

Forecasting the black Sigatoka development rate: A machine learning methods comparison

Luis-Alexander Calvo-Valverde^{a,1}, Mauricio Guzmán-Quesada^b, José-Antonio Guzmán-Alvarez^b, Pablo Alvarado-Moya^c

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

^b*Dirección de Investigaciones, Corporación Bananera Nacional S.A., Guápiles, Costa Rica*

^c*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 plants as the ability of photosynthesis decreases, causes a reduction in the quality of the fruit, and promotes premature maturation of bunches, which is the major cause of product losses due to this disease. 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 the development of this disease [6].

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

¹Corresponding author. (506)70104420

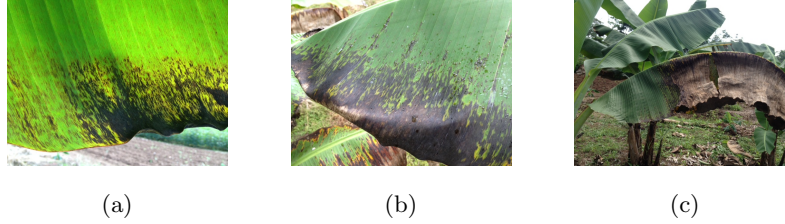


Figure 1: Examples of three disease stages of the black Sigatoka. (a) Initial stage. (b) Intermediate stage, and (c) Advanced stage.

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

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

24 There are efforts to apply machine learning methods for decision-making in
 25 agriculture, including the control of crop diseases. For example, [Camargo et al.,2012]
 26 present an intelligent systems for the assessment of crop disorders, [4] introduce
 27 a plant virus identification method based on neural networks with an evolu-
 28 tionary preprocessing stage, [5] summarize in their survey crop pests prediction
 29 methods using regression and machine learning approaches, while [7] present an
 30 intelligent agricultural forecasting system based on wireless sensor networks.

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

35 The main contribution of this work is a comparison between machine learning
 36 methods to forecast black Sigatoka development rate.

37 2. Materials and methods

38 2.1. Concepts

39 2.1.1. Biological warning system

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

50 2.1.2. Support Vector Regression (SVR)

51 From the perspective of Support Vector Regression (SVR) the regression
52 function $y = f(s)$ for a given dataset $D = (s_i, y_i)_{i=1}^n$, is represented as a
53 linear function of the form (Wei, Tao, ZhuoShu, and Zio, 2013): $f(s) = w^T s + b$
54 where w and b are respectively the weight vector and the intercept of the model,
55 and they are selected to find an optimal fit to the data available in D . For non-
56 linear cases, one proceeds by mapping the input p -dimensional vectors via a
57 nonlinear function $R^p F$, onto the feature space F . After nonlinear mapping, the
58 regression function evolves to a pervasive form: $f(s) = w^T(s) + b$ SVR uses
59 the ϵ -insensitive loss function: $l = |y - f(s)| = (0, |y - f(s)| @ |y - f(s)|_-, else)$
60 which ignores the error if the difference between the prediction value and the
61 actual value is smaller than ϵ . *ϵ -insensitivelossfunction* allows to find the
62 coefficients w and b by solving a convex optimization problem, which balances
63 the empirical error and the generalization ability. In SVR, the empirical error is
64 measured by the loss function ϵ -insensitive and the generalization ability is mea-
65 sured by the Euclidean norm of w . Then, the optimization problem to identify

66 the regression model can be formulated by (Wei, Tao, ZhuoShu, and Zio, 2013):

67 $minimize J(w, i, i^*) = 1/2 ||w||^2 + C(i = 1)^n (i, i^*) subject to (y_i - w^T(s) - b + i @ w^T(s) + b - y_i + i^* i = 1, , n @ i, i^*)$

68 where C denotes the penalty parameter between empirical and generalization

69 errors, and i, i^* are slack variables, as shown in Fig 2.

Fig 2: -insensitive loss function (Wei, Tao, ZhuoShu, and Zio, 2013) The solution of this optimization problem by the Lagrange method is: $f(s) = w^T(s) + b =$

$(i = 1)^n (i - i^*) K(s, s_i) + b$ where i, 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: $K(s_i, s_j) = (s_i)(s_j)$ 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 f . An inner product

70 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): $(\min_w) Xw - y_2^2$

71 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): $(\min_w) Xw - y_2^2 + w_2^2$ Here, 0 is a complexity parameter that controls the amount of shrinkage :

the larger the value of, the greater the amount of shrinkage and thus the coefficients become more robust to collinearity

72 2.1.5. Echo State Networks (ESN)

Recurrent Neural Networks (RNN) are useful for temporal patterns, but

when they are trained with backpropagation method, they are very slow. Echo

State Network (ESN) is an alternative training method to solve that problem.

