

# Forecasting the black Sigatoka development rate: A comparison of machine learning techniques

Luis-Alexander Calvo-Valverde<sup>a,1</sup>, Mauricio Guzmán-Quesada<sup>b</sup>, José-Antonio Guzmán-Alvarez<sup>b</sup>, Pablo Alvarado-Moya<sup>c</sup>

<sup>a</sup>*DOCINADE, Instituto Tecnológico de Costa Rica, Computer Research Center, Multidisciplinar program eScience, CNCA/CeNAT, Cartago, Costa Rica*

<sup>b</sup>*Dirección de Investigaciones, Corporación Bananera Nacional S.A., Guápiles, Costa Rica*

<sup>c</sup>*DOCINADE, Instituto Tecnológico de Costa Rica, Cartago, Costa Rica*

---

## Abstract

Pending.

*Keywords:* Machine learning, Black Sigatoka, Support vector regression, Banana disease prediction, Biological warning system

---

## 1. Introduction

The black Sigatoka disease caused by the fungus *Mycosphaerella fijiensis* Morelet is the major pathological problem of banana and plantain crops in Central America, Panama, Colombia and Ecuador, as in many parts of Africa and Asia [6].

This disease attacks the plant leaves producing a rapid deterioration of the leaf area, affects the growth and productivity of the plants due to the impairment of their photosynthetic ability causes a reduction in the quality of the fruit, and promotes premature maturation of bunches, which is the major cause of product losses associated with the black Sigatoka. Figure.1 shows three stages of this disease.

Phytopathological studies point out that precipitation, temperature, relative humidity and wind are the main climatic variables that affect its development

---

*Email address:* lualcava.sa@gmail.com (Luis-Alexander Calvo-Valverde)

<sup>1</sup>Corresponding author. (506)70104420

14 [6].

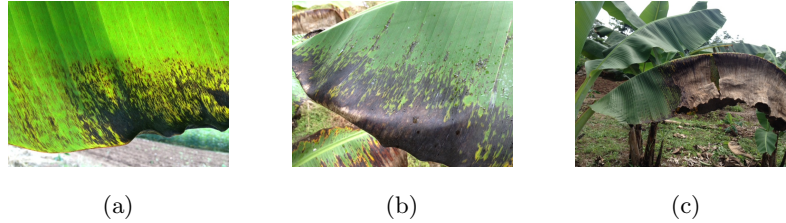


Figure 1: Examples of three disease stages of the black Sigatoka. (a) Initial stage. (b) Intermediate stage, and (c) Advanced stage.

15 According to studies by the National Banana Corporation of Costa Rica  
16 (Corbana) made in 2013, considering on average between 53 thru 57 cycles of  
17 fungicide applications per farm, the cost per hectare per year ranged between  
18 \$1800 USD and \$1900 USD. This represents about 0.76 cents of the price of  
19 a box of 18.14 kilograms. Overall, this represents 10% to 12% of the total  
20 production cost Brescani [1].

21 The past and present rates of disease development can in principle be used  
22 to predict its future behavior and to determine whether particular fungicide  
23 spray schedules will be able to effectively and economically control the disease  
24 Chuang and Jeger [3].

25 There are efforts to apply machine learning methods for decision-making in  
26 agriculture, including the control of crop diseases. For example, [Camargo et al.,2012]  
27 present an intelligent system for the assessment of crop disorders, [17] introduce  
28 a plant virus identification method based on neural networks with an evolu-  
29 tionary preprocessing stage, [5] summarize in their survey crop pests prediction  
30 methods using regression and machine learning approaches, while [7] present an  
31 intelligent agricultural forecasting system based on wireless sensor networks.

32 In this work, we compare four machine learning techniques (support vector  
33 regression (SVR), echo state networks (ESN), ridge regression and ordinary least  
34 squares linear regression) to predict the development rate of the black Sigatoka  
35 disease.

