

Network Analysis of Cropping Practices and Injury Profiles in Irrigated Rice Agroecosystems

Sith Jaisong

September 12, 2015

Abstract

Here is the abstract.....

Introduction

The use of in-field surveys is a useful tool to develop ground-truth databases that allow one identify actual constraints due to pests in an agricultural productions system. These sorts of databases provide an overview of the complex relationships between the crop, its management, pest injuries, yields. Understanding theses relationships may lead to better management, and guide researchers the new research hypotheses (Mew et al., 2004; Savary et al., 2006).

Several previous studies (Savary et al., 2000b,a, 2005; Dong et al., 2010; Reddy et al., 2011) involved surveys that have been used to identify relationships in an individual production situation (a set of factors that determine agricultural production) and the injury profiles (combination of disease and pest injuries that may occur in a given farmer's field) using nonparametric multivariate analysis such as cluster analysis, correspondence analysis, multiple correspondence analysis. Performing correspondence analysis (Savary et al., 1997), they characterized the relationships between categorized levels of variables: actual yield, production situations and injuries profiles. Their results led to the conclusions that observed injuries profiles were strongly associated with production situations, and the level of actual yields.

The components of production situation and injury profiles are biologically related. For example, the excess amount of fertilizers applied in the rice files increases the susceptibility of rice to blast, and directed seeded flooded rice fields with high seedling rate is favorable condition for sheath blight (Ou, 1985). The relationships will be more complex when the number of their components increased. A way to systemically model and intuitively interpret such relationships is the depiction as a graph or network. This approach has been widely used and proven very useful in biological studies (Moslonka-Lefebvre et al., 2011). Networks typically consist of nodes, usually representing components, while links between the nodes depict their interactions (Proulx et al., 2005). A correlation network is a type of network in which two nodes are connected if their respective correlation lies above a certain threshold. The construction of this network is obtained from pairwise correlation methods (Toubiana et al.,

2013). By using appropriate correlation measure, correlation networks can capture biologically meaningful relationships, and discover valuable information in crop health surveys.

The main objective of this study is to select a suitable association method for network construction. Selecting the suitable measure is important because the method should be able to capture the relationships with true concordance often determined the type and amount of knowledge we can gain from survey data, moreover it will affect the topological structure of network (the patterns of pairwise relationships between variables).

Materials and Methods

The limitation of each measure are difference assumption and detach different patterns. In this study, we evaluated four association methods, Pearson, Spearman rank correlation, Kendall, Biweight midcorrelation to be suitable for the properties of the survey data (*i.e.* type of variables, pattern of distribution, normality) based on their assumptions. The other evaluation is to examine which methods can associate those variables that are know to present the biological relationships.

Next, we inferred correlation network from surveys comprising five countries (India, Indonesia, Philippines, Thailand, and Vietnam), 420 lowland farmers' fields. We determine the correlation patterns among the incidence of injuries caused by animal pests and diseases and the cropping practices, potentially indicative of their occurrence relations. We then constructed the network from these pairwise correlations.

Survey datasets

Crop health survey data were collected through surveys comprising 420 farmers' fields from 2010 to 2012 for wet and dry seasons in different production environments across South and South East Asia. The survey protocol described in the IRRI publication, "A survey portfolio to characterize yield-reducing factors in rice", was used for data collection (Savary and Castilla, 2009). The variables collected included patterns of cropping practices, crop growth measurement and crop management status assessments, measurements of levels of injuries caused by pests, and direct measurements of actual yields from crop cuts. The data collected can be classified into three groups: cropping practices, injuries, and actual yield measurements.

Evaluation

One: Data exploration

There are three main properties to be determined before deciding the appropriate correlation measure for use in constructing the network.

Check data distribution This test can be achieved by significance test and visual methods. Each variable in survey dataset was tested normality using the Shapiro-Wilk test (Ghasemi and Zahediasl, 2012). The Shapiro-Wilk test is based on the correlation between the data and the corresponding normal scores.

H_0 : sample distribution is normal.

H_a : sample distribution is not normal.

