

# Yield Prediction by Machine Learning From UAS-Based Multi-Sensor Data Fusion in Soybean

Monica Herrero-Huerta (✉ [monicaherrero@usal.es](mailto:monicaherrero@usal.es))

Purdue University <https://orcid.org/0000-0002-4134-557X>

Pablo Rodriguez-Gonzalvez

Universidad de Leon

Katy M. Rainey

Purdue University

---

## Research

**Keywords:** Unmanned Aircraft System (UAS), High Throughput Phenotyping, Soybean, Structure from Motion (SfM), Machine Learning (ML), Yield, Point Clouds

**Posted Date:** May 13th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-16958/v2>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Plant Methods on June 1st, 2020. See the published version at <https://doi.org/10.1186/s13007-020-00620-6>.

# YIELD PREDICTION BY MACHINE LEARNING FROM UAS-BASED MULTI- SENSOR DATA FUSION IN SOYBEAN

Monica Herrero-Huerta <sup>\*</sup>,<sup>1</sup>, Pablo Rodriguez-Gonzalvez <sup>2</sup> and Katy M. Rainey <sup>1</sup>

<sup>1</sup> Department of Agronomy, Purdue University, IN 47906, USA – (mherrero, krainey@purdue.edu)

<sup>2</sup> Department of Mining Technology, Topography and Structures, Universidad de Leon, Ponferrada, Spain –  
p.rodriguez@unileon.es

## ABSTRACT:

### Background

Nowadays, automated phenotyping of plants is essential for precise and cost-effective improvement in the efficiency of crop genetics. In recent years, machine learning (ML) techniques have shown great success in the classification and modelling of crop parameters. In this research, we consider the capability of ML to perform grain yield prediction in soybeans by combining data from different optical sensors via RF (Random Forest) and XGBoost (eXtreme Gradient Boosting). During the 2018 growing season, a panel of 382 soybean recombinant inbred lines were evaluated in a yield trial at the Agronomy Center for Research and Education (ACRE) in West Lafayette (Indiana, USA). Images were acquired by the Parrot Sequoia Multispectral Sensor and the S.O.D.A. compact digital camera on board a senseFly eBee UAS (Unnamed Aircraft System) solution at R4 and early R5 growth stages. Next, a standard photogrammetric pipeline was carried out by SfM (Structure from Motion). Multispectral imagery serves to analyse the spectral response of the soybean end-member in 2D. In addition, RGB images were used to reconstruct the study area in 3D, evaluating the physiological growth dynamics per plot via height variations and crop volume estimations. As ground truth, destructive grain yield measurements were taken at the end of the growing season.

### Results

Algorithms and feature extraction techniques were combined to develop a regression model to predict final yield from imagery, achieving an accuracy of over **90.72%** by RF and **91.36%** by XGBoost.

### Conclusions

Results provide practical information for the selection of phenotypes for breeding coming from UAS data as a decision support tool, affording constant operational improvement and proactive management for high spatial precision.

**KEY WORDS:** Unmanned Aircraft System (UAS), High Throughput Phenotyping, Soybean, Structure from Motion (SfM), Machine Learning (ML), Yield, Point Clouds.

---

\* Corresponding author

## 31 1. BACKGROUND

32 Estimating morphological plant variables and the non-destructive characterization of traits with high accuracy and cost-  
33 effectiveness is imperative for high-throughput phenotyping in precision agriculture [1]. Recent advances in sensor technology  
34 provide great opportunities for the use of UAS (Unmanned Aircraft Systems) as a low-cost platform to derive high throughput  
35 and precise quantitative phenotyping datasets [2]. This technology offers images at high spatial, temporal and spectral  
36 resolution containing precise information about interactions from canopy and solar radiation [3]. Due to the increasing use of  
37 UAS, the development of software tools and methodologies to automatically phenotype crops is urgently required.  
38 Photogrammetric sensors on board the UAS allow for the application of digital image analysis of cover plant height estimation  
39 [4], yield estimation [5], early emergence, senescence rate [6], disease detection [7], quality evaluation [8] and canopy  
40 architecture [9]. RGB images have been used to accurately estimate vegetation index by deep neural network [10], while  
41 thermal sensors ability to capture canopy temperature has been used to detect water stress [11].

