
Random Forest Regression Model Feasibility for Plant Data Analysis

Vaaranan Yogalingam
School of Computer Science
University of Waterloo
Waterloo, ON, N2L 3G1
vyogalin@uwaterloo.ca
report due: August 13

Abstract

1 This report will discuss the results of experimenting with various models and
2 strategies in determining various plant traits given their images and ancillary data.
3 Ultimately, a random forest regression model led to the best results, seen in this
4 report. The code referenced in this report can be found here.

5 1 Introduction

6 In this report, several methods of predicting the six plant properties (X4,X11,X18,X50,X26,X3112)
7 were explored. Eventually, the highest performing strategy was to use a Random Forest Regressor
8 Model, courtesy of Sci-Kit Learn, trained over the plant trait data, to make the predictions. It was
9 later found that extracting feature vectors from the plant images, using a VGG11 network, courtesy
10 of Pytorch, which combined with the plant trait feature vectors, led to an similarly performing model
11 after a certain number of estimators were used. Such improvement in prediction (as measured by
12 the R^2 value) was minimal, with or without image data. This will be explored later on.

13 2 Related Works

14 Research from multiple works and papers supports the use of a Random Forest Regression Model,
15 for multiple outputs. In this paper on multi-output random forests, two methods of multi-output
16 regression are described. One method would be to use multiple single output regression algo-
17 rithms, and the other reform the algorithm for multiple outputs. The Random Forest Regressor
18 from SKLearn extends to multi-output problems by leveraging its use of decision trees and aver-
19 aging over them. There are real world examples of using Random Forests in regression problems,
20 such as in a paper by Hoffman et al. that used Random forest over numerical dust particle informa-
21 tion (i.e. size) to predict multiple numerical features of early planet formation. Another example
22 includes a paper by Jog et al., where Random Forests were used to analyze MRI images to address
23 issues of inconsistencies. Unlike these papers that used numerical data and images data separately,
24 this report will briefly explore combining both types of data in the Random Forest model, using a
25 VGG11 network to extract image features (a concept discussed in assignments/lectures).

26 3 Main Results

27 Two versions of Random Forest were used for experimentation. The first model was fitted solely on
28 the train.csv data provided. In other words no image data was used. In the second model, the image
29 data was incorporated. This was done using feature extraction through the VGG11 model. Below

30 is a table, displaying the score of each random forest model based on the number of estimators and
 31 whether or not the model used the image data. Note that the score is based on the coefficient of
 32 determination, or R^2 value calculated based on 55% of the data in the public tests.

Table 1: Model Comparison of Random Forest Regressor including and excluding Image Data Respectively (*Note: – represent the points that could not be recorded due to lack of time for training)

Number of Estimators	Random Forest without Image Data	Random Forest with Image Data
100	0.19098	0.18444
200	0.19529	0.18980
300	0.19556	0.19443
400	0.19648	–
500	0.19645	–
1000	0.19714	0.19714

33 This data can also be visualized in a graph like so, in figure one, where the convergence between the
 34 model using image data, and the model not using it, is much clearer:

