



A comparison of 4 different machine learning algorithms to predict lactoferrin content in bovine milk from mid-infrared spectra.

Journal:	<i>Journal of Dairy Science</i>
Manuscript ID	JDS.2020-18870.R1
Article Type:	Research
Date Submitted by the Author:	26-Jul-2020
Complete List of Authors:	Soyeurt, H�el�ene; University of Li�ge, Gembloux Agro-Bio Tech, AGROBIOCHEM dept; Statistics, Informatics and Applied modelling Unit Grelet, Cl�ement; Walloon Research Centre, Valorisation of Agricultural Products McParland, Sin�ead; Teagasc, Animal and Grassland Research and Innovation Centre, Animal and Bioscience Research Department Coffey, Mike; SRUC, Animal & Veterinary Sciences Group Calmels, Marion; Seenovia Tedde, Anthony; University of Li�ge, Gembloux Agro-Bio Tech, AGROBIOCHEM dept; Statistics, Informatics and Applied modelling Unit Delhez, Pauline; Universit� de Li�ge Gembloux Agro-Bio Tech, Terra Teaching and Research Centre Dehareng, Fr�d�eric; Walloon Agricultural Research Centre Gengler, Nicolas; ULg - Gembloux Agro-Bio Tech (GxABT), Animal Science Unit
Key Words:	milk, lactoferrin, mid-infrared, machine learning

SCHOLARONE™
Manuscripts

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

1 **A comparison of 4 different machine learning algorithms to**
2 **predict lactoferrin content in bovine milk from mid-infrared**
3 **spectra.**

4

5 H. Soyeurt^{1,*}, C. Grelet², S. McParland³, M. Calmels⁴, M. Coffey⁵, A. Tedde¹, P. Delhez^{1,6}, F.
6 Dehareng², N. Gengler¹

7

8 ¹ TERRA research and teaching centre, Gembloux Agro-Bio Tech, University of Liège, Gembloux,
9 Belgium

10 ² Valorisation of agricultural products, Walloon Research Centre, Gembloux, Belgium

11 ³ Animal & Grassland Research and Innovation Centre, Teagasc, Moorepark, Fermoy, Co. Cork,
12 Ireland

13 ⁴ Research and development, Seenovia, Saint-Berthevin, France

14 ⁵ Livestock Breeding, Animal and Veterinary Sciences, Scotland's Rural College, Midlothian, UK

15 ⁶ National fund for Scientific Research, Brussels, Belgium

16

17 **Corresponding author:** Prof. H el ene Soyeurt

18 Gembloux Agro-Bio Tech

19 Passage des D eport es, 2

20 5030 Gembloux

21 Belgium

22 hsoyeurt@uliege.be

23 +32/81/62.25.35

24

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

ABSTRACT

25
26 Lactoferrin (LF) is a glycoprotein naturally present in milk. Its content varies throughout the lactation
27 but also with mastitis, therefore potentially being an additional indicator of udder health beyond somatic
28 cell count. Therefore, there is an interest in quantifying this biomolecule routinely. First prediction
29 equations proposed in the literature to predict the content in milk using milk mid-infrared (MIR)
30 spectrometry were built using Partial Least Square regression (PLSR) due to the limited size of the
31 dataset. Thanks to a large dataset, the current study aimed to test fourth different machine learning
32 algorithms using a large dataset comprising 6,619 records collected across different herds, breeds and
33 countries. The first algorithm was a PLSR as used in past investigations. The second and third algorithms
34 used PLS factors combined with a linear and polynomial Support Vector regression (PLS + SVR). The
35 fourth algorithm also used PLS factors but included in an artificial neural network having one hidden
36 layer (PLS + ANN). The training and validation sets comprised 5,541 and 836 records, respectively.
37 Even if the calibration prediction performances were the best for PLS + polynomial SVR, their
38 validation prediction performances were the worse. The three other algorithms had similar validation
39 performances. Indeed, the validation root mean squared error (RMSE_v) ranged between 162.17 and
40 166.75 mg/L of milk. However, the lower standard deviation of cross-validation RMSE and the better
41 normality of the residual distribution observed for PLS + ANN suggest that this modeling was the more
42 suitable to predict the LF content in milk from milk MIR spectra ($R^2_v=0.60$ and $RMSE_v=162.17$ mg/L
43 of milk). This PLS+ANN model was then applied to almost 6 million spectral records. The predicted

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

44 LF showed the expected relationships with milk yield, somatic cell score, somatic cell count and stage
45 of lactation. The model tended to underestimate high LF values (higher than 600 mg/L of milk).
46 However, if the prediction threshold was set to 500 mg/L, 82% of samples from the validation having a
47 content of LF higher than 600 mg/L were detected. Future research should aim to increase the number
48 of those extremely high LF records in the calibration set.

49 Keywords: milk, lactoferrin, mid infrared, machine learning

For Peer Review

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

INTRODUCTION

50
51

52 Lactoferrin (LF), a 80-kDa glycoprotein naturally present in milk, is synthesized by the
53 mammary gland epithelial cells (Molenaar et al., 1996) and has antibacterial, antiviral and antifungal
54 activities potentially interesting to improve the cow's disease resistance. More details about the
55 nutraceutical and pharmaceutical properties of milk bovine LF can be found in the review of Giansanti
56 et al. (2016). These immune effects explain why the synthesis of LF increases in the presence of mastitis
57 infection (Gaunt et al., 1980; Hagiwara et al., 2003) but the responses can be different following the
58 incriminated pathogens (Kawai et al., 1999; Chaneton et al., 2008). Moreover, the content of LF in milk
59 can also vary naturally depending on parity, age, lactation stage (Gaunt et al., 1980; Hagiwara et al.,
60 2003), and breed (Król et al., 2010). Consequently, there is an interest to quantify the LF content in milk
61 at the individual level for different issues related to animal welfare (i.e., early detection of infection) and
62 human health due to the presence of this active bio-molecule in milk. Indeed, humans can take also
63 profit to administer orally LF due to its anti-infective, anti-cancer and anti-inflammatory properties
64 (Wakabayashi et al., 2006).

65 LF content can be quantified using an immunodiffusion method (Hagiwara et al., 2003) but its
66 quantification is more often based on enzyme-linked immunosorbent assay (ELISA) (Chen and Mao,
67 2004; Chaneton et al., 2013). Unfortunately, those methods are too labour intensive for routine screening
68 of the cow population at the individual scale as desired. Therefore, alternative methods offering a

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

69 quantification at a low cost must be found. In 2007, a first study was published about the prediction of
70 LF content in milk using milk mid-infrared (**MIR**) spectrometry (Soyeurt et al., 2007), a technology
71 largely implemented in most milk laboratories. This first study measured the content of LF in a limited
72 number of samples (i.e., 69 records). However, **this first study** and a follow-on one based on a bigger
73 dataset (i.e., 2,499 records) (Soyeurt et al., 2012) validated the potential of MIR to provide a relevant
74 indicator of LF. Indeed, **the cross- and external validation coefficients of determination (R^2) and root**
75 **mean square error (RMSE)** obtained from this second study were of 0.71 and 0.60, and of 50.55 and
76 **58.98 mg/L of milk, respectively.** Recently, the European Milk Recording network (EMR,
77 **www.milkrecording.eu**) has developed also its own equation from more than 2,000 records and offers
78 **the prediction service of this biomolecule to its members.** By combining all of those datasets, new
79 **perspectives are opening to try different machine learning algorithms which could maybe improve the**
80 **current accuracy of LF prediction.** Indeed, all LF prediction equations were built using Partial Least
81 **Squares regressions (PLSR).**