Thus if the p -value is less than the chosen alpha level, the null hypothesis is rejected and there is evidence that the data tested are not from a normally distributed population. In other words, the data are not normal. On the contrary, if the p -value is greater than the chosen alpha level, then the null hypothesis that the data came from a normally distributed sample cannot be rejected. However, for small sample sizes, normality tests have little power to reject the null hypothesis, so a QQ (quantile–quantile plot) plot and the frequency distribution (histogram) are required for verification in addition to check normality visually.

The R function for Shapiro-Wilk Normality test is `shapiro.test` (package stats), which is (R Core Team, 2014).

Check for independence Performing the distance correlation t -test (Székely et al., 2007) to check independence aims to select the pair of variables, which are able to be detected correlation. The distance correlation of two random variables is obtained by dividing their distance covariance by the product of their distance standard deviations. The distance correlation is

$$\text{dCor}(X, Y) = \frac{\text{dCov}(X, Y)}{\sqrt{\text{dVar}(X) \text{dVar}(Y)}} \quad (1)$$

- $0 \leq \text{dCor}_n(X, Y) \leq 1$ and $0 \leq \text{dCor}(X, Y) \leq 1$
- $\text{dCor}(X, Y) = 0$ if and only if X and Y are independent.
- $\text{dCor}_n(X, Y) = 1$ implies that dimensions of the linear subspaces spanned by X and Y samples respectively are almost surely equal and if we assume that these subspaces are equal, then in this subspace $Y = A + b\mathbf{C}X$ for some vector A , scalar b , and orthonormal matrix \mathbf{C} .

This test was performed using the `fda.usp` package (Febrero-Bande and Oviedo de la Fuente, 2012).

Step two: identify the most appropriate method

Pearson's product-moment correlation coefficient

The Pearson's product-moment correlation or simply Pearson's correlation is a measure of linear dependence, as the slope obtained by the linear regression of Y by X is Pearson's correlation multiplied by that ratio of standard deviations. Let $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ and $\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$ be the means of X and Y , respectively, then the Pearson's correlation coefficient ρ_{Pearson} is defined as follows:

$$\rho_{\text{Pearson}}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (2)$$

For joint normal distributions, Pearson's correlation coefficient under H_0 follows a Student's t -distribution with $n - 2$ degrees of freedom. The t statistic is as follows:

$$t = \frac{\rho_{\text{Pearson}}(X, Y) \sqrt{n - 2}}{\sqrt{1 - \rho_{\text{Pearson}}^2(X, Y)}} \quad (3)$$

When the random variables are not jointly normally distributed, the Fisher's transformation is used to get an asymptotic normal distribution.

In the case of perfect linear dependence, we have $\rho_{Pearson} = \pm 1$. The Pearson correlation is +1 in the case of a perfect positive (increasing) linear relationship and -1 in the case of a perfect negative (decreasing) linear relationship. In the case of linearly independent random variables, $\rho_{Pearson} = 0$, and in the case of imperfect linear dependence, $-1 < \rho_{Pearson} < 1$. These last two cases are the ones for which misinterpretations of correlation are possible because it is usually assumed that non-correlated X and Y means independent variables, whereas in fact, they may be associated in a non-linear fashion that Pearson's correlation coefficient is not able to identify.

The R function for Pearson's test is `cor.test` with parameter method 'Pearson' (package stats). The stats package can be downloaded from the R web page (<http://www.r-project.org>).

The Spearman correlation coefficient The Spearman correlation coefficient is defined as the Pearson correlation coefficient between the ranked variables. For a sample of size n, the n raw scores X_i, Y_i are converted to ranks x_i, y_i , and ρ is computed from:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (4)$$

where $d_i = x_i - y_i$, is the difference between ranks. See the example below. Identical values (rank ties or value duplicates) are assigned a rank equal to the average of their positions in the ascending order of the values. In the table below, notice how the rank of values that are the same is the mean of what their ranks would otherwise be:

In applications where duplicate values are known to be absent, a simpler procedure can be used to calculate .

This method should not be used in cases where the data set is truncated; that is, when the Spearman correlation coefficient is desired for the top X records (whether by pre-change rank or post-change rank, or both), the user should use the Pearson correlation coefficient formula given above.