36 The main contribution of this work is a comparison between machine learning

37 methods to forecast black Sigatoka development rate. (FALTA COMPLETAR)

## 38 **2. Materials and methods**

### 39 *2.1. Concepts*

#### 40 *2.1.1. Biological warning system*

41 This system measures the disease development state to determine when to  
42 apply fungicides [6]. This system is based on two components: a climate com-  
43 ponent, which is given by the Piche evaporation and a biological component,  
44 given by the stage of progress or the rate of disease development. Originally,  
45 this system was designed to work with young plants. One selected plant must  
46 exhibit a normal growth and be in a place that enforces a healthy development.  
47 The plant must start with 5 to 6 true leaves. The assessments are made at  
48 fixed intervals of seven days as long as possible, on the same plant. The first  
49 observations should consider the leaf emission, also the level of infection on the  
50 leaves should be evaluated considering the stages of development [6].

#### 51 *2.1.2. Support Vector Regression (SVR)*

From the perspective of Support Vector Regression (SVR) the regression  
function  $y = f(s)$  for a given dataset  $D = \{(s_i, y_i)\}_{i=1}^n$ , is represented as a  
linear function of the form [8]:

$$f(s) = w^T s + b$$

52 where  $w$  and  $b$  are respectively the weight vector and the intercept of the model,  
53 and they are selected to find an optimal fit to the data available in  $D$ .

54 For nonlinear cases, one proceeds by mapping the input  $p$ -dimensional vec-  
55 tors via a nonlinear function  $\phi : R^p \rightarrow F$ , onto the feature space  $F$ . After  
56 nonlinear mapping, the regression function evolves to a pervasive form:

$$f(s) = w^T \phi(s) + b$$

57 SVR uses the  $\epsilon - insensitive$  loss function:

$$l = |y - f(s)|_{\epsilon} = \begin{cases} 0 & |y - f(s)| \leq \epsilon \\ |y - f(s)| - \epsilon & \text{else} \end{cases}$$

58 which ignores the error if the difference between the prediction value and  
 59 the actual value is smaller than  $\epsilon$ . The  $\epsilon$  - *insensitive loss function* allows to  
 60 find the coefficients  $w$  and  $b$  by solving a convex optimization problem, which  
 61 balances the empirical error and the generalization ability. In SVR, the empirical  
 62 error is measured by the loss function -insensitive and the generalization ability  
 63 is measured by the Euclidean norm of  $w$  [9]. Then, the optimization problem  
 64 to identify the regression model can be formulated by [8]:

$$\begin{aligned} \text{miimize} \quad & J(w, \xi_i, \xi_i^*) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i, \xi_i^*) \\ & y_i - w^T \phi(s) - b \leq \epsilon + \xi_i \\ \text{subject to} \quad & w^T \phi(s) + b - y_i \leq \epsilon + \xi_i^* \quad i = 1, 2, \dots, n \\ & \xi_i, \xi_i^* \geq 0 \end{aligned} \tag{1}$$

65 where  $C$  denotes the penalty parameter between empirical and generalization  
 66 errors, and  $\xi_i, \xi_i^*$  are slack variables. Figure.2 shows this situation.

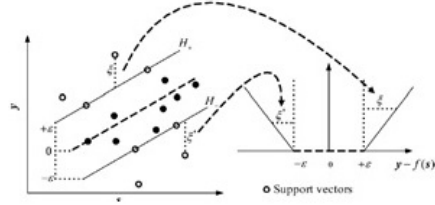


Figure 2:  $\epsilon$  - *insensitive loss function*[8].

The solution of this optimization problem by the Lagrange method is given by:

$$f(s) = w^T \phi(s) + b = \sum_{i=1}^n (\alpha_i - \alpha_i^*) K(s, s_i) + b$$

where  $\alpha_i - \alpha_i^*$  are the Lagrange multipliers of the optimization problems dual form and  $K(s_i, s_j)$  is the kernel function satisfying the Mercer condition, and holds:

$$K(s_i, s_j) = \langle \phi(s_i), \phi(s_j) \rangle$$

Operations in the kernel function  $K(s, s_i)$  are performed in the input space rather than in the potentially high dimensional feature space of  $\phi$  [10].

### 2.1.3. Ordinary least squares regression

This method fits a linear model with coefficients  $w = (w_1, \dots, w_p)$  to minimize the residual sum of squares between the observed responses in the dataset, and the responses predicted by the linear approximation. Mathematically it solves a problem of the form [? ]:

$$\min_w \|Xw - y\|_2^2$$

where  $X$  denotes the features matrix.