82 **Since the nineties, several pieces of research were conducted in dairy science using artificial**
83 **neural network (ANN), for instance, to analyze breeding dairy patterns (Finn et al., 1996), or to predict**
84 **the incidence of clinical mastitis (Yang et al., 2000) or milk yield (Grzesiak et al., 2006).** At our best
85 **knowledge, only 3 articles used mainly or partly the milk MIR spectra as predictors.** Those studies
86 **concerned the prediction of conception success for a given insemination (Hempstalk et al., 2015),**

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

87 content of blood β -hydroxybutyrate (Pralle et al., 2018) and feed intake (Dórea et al., 2018). Several
88 reasons could explain potentially the large use of PLSR to model milk MIR spectral data. Absorbance
89 values of consecutive spectral data points are highly correlated. Therefore, collinearity problems are
90 present if conventional simple linear regression is employed using all spectral data points as explanatory
91 variables. Fortunately, some solutions exist to solve this problem. The first one consists in selecting a
92 limited number of spectral points by trying to keep the most relevant information in the dataset, while
93 limiting the correlations between them. Then, those low correlated spectral points can be included in a
94 multivariate regression. The second possibility is to reduce the dimensionality of the spectral X matrix
95 by using a principal component analysis (PCA). The PCA latent variable (LV) can be then used in a
96 multivariate regression, commonly named principal component regression. However, the PCA
97 methodology used to define those LVs takes into account only the spectral variability and not the
98 variability of the trait to be predicted. This could lead to a lack of relevant spectral information to predict
99 the trait of interest. The PLS method solves this problem by defining LVs considering simultaneously
100 the variabilities of X and Y (Despaigne et al., 2000). This explains why this methodology was and is still
101 mainly used to develop milk MIR models to predict traits related to milk quality like fatty acids (Soyeurt
102 et al., 2011), cheese making properties (De Marchi et al., 2009), body weight (Soyeurt et al., 2019),
103 fertility (Delhez et al., 2020), traits related to animal welfare (Grelet et al., 2016) and environmental
104 issues (Vanlierde et al., 2016). Unfortunately, PLS is able to consider weakly non-linear relationships
105 by adding LVs, potentially leading to over-fit the developed prediction model (Thissen et al., 2004).

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

106 Other machine learning algorithms such as Support Vector Regression (SVR) and ANN have the ability
107 to model the non-linear relations (Thissen et al., 2004). Compared to ANN, SVR can deal with a high
108 number of input variables (Thissen et al. 2004). However, after a feature selection, ANN can be also
109 efficient. Moreover, using the priors obtained from a past calibration dataset, the weights used in an
110 ANN network can be updated using a new calibration dataset. This is particularly interesting when a
111 large number of phenotypes useful to predict the trait of interest is available.

112 The computational methodology differs between SVR and ANN. SVR was created by Vapnik
113 (Thissen et al., 2004) and consists in defining a classification boundary between records in order to
114 minimize the distance between the records and the boundary by taking into account a certain limit of
115 detection (called epsilon). More specifically, SVR is a method that selects a reduced number of samples,
116 the support vectors, defining the best sparse deterministic regression relationship between the MIR data
117 and the reference values. To ensure a global solution, a penalty (called C penalty) is used during the
118 computation. Different kernels can be used to compute the boundary like linear, polynomial or radial
119 kernel; each kernel has its own parameters to be optimized. More details about this method are given by
120 Thissen et al. (2004). ANN, initially introduced by McCulloch and Pitts (1943), is the basis of deep
121 learning which tries to mimic neuronal brain activity (i.e., the computer learns by experience). ANN is
122 composed of different layers including units: one input layer, one output layer and a certain number of
123 hidden layers. The number of hidden layers and its corresponding units must be defined by the user. The

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

124 higher the number of hidden layers and their corresponding units, the higher the complexity of the model
125 and therefore the higher the potential to over-fit the prediction model. The ANN algorithm aims to
126 estimate the weights of each relation among units by using, for instance, the back propagation
127 methodology. The cost function is related to the minimization of the residual error. The ANN model
128 will provide a response of one or more variables given many explanatory variables (i.e., units) (Beck,
129 2018). In conclusion, SVR and ANN require the optimization of several parameters. To achieve this
130 objective and to get a global solution, it is important to have a large dataset, explaining potentially why
131 those methods are not often used in milk MIR spectrometry as the size of the dataset is often limited.
132 However, compared to PLSR, SVR and ANN in themselves do not solve the issue of collinearity
133 between spectral data points. So, there is an interest to combine the dimension reduction obtained by
134 PLS with SVR or ANN algorithms. Therefore, the objective of the current study was to compare the
135 accuracy of predictions of milk LF content from milk MIR spectra using 4 different machine learning
136 algorithms: PLSR, linear and polynomial SVR coupled with PLS LVs and ANN coupled with PLS LVs.

137 MATERIALS AND METHODS

138 All analyses were performed using R software (version 4.0.1.; <https://www.r-project.org/>).

139 **Data**

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

140 The first dataset comprised 3,965 milk samples (50% of morning and 50% of evening milk)
141 preserved with bronopol collected between April 2005 and April 2006 in Belgium, between April 2009
142 and August 2009 in Ireland, and during August 2009 in Scotland. Part of those samples was also used
143 in a previous study (Soyeurt et al., 2012). The Belgian samples (N=549) were analyzed using one
144 MilkoScan FT6000 spectrometer (Foss, Hillerod, Denmark) located in the milk laboratory “Comité du
145 Lait” (Battice, Belgium). The Irish and Scottish samples (N=3,416) were also analyzed on the same
146 brand of spectrometer at the Animal and Grassland Research and Innovation Centre, Teagasc
147 Moorepark (Fermoy, Co. Cork, Ireland). The spectral data of each sample contained 1,060
148 wavenumbers. The second dataset comprised 2,654 milk samples (50% of morning milk and 50% of
149 evening milk) collected by the EMR in France (N=1,333), Luxembourg (N=246), England (N=500) and
150 Germany (N=575) between June 2016 and January 2017. The samples in this second dataset were
151 selected based on their LF content predicted using the equation developed by Soyeurt et al. (2012) to
152 increase the variability over what was present in the samples of the first dataset. All samples were
153 analyzed using either FT+ MilkoScan spectrometers (Foss, Hillerod, Denmark) or Bentley
154 spectrometers (Chaska, MN, United States). The spectral data were then standardized based on the
155 procedure explained by Grelet et al. (2017). All aspects related to this standardization was managed by
156 the European Milk recording network and resulted in all samples having 1,060 harmonized
157 wavenumbers and absorbance values available for analysis. A single milk sample per cow was selected
158 from animals of different breeds across several herds and countries right across lactation.

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

159 **Lactoferrin Quantification**

160 LF concentration was quantified from the milk samples already analyzed by infrared
161 spectroscopy using commercial ELISA kits: Bovine Lactoferrin ELISA Quantification kit from Bethyl
162 Laboratories Inc. (Montgomery, TX, USA) for the first dataset and e_bLF_01 kit from IDBiotech
163 (Issoire, France) for the second dataset. The Belgian samples were analyzed by Gembloux Agro-Bio
164 Tech – University of Liège (Gembloux, Belgium). The ELISA analyses of the Irish and Scottish samples
165 were conducted by Enfer Laboratories (Naas Co. Kildare, Ireland). The ELISA analysis of the second
166 dataset was conducted at Seenovia (Saint-Berthevin, France). The samples were diluted 1:1000; 1:2000;
167 1:4000; 1:6000; 1:8000 or 1:10000 in sample buffer. The LF concentrations used for the calibration
168 were the average of at least two ELISA measures taken on the same milk sample.