ESN is based on the observation that if a random RNN possesses certain algebraic

properties, training only a linear readout from it is often sufficient to

achieve excellent performance in practical applications (Lukoeviius and Jaeger,

2009). For a given training input signal $u(n)$ $R(N_u)$ a desired target output signal $y^{(target(n))}$ $R(N_y)$ is known. Here $1, \dots, T$ is the discrete time and T is the number of data points in the training dataset. The task is to learn a model without

73 Figure 3: An echo state network (Lukosevicius, 2012) The connections be-
74 tween the different elements of an Echo State Network have weights randomly
75 generated. The weights of the internal connections of the reservoir (W) as well
76 as the weights of the input layer (W_{in}), after being generated are set statically
77 during all stages of implementation of the algorithm. The weights between the
78 reservoir and the output layer (W_{out}) are subject to changes of a supervised
79 learning algorithm to correct the degree of error generated by the entire system
80 (Lukoievius M. , 2012). Related works A related work, no machine learning
81 approach, was performed by Romero (1995) who in the third chapter of his
82 doctoral thesis in the field of plant pathology, proposed regression models using
83 stepwise procedure to predict incubation and latency times of black Sigatoka.
84 The author performed experiments on two farms located in Costa Rica (Rita
85 and Waldeck, the same as those used in this study but with different names).
86 The time intervals used for that study were: December 1993 thru August 1995.
87 Romero concluded that the model to predict the incubation period accounted
88 a R^2 of 69% in his observed data but it was not a good predictor when it was
89 validated against an independent dataset (cross validation). For latency, he de-
90 veloped two models that accounted a R^2 of 78% in the observed data, however,
91 when validated against an independent dataset (cross validation), the model
92 was incorrect for Weldeck, and for Rita obtained an adjusted R^2 of 82%. A ma-
93 chine learning method was proposed by Glezakos, Moschopoulou, Tsiligiridis,
94 Kintzios, and Yialouris (2010), who presented a genetic algorithm as to smooth
95 out the initial information while, the so produced meta-data sets were used in
96 the training and testing of the applied neural network, producing fitter training
97 data. Given the features of the acquired virus time-series signals of the problem
98 under study, an evolutionary method was proposed in order to produce meta-
99 data from the original time-series initial information, reduce the dimensionality
100 of the input data space, and eliminating the noise inherent in the initial raw
101 information The method was tested against some of the most commonly used

classifiers in machine learning (Bayes, Trees and k-NN) via cross-validation and
proved its potential towards assisting virus identification. They made their test
with CGMM and TR viruses. In agricultural area, Alves, de Carvalho, Pozza,
Sanches, and Mai (2011) selected the zones that are potentially favorable to
coffee, soybean and banana diseases in Brazil according to the spatial-temporal
variability of climatic variables and the geographical distribution of hosts. Their
study applied methodology enabled the visualization of the variation of areas
favorable to epidemics under future scenarios of climate change. The geosci-
entific and statistical modeling techniques developed in that study enabled the
development of predictive models and the characterization of risk areas for soy-
bean rust, coffee rust and black Sigatoka disease of banana. There have been
attempts to generate software tools, Camargo, Molina, Cadena-Torres, Jimnez,
and Kim (2012) presented an information system for the assessment of plant
disorders (Isacrodi). They proposed that experts will attain a much better
accuracy than the Isacrodi classifier, particularly when provided with samples
from the affected crop. However, where such expertise is not available, they
suggest that Isacrodi can provide valuable support to farmers. Isacordi includes
15 crop disorders, but the black Sigatoka no is one. The prediction process is
based on multi-class Support Vector Machines. Regarding black Sigatoka with
machine learning methods, Bendini, Moraes, da Silva, Tezuka, and Cruvinel
(2013) presented a study about the risk analysis of black Sigatoka occurrence
based on polynomial models. A case study was developed in a commercial ba-
nana plantation located in Jacupiranga, Brazil, it was monitored weekly during
the period from February to December 2005. Data were the weekly monitoring
of the diseases evolution stage, time series of meteorological data and remote
sensing data. They obtained a model to estimate the evolution of the disease
from satellite imagery. This model relates gray levels (NC) of the corresponding
image, band 2 of the Landsat-5 satellite, with the progress status or disease
severity (EE): Authors express have reach an R2 of 90Also there are research
related to banana fruit, Soares, Pasqual, Lacerda, Silva, and Donato (2014)
show in their study that to the analyses, the neural network proved to be more