According Pedregosa et al. [11] the coefficient estimates for Ordinary Least Squares rely on the independence of the model terms. When terms are correlated and the columns of the design matrix  $X$  have an approximate linear dependence, the design matrix becomes close to singular and as a result, the least-squares estimate becomes highly sensitive to random errors in the observed response, producing a large variance. This situation of multicollinearity can arise, for example, when data are collected without an experimental design

### 2.1.4. Ridge regression

The ridge regression addresses some of the problems of ordinary least squares regression by imposing a penalty on the size of the coefficients. The ridge coefficients minimize a penalized residual sum of squares [11]:

$$\min_w \|Xw - y\|_2^2 + \alpha \|w\|_2^2$$

Here,  $\alpha > 0$  is a complexity parameter that controls the amount of shrinkage: the larger the value of  $\alpha$ , the greater the amount of shrinkage and thus the coefficients become more robust to collinearity.

86 2.1.5. Echo State Networks (ESN)

87 Recurrent Neural Networks (RNN) are useful for temporal patterns, but  
 88 when they are trained with backpropagation methods, they are very slow. Echo  
 89 State Network (ESN) is an alternative training method to solve that problem.  
 90 ESN is based on the observation that if a random RNN possesses certain al-  
 91 gebraic properties, training only a linear readout from it is often sufficient to  
 92 achieve excellent performance in practical applications [12]. For a given train-  
 93 ing input signal  $u(n) \in R^{N_u}$  a desired target output signal  $y^{target}(n) \in R^{N_y}$  is  
 94 known. Here  $n = 1, \dots, T$  is the discrete time and  $T$  is the number of data points  
 95 in the training dataset. The task is to learn a model with output  $y(n) \in R^{N_y}$ ,  
 96 where  $y(n)$  matches  $y^{target}(n)$  as well as possible, minimizing an error measure  
 97  $E(y, y^{target})$ , and, more importantly, generalizes well to unseen data. The un-  
 98 trained RNN part of an ESN is called a dynamical reservoir, and the resulting  
 99 states  $x(n)$  are termed echoes of its input history [13]. Finally, these signals are  
 sent to an output layer as shown in the Figure.3.

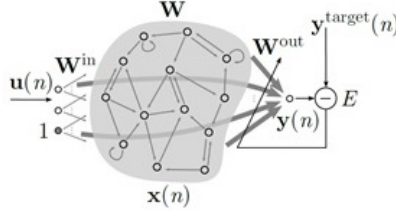


Figure 3: An echo state network [13].

100

101 The connections between the different elements of an Echo State Network  
 102 have weights randomly generated. The weights of the internal connections of  
 103 the reservoir ( $W$ ) as well as the weights of the input layer ( $W_{in}$ ), after being  
 104 generated are set statically during all stages of implementation of the algorithm.  
 105 The weights between the reservoir and the output layer ( $W_{out}$ ) are subject  
 106 to changes of a supervised learning algorithm to correct the degree of error  
 107 generated by the entire system [13].

## 108 2.2. Related works

109 Huang et al. [17] surveyed the development of soft computing techniques  
110 in agricultural and biological engineering, including fuzzy logic, artificial neural  
111 networks, genetic algorithms, bayesian inference and decision trees.

112 A related work, proposed by Romero [14] relies on regression models using  
113 a stepwise procedure to predict incubation and latency times of black Sigatoka.  
114 The author performed experiments on two farms located in Costa Rica (Rita  
115 and Waldeck, the same as those used in this study but with different names).  
116 The study used data from: December 1993 to August 1995. Romero concluded  
117 that the model to predict the incubation period accounted a  $R^2$  of 69% in his  
118 observed data but it was not a good predictor when it was validated against an  
119 independent dataset (cross validation). For latency, he developed two models  
120 that accounted a  $R^2$  of 78% PONER EL VALOR OBTENIDO in the observed  
121 data, however, when validated against an independent dataset (cross validation),  
122 the model was incorrect PONER EL VALOR OBTENIDO for Weldeck, and for  
123 Rita obtained an adjusted  $R^2$  of 82%.

124 Glezakos et al. [15] proposed to use Genetic Algorithms (GA) and Neu-  
125 ral Networks (NN) to identify plant virus (Tobacco Rattle Virus (TRV) and  
126 the Cucumber Green Mottle Mosaic Virus (CGMMV)). This is achieved by  
127 the development of analytical tools of evolutionary adaptive width, propelled  
128 by Genetic Algorithms (GAs) and Neural Networks (NNs). The method was  
129 tested against some of the most commonly used classifiers in machine learning  
130 (Bayes, Trees and k-NN) via cross-validation and proved its potential towards  
131 the identification.

