

1 **Accepted version** of manuscript *From spreadsheets to sugar content modeling: a*
2 *data mining approach* This version of the work has been peer-reviewed.

3 Published version of this work is available at:

4 <https://doi.org/10.1016/j.compag.2016.11.012>

5

6 Authors' emails were updated and do not necessarily correspond to the ones in the
7 published work.

8

9

10 **From spreadsheets to sugar content modeling: a data mining approach**

11 Monique Pires Gravina de Oliveira ^{a,b}, Felipe Ferreira Bocca ^{a,c}, Luiz Henrique

12 Antunes Rodrigues ^{a,d}

13 ^a School of Agricultural Engineering, University of Campinas

14 ^b moniquepgoliveira@gmail.com

15 ^c -

16 ^d lique@feagri.unicamp.br - Corresponding author

17

18 **Abstract**

19 Sugarcane mills need sugar content estimates in advance to establish their
20 commercial strategy. To obtain these estimates, mills rely on historical averages
21 or maturation curves. Crop models have also been developed to provide those
22 estimates. Leveraging mill data about fields and sugar content at harvest, we
23 developed empirical models using different data mining techniques along with the
24 RReliefF algorithm for feature selection. The best model was attained with
25 Random Forest with features selected by RReliefF, having a mean absolute error
26 of 2.02 kg Mg⁻¹. This model outperformed Support Vector Regression and
27 Regression Trees with and without feature selection. Models were also evaluated
28 by the Regression Error Characteristic Curves, which showed that the best model
29 was able to predict 90% of the observations within a precision of 5.40 kg Mg⁻¹.

30

31 **Keywords**

32 Sugarcane Ripening; Machine Learning; Empirical Modeling; Total Recoverable
33 Sugar; Crop Modeling

34

35 **1. Introduction**

36 Planning in the sugar industry requires estimates about the amount of sugar
37 that will be produced in the following cycle. This information is then used in
38 forward selling, forward pricing, and managing storage and shipping schedules
39 (Everingham et al., 2007). In Brazil, where harvests mostly occur from April until
40 November, the commercial strategy for the following season starts being
41 established in August of the current season (Bocca et al., 2015). Overestimates
42 could compromise previous selling commitments while underestimates could lead
43 to difficulties in storing and shipping (Everingham et al., 2003). Sugar estimates
44 are also useful for operational level plans, such as prioritizing harvesting areas
45 (Scarpari and Beauclair, 2004).

46 Two forecasts are required to achieve such estimates: sugarcane (*Saccharum*
47 spp.) fresh mass yield and sugar content in sugarcane stalk (Alvarez et al., 1982;
48 Bocca et al., 2015). The former has been addressed by Everingham et al. (2009,
49 2016) and by Bocca and Rodrigues (2016). For the latter, industries use either
50 averages from the previous years or variety-specific maturity curves (Scarpari and
51 Beauclair, 2004, 2009). Both approaches, however, do not allow for the inclusion
52 of factors that favor sucrose storage in sugarcane stalks, e.g. weather variability
53 and management practices (van Heerden et al., 2013). Particularly for the case of
54 weather, the increase in weather variability leads to the need of tools to assess the
55 effect of weather uncertainty in production. The urge for climate risk assessment
56 is increasing among companies as an effect of climate change (Surminski, 2013).

57 To take weather variability and management practices into account, crop
58 models could also be explored. Crop yield models gather knowledge about crop

59 growth and development and are able to predict its behavior (Boote et al., 1996;
60 Monteith, 1996).

61 There are mainly two approaches to modeling: to deepen understanding and
62 knowledge of a topic and to make accurate predictions for specific decisions.
63 Frequently, different levels of both can be found in most models (Shmueli, 2010;
64 Singels, 2013). The first approach is seen in models that simulate sugarcane
65 phenological and physiological processes. These models try and describe plants'
66 processes and deepen the understanding of plant physiology and its interactions
67 with the environment (Passioura, 1996; O'Leary, 2000; Singels, 2013). To
68 achieve higher prediction accuracy, aiming at production planning, one could use
69 empirical models, which are independent of the simulations aforementioned.

70 Empirical models are conceptually simpler models and are based on
71 relationships between crop outputs and its driving factors, e. g. water availability,
72 weather conditions, and agricultural practices (Monteith, 1996; Passioura, 1996;
73 Surendran Nair et al., 2012; Singels, 2013). The relationships explored in these
74 models vary from proxies to direct effects, such as the distance from the lake
75 feature used by Alvarez et al. (1982) and variety, respectively.

76 Scarpari and Beauclair (2004, 2009) developed empirical models to predict
77 total recoverable sugar by using stepwise regression. In the paper published in
78 2004, the only variables used by the authors were negative degree-days and
79 available water content during crop development. In 2009, they added another
80 one, concerning photoassimilate production. Despite aiming not to make
81 predictions but to describe the relationship between variables, Lawes et al. (2002)
82 modeled commercial cane sugar by using linear mixed models. Their final model
83 included the year, month of harvest, farm of origin, variety and an interaction

84 between the month of harvest and year of harvest. More recently, Cardozo et al.
85 (2015) established an exponential relationship between total recoverable sugar
86 and accumulated rainfall in the 120 days before the harvest.