42 Plant height is a crucial variable connected to stability, yield potential and lodging resistance. This variable has been assessed  
43 by UAS as a Structure from Motion (SfM), obtaining high correlations with ground reference measurements for barley [12],  
44 wheat [13], poppy [14] and sorghum [15]. In addition, Light Detection and Ranging (LiDAR) is capable of providing 3D data  
45 including height and vegetation density areas on canopy structure [16]. It has been used to derive canopy height, fractional  
46 cover and above ground biomass [17].

47 Lately, machine learning (ML) models have been used to model plant traits based on image data. These methods employ  
48 sophisticated statistical techniques, being able to approximate complex non-linear functions between image features and  
49 biophysical parameters. Concretely, deep learning has been used for temporal phenotype/genotype classification [18].  
50 Moreover, [19] use  $k$ -NN as a classification method to analyse images of diverse germination phenotypes as well as to detect  
51 single seed germination. In addition, geometric parameters such as leaf counting have been addressed through plant models  
52 by [20]. The best characteristic of ML is the limited prior information necessary for it to be applied. This is due to these  
53 model's ability to capture assumptions and essential distributions directly from the training dataset [21]. Thus, the effect of the  
54 unknown variability is significantly reduced. As a disadvantage, the over-fitting of the models is a continuing problem that is  
55 difficult to mitigate [22]. Another weak point common in ML is the necessity for a similar distribution between training and  
56 testing datasets so that the model has the ability to properly predict variables; even for extensive training data. When  
57 distribution differences between both datasets exist, two related common errors appear, so-called *covariance shift* (the  
58 distribution changes between trained and testing data) and *dataset shift* (different distribution of the outputs and inputs from  
59 the test dataset regarding the distribution from the training dataset) [23]. Moreover, many ML approaches hold huge  
60 computational complexity, such as tuning learning parameters that may affect the model's robustness.

61 In this research, senseFly eBee was chosen as a UAS platform to automate the mapping at high spatial resolution using an  
62 onboard narrowband Parrot Sequoia Multispectral sensor and the the senseFly's S.O.D.A. compact digital camera. The images  
63 were separately processed through an end-to-end photogrammetric pipeline by computing the view of each image and,

64 subsequently, the generation of a dense and scaled 3D model of the crop and orthomosaic production. Next, the plot extraction  
 65 is carried out in 2D for the multispectral imagery and in 3D point clouds for the RGB data. The multi-spectral imagery (MSI)  
 66 features per plot are calculated applying the ‘Triple S’ pipeline (Statistical computing of Segmented Soybean multispectral  
 67 imagery) by statistically analysing the pixel values of soybean end-members by filtering the image through k-means clustering.  
 68 For RGB data, algorithms were employed to analyse height variations per plot and mesh calculations were applied to quantify  
 69 canopy volume using point clouds as a photogrammetric product. Features coming from both optical sensors are extracted to  
 70 perform a ML model by RF (Random Forest) and XGBoost (eXtreme Gradient Boosting), training the learning process and  
 71 validating it with grain yield field measurement. Therefore, the main goal is to predict the final yield based on imagery data  
 72 that will allow the selection of phenotypes for practical breeding, affording constant operational improvement and proactive  
 73 management with high spatial precision.  
 74 After this brief introduction, the employed materials and the proposed methodology will be described, followed by the  
 75 experimental results and analysis. To finalize, the conclusions and further studies are summarized.

## 76 2. MATERIALS

### 77 2.1. Materials

78 The materials used for the data acquisition are described below:

- 79 • A GNSS device from TopCon to georeference the Ground Control Points (GCPs), Hiper V receiver. The  
 80 topographic surveying was done using Real-Time Kinematic (RTK).
- 81 • A general purpose GER 1500 spectroradiometer to acquire spectral measurements of the calibration targets.
- 82 • A senseFly’s S.O.D.A. Digital Camera as an RGB photogrammetric sensor, with the following technical  
 83 specifications (Table 1):

84 Table 1. Technical specifications of the senseFly’s S.O.D.A. Digital Camera.

Parameter	Value
Optical Sensor Size	116.2 mm <sup>2</sup>
Image Size	5742*3648 pixels
Focal length	10.6 mm
Pixel size	3 μm

85

- 86 • A four narrowband passive sensor (Green, Red, Red-edge and Near infrared): Parrot Sequoia Multispectral sensor.  
 87 The camera specifications are detailed in Table 2. It has a global shutter to avoid problems in data processing [24]  
 88 and it is self-calibrating, using the incorporated Sunshine sensor.

