

Project Part 2.2: Modelling Plan

Students: Bahadir Bagci, Hassan Majeed, Zoheb Rahman

Team Name: Spectral Data Soft Sensor (B)

1. Introduction

In the last part, we completed initial data processing and number of principal component selection steps. From now on, we will be working on developing a model in order to predict selected five plant traits: carbon, chlorophyll, equivalent water thickness, leaf mass per area, and nitrogen, those were chosen because they have fewer missing values compared to the others.

2. Methodology

The data is three-dimensional, with spectral wavelengths (X), pixels (samples), and plant traits (Y). Our aim is not just to reduce the number of predictors, but also to ensure that the variation in X is linked with the variation in Y. Hence, Partial Least Squares (PLS) regression is applied.

2.1.Tools and Path

We will be using libraries scikit-learn(for PLS regression) and matplotlib(for visualization) in Python

2.2.Model Calibration Strategy

Based on the nature of this dataset, missing trait values do not imply that the corresponding traits are absent or zero. Rather, they indicate that those traits were not measured for that particular sample. Since even the least missing trait has over 30% missing values, we did not apply imputation. Filling in such a large amount of missing data without clear patterns could distort the relationships the model tries to learn. Instead, for each selected trait, we filtered out the rows where the target value is missing. We made sure to use only the real, available values for each trait, instead of filling in missing ones.

First of all, we split the dataset into training and test sets.(%70-%30) Then, for generalization, we checked how well the models perform during training. In order to do that, we used 5-fold cross validation on the training set. In brief, this method splits the data into five parts, trains the model on four parts, and validates it on the remaining part—repeating this process five times to help us understand the consistency of the models.