The standard error of the coefficient (σ) was determined as

$$\sigma = \frac{0.6325}{\sqrt{n-1}} \quad (5)$$

Kendall rank correlation : Kendall rank correlation is a non-parametric test that measures the strength of dependence between two variables. If we consider two samples, a and b, where each sample size is n, we know that the total number of pairings with a b is $n(n-1)/2$. The following formula is used to calculate the value of Kendall rank correlation:

The Kendall τ coefficient is defined as:

$$\tau = \frac{(N_c) - (N_d)}{\frac{1}{2}n(n-1)} \quad (6)$$

Where:

N_c = Number of concordant pairs

N_d = Number of discordant pairs

The Kendall τ coefficients range $-1 \leq \tau \leq 1$.

If the agreement between the two rankings is perfect (*i.e.*, the two rankings are the same) the coefficient has value 1. If the disagreement between the two rankings is perfect (*i.e.*, one ranking is the reverse of the other) the coefficient has value -1. If X and Y are independent, then we would expect the coefficient to be approximately zero.

Biweight midcorrelation (also called bicor) is a measure of similarity between samples. It is median-based, rather than mean-based, thus is less sensitive to outliers, and can be a robust alternative to other similarity metrics, such as Pearson correlation or mutual information. In order to define the biweight midcorrelation (?) of two vectors x and y , with $i = 1, 2, \dots, m$ items, representing each item in the vector as x_1, x_2, \dots, x_m and y_1, y_2, \dots, y_m . First, we define $\text{med}(x)$ as the median of a vector x and $\text{mad}(x)$ as the median absolute deviation (mad), then define u_i and v_i as,

$$\begin{aligned} u_i &= \frac{x_i - \text{med}(x)}{9 \text{mad}(x)}, \\ v_i &= \frac{y_i - \text{med}(y)}{9 \text{mad}(y)}. \end{aligned} \quad (7)$$

The weights $w_i^{(x)}$ and $w_i^{(y)}$ are defined as,

$$\begin{aligned} w_i^{(x)} &= (1 - u_i^2)^2 I(1 - |u_i|) \\ w_i^{(y)} &= (1 - v_i^2)^2 I(1 - |v_i|) \end{aligned} \quad (8)$$

where I is the identity function where,

$$I(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

Then we normalize so that the sum of the weights is 1:

$$\begin{aligned} \tilde{x}_i &= \frac{(x_i - \text{med}(x)) w_i^{(x)}}{\sum_{j=1}^m [(x_j - \text{med}(x)) w_j^{(x)}]^2} \\ \tilde{y}_i &= \frac{(y_i - \text{med}(y)) w_i^{(y)}}{\sum_{j=1}^m [(y_j - \text{med}(y)) w_j^{(y)}]^2}. \end{aligned} \quad (10)$$

Finally, biweight midcorrelation is defined as,

$$\text{bicor}(x, y) = \sum_{i=1}^m \tilde{x}_i \tilde{y}_i \quad (11)$$

The value of bicor ranges from -1 to 1, where -1 represents the maximum negative correlation and 1 represents the maximum positive correlation. 0 represents irrelevant.

Network Construction

Correlation network construction

The matrix can be viewed as an adjacency matrix of a weighted network. The matrix contains the correlation coefficient between each node (i.e., the variable). Thus the matrix can be thought of as the population average of the network structure. Because we are looking at several specific links, we control for multiple testing by controlling the False Discovery Rate (FDR method) at 5%. The generated network structure can be visualized through the R package qgraph (?). Only connections that surpass the significance threshold are shown in the visual representation.

Network analysis

Important information about a network can be gained by analyzing its global structure, for example by looking at the relative centrality of different nodes. In a centrality analysis, nodes are ordered in terms of the degree to which they occupy a central place in the network. Global descriptors of the modules were obtained using package qgraph in R. The neighborhood of a given node n is the set of its neighbors. The connectivity is the size of its neighborhood. The average number of neighbors indicates the average connectivity of a node in the network. A normalized version of this parameter is the network density. Density ranges between 0 and 1. It shows how densely the network is populated with edges. A network, which contains no edges and solely isolated nodes has a density of 0.