169 **Spectral Pre-treatment**

170 The spectral data coming from the first dataset were not standardized as this procedure did not
171 exist when the samples were collected. Therefore, in order to correct for a potential baseline drift, the
172 first derivative was applied to the recorded spectra for the dataset 1 and standardized spectra for the
173 dataset 2 using the formula:

$$174 \quad wavenumber'_i = wavenumber_i - wavenumber_{i+gap}$$

175 where $wavenumber'_i$ represents the first derivative value of the i th wavenumber, $wavenumber_i$ is the
176 raw value observed for the i th wavenumber and the gap is the windows chosen for the derivation and

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

177 was equal to 5. Then, the wavenumbers located in the most informative regions were selected. So, a
178 total of 277 wavenumbers were kept for this study and were located from 950 to 1,580 cm^{-1} ; from 1,720
179 to 1,770 cm^{-1} , from 1,780 to 1,850 cm^{-1} , and from 2,800 to 2,970 cm^{-1} .

180 The presence of potential spectral outliers was assessed by estimating the standardized
181 Mahalanobis distance (also called GH distance) for all recorded spectra. In order to allow the inversion
182 of the matrix needed to calculate the Mahalanobis distance due to the high collinearity existing between
183 some spectral points, a PCA was performed using FactoMineR package (version 1.42; Lê et al., 2008),
184 defining 22 uncorrelated principal components (PC) which explained 99.04% of the spectral variability.
185 The formula used to calculate the GH distance was:

$$186 \quad GH = \sqrt{(x - \mu)^T S^{-1} (x - \mu)} / nPC$$

187 where, x is the PC scores of a specific spectrum; μ is the mean of PC scores estimated from the entire
188 dataset; S corresponds to the (co-)variance matrix between PC scores estimated from the entire dataset;
189 nPC is the number of PC used in the calculation (i.e., 22 in our case). The PC analysis was performed
190 on a combined dataset containing all the records to improve the certainty of spectral outlier detection.
191 A total of 86 records with a GH higher than 5 were discarded from the dataset. The cleaned dataset
192 contained finally 6,533 records (i.e., 3,931 and 2,602 records of the first and second datasets,
193 respectively).

194 Data Splitting

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

195 To perform a complete external validation, the data coming from 2 different DHI organizations
196 (one in Germany and one in England) were not used to calibrate the model. This external validation
197 dataset represented 836 samples. The remaining samples (N=5,541) were used to calibrate the model
198 and performed two different cross-validations: 10-fold cross-validation where the samples were chosen
199 randomly in the calibration dataset and a leave-one DHI out cross-validation. The leave-one DHI out
200 cross-validation allows to evaluate the models with samples coming from the same context (countries,
201 diets, breeds) as mentioned by Prekopcsak et al. (2010). So, we supposed that the first dataset contained
202 9 DHI which corresponded in this case to 9 herds. The second dataset contained records coming from 6
203 different DHI but records coming from 2 DHI were kept for the external validation as explained
204 previously. Therefore, the leave-one DHI out cross-validation procedure considered 13 groups. Two
205 cross-validations were tested and compared in this study as a random N-fold cross-validation could
206 provide over-optimistic prediction performances (Wang and Bovenhuis, 2019).

207 Machine Learning Algorithms

208 All machine learning algorithms used in this study were implemented using the CARET
209 package version 6.0-86 (Kuhn, 2008). For all models, the spectral data were scaled and centered before
210 computation.

211 PLSR was performed on the 277 selected wavenumbers using the method="pls" as an argument
212 in the train function of CARET package. The maximum number of PLS latent variables was set to 50.

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

213 The optimized number of factors was chosen using the selectionFunction="oneSE" and "best" as
214 arguments in the train function of CARET package. The "best" selection function defines the optimal
215 model as the one having the lowest RMSEcv. The "oneSE" selection function allows to select a simpler
216 model having a RMSE lower within one standard error from the lowest obtained RMSEcv. This simpler
217 model is assumed to have a better generalization.

218 The computation of SVR was based on linear and polynomial kernels and was implemented
219 using the method="svmLinear" or "svmPoly" as arguments in the train function of the CARET package.
220 For both kernels, the expand.grid function was used to test different values to optimize the required
221 parameters. So, for "svmLinear", the tested C values were 0.01, 0.05, 0.1, 0.25, 0.5, 0.75, 1, 1.25, 1.5,
222 1.75, 2, and 5. For the "svmPoly" kernel, the tested values were 0.25, 0.5, 1 and 2 for C; 1 until 3 for
223 the polynomial degree; and 0,001, 0.01, 0.1 and 1 for the scale. For both kernels, the epsilon parameter
224 was set to 0.1. As for PLSR, the optimal parameters were chosen using the selection function "best" or
225 "oneSE". A total of 26 PLS factors explaining 99% of the spectral variability were combined in SVR to
226 limit the problem of overfitting. The interest in using PLS factors instead of other more conventional
227 selections of features is based on the fact that PLS will extract factors by considering simultaneously
228 the spectral variability and the variability of the trait to be predicted.

229 ANN seems to be more powerful when the selection of features is made before the modelling
230 (Thissen et al., 2004). So as performed for SVR, ANN included the 26 PLS latent variables instead of

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

231 the 277 initially selected wavenumbers. ANN based on one-layer perceptron was tested in this study.

232 This ANN architecture is composed of one hidden layer in order to minimize the risk of overfitting. To

233 estimate the weights related to this ANN design, a back propagation was used. This model was

234 performed using method="nnet" as argument in the train function of the CARET package. Different

235 numbers of units in the unique hidden layer (ranging from 1 to 5) were tested using the expand.grid

236 function. In order to ensure a global solution, a penalty was introduced during the computation of

237 weights (called decay). Decay values of 0, 0.0001, 0.001, 0.01, 0.1, 0.2, 0.3, 0.4, and 0.5 were tested

238 using the expand.grid function. As done for other algorithms, the optimized values for the size and decay

239 were chosen using the selection function "best" or "oneSE". The maximum iteration for the weight

240 estimations was set to 1000 in order to be sure to reach the convergence.

241 The prediction performances of the different models developed were assessed by estimating the

242 calibration, 10-fold cross-validation, leave-one DHI out cross-validation and external validation

243 coefficients of determination (R^2c , 10-fold R^2cv , DHI R^2cv and R^2v , respectively) as well as their

244 corresponding RMSE (i.e., $RMSEc$, 10-fold $RMSEcv$, DHI $RMSEcv$ and $RMSEv$, respectively).

245 Distributions of residuals were also studied for all models.

246 Prediction of LF from The Walloon Milk Recording Database

247 As the calibration and validation sets were composed of a limited number of records and were

248 not representative of the studied dairy population due to the sampling procedure used, there is an interest

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

249 in observing the behavior of the prediction on a large scale spectral database. Indeed, this allows
250 observing the behavior of the predictions according to known sources of variation like the stage of
251 lactation, parity, breed, and season. So, the machine learning algorithm chosen as the best based on the
252 validation prediction performances was applied on **the first derived milk MIR spectra**. The spectral
253 database is managed by the Walloon Breeding Association (Awé, Ciney, Belgium). This database is
254 related to the milk recording. A total of **5,651,470** records were collected between January 2007 and
255 **March 2020** from **349,396 cows in 1,963** herds. The average values for the predicted LF were estimated
256 according to the stage of lactation, the milk yield and the somatic cell score (SCS; $\log_2(\text{somatic cell}$
257 $\text{count (SCC)/100000}+3)$). The correlations between the predicted LF and milk yield, fat and protein
258 contents as well as SCC and SCS were also estimated.

