

STAT 427 CONSULTING PROJECT

TERRAREF PREDICTION METRICS

MEASURES OF PREDICTION AND RANKING ACCURACY FOR
THE TERRAREF PROJECT

Title

Author:
Kyle PAYNE

Manze QIN

February 25, 2016



Abstract

Your abstract.

1 Introduction

The following report consists of the Author's Recommendation for Measures of Prediction Quality for the TerraRef Project at the University of Illinois. In this preliminary report we will focus on two particular problems that have been addressed so far:

- How to determine the accuracy of a continuous prediction on a continuous target value (e.g. phenotype).
- How to Score Predicted Rankings for some subset of lines (e.g. genotypes).

While we only address these problems in a relatively closed sense, the measures that we propose may be applicable to other settings as well. We will define our measures, describe their respective numerical and statistical properties, make recommendations for using these measures in practice, and produce functions in both the Python and R languages for their implementation.

2 Measures

2.1 Continuous Phenotype Prediction

An example of the types of problems that fall under 'How to determine the accuracy of a continuous prediction on a continuous target value (e.g. phenotype)' would be to demonstrate that predicted values are within 20 percent of ground truth values. Thus, we need a measure that accounts for the difference between the predicted and ground truth (the true observed phenotype values), while also accounting for the *relative degree* by which the predictions are different from the ground truth values. One metric that appears in the literature (citation) is the Relative Mean Squared Prediction Error, or the RMSPE, which we define in equation (*)

Let Y_1, \dots, Y_n be a set of ground truth phenotype values, and let $\hat{Y}_1, \dots, \hat{Y}_n$ be the set of corresponding predictions for the ground truth phenotype values, then the $RMSPE$ is defined as.

$$RMSPE = \frac{\sqrt{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}}{\sqrt{\sum_{i=1}^n Y_i^2}} \quad (1)$$

The square of the numerator $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ is a well-studied function within the machine learning and statistics community, known as the Mean Squared Prediction Error, or (MSPE). This is an easily computed, and numerically stable quantity that provides several desirable large sample properties. The denominator of $RMSPE$ can be viewed as the difference of the continuous ground-truth phenotype values from 0. Thus, $RMSPE$ acts a relative measure of the difference between the ground truth continuous phenotype values, expressed in units of the continuous phenotype. Using this measure, an experimenter could make a statment such as, "I demonstrated that the predicted values are within 0.2 of the ground truth values". This type of measure could be helpful for determining prediction accuracy for both Terminal Biomass and 3D Structural Models. Moreover, the $RMSPE$ could also be used to compare the prediction accuracy between competing models, such as comparing algorithm predictions v.s. the Lemnatec Software.

Examples

If a prediction model produces a set of predictions $\hat{Y}_1, \dots, \hat{Y}_n$, for a set of ground truth continuous phenotype values Y_1, \dots, Y_n , then let's state that

$$RMSPE \leq 0.20 \Rightarrow \sqrt{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2} \leq 0.2 \sqrt{\sum_{i=1}^n Y_i^2} \quad (2)$$

The equation above describes the case where the $RMSPE$ being less than or equal to 20 percent is equivalent to stating that the mean squared prediction error is bounded by 0.2 the size of the ground truth phenotype values. The quantity on the left-hand side of the inequality is an example of a commonly used measure of distance in mathematics, engineering, statistics, computer science, etc. known as the *Norm* (Weisstein, Eric W. "Norm." From MathWorld—A Wolfram Web Resource. <http://mathworld.wolfram.com/Norm.html>).

Let $RMSPE_L$ be the relative mean squared prediction error of Lemnatec software predictions on the ground truth continuous phenotype values Y_1, \dots, Y_n , which we will denote as Y_1^*, \dots, Y_n^* . Let $RMSPE_M$ be the relative mean squared prediction error of some prediction model or algorithm. Then, making a determination such as 'algorithm predictions are no less accurate than values predicted via LemnaTec software', would require comparing $RMSPE_L$ and $RMSPE_M$, thus

$$RMSPE_M \leq RMSPE_L \quad (3)$$

Let's assume that the predictions are for the same set of ground truth continuous phenotype values Y_1, \dots, Y_n , then the preceding equation is equivalent to

$$RMSPE_M \leq RMSPE_L \Rightarrow \frac{\sqrt{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}}{\sqrt{\sum_{i=1}^n Y_i^2}} \leq \frac{\sqrt{\sum_{i=1}^n (Y_i - Y_i^*)^2}}{\sqrt{\sum_{i=1}^n Y_i^2}} \Rightarrow \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \leq \sum_{i=1}^n (Y_i - Y_i^*)^2 \quad (4)$$

Thus, if we compare two prediction models in $RMSPE$ on the same set of data, the comparison is equivalent to just comparing the average squared difference between the predictions and the ground truth values.