132 In the agricultural context, Alves et al. [16] used geoinformation techniques  
133 to develop predictive models to study the areas of risk to soybean rust in soy-  
134 bean, coffee leaf rust in coffee, and black Sigatoka in banana, considering Brazils  
135 climatic characterization and the distribution of soybean, coffee and banana  
136 crops. Temperature and rainfall data were obtained for the period from 1950  
137 to 2000, and of simulations for 2020, 2050 and 2080 using the SRES A2 cli-  
138 mate change scenarios. Using principal components analysis, a single variable

139 was generated based on 57 variables, in order to determine an index explain-  
140 ing 87%, 88% and 90% of the variability of soybean, coffee and banana crops,  
141 respectively, in municipal districts across Brazil. The climatic model was used  
142 to generate the zoning of the three plant diseases, using temperature and leaf  
143 wetness as input. Areas of favorability for the diseases were plotted against the  
144 main coffee, soybean and banana growing areas in Brazil. This methodology  
145 enabled the visualization of the changes in areas favorable for epidemics under  
146 possible future scenarios of climate change.

147 Other applications of machine learning methods in precision agriculture in-  
148 clude the use of support vector regression to predict carcass weight in beef cattle  
149 in advance to the slaughter [10], machine learning assessments of soil drying for  
150 agricultural planning [18], and early detection and classification of plant diseases  
151 with support vector machines based on hyperspectral reflectance [19].

152 Furthermore, there have been attempts to generate software tools. Camargo  
153 et al. [Camargo et al.,2012] presented an information system for the assessment  
154 of plant disorders (Isacrodi). They proposed that experts will attain a much  
155 better accuracy than the Isacrodi classifier, particularly when provided with  
156 samples from the affected crop. However, those cases where such expertise is  
157 not available, they suggest that Isacrodi can provide valuable support to farmers.  
158 Isacrodi includes 15 crop disorders, but the black Sigatoka no is one of them.  
159 The prediction process is based on multi-class support vector machines.

160 Regarding the prediction of the development of the black Sigatoka with ma-  
161 chine learning methods, Bendini et al. [20] presented a study about the risk  
162 analysis of black Sigatoka occurrence based on polynomial models. A case study  
163 was developed in a commercial banana plantation located in Jacupiranga, Brazil.  
164 It was monitored weekly during the period from February to December 2005.  
165 Data included the weekly monitoring of the diseases evolution stage, time series  
166 of meteorological data and remote sensing data. They obtained a model to esti-  
167 mate the evolution of the disease from satellite imagery. This model relates gray  
168 levels (NC) of the band 2 images of the Landsat-5 satellite, with the progress  
169 status or disease severity (EE). The authors claim to reach an  $R^2$  of 90%.



Also there are works related to banana fruit. Soares et al. [21] apply two techniques: artificial neural networks (ANNs) and multiple linear regression (MLR) in banana plant to predict the yield, their results show that the neural network proved to be more accurate in forecasting the weight of the bunch in comparison to the multiple linear regressions in terms of the mean prediction-error ( $MPE = 1.40$ ), mean square deviation ( $MSD = 2.29$ ) and coefficient of determination ( $R^2 = 91\%$ ).

In general, the machine learning methods applied to predict the evolution of plant diseases, can be classified in two main approaches: 1) Those whose main inputs are images, and 2) Those whose main inputs are environmental and biological variables. Our study focuses in the second case.

### 2.3. Data

In this work we used data acquired in two research farms of Corbana in Costa Rica: 1) 28 Millas (located at Matina) and La Rita (located at Pococ), both in the province of Limn, Costa Rica. The banana type is Musa AAA, subgroup Cavendish, cv. Grande Naine. Table 1 shows the variables considered initially. Table 1 Variables used in the study Variable Meaning