259 RESULTS AND DISCUSSION

260 Descriptive Statistics and Data Cleaning

261 The LF content measured in the samples included in the first dataset (N=3,931) ranged from 3
262 to 2,038 mg/L of milk, with an average of 202 ± 170 mg/L of milk. LF content in the second dataset
263 (N=2,602) varied from 6 to 1,299 mg/L of milk, with an average of **325 ± 257 mg/L of milk**. The average
264 content observed in the first dataset is within the expected range compared to other published articles.
265 For instance, Gaunt et al. (1980) found an average content of LF of 266 ± 136 mg/L of milk from a first
266 set of 4 herds and $228 \text{ mg} \pm 112 \text{ mg/L}$ of milk from a second set of 4 herds. Cheng et al. (2008) found

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

267 a slightly lower content of LF in bovine milk (177 ± 120 mg/L) from samples collected on cows without
268 mastitis infection. In a previous article using a part of the first dataset (N=2,499), Soyeurt et al. (2012)
269 found an average content of LF equal to 163 ± 103 mg/L of milk. The content observed in the second
270 dataset seemed to be high compared to the literature. Moreover, the standard deviation was also higher
271 compared to the first dataset. This can be related to the sample selection. Indeed, the samples included
272 in the second dataset did not come from entire herds as they were selected based on a past LF MIR
273 predictive model in order to cover as much as possible the LF content and spectral variation. However,
274 the ranges of variation observed in both datasets were very high, especially extreme high LF
275 measurements were obtained. This could be related to the fact that some samples could be collected
276 from cows having subclinical mastitis. Indeed, some authors found a positive relationship between the
277 content of LF and the presence of mastitis (Kawai et al., 1999) even if the response differs following the
278 incriminated pathogens (Chaneton et al., 2008). For instance, Gaunt et al. (1980) measured average LF
279 content of 222 ± 168 mg/L of milk for healthy cows to 640 ± 250 mg/L of milk for cows presenting
280 mastitis. Cheng et al. (2008) obtained similar values (i.e., 742 ± 374 mg/L of milk for cows suspected
281 of having mastitis on the basis of the SCC of the milk). The distribution of LF (Figure 1) and spectra
282 (data not shown) from the 2 datasets were complementary. This was expected since samples of the
283 dataset 2 were selected to complement the dataset 1.

Lactoferrin Predictions Using Milk MIR Spectrometry

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

285 Two kinds of cross-validation procedures were tested **in this study** to fix the model parameters
286 **(i.e., number of LVs for PLS regressions; C value for linear SVR; C, scale and degree for polynomial**
287 **SVR; and size and decay for ANN). The leave-one DHI out cross-validation** leads to the model being
288 under-fit resulting a higher prediction error. Indeed, RMSE_{cv} values were higher than 175 mg/L of milk
289 and were always greater for models developed using the **leave-one DHI out cross-validation** when
290 compared to models built using the 10-fold cross-validation (Table 1). Moreover, the high RMSE_{cv} SD
291 for models building using **the leave-one DHI out cross-validation** (i.e., values higher than 90 mg/L of
292 milk) confirmed the low robustness of the developed models. This suggests that too many informative
293 samples were taken out from the calibration set. **For instance, during a 10-fold cross-validation, samples**
294 **coming from the same herd can be in the training and validation set involving relevant information to**
295 **provide a better prediction.** Consequently, the use of a 10-fold cross-validation to parametrize a model
296 is still relevant to limit the under-fitting but as mentioned by **Wang and Bovenhuis (2019)**, this procedure
297 leads to be over-optimistic concerning the prediction performances. Indeed, the observed RMSE_{cv} were
298 always lower than the one observed for the validation set (Table 1). **R²_{cv}** values obtained from models
299 developed using 10-fold cross-validation were similar between models used and ranged from 0.51 to
300 0.56. Similarly, the observed RMSE_{cv} were also globally the same and ranged from 138.40 to 144.60
301 mg/L (Table 1). This suggests similar prediction performances. **However, RMSE_{cv} SD was higher for**
302 **PLS + polynomial SVR compared to other tested algorithms.** Based on the external validation, PLSR,
303 PLS + polynomial SVR and PLS + ANN showed similar validation prediction performances with

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

304 respective RMSE values of 163.76, 166.75 and 162.17 mg/L of milk. However, the correlation values
305 between predictions on the validation set (N=836) suggested some differences. Indeed, higher
306 correlations were observed between the predictions given by the PLSR and PLS + linear SVR models
307 (0.99) compared to PLS + polynomial SVR or PLS + ANN (0.95 for both algorithms). The correlation
308 between PLS + ANN and PLS + polynomial SVR was 0.94. From Figure 2, it is clear that the
309 relationships between the predictions made from PLSR and PLS + linear SVR models is strong.
310 However, the relationship of those models with other tested ones was not linear. There appears to be a
311 saturation for low and high prediction values (i.e., S shape).

312 Therefore, even though the predictions made by the 4 models were highly correlated (i.e., higher
313 than 0.94), the low and high values behaved differently. Moreover, the range of predictions is really
314 different, with PLS + linear SVM and PLS + ANN having a reduced range compared to PLSR and PLS
315 + polynomial SVR. Moreover, except PLS + ANN, all other tested algorithms had the tendency to
316 predict negative values (Table 2). The correlation between residuals and predicted content of LF ranged
317 from 0.64 for PLS + ANN to 0.77 for PLS + linear SVM based on the training set. From the validation
318 set, these correlation values were comprised between 0.60 for PLS + ANN to 0.83 for PLS + linear
319 SVM. These correlation values were lower with the squared residuals (from 0.49 for PLS + ANN to
320 0.61 for PLS + linear SVR and from 0.30 for PLS + ANN to 0.68 for PLS + linear SVR based on the
321 training and validation sets, respectively). This suggests that higher errors were made for samples having

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

322 a high content of LF. The validation prediction performances, the robustness (low RMSE_{cv} SD) (Table
323 1), the prediction of positive values (Table 2) and the lowest correlation between squared residuals and
324 LF content observed for PLS + ANN suggest that this modeling is the most relevant to predict daily LF
325 content in milk from milk MIR spectrometry. For another application dedicated to dairy science, Dórea
326 et al. (2018) obtained also better prediction performances using ANN including one hidden layer after
327 a selection of input variables compared to PLSR to predict feed intake. Pralle et al. (2018) obtained
328 similar performances for PLSR and ANN including also one hidden layer to predict blood β -
329 hydroxybutyrate.

330 Compared to the previous studies published by our team on the same topic, the prediction error
331 observed in the current study is higher than the one observed in the past (RMSE_v = 77.26 mg/L of milk
332 in Soyeurt et al. (2012) vs. 163.76 mg/L in the current study based on PLSR or 162.17 mg/L of milk for
333 PLS+ANN). However, this is difficult to compare as the validation dataset was not the same and all
334 spectra were not standardized. If the equation published in 2012 (Soyeurt et al., 2012) is applied on the
335 current validation set, the validation prediction error is of 462 mg/dL with a R²_v equal to 0.02. Even if
336 this old equation was built from a part of the dataset 1 (2499 samples), the variability in the current
337 datasets is higher and the past prediction equation is not suitable to predict those records. This could
338 also be related to the fact that 2 different ELISA kits were used. However, the residual distributions
339 obtained after using the PLS + ANN model were similar for datasets 1 and 2 (data not shown). Moreover,

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

340 an analysis of variance also confirmed that the differences between residuals observed from datasets 1
341 and 2 were not significant. It is interesting to note that 50% of records had a residual error between –
342 69.06 and 51.66 mg/L of milk and between -85.78 and 91.37 mg/L of milk for the training and validation
343 set, respectively.