87 In 1982, Alvarez et al. (1982) had already highlighted not only the vast
88 number of variables that could affect sugarcane yield, but also the complexity of
89 the relations between them. Different approaches have been used to address these
90 issues. Scarpari and Beauclair (2009) generated a set of models: one model was
91 fitted for each combination of variety, number of cuts and type of management
92 zone, for early, mid and late period of harvest during the season. Lawes et al.
93 (2002), in turn, used pairwise combinations of variables, while Cardozo et al.
94 (2015) selected one variable most correlated to sugar content to be included in
95 their three models, for each ripening pattern.

96 These examples draw attention to the limitations of the methods being used
97 to model sugar content: they either assume linearity, do not extensively account
98 for interactions or both. Also, they should not be directly used for non-normal
99 data with auto-correlated features, which underlines the need for other techniques,
100 such as those highlighted by Breiman (2001), which he called algorithmic models,
101 referring to the models obtained by data mining or machine learning techniques.

102 Data mining techniques have been long applied in agriculture-related
103 problems, e.g. prediction of wine-fermentation results, evaluation of imperfections
104 in fruits, both with images and X-ray, classification of sounds from pigs and birds,
105 meat analysis and the use of energy in agriculture (Mucherino et al., 2009). The
106 successful application of these techniques is due to their capacity to deal with the
107 previously mentioned aspects of agricultural data.

108 One further reason to use these other techniques is the availability of data.
109 Lawes et al. (2002) stated that for the Australian production context, some sugar
110 mills collect block-productivity data such as cane yield and commercial cane
111 sugar from every block or paddock harvested during the season, as well as
112 information on the block size, the cane variety, the time of harvest and how many
113 ratoons the cane has. Data collection for Brazilian context is not only similar but
114 also enhanced by the fact that the mill is either owner or responsible for the
115 production (Bocca et al., 2015) and therefore has additional information regarding
116 soil analysis and agricultural practices.

117 Furthermore, the use of data mining techniques allows for more accurate
118 models since they can identify new and unknown patterns in large datasets
119 (Witten et al., 2011). An attempt in this direction has already been made by
120 Everingham and Sexton (2011), although still with a limited number of variables.
121 It is possible to achieve better estimates by exploring more variables, and by
122 looking at further available algorithms.

123 Bocca et al. (2015) suggested the use of yield models associated to climate
124 forecasts and production data in an integrated system in order to obtain yield
125 forecasts. In this study, we present the development of an element of this system:
126 a sugar content model that could be used in conjunction with both weather
127 forecasts and production data. To model sugar content (Total Recoverable Sugar -
128 TRS), we use a commercial sugarcane production database and the data mining
129 framework (feature selection, parameter tuning, modeling and validation in an
130 independent set).

131

132 **2. Materials and methods**

133 **2.1. Dataset**

134 Data used in this study were supplied by Alcídia mill, operated by
135 Odebrecht Agroindustrial, located in the city of Teodoro Sampaio, state of São
136 Paulo (SP), Brazil. The mill annual production area is almost 25 thousand hectares
137 of land and its production reaches 1.6 million tons of sugarcane. Harvests that
138 happened in 2011 and 2012 provided 2,102 observations, with each observation
139 referring to one block in the farms in each year. The 53 variables of the dataset
140 belong to four categories: soil physics and soil chemistry, weather, agricultural
141 practices, and those related to the crop (Table 1).

142 It is worth noting that some variables were created, particularly regarding
143 the developmental stages of the crop, based on the planting dates. Plant cycle was
144 simplified into four stages: (1) sprouting, (2) tillering, (3) growth and (4) maturity.
145 With this approach, we could group weather and phenological information,
146 providing estimates of the weather in each of the plant's stage, rather than
147 averages for the whole cycle.

148 Variables that delve too much into the cycle, i.e. that are intrinsically related
149 to harvest, such as the occurrence of pests, that is only verified by harvesting time
150 and cannot be predicted in advance, were removed. The remaining variables are
151 either defined in the beginning of the cycle, as is the case for fertilization, or refer
152 to the weather and can be estimated by weather forecasts.

153 Two scenarios were modeled: (a) using all available features, and (b)
154 performing feature selection and using only the selected features in the modeling
155 process. Feature selection will be further explained in section 2.2.2.

156

157 **2.2. Model development**

158 **2.2.1. Algorithms deployed**

159 In data mining, the prediction of a continuous variable, such as Total
160 Recoverable Sugar, is known as a regression problem. In this paper, three
161 techniques were used to tackle this problem: Support Vector Regression (SVR),
162 Random Forests (RF) and Regression Trees (RT). Statistical software R, version
163 3.1.1 (R Core Team, 2015) was used in the modeling process with packages
164 *e1071* (Meyer et al., 2014), *randomForest* (Liaw and Wiener, 2002) and *rpart*
165 (Therneau et al., 2014).

166

167 **2.2.2. Feature selection**

168 Feature selection was only performed on the part of the dataset reserved for
169 training, with 1,402 observations. The algorithm that was chosen to perform
170 feature selection was RReliefF (Robnik-Šikonja and Kononenko, 2003) as it is
171 able to estimate the quality of attributes in problems with strong dependencies
172 between attributes. Since the dataset is comprised of weather and edaphic
173 features, this has turned out to be an important characteristic. The parameters
174 number of neighbors and number of iterations were kept as suggested in Robnik-
175 Šikonja and Kononenko (2003) and Robnik-Sikonja and Kononenko (1997): 10
176 and 100, respectively. Importance values provided by the algorithm were
177 averaged after 10 repetitions. The RReliefF algorithm implemented in the
178 CORElearn package was used (Robnik-Šikonja et al., 2015).