2.2 Performance

The Relative Mean Squared Prediction Error Performs well in situations in which there is additional noise in the continuous phenotype values Y_1, \dots, Y_n . For the following example, let's assume that some prediction algorithm has been fitted to a set of training data. In this case, we chose a Random Forest Model to predict Stem Biomass using the plot identifier and the precipitation data on a sub-sample of simulated data. We trained the random forest model on a sub-sample of simulated data. Predictions were then made using a test sample of the data, with additional mean 0 gaussian random noise applied to the Stem Biomass data. The RMSPE measure increases like a polynomial with increasingly variable noise. However, the RMSPE remains relatively robust to deviations from the true Prediction Error, and only increases to very large values as the

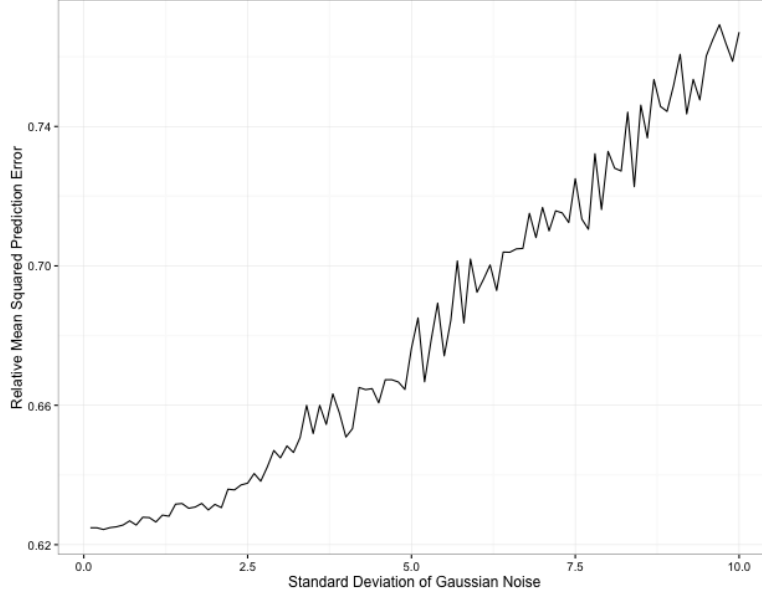


Figure 1: Robustness of RMSPE to Noise in Continuous Phenotype

3 Ranking Scores

Another problem of interest is assigning a score to rank predictions, e.g. given a set of lines of sorghum, and we score the predicted rank orders from some prediction algorithm? The literature is replete with measures of rank ordering. We will focus here two measures of rank ordering correlation, Kendall's Tau and the Normalized Discounted Cumulative Gain.

4 Kendall's Tau

Kendall's Tau is a measure of the correlation between two sets of rankings. Let Y_1, \dots, Y_n be a set of observations and let $\hat{Y}_1, \dots, \hat{Y}_n$ be a set of predicted observations. Then Any pair of observations (Y_i, \hat{Y}_i) and (Y_j, \hat{Y}_j) , where $i \neq j$, are said to be concordant if the ranks for both elements agree: that is, if both $Y_i > Y_j$ and $\hat{Y}_i > \hat{Y}_j$ or if both $Y_i < Y_j$ and $\hat{Y}_i < \hat{Y}_j$. They are said to be discordant, if $Y_i > Y_j$ and $\hat{Y}_i < \hat{Y}_j$ and $Y_i < Y_j$ and $\hat{Y}_i > \hat{Y}_j$. If $Y_i = Y_j$ and $\hat{Y}_i = \hat{Y}_j$, the pair is neither concordant nor discordant.

Kendall's Tau is defined as:

$$\tau = \frac{(\textit{number of concordant pairs}) - (\textit{number of discordant pairs})}{1/2n(n-1)} \quad (5)$$

4.1 Tables and Figures