In correlation (undirected) networks, the clustering coefficient is the number of connected pairs between all neighbors of the network. The clustering coefficient of a node is always a number between 0 and 1. The network clustering coefficient is the average of the clustering coefficients for all nodes in the network. Nodes with less than two neighbors are assumed to have a clustering coefficient of 0. We then determined network centralities on the modules obtained from network analysis. Centralities were assessed using qgraph package in R. We calculated Degree centrality and Betweenness centrality.

Results

Discussion

References

- Dong, K., Chen, B., Li, Z., Dong, Y., and Wang, H. (2010). A characterization of rice pests and quantification of yield losses in the japonica rice zone of yunnan, china. *Crop Protection*, 29(6):603–611.
- Febrero-Bande, M. and Oviedo de la Fuente, M. (2012). Statistical computing in functional data analysis: The R package fda.usc. *Journal of Statistical Software*, 51(4):1–28.

- Ghasemi, A. and Zahediasl, S. (2012). Normality tests for statistical analysis: a guide for non-statisticians. *International journal of endocrinology and metabolism*, 10(2):486.
- Mew, T. W., Leung, H., Savary, S., Vera Cruz, C. M., and Leach, J. E. (2004). Looking ahead in rice disease research and management. *Critical Reviews in Plant Sciences*, 23(2):103–127.
- Moslonka-Lefebvre, M., Finley, A., Dorigatti, I., Dehnen-Schmutz, K., Harwood, T., Jeger, M. J., Xu, X., Holdenrieder, O., and Pautasso, M. (2011). Networks in plant epidemiology: from genes to landscapes, countries, and continents. *Phytopathology*, 101(4):392–403.
- Ou, S. H. (1985). *Rice diseases*. International Rice Research Institute (IRRI).
- Proulx, S., Promislow, D., and Phillips, P. (2005). Network thinking in ecology and evolution. *Trends in Ecology & Evolution*, 20(6):345–353.
- R Core Team (2014). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Reddy, C. S., Laha, G. S., Prasad, M. S., and Krishnaveni, D. (2011). Characterizing multiple linkages between individual diseases, crop health syndromes, germplasm deployment, and rice production situations in India. *Field Crops*.
- Savary, S. and Castilla, N. (2009). A survey portfolio to characterize yield-reducing factors in rice. *IRRI Discussion Paper No 18*, page 32.
- Savary, S., Castilla, N. P., Elazegui, F., and Teng, P. S. (2005). Multiple effects of two drivers of agricultural change, labour shortage and water scarcity, on rice pest profiles in tropical asia. *Field Crops Research*, 91(2):263–271.
- Savary, S., Elazegui, F., Pinnschmidt, H., Castilla, N., and Teng, P. (1997). A new approach to quantify crop losses due to rice pests in varying production situations. *IRRI discusion paper series no20. International Rice Research Institute, PO Box*, 933.
- Savary, S., Teng, P. S., Willocquet, L., and Nutter, F. W. (2006). Quantification and modeling of crop losses: a review of purposes. *Annual Review of Phytopathology*, 44(1):89–112.
- Savary, S., Willocquet, L., Elazegui, F. A., Castilla, N. P., and Teng, P. S. (2000a). Rice pest constraints in tropical Asia: quantification of yield losses due to rice pests in a range of production situations. *Plant disease*, 84(3):357–369.
- Savary, S., Willocquet, L., Elazegui, F. A., Teng, P. S., Van Du, P., Zhu, D., Tang, Q., Huang, S., Lin, X., and Singh, H. M. (2000b). Rice pest constraints in tropical Asia: characterization of injury profiles in relation to production situations. *Plant Disease*, 84(3):341–356.
- Székely, G. J., Rizzo, M. L., Bakirov, N. K., et al. (2007). Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35(6):2769–2794.

Toubiana, D., Fernie, A. R., Nikoloski, Z., and Fait, A. (2013). Network analysis: tackling complex data to study plant metabolism. *Trends in Biotechnology*, 31(1):29–36.