179 Robnik-Šikonja and Kononenko (2003) state that importance value is
180 analogous to the percentage of explained variance if scaled to sum 1. Based on
181 that, the criterion to limit the number of features was to select only the best-

182 ranked attributes that accounted for 0.9 – 90% – of the RReliefF explained
183 variance.

184 We chose not to perform the variable importance for Regression Tree and
185 Random Forest so that we could use the same variables, chosen by the criteria
186 proposed by RReliefF, in all models.

187

188 **2.2.3. Parameter tuning**

189 To achieve better results, parameters were tuned using a two-stage grid
190 search. The second stage used smaller step sizes for the grid search in the region
191 close to the best result found in the first stage. Different parameters were searched
192 for the different feature sets used. Details of the ranges and step size for each
193 parameter can be found in the Appendix. Optimal values obtained can be seen in
194 Table 2. Tuning and modeling processes used only data available in the training
195 set, in the two situations modeled, after the groups of features were defined.

196

197 **2.2.4. Validation**

198 Models were trained in a dataset constituted by two-thirds of the records and
199 validation was achieved through testing in the remaining third (700 observations),
200 still unseen by the model. This strategy is also known as hold-out validation.

201 Three metrics were computed to report model quality: root mean squared
202 error (RMSE), mean absolute error (MAE) and coefficient of determination
203 between model output and observed values (R^2). The REC Curve, developed by
204 Bi and Bennet (2003), was also used with this objective. All metrics of evaluation
205 were reported for the hold-out data.

206

207 **3. Results and discussion**

208 **3.1. Feature Selection**

209 As expected, Algorithm RReliefF was able to identify those variables more
210 often associated with sugar accumulation, e.g. the occurrence of precipitation
211 events and low temperatures close to harvesting (van Heerden et al., 2013). The
212 algorithm was also able to distinguish other variables that could bring information
213 to the model, even though they are less obvious, such as Silt in the intermediate
214 layer. RReliefF selected 31 features from the 53 in the original set (Figure 1).

215 One can see how weather and crop variables predominate, as not only all of
216 the available ones were selected, but also all of the 13 most important belong to
217 one of those two categories (up to Mean Min Temp 1 - Figure 1). Soil variables
218 were already deemed less important and only 7 remained from the original 24.
219 Even though soil variables are usually autocorrelated, which could have lowered
220 their importance (Robnik-Šikonja and Kononenko, 2003), weather variables are as
221 well. Despite that, these variables had the highest values of importance,
222 highlighting how much they influenced sugar production and accumulation
223 processes.

224 We do not aim to try and explain every feature chosen by the algorithm, but
225 some will draw attention to other important points. In the already mentioned case
226 of Silt, its selection could have been caused by autocorrelation among features.
227 Along with Texture, they are probably proxies for other soil features and, due to
228 the aforementioned behavior of RReliefF, only Silt and Texture showed enough
229 importance.

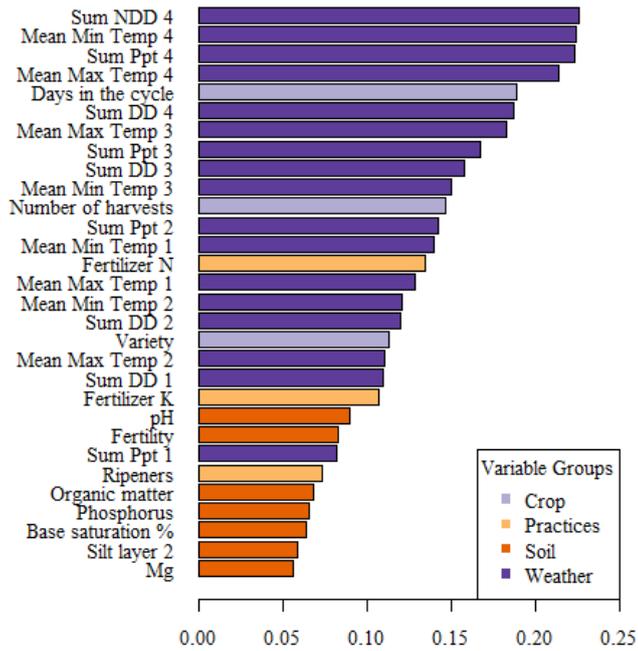
230 The use of soil fractions as features for the models requires a brief
231 discussion. Since soil fractions are compositional data, this characteristic must be

232 addressed for the use in models such as least squares linear regressions in which
233 independence between features is important and compositional data explicitly
234 violate this assumption. Empirical investigations seem to indicate that such
235 transforms are not required (Ranganathan and Borges, 2011), or that they can be
236 alleviated by the use of GLM models, in which one can specify different error
237 distribution and link functions, or by the use of random forests (Lopatin et al.,
238 2016). In a practical sense, tree methods are invariant to monotonic
239 transformations in the data, meaning these transformations would not change the
240 models. For Support Vector Machines/Regression, if transformations do not
241 change the relative distances between samples, relation between inputs and
242 outputs would not change, but some changes in the internal structure could
243 happen.

244 Reproducing this methodology with another dataset would probably not
245 present the same outcome, i.e. instead of Silt, Clay could have been chosen. Other
246 ways of engineering the data, such as including latent variables, could also have
247 lead to different outcomes. Despite that, the most important variables are expected
248 to be quite similar.