The value to be predicted in all cases was ES, that is the total measure of the biological warning system. The data on the biological warning system are collected once a week. Although Corbana has meteorological stations that take data every five minutes, for these experiments, weekly averages generated by nearby stations to each of the farms were used. The time intervals used for this study were: La Rita, week 48 of 2002 to week 17 of the 2015 (647 weeks) and for 28 Miles, week 37 of 2003 to week 18 of 2015 (605 weeks). Data preprocessing In 28 Miles farm we detect that 1% of the data were missing, while in La Rita 2.25% of the data were missing. To complete the missing value we use spline interpolation. The data collected did not exhibit outliers. Due the fact that the variables measure meteorological or biological process, they are discretized in order to reflect trends in the data rather than the specific continuous values. The coefficient of variation  $C_v(x)$  of each variable x was used to determine the

numbers of discretizations with:

$$n = 100C_v(x)$$

Each discretization range was uniformly partitioned. Besides enabling the capture of tendencies, the discretization removes the effect of small variations in the data collection, either by inaccuracies of the instruments (meteorological variables) or by subjective bias introduced by the human who collects the data (biological warning system). Each feature was scaled in a range 0 between 1. The variable to be predicted was not scaled.

#### 2.4. Evaluation criteria

Although there are many types of indicators to assess the quality of the prediction, we selected the root mean square error (RMSE) and the determination coefficient (R2). This decision is supported by the widespread use in machine learning and agriculture areas (Soares, Pasqual and Lacerda (2013); Soares, Pasqual and Lacerda (2014); Ibrahim and Wibowo (2014) and Demir and Bruzzzone (2014)).

#### 2.5. Methods

This research had two phases.

##### **Phase one**

In the phase one, we did ten-fold-cross-validation and did a lot of proofs with different machine learning methods and different configuration. We proved with several combinations: Patterns: From one week of observed data to predict the next week until nine weeks before to predict two weeks later. Algorithms: Support vector regression with different kernel functions: linear, RBF (Gaussian) and sigmoid; echo state networks; ordinary least squares linear regression and ridge regression. Variables included in the model. We proved the following combinations: All variables. Only variables that according to expert judgment have more impact on the black Sigatoka development: humidity, precipitation,

213 temperature and wind speed (Marin Vargas and Romero Caldern, 1995). From  
214 the four variables listed in the previous paragraph, runs were conducted using  
215 each of the variables separately, and combining other runs all the possible pairs  
216 of those four variables.

## 217 **Phase two**

218 In the second phase, we used the best configurations obtained en la phase one  
219 and did validation with the last 52 and 102 weeks. This second phase pretended  
220 to show how these methods behaved on a time of important climate change how  
221 are 2014 and 2015 years.

222       Programming environment We use python programming language with the  
223 Integrated Development Environment (IDE) Spyder, particularly with libraries:  
224 pandas (Comunity, 2014); numpy (numpy.org, 2013); for SVR, ridge and ordi-  
225 nary least squares, we used sklearn (Pedregosa, et al., 2011); and for ESN the  
226 python-based code used belongs to Dr. Mantas Lukoeviius (2012) from which  
227 we made the necessary adjustments for the experiments of this research. The  
228 computer was a Lenovo ThinkPad, processor Intel(R) Core i7-4800MQ CPU @  
229 2.70GHz, 16.0 GB RAM, running Windows 8 Pro.

## 230 **3. Results**

## 231 **4. Discussion and conclusions**

## 232 **5. References**

- 233 [1] Brescani, XXXXX.
- 234 [Camargo et al.,2012] Camargo, A., Molina, J., Cadena-Torres, J.,  
235 Jiménez, N., Kim, J. 2012. Intelligent systems for the assessment of  
236 crop disorders. Computers and Electronics in Agriculture(85), 1-7.  
237 doi:10.1016/j.compag.2012.02.017.
- 238 [3] Chuang, T., Jeger, M. 1987. Predicting the Rate of Development of Black  
239 Sigatoka ( Mycosphaerella fijiensis var. difformis ) Disease in Southern Tai-  
240 wan. Phytopathology, 77, 1542-1547.

