

# 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].

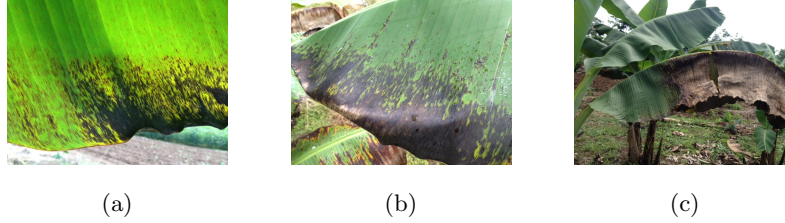


Figure 1: Examples of three disease stages of the black Sigatoka. (a) Initial stage. (b) Intermediate stage, and (c) Advanced stage.

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, [4] 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 The 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 (Wei, Tao, ZhuoShu, and Zio, 2013):  $f(s) = w^T s + b$   
where  $w$  and  $b$  are respectively the weight vector and the intercept of the model,  
and they are selected to find an optimal fit to the data available in  $D$ . For  
nonlinear cases, one proceeds by mapping the input p-dimensional vectors via  
a nonlinear function  $R^p F$ , onto the feature space  $F$ . After nonlinear mapping,  
the regression function evolves to a pervasive form:

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

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

53 which ignores the error if the difference between the prediction value and  
54 the actual value is smaller than  $\epsilon$ .  $\epsilon$ -insensitive loss function allows to find the

coefficients  $w$  and  $b$  by solving a convex optimization problem, which balances the empirical error and the generalization ability. In SVR, the empirical error is measured by the loss function  $\epsilon$ -insensitive and the generalization ability is measured by the Euclidean norm of  $w$ . Then, the optimization problem to identify the regression model can be formulated by (Wei, Tao, ZhuoShu, and Zio, 2013): where  $C$  denotes the penalty parameter between empirical and generalization errors, and  $\xi_i, \xi_i^*$  are slack variables, as shown in Fig 2.

Fig 2:  $\epsilon$ -insensitive loss function (Wei, Tao, ZhuoShu, and Zio, 2013) The solution of this optimization problem by the Lagrange method is: 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 Mercer condition, and can be described by: Common kernel functions are: linear, polynomial and sigmoid. 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  $\Phi$ . An inner product in the feature space has an equivalent kernel in the input space (Alonso, Rodriguez Castan, and Bahamonde, 2013).

### 2.1.3. Ordinary least square

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 (scikit-learn developer, 2014):

### 2.1.4. Ridge regression

This addresses some of the problems of Ordinary Least Squares by imposing a penalty on the size of coefficients. The ridge coefficients minimize a penalized residual sum of squares (scikit-learn developer, 2014): Here,  $\lambda$  is a complexity parameter that controls the amount of shrinkage: the larger the value of  $\lambda$ , the greater the amount of shrinkage and thus the coefficients become more robust to collinearity.

### 83 2.1.5. Echo State Networks (ESN)

84 Recurrent Neural Networks (RNN) are useful for temporal patterns, but  
85 when they are trained with backpropagation method, they are very slow. Echo  
86 State Network (ESN) is an alternative training method to solve that problem.  
87 ESN is based on the observation that if a random RNN possesses certain al-  
88 gebraic properties, training only a linear readout from it is often sufficient to  
89 achieve excellent performance in practical applications (Lukoeviius and Jaeger,  
90 2009). For a given training input signal  $u(n)R^{(N_u)}$  a desired target output  
91 signal is known. Here  $n = 1, \dots, T$  is the discrete time and  $T$  is the number of  
92 data points in the training dataset. The task is to learn a model with output  
93  $y(n)R^{(N_y)}$ , where  $y(n)$  matches  $y^{target}(n)$  ( $n$ ) as well as possible, minimizing  
94 an error measure  $E(y, y^{target})$ , and, more importantly, generalizes well to un-  
95 seen data. The untrained RNN part of an ESN is called a dynamical reservoir,  
96 and the resulting states  $x(n)$  are termed echoes of its input history (Lukoeviius  
97 M. , 2012). Finally, these signals are sent to an output layer as shown in the  
98 following Fig.