249 Another noteworthy result in this selection is the low score presented by the
250 Ripeners feature, which may have occurred due to bias introduced in the dataset.
251 In both cases, when ripeners were applied and when they were not, TRS would
252 never present an excessively low value, or else the area would not have been
253 harvested. Since low TRS fields were not harvested, they did not make it into the
254 commercial dataset. This illustrates how biases might often be unaccounted for in
255 the modeling process.

256



257

258 Figure 1. Ranked feature importance of features selected by the RReliefF
 259 algorithm. Weather features names follow the pattern of Mean/Sum, variable
 260 studied and period for the calculation (1 for sprouting, 2 for tillering, 3 for growth
 261 and 4 for maturity), e.g. Sum DD 4 is the sum of degree days in the maturity
 262 period and Mean Max Temp 1 is the mean value of the daily maximum air
 263 temperature for the sprouting period.

264

265 This unexpectedness of some selected features, nevertheless, creates
 266 opportunities for investigating other unknown variables that play a role in sucrose
 267 accumulation in sugarcane stalks.

268

269 3.2. Modeling

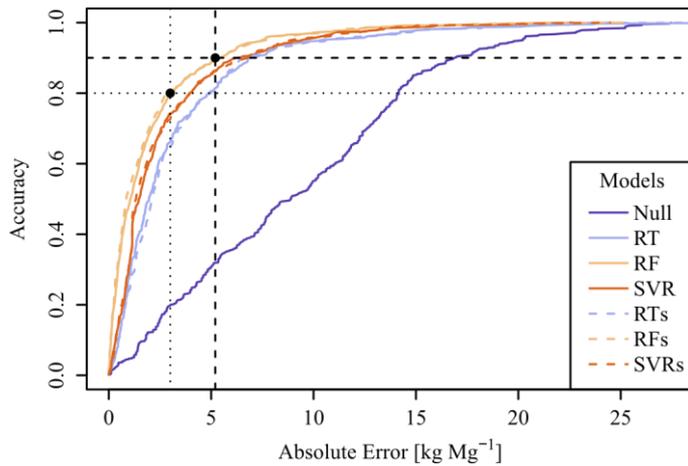
270 On average, developed models achieved mean absolute error values that
 271 represent about 2.5% of the smallest value in the test set. This result means that

272 decision makers are prone to forecast errors that depend almost exclusively on the
273 quality of the weather forecast models available.

274 Based on MAE or RMSE values, one could affirm the best prediction
275 outcome came from algorithm Random Forest both when feature selection was
276 performed and when it was not (Table 3). However, the REC curve allows for
277 further understanding of the prediction errors (Figure 2). We observe that once the
278 acceptable error is larger than 15 kg Mg⁻¹, all models present similar performance.
279 This means criteria other than the accuracy of the predictions might be used, e.g.
280 processing speed and interpretability. Regression trees often have an advantage of
281 interpretability (Breiman, 2001), but taking into account the complexity parameter
282 achieved through tuning, this is not the case for the models obtained. The trees
283 turned out to have more than 150 nodes.

284 One should note that not only the Random Forest mean absolute error is
285 2.02 kg Mg⁻¹, but that 80% of the predicted values have a mean absolute error
286 lower than 3.00 kg Mg⁻¹ and that 90% have an error lower than 5.40 kg Mg⁻¹.
287 These limits – if acceptable to decision makers – reinforce the model’s utility in
288 the decision support process.

289



290

291 Figure 2. REC Curve. Dotted line indicates an error threshold of 3.00 kg Mg⁻¹.

292 Dashed line indicates an error threshold of 5.40 kg Mg⁻¹. Letter “s” indicates

293 feature selection.

294

295 Even though the coefficient of determination between model output and

296 observed values should not be used as a validation metric (Harrison, 1990;

297 Mitchell, 1997), we do so (Figure 3) as this is the only reported metric regarding

298 the validation process in some of previous research papers.

299 We stress that comparisons of the present models with previous works

300 should consider the differences in the validation strategy. Everingham and Sexton

301 (2011) used a cross-validation approach and obtained R² of 0.55, while Scarpari

302 and Beauclair (2004) reported R² of the training set of 0.70.

303 Likewise, we believe our results should not be compared to the ones

304 obtained by Scarpari and Beauclair (2009) and by Cardozo et al. (2015) due to the