344 Besides the interest to predict the quantity of lactoferrin in milk, it could be useful to know if
345 the models were able to detect extreme values. Indeed, even if a high prediction error exists for the
346 records having a high content of lactoferrin, the prediction increased as expected but with a lower
347 intensity. PLS + ANN model gave predictions allowing detection of 65% of the records with a content
348 of LF higher than 600 mg/L of milk in the training set if we fixed the prediction limit to 500 mg/L of
349 milk. This proportion reached to 82% for the samples in the validation set. The threshold of 600 mg/L
350 was used as the RMSE started to be related to the content of LF from this content and because authors
351 like Gaunt et al. (1980) and Kawai et al. (1999) mentioned that a such high content is potentially
352 related to cows having mastitis.

353 Due to the distribution of LF observed in Figure 1, we have also tested the log transformation
354 but the results were not better (data not shown).

355 PLS + ANN Model Applied to a Large Spectral Database

356 PLS + ANN model was applied to 5,651,470 records from cows within the first 365 days in
357 milk. The obtained average prediction was 307.80 mg/L of milk with a SD of 209.17 mg/L of milk.

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

358 The minimum and maximum values were 19.74 and 1121.00 mg/L, respectively. So as observed on
359 the training and validation datasets, no negative predictions were observed.

360 The LF content predicted using MIR varied according to the stage of lactation (Figure 3). This
361 variation was already observed by Gaunt et al. (1980). The differences per lactation stages of the log
362 transformed LF contents obtained by predictions and found by Cheng et al. (2008) and Hagiwara et al.
363 (2003) were very similar (Table 3) even if the contents observed in this study were closer to the ones
364 obtained by Hagiwara et al. (2003).

365 We also observed a negative correlation with milk yield (-0.24) but positive with fat and
366 protein contents (0.11 and 0.28, respectively). Cheng et al. (2008) also found a strongly positive
367 relationship with protein ($r=0.48$) but they mentioned that the correlation with fat content was not
368 significantly different from 0. However, the negative relationship between LF and milk yield was
369 stronger for Cheng et al. (2008) ($r=-0.47$). The difference of LF after log transformation and observed
370 by the level of milk yield was similar even if the contents found in this study were higher (Table 4).

371 As expected, positive correlations with predicted LF were observed for SCC and SCS (0.21 and 0.30,
372 respectively; $N=5,477,197$). Cheng et al. mentioned that the correlation between LF and SCC was not
373 significantly different than 0 but the correlation found by these authors for SCS ($r=0.37$) was similar to
374 the one estimated in the current study. The evolution of log-transformed predicted LF was also in
375 agreement with the results found by Hagiwara et al. (2003) and Cheng et al. (2008) (Table 5).

CONCLUSIONS

376

377 This study tried 4 examined machine learning algorithms to predict the daily content of
378 lactoferrin in cow's milk from milk MIR spectral data. It found, based on the validation prediction
379 performances, that PLS, PLS + polynomial SVR and PLS + ANN provided similar results, but that the
380 model using PLS factors combined with an ANN was the best. This model was then applied to the
381 Walloon milk recording spectral database to observe the relationships between predicted LF content and
382 the main milk components as well as SCC and SCS which were in line with the literature. However, the
383 model still had some difficulties in predicting extremely high values. Indeed, the observed RMSE
384 increased strongly once LF content exceeded 600 mg/L of milk; 12 percent of records coming from the
385 Walloon dairy cow population reached this level of LF production. Including extreme values of milk LF
386 content to the calibration set could help to improve the prediction models. We now have the possibility
387 to directly predict the content of LF in the milk lab allowing identification of those specific samples.
388 Moreover, as the quantity of milk required for the ELISA analysis is quite low, the same sample as the
389 one analyzed for the routine milk recording could be used. This could be an appropriate task to improve
390 in the future the ability of the ANN network to discriminate low and high LF samples. Until now, no
391 implementation of this LF prediction is done by DHI. However, there is an interest for them to use this
392 molecule information to improve the detection of subclinical mastitis. The inclusion of the LF trait into

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 4 MACHINE LEARNING ALGORITHMS

393 breeding program to improve the cow robustness or the milk nutritional quality has not been investigated
394 yet.

395 ACKNOWLEDGMENTS

396 This research received a financial support from the European Commission, Directorate-General for
397 Agriculture and Rural Development, under Grant Agreement 211708 and from the Commission of the
398 European Communities, FP7, KBBE-2007-1. This paper does not necessarily reflect the view of these
399 institutions and in no way anticipates the Commission's future policy in this area. The Ministry of
400 Agriculture of the Walloon Region of Belgium [Service public de Wallonie, Direction générale de
401 l'Agriculture, des Ressources naturelles et de l'Environnement, Direction de la Recherche] is
402 acknowledged for their financial support through the project NovaUdderHealth (D31-1207). Scottish
403 Agricultural College receives funding from the Scottish Government and they are acknowledged for
404 funding the long-term selection experiment producing milk samples used in this analysis. The work of
405 Sinead McParland was funded by a Science Foundation Ireland Starting Investigator Research Grant
406 (18/SIRG/5562).

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

407References

408

- 409Beck, M.W. 2018. NeuralNetTools: Visualization and analysis tools for neural networks. *J. Stat. Softw.*
410 85:1–20. doi:10.18637/jss.v085.i11.
- 411Chaneton, L., M. Bontá, M. Pol, L. Tirante, and L.E. Bussmann. 2013. Milk lactoferrin in heifers:
412 Influence of health status and stage of lactation. *J. Dairy Sci.* 96:4977–4982. doi:10.3168/jds.2012-
413 6028.
- 414Chaneton, L., L. Tirante, J. Maito, J. Chaves, and L.E. Bussmann. 2008. Relationship between milk
415 lactoferrin and etiological agent in the mastitic bovine mammary gland. *J. Dairy Sci.* 91:1865–1873.
416 doi:10.3168/jds.2007-0732.
- 417Chen, P.W., and F.C. Mao. 2004. Detection of lactoferrin in bovine and goat milk by enzyme-linked
418 immunosorbent assay. *J. Food Drug Anal.* 12:133–139.
- 419Cheng, J.B., J.Q. Wang, D.P. Bu, G.L. Liu, C.G. Zhang, H.Y. Wei, L.Y. Zhou, and J.Z. Wang. 2008.
420 Factors Affecting the Lactoferrin Concentration in Bovine Milk. *J. Dairy Sci.* 91:970–976.
421 doi:10.3168/JDS.2007-0689.
- 422Delhez, P., P.N. Ho, N. Gengler, H. Soyeurt, and J.E. Pryce. 2020. Diagnosing the pregnancy status of
423 dairy cows: How useful is milk mid-infrared spectroscopy?. *J. Dairy Sci.* 103:3264–3274.
424 doi:10.3168/jds.2019-17473.
- 425Despaigne, F., D. Luc Massart, and P. Chabot. 2000. Development of a robust calibration model for
426 nonlinear in-line process data. *Anal. Chem.* 72:1657–1665. doi:10.1021/ac991076k.
- 427Dórea, J.R.R., G.J.M. Rosa, K.A. Weld, and L.E. Armentano. 2018. Mining data from milk infrared
428 spectroscopy to improve feed intake predictions in lactating dairy cows. *J. Dairy Sci.* 101:5878-
429 5889.
- 430Finn, G.D., R. Lister, T. Szabo, D. Simonetta, H. Mulder, and R. Young. 1996. Neural Networks Applied
431 to a Large Biological Database to Analyse Dairy Breeding Patterns. *Neural Comput. & Applic.*
432 4:237-253.
- 433Gaunt, S.N., N. Raffio, E.T. Kingsbury, R.A. Damon, W.H. Johnson, and B.A. Mitchell. 1980. Variation
434 of Lactoferrin and Mastitis and Their Heritabilities. *J. Dairy Sci.* 63:1874–1880.
435 doi:10.3168/jds.S0022-0302(80)83154-7.
- 436Giansanti, F., G. Panella, L. Leboffe, and G. Antonini. 2016. Lactoferrin from milk: Nutraceutical and
437 pharmacological properties. *Pharmaceuticals* 9:1–15. doi:10.3390/ph9040061.
- 438Grelet, C., C. Bastin, M. Gelé, J.B. Davière, M. Johan, A. Werner, R. Reding, J.A. Fernandez Pierna, F.G.
439 Colinet, P. Dardenne, N. Gengler, H. Soyeurt, and F. Dehareng. 2016. Development of Fourier
440 transform mid-infrared calibrations to predict acetone, β -hydroxybutyrate, and citrate contents in
441 bovine milk through a European dairy network. *J. Dairy Sci.* 99:4816–4825. doi:10.3168/jds.2015-
442 10477.
- 443Grelet, C., J.A.F. Pierna, P. Dardenne, H. Soyeurt, A. Vanlierde, F. Colinet, C. Bastin, N. Gengler, V.
444 Baeten, and F. Dehareng. 2017. Standardization of milk mid-infrared spectrometers for the transfer
445 and use of multiple models. *J. Dairy Sci.* 100:7910–7921. doi:10.3168/jds.2017-12720.
- 446