- [17] Huang, Y., Lan, Y., Thomson, S., Fang, A., Hoffmann, W., Lacey, R. 2010. Development of soft computing and applications in agricultural and biological engineering. *Computers and Electronics in Agriculture*,(71(2)), 107127. doi:10.1016/j.compag.
- [5] Kim, Y., Yoo, S., Gu, Y., Lim, J., Han, D., Baik, S. 2014. Crop Pests Prediction Method Using Regression and Machine Learning Technology: Survey. *IERI Procedia*(6), 5256. doi:10.1016/j.ieri.2014.03.009.
- [6] Marin Vargas, D., Romero Caldern, R. 1995. El combate de la Sigatoka Negra. *Boletín Departamento de Investigaciones, Corbana Costa Rica*.
- [7] Zhao, L., He, L., Harry, W., Jin, X. 2013. Intelligent Agricultural Forecasting System Based on Wireless Sensor. *Journal of Networks*(8), 18171824. doi:10.4304/jnw.8.8.1817-1824.
- [8] Wei, Z., Tao, T., ZhuoShu, D., Zio, E. (2013). A dynamic particle filter-support vector regression method for reliability prediction. *Reliability Engineering & System Safety*, 109116. doi:10.1016/j.res.2013.05.021.
- [9] Libro SVM XXXX.
- [10] Alonso, J., Rodriguez Castan, ., Bahamonde, A. (2013). Support Vector Regression to predict carcass weight in beef cattle in advance of the slaughter. *Computers and Electronics in Agriculture*, 116-120.
- [11] Pedregosa, F. and Varoquaux, G. and Gramfort, A. and Michel, V. and Thirion, B. and Grisel, O. and Blondel, M. and Prettenhofer, P. and Weiss, R. and Dubourg, V. and Vanderplas, J. and Passos, A. and Cournapeau, D. and Brucher, M. and Perrot, M. and Duchesnay, E. (2011) Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, vol. 12, 2825–2830.
- [12] Lukosevicius, M. and Jaeger, H. (2009). Reservoir computing approaches to recurrent neural network training. *Computer Science Review*(3), 127149. doi:10.1016/j.cosrev.2009.03.005.

- 269 [13] Lukosevicius, M. (2012). A Practical Guide to Applying Echo State Net-  
270 works. *Neural Networks: Tricks of the Trade*. 1-20
- 271 [14] Romero, R. (1995). Dynamics of fungicide resistant populations of *My-*  
272 *cosphaerella fijiensis* and Epidemiology of black Sigatoka of banana. Costa  
273 Rica: North Carolina State University.
- 274 [15] Glezakos, T., Moschopoulou, G., Tsiligrdis, T., Kintzios, S., Yialouris, C.  
275 (2010). Plant virus identification based on neural networks with evolutionary  
276 preprocessing. *Computers and Electronics in Agriculture*, 70, 263-275.
- 277 [16] Alves, M., de Carvalho, L., Pozza, E., Sanches, L., Maia, J. (2011). Eco-  
278 logical zoning of soybean rust, coffee rust and banana black sigatoka based  
279 on Brazilian climate changes. *Procedia Environmental Sciences*, 6, 35-49.
- 280 [17] Huang, Y., Lan, Y., Thomson, S., Fang, A., Hoffmann, W., Lacey, R.  
281 (2010). Development of soft computing and applications in agricultural and  
282 biological engineering. *Computers and Electronics in Agriculture*, (71(2)),  
283 1071-127. doi:10.1016/j.compag.2010.06.009
- 284 [18] Coopersmith, E. J., Minsker, B. S., Wenzel, C. E., Gilmore, B.  
285 J. (2014). Machine learning assessments of soil drying for agricul-  
286 tural planning. *Computers and Electronics in Agriculture*, 104, 93-104.  
287 <http://doi.org/10.1016/j.compag.2014.04.004>
- 288 [19] Rumpf, T., Mahlein, a.-K., Steiner, U., Oerke, E.-C., Dehne, H.-  
289 W., Plmer, L. (2010). Early detection and classification of plant  
290 diseases with Support Vector Machines based on hyperspectral re-  
291 flectance. *Computers and Electronics in Agriculture*, 74(1), 91-99.  
292 <http://doi.org/10.1016/j.compag.2010.06.009>
- 293 [20] Bendini, H., Moraes, W., da Silva, S., Tezuka, E., Cruvinel, P. (2013). An-  
294alise de risco da ocorrência de Sigatoka-negra baseada em modelos polinomiais:  
295 um estudo de caso. *Tropical Plant Pathology*, 38, 035-043.

296 [21] Soares, J., Pasqual, M., Lacerda, W., Silva, S., Donato, S. (2014). Compar-  
297 ison of techniques used in the prediction of yield in banana plants. *Scientia*  
298 *Horticulturae* journal, 167, 84-90.