133 accurate in forecasting the weight of the bunch in comparison to the multiple
134 linear regressions in terms of the mean prediction-error ($MPE = 1.40$), mean
135 square deviation ($MSD = 2.29$) and coefficient of determination ($R^2 = 91$). In
136 general, machine learning methods applied to prediction plant diseases can be
137 classified in two main approaches: 1) Those that their main inputs are images,
138 and 2) Those that their main inputs are environmental and biological variables.
139 Our study is focus in the second one.

140 2.1.6. 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 $T_{(a_{max})}$ Max air temperature $T_{(a_{min})}$ Min air temperature
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)
 $n = 100$ $C_v(x)$ Each discretization range was uniformly partitioned. Besides enabling the capture of tendencies, the

141 2.1.7. Evaluation criteria

142 Although there are many types of indicators to assess the quality of the
143 prediction, we selected the root mean square error (RMSE) and the determination coefficient (R^2). This decision is supported by the widespread use in
144 machine learning and agriculture areas (Soares, Pasqual and Lacerda (2013);
145 Soares, Pasqual and Lacerda (2014); Ibrahim and Wibowo (2014) and Demir
146 and Bruzzone (2014)).

148 2.1.8. Methods

149 This research had two phases. **Phase one**

150 In the phase one, we did ten-fold-cross-validation and did a lot of proofs with
151 different machine learning methods and different configuration. We proved with
152 several combinations: Patterns: From one week of observed data to predict the
153 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, temperature and wind speed (Marin Vargas and Romero Caldern, 1995). From the four variables listed in the previous paragraph, runs were conducted using each of the variables separately, and combining other runs all the possible pairs of those four variables. Phase two In the second phase, we used the best configurations obtained en la phase one and did validation with the last 52 and 102 weeks. This second phase pretended to show how these methods behaved on a time of important climate change how are 2014 and 2015 years. Programming environment We use python programming language with the Integrated Development Environment (IDE) Spyder, particularly with libraries: pandas (Comunity, 2014); numpy (numpy.org, 2013); for SVR, ridge and ordinary least squares, we used sklearn (Pedregosa, et al., 2011); and for ESN the python-based code used belongs to Dr. Mantas Lukoeviius (2012) from which we made the necessary adjustments for the experiments of this research. The computer was a Lenovo ThinkPad, processor Intel(R) Core i7-4800MQ CPU @ 2.70GHz, 16.0 GB RAM, running Windows 8 Pro.

3. Results

4. Discussion and conclusions

5. References

- [1] Brescani, XXXXX.
- [Camargo et al.,2012] Camargo, A., Molina, J., Cadena-Torres, J., Jiménez, N., Kim, J. 2012. Intelligent systems for the assessment of crop disorders. Computers and Electronics in Agriculture(85), 1-7. doi:10.1016/j.compag.2012.02.017.

- 182 [3] Chuang, T., Jeger, M. 1987. Predicting the Rate of Development of Black
183 Sigatoka (*Mycosphaerella fijiensis* var. *difformis*) Disease in Southern Tai-
184 wan. *Phytopathology*, 77, 1542-1547.
- 185 [4] Huang, Y., Lan, Y., Thomson, S., Fang, A., Hoffmann, W., Lacey, R. 2010.
186 Development of soft computing and applications in agricultural and biolog-
187 ical engineering. *Computers and Electronics in Agriculture*, (71(2)), 107127.
188 doi:10.1016/j.compag.
- 189 [5] Kim, Y., Yoo, S., Gu, Y., Lim, J., Han, D., Baik, S. 2014. Crop Pests Predic-
190 tion Method Using Regression and Machine Learning Technology: Survey.
191 *IERI Procedia*(6), 5256. doi:10.1016/j.ieri.2014.03.009.
- 192 [6] Marin Vargas, D., Romero Caldern, R. 1995. El combate de la Sigatoka
193 Negra. *Boletín Departamento de Investigaciones, Corbana Costa Rica*.
- 194 [7] Zhao, L., He, L., Harry, W., Jin, X. 2013. Intelligent Agricultural Forecast-
195 ing System Based on Wireless Sensor. *Journal of Networks*(8), 18171824.
196 doi:10.4304/jnw.8.8.1817-1824.