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

- 447 Grzesiak, W., P. Blaszczyk, and R. Lacroix. 2006. Methods of predicting milk yield in dairy cows—
448 Predictive capabilities of Wood's lactation curve and artificial neural networks (ANNs). *Computers*
449 *and Electronics in Agriculture* 54:69–83.
- 450 Hagiwara, S.I., K. Kawai, A. Anri, and H. Nagahata. 2003. Lactoferrin concentrations in milk from
451 normal and subclinical mastitic cows. *J. Vet. Med. Sci.* 65:319–323. doi:10.1292/jvms.65.319.
- 452 Hempstalk, K., S. McParland, and D.P. Berry. 2015. Machine learning algorithms for the prediction of
453 conception success to a given insemination in lactating dairy cows. *J. Dairy Sci.* 98:5262-5273.
- 454 Kawai, K., S. Hagiwara, A. Anri, and H. Nagahata. 1999. Lactoferrin Concentration in Milk of Bovine
455 Clinical Mastitis. *Vet. Res. Commun.* 23:391–398. doi:10.1023/A:1006347423426.
- 456 Król, J., Z. Litwińczuk, A. Brodziak, and J. Barłowska. 2010. Lactoferrin, lysozyme and immunoglobulin
457 G content in milk of four breeds of cows managed under intensive production system. *Pol. J. Vet.*
458 *Sci.* 13:357–361.
- 459 Kuhn, M. 2008. caret Package. *J. Stat. Softw.* 28:1–26.
- 460 Lê, S., J. Josse, and F. Husson. 2008. FactoMineR: An R package for multivariate analysis. *J. Stat. Softw.*
461 25:1–18. doi:10.18637/jss.v025.i01.
- 462 De Marchi, M., C.C. Fagan, C.P. O'Donnell, A. Cecchinato, R. Dal Zotto, M. Cassandro, M. Penasa, and
463 G. Bittante. 2009. Prediction of coagulation properties, titratable acidity, and pH of bovine milk
464 using mid-infrared spectroscopy. *J. Dairy Sci.* 92:423–432. doi:10.3168/jds.2008-1163.
- 465 McCulloch, W.S., and W. Pitts. 1943. A logical calculus of the ideas imminent in nervous activity.
466 *Bulletin of Mathematical Biophysics* 5:115-133. doi:10.1007/bf02478259.
- 467 Molenaar, A.J., Y.M. Kuys, S.R. Davis, R.J. Wilkins, P.E. Mead, and J.W. Tweedie. 1996. Elevation of
468 Lactoferrin Gene Expression in Developing, Ductal, Resting Regressing Parenchymal Epithelium of
469 the Ruminant Mammary Gland. *J. Dairy Sci.* 79:1198–1208. doi:10.3168/jds.S0022-0302(96)76473-
470 1.
- 471 Pralle, R.S., K.W. Weigel, and H.M. White. 2018. Predicting blood β -hydroxybutyrate using milk Fourier
472 transform infrared spectrum, milk composition, and producer-reported variables with multiple linear
473 regression, partial least squares regression, and artificial neural network. *J. Dairy Sci.* 101:4378-
474 4387.
- 475 Prekopcsak, Z., T. Henk, and C. Gaspar-Papanek. 2010. Cross-validation : the illusion of reliable
476 performance estimation. RCOMM RapidMiner Community Meet. Convergence 1–6.
- 477 Soyeurt, H., C. Bastin, F.G. Colinet, V.M.R. Arnould, D.P. Berry, E. Wall, F. Dehareng, H.N. Nguyen, P.
478 Dardenne, J. Schefers, J. Vandenplas, K. Weigel, M. Coffey, L. Thé Ron, J. Detilleux, E. Reding, N.
479 Gengler, and S. McParland. 2012. Mid-infrared prediction of lactoferrin content in bovine milk:
480 Potential indicator of mastitis. *Animal* 6:1830–1838. doi:10.1017/S1751731112000791.
- 481 Soyeurt, H., F.G. Colinet, V.M.R. Arnould, P. Dardenne, C. Bertozzi, R. Renaville, D. Portetelle, and N.
482 Gengler. 2007. Genetic variability of lactoferrin content estimated by mid-infrared spectrometry in
483 bovine milk. *J. Dairy Sci.* 90:4443–4450. doi:10.3168/jds.2006-827.
- 484 Soyeurt, H., F. Dehareng, N. Gengler, S. McParland, E. Wall, D.P. Berry, M. Coffey, and P. Dardenne.
485 2011. Mid-infrared prediction of bovine milk fatty acids across multiple breeds, production systems,
486 and countries. *J. Dairy Sci.* 94:1657–1667. doi:10.3168/jds.2010-3408.
- 487 Soyeurt, H., E. Froidmont, I. Dufrasne, D. Hailemariam, Z. Wang, C. Bertozzi, F.G. Colinet, F. Dehareng,

MID-IRRED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

- 488 and N. Gengler. 2019. Contribution of milk mid-infrared spectrum to improve the accuracy of test-
489 day body weight predicted from stage, lactation number, month of test and milk yield. *Livest. Sci.*
490 227:82–89. doi:10.1016/j.livsci.2019.07.007.
- 491Thissen, U., M. Pepers, B. Üstün, W.J. Melssen, and L.M.C. Buydens. 2004. Comparing support vector
492 machines to PLS for spectral regression applications. *Chemom. Intell. Lab. Syst.* 73:169–179.
493 doi:10.1016/j.chemolab.2004.01.002.
- 494Vanlierde, A., M.L. Vanrobays, N. Gengler, P. Dardenne, E. Froidmont, H. Soyeurt, S. McParland, E.
495 Lewis, M.H. Deighton, M. Mathot, and F. Dehareng. 2016. Milk mid-infrared spectra enable
496 prediction of lactation-stage-dependent methane emissions of dairy cattle within routine population-
497 scale milk recording schemes. *Anim. Prod. Sci.* 56:258–264. doi:10.1071/AN15590.
- 498Wang, Q., and H. Bovenhuis. 2019. Validation strategy can result in an overoptimistic view of the ability
499 of milk infrared spectra to predict methane emission of dairy cattle. *J. Dairy Sci.* 102:6288–6295.
500 doi:10.3168/jds.2018-15684.
- 501Wakabayashi, H., K. Yamauchi, and M. Takase. 2006. Lactoferrin research, technology and applications.
502 *Int. Dairy J.* 16:1241-1251.
- 503Yang, X.Z., R. Lacroix, and K.M. Wade. 2000. Investigation into the production and conformation traits
504 associated with clinical mastitis using artificial neural networks. *Can. J. Anim. Sci.* 80: 415–426.
505
506

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

507 Table 1. Ten-fold cross-validation and external validation performances for predicting lactoferrin content

508 in milk using four different machine learning algorithms.