305 different strategies of validation. In the former paper, several models were

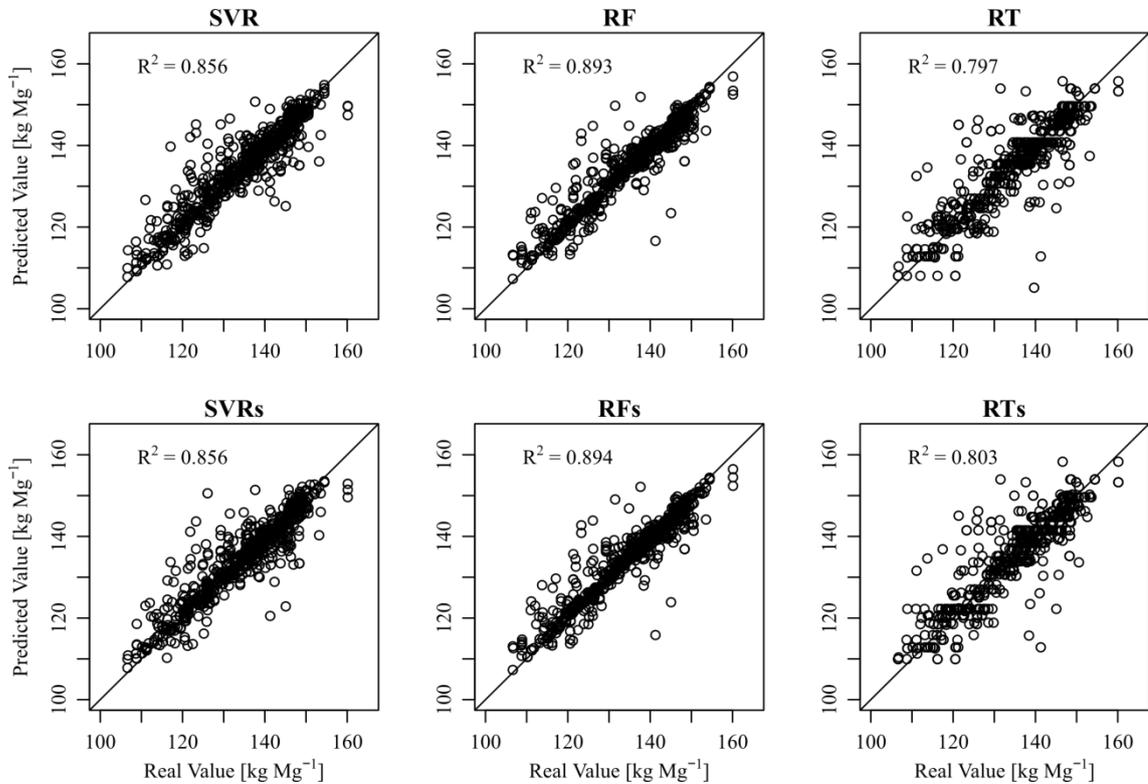
306 developed and different coefficients are presented for each model. In the latter,

307 measured values of TRS are averages of all of the blocks in each city, weighted by

308 the proportion of each group of cultivars, for each month. The average value is

309 then compared to the predicted TRS for each month of harvest. Due to this

310 aggregation, observed data variability is reduced, diminishing residuals from
 311 comparison between predicted and observed values.
 312



313
 314 Figure 3. Real values vs. predicted values obtained by the different models.

315 Techniques used were Support Vector Regression (SVR), Random Forests (RF)
 316 and Regression Tree (RT). Straight line indicates the 1:1 reference. Letter “s”
 317 indicates feature selection.

318
 319 A data mining approach is different from those used in previous modeling,
 320 not only regarding the choice of which variables would be included in the model,
 321 but also because it allows for variables that are less generic than Farm and Year.
 322 Lawes et al. (2004) pointed out that the Farm and Year variables capture the
 323 influence of management and local variations in the environment. In variables
 324 such as month or year there is an implicit consideration of past weather, which

325 happened during the months of growth of the observations in the training set.
326 Breaking these variables into several others allows the algorithms to extract the
327 knowledge that each individually brings to the pool, as well as prevent inclusion
328 in the model of information that would not be useful. Also, when variables are
329 explicitly considered, they allow for including estimates, such as the amount of
330 nitrogen applied for the cycle.

331 We acknowledge that the modeling errors do not correspond to forecasting
332 errors, since weather data used correspond to past data. These models, when used
333 with forecasting intentions, have to be used in combination with weather
334 forecasting technologies (Everingham et al., 2007). In the system proposed by
335 Bocca et al. (2015), in what concerns production data, some estimates should be
336 included, such as amount of fertilizer applied, ripener application, etc.

337 Therefore, the highest importance of weather variables in the last period has
338 both a positive and negative impact on the potential use of the models for
339 operational decision-making. Concerning decisions regarding harvesting dates –
340 which can be changed late in the crop cycle – weather forecasts will be needed
341 only a few months in advance, leading to better estimates. On the other hand,
342 decisions that depend on long-range forecasts will have larger errors for the latest
343 months.

344 As for feature selection, one can note how similar the results were when the
345 technique was performed and when it was not. We expected improvements in the
346 results, since the technique was developed to deal with performance degradation
347 both in speed and in performance accuracy due to high dimensionality (Kira and
348 Rendell, 1992). Even though there was no pronounced improvement, the similar

349 performance despite the lower number of variables implies that the developed
350 models were more robust.

351 We underline that there are also other well-known techniques that could
352 have been explored, e.g. neural networks, boosted models or even creating an
353 ensemble of the models. However, due to its iterative nature, data-mining could
354 be indefinitely performed for improved performance (more data, more features,
355 refined tuning, etc.), and we have already shown the better performance of data
356 mining techniques.

357 Finally, we must add that due to the scope of the data used to fit our models,
358 their use should be limited to the mill that provided the data. At most, they could
359 be used by mills with similar conditions of soil and weather. If the conditions
360 differ, our workflow should be repeated with data from the different mill.

361

362 **4. Conclusions**

363 Data mining techniques not only showed potential in Total Recoverable
364 Sugar prediction at block level, but also showed an improvement when compared
365 to previous modeling attempts. We attributed this improvement to the inclusion of
366 more detailed data and to the refinement of these techniques. Developed models
367 were allowed to explore relationships other than linear or first degree interactions.

368 By exploring a feature selection technique, we reinforced the importance of
369 weather variables and we also could detect the introduction of bias in the dataset.
370 If the technique is to be used, said bias should be removed by either not including
371 blocks in which ripeners were applied or by including data before and after
372 ripeners were applied.