99 Figure 3: An echo state network (Lukosevicius, 2012) The connections be-  
100 tween the different elements of an Echo State Network have weights randomly  
101 generated. The weights of the internal connections of the reservoir ( $W$ ) as well  
102 as the weights of the input layer ( $W_{in}$ ), after being generated are set statically  
103 during all stages of implementation of the algorithm. The weights between the  
104 reservoir and the output layer ( $W_{out}$ ) are subject to changes of a supervised  
105 learning algorithm to correct the degree of error generated by the entire system  
106 (Lukoeviius M. , 2012). Related works A related work, no machine learning  
107 approach, was performed by Romero (1995) who in the third chapter of his  
108 doctoral thesis in the field of plant pathology, proposed regression models using  
109 stepwise procedure to predict incubation and latency times of black Sigatoka.  
110 The author performed experiments on two farms located in Costa Rica (Rita  
111 and Waldeck, the same as those used in this study but with different names).  
112 The time intervals used for that study were: December 1993 thru August 1995.

113 Romero concluded that the model to predict the incubation period accounted  
 114 a R2 of 69% in his observed data but it was not a good predictor when it was  
 115 validated against an independent dataset (cross validation). For latency, he de-  
 116 veloped two models that accounted a R2 of 78% in the observed data, however,  
 117 when validated against an independent dataset (cross validation), the model  
 118 was incorrect for Weldeck, and for Rita obtained an adjusted R2 of 82%. A ma-  
 119 chine learning method was proposed by Glezakos, Moschopoulou, Tsiligiridis,  
 120 Kintzios, and Yialouris (2010), who presented a genetic algorithm as to smooth  
 121 out the initial information while, the so produced meta-data sets were used in  
 122 the training and testing of the applied neural network, producing fitter training  
 123 data. Given the features of the acquired virus time-series signals of the problem  
 124 under study, an evolutionary method was proposed in order to produce meta-  
 125 data from the original time-series initial information, reduce the dimensionality  
 126 of the input data space, and eliminating the noise inherent in the initial raw  
 127 information The method was tested against some of the most commonly used  
 128 classifiers in machine learning (Bayes, Trees and k-NN) via cross-validation and  
 129 proved its potential towards assisting virus identification. They made their test  
 130 with CGMM and TR viruses. In agricultural area, Alves, de Carvalho, Pozza,  
 131 Sanches, and Mai (2011) selected the zones that are potentially favorable to  
 132 coffee, soybean and banana diseases in Brazil according to the spatial-temporal  
 133 variability of climatic variables and the geographical distribution of hosts. Their  
 134 study applied methodology enabled the visualization of the variation of areas  
 135 favorable to epidemics under future scenarios of climate change. The geosci-  
 136 entific and statistical modeling techniques developed in that study enabled the  
 137 development of predictive models and the characterization of risk areas for soy-  
 138 bean rust, coffee rust and black Sigatoka disease of banana. There have been  
 139 attempts to generate software tools, Camargo, Molina, Cadena-Torres, Jimnez,  
 140 and Kim (2012) presented an information system for the assessment of plant  
 141 disorders (Isacrodi). They proposed that experts will attain a much better  
 142 accuracy than the Isacrodi classifier, particularly when provided with samples  
 143 from the affected crop. However, where such expertise is not available, they

144 suggest that Isacrodi can provide valuable support to farmers. Isacordi includes  
 145 15 crop disorders, but the black Sigatoka no is one. The prediction process is  
 146 based on multi-class Support Vector Machines. Regarding black Sigatoka with  
 147 machine learning methods, Bendini, Moraes, da Silva, Tezuka, and Cruvinel  
 148 (2013) presented a study about the risk analysis of black Sigatoka occurrence  
 149 based on polynomial models. A case study was developed in a commercial ba-  
 150 nana plantation located in Jacupiranga, Brazil, it was monitored weekly during  
 151 the period from February to December 2005. Data were the weekly monitoring  
 152 of the diseases evolution stage, time series of meteorological data and remote  
 153 sensing data. They obtained a model to estimate the evolution of the disease  
 154 from satellite imagery. This model relates gray levels (NC) of the corresponding  
 155 image, band 2 of the Landsat-5 satellite, with the progress status or disease  
 156 severity (EE): Authors express have reach an R2 of 90Also there are research  
 157 related to banana fruit, Soares, Pasqual, Lacerda, Silva, and Donato (2014)  
 158 show in their study that to the analyses, the neural network proved to be more  
 159 accurate in forecasting the weight of the bunch in comparison to the multiple  
 160 linear regressions in terms of the mean prediction-error ( $MPE = 1.40$ ), mean  
 161 square deviation ( $MSD = 2.29$ ) and coefficient of determination ( $R^2 = 91$ In  
 162 general, machine learning methods applied to prediction plant diseases can be  
 163 classified in two main approaches: 1) Those that their main inputs are images,  
 164 and 2) Those that their main inputs are environmental and biological variables.  
 165 Our study is focus in the second one.