Selection function		PLSR	PLS + Linear SVR	PLS + Polynomial SVR	PLS + ANN
		oneSE	oneSE	best	best = oneSE*
Calibration (N=5541)	parameters	nLV=23	C=5	degree=3; scale=0.01; C=1	size=4; decay=0.5
	R ² c	0.53	0.53	0.64	0.60
	RMSEc	140.94	144.32	125.89	130.59
Cross-validation	R ² cv	0.51	0.53	0.56	0.55
	R ² cv SD	0.03	0.03	0.03	0.03
	RMSEcv	144.31	144.60	138.40	139.01
	RMSEcv SD	5.77	5.61	8.08	5.05
	RPD	1.43	1.42	1.49	1.48
External validation (N=836)	R ² v	0.61	0.63	0.62	0.60
	RMSEv	163.76	174.92	166.75	162.17

509 PLSR = Partial least squares regression; PLS + Linear SVR = Linear Support Vector Regression based on 26 PLS latent
510 variables; PLS + Polynomial SVR = Linear SVR based on 26 PLS latent variables; PLS + ANN = modeling based on
511 artificial neural network including 26 PLS latent variables in the input layer and one hidden layer; nLV = number of
512 PLS latent variables; C = cost penalty used for SVR; * = the selection function 'best' and 'oneSE' provided the same
513 results.

514

515

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

516 Table 2. Data distribution of the reference lactoferrin contents as well as the predictions obtained from the
 517 developed models after a 10-fold cross-validation.

		Data distribution (mg/L of milk)								
		0%	5%	10%	25%	50%	75%	90%	95%	100%
Training set	Lactoferrin	3.01	34.00	57.22	99.84	179.72	302.69	532.39	669.97	2038.27
	PLSR	-277.95	35.10	77.14	142.16	218.59	308.55	430.42	521.86	1244.38
	PLS + linSVR	-195.05	37.97	74.54	128.68	194.55	271.75	385.47	466.19	1060.68
	PLS + polSVR	-47.40	51.77	74.18	119.48	187.15	271.52	413.47	532.49	1423.58
	PLS + ANN	31.27	77.53	95.71	135.09	188.21	289.78	472.00	605.24	1053.65
Validation set	Lactoferrin	6.00	32.32	47.46	122.30	280.58	476.69	707.86	860.52	1286.86
	PLSR	-162.60	13.98	70.29	165.56	320.69	420.68	518.84	591.00	912.53
	PLS + linSVR	-129.09	22.84	68.17	150.40	290.59	367.41	455.30	520.24	795.47
	PLS + polSVR	-72.59	27.30	55.13	123.53	263.11	381.74	533.62	626.71	1200.00
	PLS + ANN	36.67	67.71	85.11	127.48	313.83	483.87	633.11	682.23	851.38

518 PLSR = Partial least squares regression; PLS + linSVR = linear Support Vector Regression based on 26 PLS latent
 519 variables; PLS + polSVR = polynomial SVR based on 26 PLS latent variables; PLS + ANN = modeling based on an
 520 artificial neural network including 26 PLS latent variables in the input layer and one hidden layer.

521

522

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

523 Table 3. Evolution of lactoferrin content predicted by MIR following the stage of lactation and comparison
 524 with the literature.
 525

	N	MIR lactoferrin mg/L of milk		log(lactoferrin)		log(lactoferrin) [Cheng et al.,2008]			log(lactoferrin) [Hagiwara et al., 2003]*		
		Mean	SD	Mean	SD	N	Mean	SD	N	Mean	SD
DIM ≤ 20	317,322	229.75	181.22	2.23	0.33						
20 > DIM ≤ 100	1,486,501	245.45	176.28	2.29	0.30	49	1.90	0.12	8	2.06	0.43
100 > DIM ≤ 200	1,736,757	297.04	198.64	2.37	0.30	45	2.03	0.28	59	2.23	0.38
200 > DIM ≤ 365	2,110,890	372.31	223.49	2.48	0.29	28	2.20	0.20	32	2.30	0.45

526 DIM = days in milk; * The range of DIM was slightly different

527

528

529

MID-IRRED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

530Table 4. Evolution of lactoferrin content predicted by MIR following the milk yield (kg/day) and
531comparison with the literature.

	N	MIR lactoferrin mg/L of milk		log(lactoferrin)		log(lactoferrin) [Cheng et al.,2008]		
		Mean	SD	Mean	SD	N	Mean	SD
milk yield < 20 kg	1,611,617	374.43	228.40	2.48	0.30	36	2.18	0.25
20 kg ≤ milk yield < 25 kg	1,333,645	308.94	204.31	2.39	0.30	34	1.99	0.23
25 kg ≤ milk yield < 30 kg	1,201,798	283.77	194.82	2.35	0.30	33	1.93	0.16
30 kg ≥ milk yield	1,504,410	254.63	181.96	2.30	0.30	19	1.89	0.18

532

533

534

535

536

For Peer Review

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

537Table 5. Evolution of lactoferrin content predicted by MIR following the somatic cell score (SCS) and
538comparison with the literature.

SCS	N	MIR lactoferrin		log(lactoferrin)		log(lactoferrin) [Cheng et al.,2008]			log(lactoferrin) [Hagiwara et al., 2003]		
		Mean	SD	Mean	SD	N	Mean	SD	N	Mean	SD
0	163,716	236.23	166.91	2.28	0.29	12	1.91	0.14	36	2.18	0.19
1	996,710	233.82	169.95	2.27	0.29	20	2.02	0.17	28	2.16	0.42
2	1,365,672	267.61	184.44	2.33	0.30	50	1.98	0.19	39	2.27	0.51
3	1,164,258	305.80	199.50	2.39	0.30	40	2.06	0.26			
4	860,534	347.21	214.41	2.45	0.30	34	2.10	0.25			
5	510,941	382.85	227.91	2.49	0.30	20	2.26	0.32			
6	292,339	403.68	234.51	2.52	0.30	22	2.28	0.30			
7	160,828	421.79	234.70	2.54	0.30						
8	83,575	461.05	235.88	2.59	0.28						
9	42,340	569.97	244.31	2.70	0.25						

539* 10,557 records were deleted because the SCS had a negative value.