373 We evaluated models according to several metrics, and all of them
374 supported the use of these techniques in decision making in the sugar industry,
375 considering the availability of appropriate weather data. Particularly, the use of
376 the REC curve provided ranges of accuracy that can also be used by decision
377 makers.

378

379 **5. Acknowledgements**

380 This work was supported by Bioen/Fapesp and Odebrecht Agro-industrial
381 (Process # 12/50049-3 - Link: [http://www.bv.fapesp.br/en/auxilios/55622/data-](http://www.bv.fapesp.br/en/auxilios/55622/data-mining-techniques-applied-to-the-analysis-and-prediction-of-sugarcane-yield/)
382 [mining-techniques-applied-to-the-analysis-and-prediction-of-sugarcane-yield/](http://www.bv.fapesp.br/en/auxilios/55622/data-mining-techniques-applied-to-the-analysis-and-prediction-of-sugarcane-yield/));
383 and CAPES (Coordination for the Improvement of Higher Education Personnel).

384

385 **6. References**

386 Alvarez, J., Crane, D.R., Spreen, T.H., Kidder, G., 1982. A yield prediction model
387 for Florida sugarcane. *Agr. Syst.* 9, 161–179. doi:10.1016/0308-
388 521X(82)90018-X

389 Bi, J., Bennet, K.P., 2003. Regression error characteristic curves, in: Proceedings
390 of the Twentieth International Conference on Machine Learning (ICML-
391 2003). Washington DC, pp. 43–50.

392 Bocca, F.F., Rodrigues, L.H.A., 2016. The effect of tuning, feature engineering,
393 and feature selection in data mining applied to rainfed sugarcane yield
394 modelling. *Comput. Electron. Agric.* 128, 67–76.
395 doi:10.1016/j.compag.2016.08.015

396 Bocca, F.F., Rodrigues, L.H.A., Arraes, N.A.M., 2015. When do I want to know
397 and why? Different demands on sugarcane yield predictions. *Agr. Syst.* 135,
398 48–56. doi:10.1016/j.agsy.2014.11.008

399 Boote, K.J., Jones, J.W., Pickering, N.B., 1996. Potential uses and limitations of

- 400 crop models. *Agron. J.* 88, 704–716.
401 doi:10.2134/agronj1996.00021962008800050005x
- 402 Breiman, L., 2001. Statistical Modeling: The Two Cultures. *Stat. Sci.* 16, 199–
403 231. doi:10.1214/ss/1009213726
- 404 Cardozo, N.P., Sentelhas, P.C., Panosso, A.R., Palhares, A.L., Ide, B.Y., 2015.
405 Modeling sugarcane ripening as a function of accumulated rainfall in
406 Southern Brazil. *Int. J. Biometeorol.* doi:10.1007/s00484-015-0998-6
- 407 Demattê, J.L.I., Demattê, J.A.M., 2009. Ambientes de produção como estratégia
408 de manejo na cultura da cana-de-açúcar. *Informações Agronômicas* 10–18.
- 409 Everingham, Y., Sexton, J., Skocaj, D., Inman-Bamber, G., 2016. Accurate
410 prediction of sugarcane yield using a random forest algorithm. *Agron.*
411 *Sustain. Dev.* 36, 27. doi:10.1007/s13593-016-0364-z
- 412 Everingham, Y.L., Inman-Bamber, N.G., Thorburn, P.J., McNeill, T.J., 2007. A
413 Bayesian modelling approach for long lead sugarcane yield forecasts for the
414 Australian sugar industry. *Aust. J. Agric. Res.* 58, 87–94.
415 doi:10.1071/AR05443
- 416 Everingham, Y.L., Muchow, R.C., Stone, R.C., Coomans, D.H., 2003. Using
417 southern oscillation index phases to forecast sugarcane yields: A case study
418 for northeastern Australia. *Int. J. Climatol.* 23, 1211–1218.
419 doi:10.1002/joc.920
- 420 Everingham, Y.L., Sexton, J., 2011. An introduction to multivariate adaptive
421 regression splines for the cane industry, in: *Proceedings of the 2011*
422 *Conference of the Australian Society of Sugar Cane Technologists.*
- 423 Everingham, Y.L., Smyth, C.W., Inman-Bamber, N.G., 2009. Ensemble data
424 mining approaches to forecast regional sugarcane crop production. *Agric.*
425 *For. Meteorol.* 149, 689–696. doi:10.1016/j.agrformet.2008.10.018
- 426 Harrison, S.R., 1990. Regression of a model on real-system output: An invalid test
427 of model validity. *Agr. Syst.* 34, 183–190. doi:10.1016/0308-

428 521X(90)90083-3

429 Kira, K., Rendell, L., 1992. A practical approach to feature selection, in:
430 Proceedings of the Ninth International Workshop on Machine Learning.
431 Morgan Kaufmann Publishers Inc., pp. 249–256.

432 Lawes, R.A., McDonald, L.M., Wegener, M.K., Basford, K.E., Lawn, R.J., 2002.
433 Factors affecting cane yield and commercial cane sugar in the Tully district.
434 *Aust. J. Exp. Agric.* 42, 473–480. doi:10.1071/EA01020

435 Lawes, R.A., Wegener, M.K., Basford, K.E., Lawn, R.J., 2004. The evaluation of
436 the spatial and temporal stability of sugarcane farm performance based on
437 yield and commercial cane sugar. *Aust. J. Agric. Res.* 55, 335–344.
438 doi:10.1071/AR03169

439 Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest. *R*
440 News 2, 18–22.