## 166 *2.2. Data*

167 In this work we used data acquired in two research farms of Corbana in  
 168 Costa Rica: 1) 28 Millas (located at Matina) and La Rita (located at Pococ),  
 169 both in the province of Limn, Costa Rica. The banana type is Musa AAA,  
 170 subgroup Cavendish, cv. Grande Naine. Table 1 shows the variables considered  
 171 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.

### 2.3. 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)).

### 2.4. Methods

This research had two phases.

#### Phase one



188

189       In the phase one, we did ten-fold-cross-validation and did a lot of proofs with  
190 different machine learning methods and different configuration. We proved with  
191 several combinations: Patterns: From one week of observed data to predict the  
192 next week until nine weeks before to predict two weeks later. Algorithms: Sup-  
193 port vector regression with different kernel functions: linear, RBF (Gaussian)  
194 and sigmoid; echo state networks; ordinary least squares linear regression and  
195 ridge regression. Variables included in the model. We proved the following  
196 combinations: All variables. Only variables that according to expert judgment  
197 have more impact on the black Sigatoka development: humidity, precipitation,  
198 temperature and wind speed (Marin Vargas and Romero Caldern, 1995). From  
199 the four variables listed in the previous paragraph, runs were conducted using  
200 each of the variables separately, and combining other runs all the possible pairs  
201 of those four variables.

## 202 **Phase two**

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

207       Programming environment We use python programming language with the  
208 Integrated Development Environment (IDE) Spyder, particularly with libraries:  
209 pandas (Comunity, 2014); numpy (numpy.org, 2013); for SVR, ridge and ordi-  
210 nary least squares, we used sklearn (Pedregosa, et al., 2011); and for ESN the  
211 python-based code used belongs to Dr. Mantas Lukoeviius (2012) from which  
212 we made the necessary adjustments for the experiments of this research. The  
213 computer was a Lenovo ThinkPad, processor Intel(R) Core i7-4800MQ CPU @  
214 2.70GHz, 16.0 GB RAM, running Windows 8 Pro.

215 **3. Results**

216 **4. Discussion and conclusions**

217 **5. References**

- 218 [1] Brescani, XXXXX.
- 219 [Camargo et al.,2012] Camargo, A., Molina, J., Cadena-Torres, J.,  
220 Jiménez, N., Kim, J. 2012. Intelligent systems for the assessment of  
221 crop disorders. Computers and Electronics in Agriculture(85), 1-7.  
222 doi:10.1016/j.compag.2012.02.017.
- 223 [3] Chuang, T., Jeger, M. 1987. Predicting the Rate of Development of Black  
224 Sigatoka ( *Mycosphaerella fijiensis* var. *difformis* ) Disease in Southern Tai-  
225 wan. Phytopathology, 77, 1542-1547.
- 226 [4] Huang, Y., Lan, Y., Thomson, S., Fang, A., Hoffmann, W., Lacey, R. 2010.  
227 Development of soft computing and applications in agricultural and biolog-  
228 ical engineering. Computers and Electronics in Agriculture,(71(2)), 107127.  
229 doi:10.1016/j.compag.
- 230 [5] Kim, Y., Yoo, S., Gu, Y., Lim, J., Han, D., Baik, S. 2014. Crop Pests Predic-  
231 tion Method Using Regression and Machine Learning Technology: Survey.  
232 IERI Procedia(6), 5256. doi:10.1016/j.ieri.2014.03.009.
- 233 [6] Marin Vargas, D., Romero Caldern, R. 1995. El combate de la Sigatoka  
234 Negra. Boletín Departamento de Investigaciones, Corbana Costa Rica.
- 235 [7] Zhao, L., He, L., Harry, W., Jin, X. 2013. Intelligent Agricultural Forecast-  
236 ing System Based on Wireless Sensor. Journal of Networks(8), 18171824.  
237 doi:10.4304/jnw.8.8.1817-1824.