To this purpose,
391 their genetic indexes will be added to the analysis, as well as the levels of consanguinity of calves that will be born from
392 ongoing pregnancies. Comparisons with other machine learning methods will be performed, to inspect better the
393 potential of GP in the zootechnical field, and to explore possibly better models.

394 In future developments, data regarding environmental conditions inside the farm will also be taken into account, such
395 as the size of the boxes and the surface available to the animals, air and water quality and the composition
396 of the food ration. These factors are usually considered as marginal. It is common to think that cow-calf problems are
397 almost exclusively induced by genetic and pathological factors associated to pregnancy and childbirth. Indeed, not
398 enough importance is given to the period after the birth, in which the cow and the calf need feeding and environmental
399 conditions, suitable for the respective postpartum and weaning phases. In this context, once again, the ability of GP to
400 automatically select features will be very important to understand if and which of these variables are influential.

401

402 ACKNOWLEDGMENT

403 This work was partially supported by FCT, Portugal, through funding of LASIGE Research Unit (UID/CEC/00408/2019)
404 and projects BINDER (PTDC/CCI-INF/29168/2017), GADgET (DSAIPA/DS/0022/2018), AICE (DSAIPA/DS/0113/2019) and
405 PREDICT (PTDC/CCI-CIF/29877/2017), and by the Slovenian Research Agency (research core funding No. P5-0410).

406

407 REFERENCES

408 [1] Bona, M., Albera, A., Bittante, G., Moretta, A., Franco, G.: L'allevamento della manza e della vacca piemontese,
409 Supplemento al n. 44 dei Quaderni della Regione Piemonte-Agricoltura, pp. 65-129. (2005).

410

411 [2] Lo svezzamento del vitello Piemontese,
412 pp. 3-5, <http://www.anaborapi.it/images/media/pdf/rivista/2012/2012-05.pdf>
413 pp. 9 -11, <http://www.anaborapi.it/images/media/pdf/rivista/2012/2012-06.pdf>

414

- 415 [3] Associazione Nazionale Allevatori Bovini Razza Piemontese, <http://www.anaborapi.it>
- 416
- 417 [4] Silva, S.: GPLAB a genetic programming toolbox for Matlab, (2007). <http://gplab.sourceforge.net/index.html>
- 418
- 419 [5] Poli, R., Langdon, W., McPhee, N.: A Field Guide to Genetic Programming. Lulu Enterprises, UK Ltd. (2008).
- 420 <https://doi.org/10.1007/s10710-008-9073-y>
- 421
- 422 [6] Relazione Tecnica e Statistiche al 31.12.2018,
- 423 <http://www.anaborapi.it/images/media/pdf/stat/relazionetecnica2018.pdf>
- 424
- 425 [7] Berckmans, D., Guarino, M., From the Editors: Precision livestock farming for the global livestock sector, *Animal*
- 426 *Frontiers*, Volume 7, Issue 1, January 2017, Pages 45. <https://doi.org/10.2527/af.2017.0101>
- 427
- 428 [8] J. B. Cole, S. Newman, F. Foertter, I. Aguilar, M. Coffey,: BREEDING AND GENETICS SYMPOSIUM: Really big data:
- 429 Processing and analysis of very large data sets, *Journal of Animal Science*, Volume 90, Issue 3, March 2012, Pages 723-
- 430 733. <https://doi.org/10.2527/jas.2011-4584>
- 431
- 432 [9] Lokhorst, C., de Mol, R.M., Kamphuis, C.: Invited review: Big Data in precision dairy farming. *Animal*.
- 433 13(7):15191528. (2019). <https://doi.org/10.1017/S1751731118003439>
- 434
- 435 [10] Morota, G., Ventura, R. V., Silva, F. F., Koyama, M., Fernando, S. C.: BIG DATA ANALYTICS AND PRECISION
- 436 ANIMAL AGRICULTURE SYMPOSIUM: Machine learning and data mining advance predictive big data analysis in precision
- 437 animal agriculture. *Journal of animal science*, 96(4), 15401550. (2018). <https://doi.org/10.1093/jas/sky014>
- 438
- 439 [11] Liakos KG, Busato P, Moshou D, Pearson S, Bochtis D. Machine Learning in Agriculture: A Review. *Sensors*
- 440 (Basel). 2018;18(8):2674. Published 2018 Aug 14. <https://doi.org/10.3390/s18082674>
- 441

- 442 [12] Yao, C., Zhu, X., & Weigel, K. A.: Semi-supervised learning for genomic prediction of novel traits with small
443 reference populations: an application to residual feed intake in dairy cattle. *Genetics, selection, evolution: GSE*, 48(1),
444 84. (2016). <https://doi.org/10.1186/s12711-016-0262-5>
445
- 446 [13] M. Craninx, V. Fievez, B. Vlaeminck, B. De Baets Artificial neural network models of the rumen fermentation
447 pattern in dairy cattle. *Comput. Electron. Agric.*, 60 (2008), pp. 226-238. <https://doi.org/10.1016/j.compag.2007.08.005>
448
- 449 [14] Williams, M.L., Parthalin, N.M., Brewer, P., James, W.P.J., Rose, M.T.: A novel behavioral model of the pasture-
450 based dairy cow from GPS data using data mining and machine learning techniques. *J Dairy Sci.*, 99(3):20632075. (2016).
451 <https://doi.org/10.3168/jds.2015-10254>
452
- 453 [15] R. Dutta, D. Smith, R. Rawnsley, G. Bishop-Hurley, J. Hills, G. Timms, D. Henry, Dynamic cattle behavioural
454 classification using supervised ensemble classifiers. *Comput. Electron. Agric.*, 18–28 (2015),
455 <https://doi.org/10.1016/j.compag.2014.12.002>
456
- 457 [16] O. Guzhva, H. Ard, A. Herlin, M. Nilsson, K. strm, C. Bergsten,: Feasibility study for the implementation of an
458 automatic system for the detection of social interactions in the waiting area of automatic milking stations by using a
459 video surveillance system. *Computers and Electronics in Agriculture*, Volume 127, Pages 506-509, ISSN 0168-1699.
460 (2016). <https://doi.org/10.1016/j.compag.2016.07.010>.
461
- 462 [17] Ortiz-Pelaez, A., Pfeiffer, D.U.: Use of data mining techniques to investigate disease risk classification as a proxy
463 for compromised biosecurity of cattle herds in Wales. *BMC Vet Res.*;4:24. (2008). [https://doi.org/10.1186/1746-6148-](https://doi.org/10.1186/1746-6148-4-24)
464 4-24
465
- 466 [18] Machado, G., Mendoza, M. R. & Corbellini, L. G.: What variables are important in predicting bovine viral
467 diarrhea virus? A random forest approach. *Vet. Res.* 46 (2015). <https://doi.org/10.1186/s13567-015-0219-7>
468

469 [19] Alonso, J., Castañón, Á.R., Bahamonde, A., Support Vector Regression to predict carcass weight in beef cattle
470 in advance of the slaughter. (2013) Computers and Electronics in Agriculture, 91, pp. 116-120.
471 <https://doi.org/10.1016/j.compag.2012.08.009>

472

473 [20] Amrine, D. E., White, B. J., & Larson, R. L.: Comparison of classification algorithms to predict outcomes of feedlot
474 cattle identified and treated for bovine respiratory disease. Computers and Electronics in Agriculture, 105, 9-19. (2014).
475 <https://doi.org/10.1016/j.compag.2014.04.009>

476

477 [21] Bovine Diseases and Resources, <http://www.cfsph.iastate.edu/Species/bovine.php>

478

479