441 Lopatin, J., Dolos, K., Hernández, H.J., Galleguillos, M., Fassnacht, F.E., 2016.
442 Comparing Generalized Linear Models and random forest to model vascular
443 plant species richness using LiDAR data in a natural forest in central Chile.
444 *Remote Sens. Environ.* 173, 200–210. doi:10.1016/j.rse.2015.11.029

445 Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., Leisch, F., 2014. e1071:
446 Misc Functions of the Department of Statistics (e1071), TU Wien.

447 Mitchell, P.L., 1997. Misuse of regression for empirical validation of models.
448 *Agr. Syst.* 54, 313–326. doi:10.1016/S0308-521X(96)00077-7

449 Monteith, J.L., 1996. The quest for balance in crop modeling. *Agron. J.* 88, 695–
450 697. doi:10.2134/agronj1996.00021962008800050003x

451 Mucherino, A., Papajorgji, P.J., Pardalos, P.M., 2009. *Data Mining in Agriculture,*
452 Springer Optimization and Its Applications. Springer New York, New York,
453 NY. doi:10.1007/978-0-387-88615-2

454 O’Leary, G.J., 2000. A review of three sugarcane simulation models with respect

- 455 to their prediction of sucrose yield. *Field Crop. Res.* 68, 97–111.
456 doi:10.1016/S0378-4290(00)00112-X
- 457 Passioura, J.B., 1996. Simulation models: Science, snake oil, education, or
458 engineering? *Agron. J.* 88, 690–694.
459 doi:10.2134/agronj1996.00021962008800050002x
- 460 R Core Team, 2015. *R: A Language and Environment for Statistical Computing.*
- 461 Ranganathan, Y., Borges, R.M., 2011. To transform or not to transform. *Plant*
462 *Signal. Behav.* 6, 113–116. doi:10.4161/psb.6.1.14191
- 463 Robnik-Šikonja, M., Kononenko, I., 2003. Theoretical and empirical analysis of
464 ReliefF and RReliefF. *Mach. Learn.* 53, 23–69.
465 doi:10.1023/A:1025667309714
- 466 Robnik-Šikonja, M., Kononenko, I., Kononenko, I., 1997. An adaptation of Relief
467 for attribute estimation in regression, in: *Proceedings of the Fourteenth*
468 *International Conference on Machine Learning - ICML '97.* Morgan
469 Kaufmann Publishers Inc., San Francisco, CA, pp. 296–304.
- 470 Robnik-Šikonja, M., Savicky, P., Alao, J.A., 2015. *CORElearn: Classification,*
471 *Regression and Feature Evaluation.*
- 472 Scarpari, M.S., Beauclair, E.G.F. de, 2009. Physiological model to estimate the
473 maturity of sugarcane. *Sci. Agric.* 66, 622–628. doi:10.1590/S0103-
474 90162009000500006
- 475 Scarpari, M.S., Beauclair, E.G.F. de, 2004. Sugarcane maturity estimation through
476 edaphic-climatic parameters. *Sci. Agric.* 61, 486–491. doi:10.1590/S0103-
477 90162004000500004
- 478 Shmueli, G., 2010. To Explain or to Predict? *Stat. Sci.* 25, 289–310.
479 doi:10.1214/10-STS330
- 480 Singels, A., 2013. *Crop Models*, in: *Sugarcane: Physiology, Biochemistry, and*
481 *Functional Biology.* John Wiley & Sons Ltd, pp. 541–577.

482 doi:10.1002/9781118771280.ch20

483 Surendran Nair, S., Kang, S., Zhang, X., Miguez, F.E., Izaurralde, R.C., Post,
484 W.M., Dietze, M.C., Lynd, L.R., Wullschleger, S.D., 2012. Bioenergy crop
485 models: Descriptions, data requirements, and future challenges. *GCB*
486 *Bioenergy* 4, 620–633. doi:10.1111/j.1757-1707.2012.01166.x

487 Surminski, S., 2013. Private-sector adaptation to climate risk. *Nat. Clim. Chang.*
488 3, 943–945. doi:10.1038/nclimate2040

489 Therneau, T., Atkinson, B., Ripley, B., 2014. *rpart: Recursive Partitioning and*
490 *Regression Trees.*

491 van Heerden, P.D.R., Eggleston, G., Donaldson, R.A., 2013. Ripening and
492 Postharvest Deterioration, in: *Sugarcane: Physiology, Biochemistry, and*
493 *Functional Biology.* John Wiley & Sons Ltd, pp. 55–84.
494 doi:10.1002/9781118771280.ch4

495 Witten, I.H., Frank, E., Hall, M.A., 2011. *Data Mining: Practical Machine*
496 *Learning Tools and Techniques, Data Mining: Practical Machine Learning*
497 *Tools and Techniques.* Elsevier. doi:10.1016/B978-0-12-374856-0.00001-8

498

Appendix

Table 4. Ranges explored in the parameter tuning process.