540

541

542

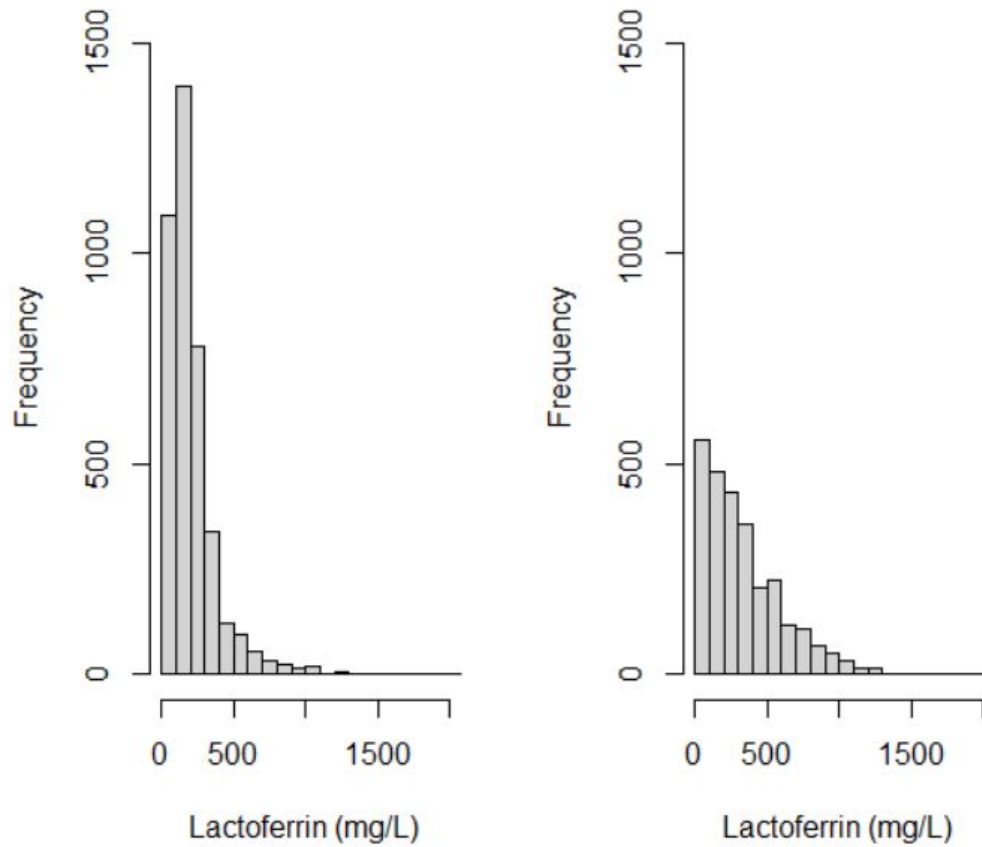
543

544

545

546

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS



547

548 Figure 1. Distribution of ELISA Lactoferrin quantifications in the first dataset on the left and the second

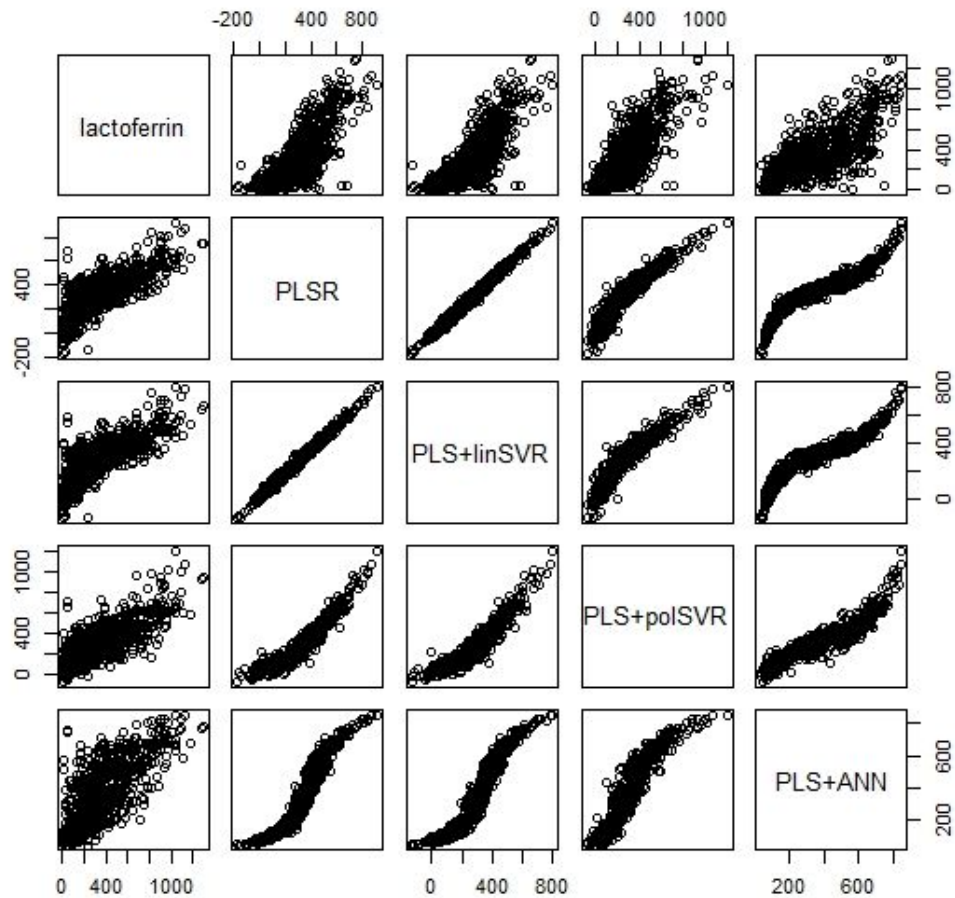
549 dataset on the right.

550

551

552

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS

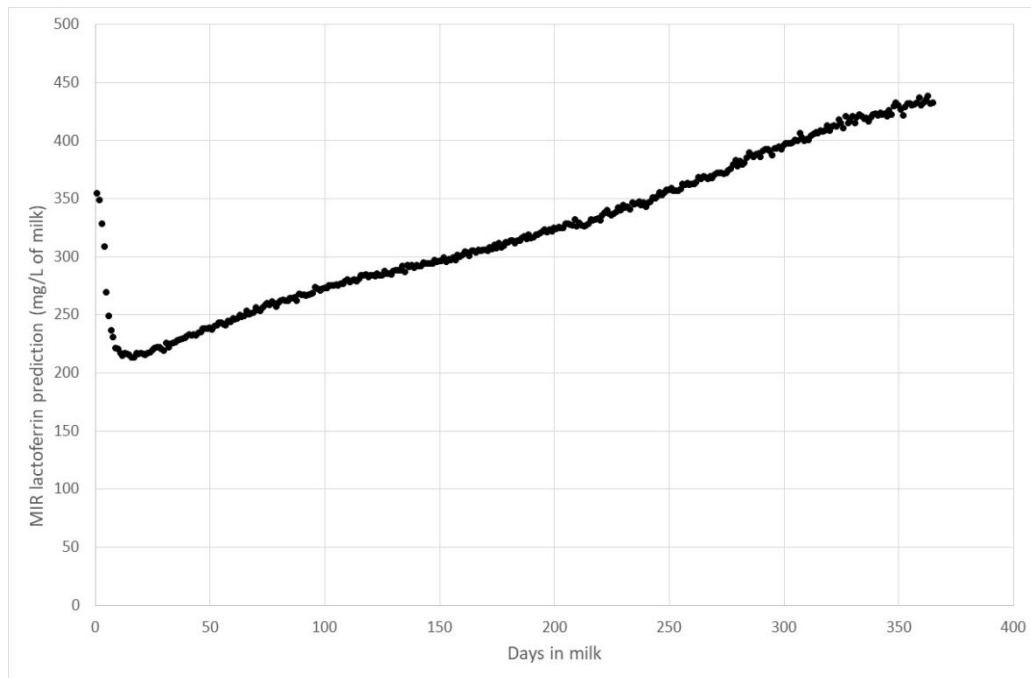


553

554 **Figure 2.** Relationships between reference lactoferrin content (mg/L of milk) and the predictions obtained
 555 using fourth different machine learning approaches applied on the validation set (PLSR = Partial Least
 556 Squares Regression; PLS+linSVR = 26 PLS factors included in a linear Support Vector Regression;
 557 PLS+polSVR = 26 PLS factors included in a polynomial SVR; PLS+ANN = 26 PLS factors included in an
 558 artificial neural network having one hidden layer).

559

MID-INFRARED LACTOFERRIN PREDICTION IN MILK THROUGH 3 MACHINE LEARNING ALGORITHMS



560

561 Figure 3. Evolution of lactoferrin content predicted by mid-infrared spectrometry following the stage of
562 lactation.