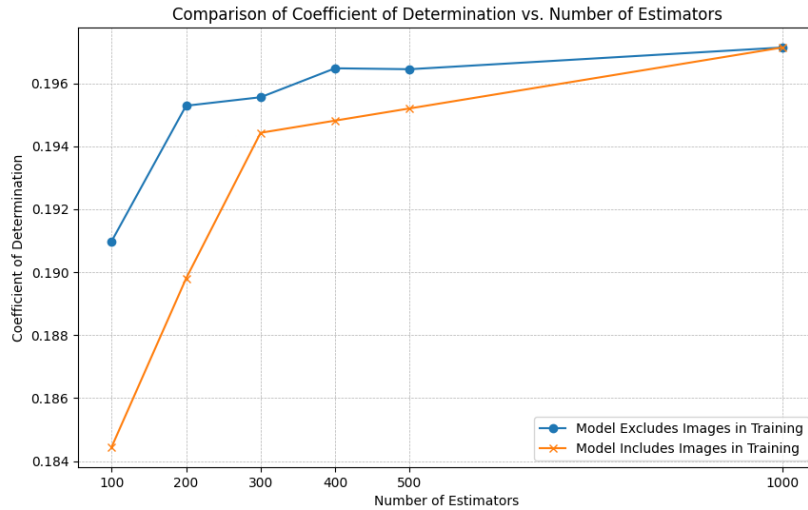


Figure 1: The progress of the coefficient of determination (R^2) as the number of estimators in a Random Forest Regressor model increases. Each line corresponds to when the model does and does not incorporate image data (incorporated via extraction using the VGG11 model). *Note: The model that incorporates the image data in prediction and training, has two interpolated points (at number of estimators being 400 and 500) due to lack of time to complete training.

35 As can be seen in these models, the greatest jump in the score is from 100 to 200 or 200 to 300
 36 estimators. In general, it would be expected that in any random forest model, that as the number of
 37 estimators increases, the measure of accuracy would increase. Increasing the number of estimators,
 38 increases the number of trees, which as a bagging algorithm, gives it more data to work off of, but
 39 prevents over-fitting, due to how simple these trees are. However, it is worth noticing that while both
 40 models eventually converge to approximately the same score (after using 1000 estimators each), the
 41 model that uses image data has a lower start, and has consistently large jumps in score between
 42 intervals of estimators. Furthermore, some spots are left blank due to a lack of time to run the
 43 models with the corresponding number of estimators (as running the model with 1000 estimators
 44 took the most time). Both these observations are important in concluding the value of including the
 45 image data in training and testing.

46 **4 Conclusion**

47 Given the little difference incorporating the image data into the Random Forest regression model
48 made, it is worth questioning when all data is too much data. While similar results were achieved
49 when compared to the model that did not use the image data, the time and energy costs of including
50 image data in the training and testing were probably not worth it. Also, it required a large number
51 of estimators before the model using image data could keep up with the model that did not need
52 that additional data. This could explain why the real world examples mentioned in the introduction
53 of this report, used either image or numerical data, but not both. It would be interesting to explore
54 other possible parameters or methods to make the image data more useful (perhaps through some
55 more complex transformations). Furthermore, perhaps a human analyses of the image would be
56 good to see if they have any relevance to the predicted values (X_4 , X_{11} , etc.). Ultimately, the
57 model that did not incorporate the image data was used to submit predictions on Kaggle (`./submis-`
58 `sions/20901584_Yogalingam_1000.jpeg` contains the file data submitted to Kaggle).

59 Acknowledgement

60 I would like to acknowledge the various sources of documentation I used and referred to when
61 developing my models, including those from Pytorch, Pandas, Pillow, and Scikit Learn, namely
62 Scikit Ensembles, where I was able to retrieve and use the RandomForestRegressor model.

63 References

- 64 Hoffman, K., J. Y. Sung, and A. Zazzera (2021). “Multi-Output Random Forest Regression to Em-
65 ulate the Earliest Stages of Planet Formation”. In: *2021 Systems and Information Engineering*
66 *Design Symposium (SIEDS)*, pp. 1–6.
- 67 Jog, A., A. Carass, S. Roy, D. L. Pham, and J. L. Prince (2017). “Random forest regression for
68 magnetic resonance image synthesis”. *Medical Image Analysis*, vol. 35, pp. 475–488.
- 69 Kinaneva, D., G. Hristov, P. Kyuchukov, G. Georgiev, P. Zahariev, and R. Daskalov (2021). “Ma-
70 chine Learning Algorithms for Regression Analysis and Predictions of Numerical Data”. In: *2021*
71 *3rd International Congress on Human-Computer Interaction, Optimization and Robotic Applica-*
72 *tions (HORA)*, pp. 1–6.
- 73 Linusson, H. (2013). “Multi-Output Random Forests”. Program: Magisterutbildning i informatik.