Algorithms	Parameters	First tuning	Second tuning
SVR	Gamma	From 2^{-9} to 2^1 . Step: 2^n	20 points linearly spaced from 2^{-6} to 2^{-3} .
	Cost	From 2^0 to 2^{10} . Step: 2^n	10 points linearly spaced from 2 to 8.
RF	Percentages of attributes in the split	15%, 22%, 33%, 44%, 55%, 66%, 77%	10 points linearly spaced from 33% to 55%.
	Minimum size of terminal nodes	1,2,5,10,20	1,2
	Number of trees to grow	100,250,500,750, 1000, 1500	20 points linearly spaced from 250 to 1500.
RT	Minimum number of objects in a node for a split	2,5,10,25,50	-
	Complexity parameter	0.001, 0.005,0.01,0.05,0.1, 0.25	-

FIGURE CAPTIONS

Figure 1. Ranked feature importance of features selected by the RReliefF algorithm.

Weather features names follow the pattern of Mean/Sum, variable studied and period for the calculation (1 for sprouting, 2 for tillering, 3 for growth and 4 for maturity), e.g. Sum DD 4 is the sum of degree days in the maturity period and Mean Max Temp 1 is the mean value of the daily maximum air temperature for the sprouting period.

Figure 2. REC Curve. Dotted line indicates an error threshold of 3.00 kg Mg^{-1} . Dashed line indicates an error threshold of 5.40 kg Mg^{-1} . Letter “s” indicates feature selection.

Figure 3. Real values vs. predicted values obtained by the different models. Techniques used were Support Vector Regression (SVR), Random Forests (RF) and Regression Tree (RT). Straight line indicates the 1:1 reference. Letter “s” indicates feature selection.

Table 1. List of variables used for modeling in the full dataset.

Variable		Details [Unit]
Soil	Sand, Clay, and Silt	Numeric variables. Each percentage is available in three depths: 0-0.25 m, 0.25-0.50 m and 0.80-1.00 m [%].
	Texture	Categorical variable with seven levels. The percentage of clay in the soil is the main criteria used to define the levels. ^(a)
	Fertility	Categorical variable with seven levels. Levels are determined based on soil percent base saturation, soil aluminum saturation, and soil pH. ^(a)
	Soil Density	Numeric variable [g cm^{-3}].
	Chemical properties	Numeric variables. Extractable aluminum [$\text{mmol}_c \text{ dm}^{-3}$], extractable calcium [$\text{mmol}_c \text{ dm}^{-3}$], extractable potassium [$\text{mmol}_c \text{ dm}^{-3}$], extractable magnesium [$\text{mmol}_c \text{ dm}^{-3}$], soil pH, phosphorus concentration [mg dm^{-3}], organic matter [g dm^{-3}].
	Derived chemical properties	Numeric variables. Cation-exchange capacity [$\text{mmol}_c \text{ dm}^{-3}$], exchangeable acidity [$\text{mmol}_c \text{ dm}^{-3}$], percent base saturation [%], aluminum saturation [%], and sum of bases [$\text{mmol}_c \text{ dm}^{-3}$].
Weather	Sum of degree days (Sum DD), Accumulated	Numeric variables. Each one is available for the four stages of plant's

	precipitation (Sum Ppt), Mean maximum air temperature (Mean Max Temp), Mean minimum air temperature (Mean Min Temp)	development: (1) sprouting, (2) tillering, (3) growth and (4) maturity.
	Sum of Negative Degree Days (Sum NDD 4)	Numeric variable. Only available in the maturity stage. ^(b)
Agricultural practices	Ripeners	Categorical variable referring to whether ripeners were applied or not
	Vinasse	Numerical variable referring to the irrigation depth using vinasse [mm]
	Cake Filter	Numerical variable referring to the amount of cake filter applied [Mg ha ⁻¹]
	Fertilization rates of Nitrogen, Phosphorus, Potassium and Molybdenum.	Rates were calculated according to applied bulk rate and formula percentages. Particularly for Potassium, rate could also have been determined by using its concentration in vinasse and how much vinasse was applied. [kg ha ⁻¹]
	Source of Potassium	Categorical variable referring to whether K source is either mineral fertilizer or vinasse.
Crop related	Variety	Categorical variable
	Days in the cycle	Numeric variable referring to the interval from planting or previous harvest to the

	following harvest
Number of harvests	Categorical variable with labels 1 to 6. Harvests that happened within 12, 15 and 18 months after planting or previous harvest are hardcoded as 1.
Total Recoverable Sugar (TRS)	Numeric variable for the mass of sugar in a metric ton of sugarcane. Expressed in [kg Mg ⁻¹].

(a) As proposed by Demattê and Demattê (2009). (b) Negative degree days as described by Scarpari and Beauclair (2004).

Table 2. Parameters defined by tuning for the full dataset and for the dataset with selected features.

Algorithms	Parameters	Full dataset	Dataset with feature selection
SVR	Gamma	0.02138158	0.06743421
	Cost	6	4
	Percentages of attributes in the split	43% (23 attributes)	43% (13 attributes)
RF	Minimum size of terminal nodes	1	1
	Number of trees to grow	579	382
RT	Minimum number of objects in a node for a split	2	2
	Complexity parameter	0.001	0.001

Table 3. Mean Absolute Errors and Root Mean Squared Errors of data mining techniques in the considered scenarios (kg Mg⁻¹).

Features used		Regression	Random	SVR	Null
		Tree	Forest		model*
MAE	All features	3.23	2.08	2.66	9.26
	Selected features	3.27	2.02	2.64	9.26
RMSE	All features	5.06	3.62	4.18	11.05
	Selected features	4.98	3.60	4.19	11.05

*Predicted values correspond to the mean TRS value in the training set.