89

Table 2. Technical specifications of the Parrot Sequoia Multispectral sensor.

Parameter	Value
Spectral range	350-2500 nm
Shooting time	0.1 s
Spectral resolution	1 nm
Field of view	25°
Pixel size	3.75 $\mu\text{m}$
Focal length	3.98 mm
Image size	1280*960 pixels

90

- 91 • The senseFly eBee, designed as a fixed wing UAS for application in precision agriculture with incorporated GPS,  
92 IMU and magnetometer. It has a weight of 700 g and a payload of 150 g. The digital camera on-board is controlled  
93 by the senseFly eBee autopilot during the flight.

## 94 2.2. Experimental setup

95 The soybean yield trial was performed at the Agronomy Center for Research and Education (ACRE) in 2018 in West Lafayette  
96 (Indiana, USA). An alpha lattice incomplete block design with 382 recombinant inbred lines, two complete replications and  
97 32 incomplete blocks per replication was planted [25]. Concretely, the panel includes lines from three classes of families: 16  
98 from elite parents, 12 with diverse pedigrees, and four that are high-yielding under drought conditions. The soybean field was  
99 on a silt loam soil with a pH of 6.5. The planting was performed at 2.5 cm depth in rows 0.76 m apart to a density of 40  
100 seed/m<sup>2</sup> on May 22nd, 2018. No fertilizers or herbicides for weed control were applied. Temperatures as measured by the on-  
101 farm weather station during the growing season averaged 20.56 °C in May, 22.68 °C in June, 22.78 °C in July, 22.57 °C in  
102 August, 20.98 °C in September and 11.75 °C in October. Monthly humidity, documented by the same weather station, was  
103 72% in May, 83% in June, 82% in July, 84% in August, 81% in September and 81% in October. The study area was  
104 282.4\*109.5 m<sup>2</sup>, consisting of 20 rainfed plots in vertical and 45 plots in horizontal, with different widths (6 and 8 rows). The  
105 photogrammetric flight configuration was with along-and across-track overlap of ca. 75%, adequate to Pix4D software  
106 processing. A flight altitude over the ground of 60 m for MSI (MultiSpectral Imagery) and 95 m for RGB was obtained by  
107 Sensefly software, given the camera focal and the required GSD (2 inches for MSI and 1 inch for RGB). A total of 114 MSI  
108 and 63 RGB images were used for the photogrammetric processing. For the RGB flight, the exposure time was fixed to 1/2000  
109 sec and the ISO was 125. 6 GCPs were placed on the ground for scaling and georeferencing purposes, identified by hand, and  
110 measured with GNSS, using RTKNAVI software [26]. GCPs are marked as dark grey rectangles and the study area was  
111 delimited by a black rectangle in Figure 2.

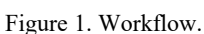
112 UAS flight performed as the planning flight was designed via autonomous flying mode on June 7<sup>th</sup> 2018 (Day After Planting  
113 (DAP) 15) with the G9X sensor to get the reference point cloud from the terrain and July 23<sup>rd</sup> 2018 (DAP 61) and August 1<sup>st</sup>  
114 (DAP 70) the Sequoia and G9X sensor for the study dataset, before the seed filling phenological period, from late R4 and early  
115 R5. All the experimental results obtained below were run on a 3.6-GHz desktop computer with an Intel CORE I7 CPU and  
116 32-GB RAM.

117 Plant height was checked against that of 5 fixed bars randomly placed over the study area for further analysis.

118 Soybean harvest was conducted on October 15<sup>th</sup>, 2018 with a small-plot research combine from Almaco. Grain Yield (GY)  
119 was performed by destructively harvesting an area of 0.5 × 0.5 m in the centre of each plot. Seed samples were processed in a  
120 drying oven at 105 °C for 48 h and later weighed. For analysis, weights were extrapolated to kg/ha and converted to 13%  
121 moisture to standardize the weight between plots. From a total of 876 plots, the mean GY value per plot was 3783.409 kg/ha  
122 with a standard deviation of 769.627 kg/ha. The minimum GY was 1915.249 kg/ha and the maximum 5422.898 kg/ha; the  
123 different quartiles reach the following value of 3442.216 kg/ha (25%), 3808.639 kg/ha (50%) and 4174.101 kg/ha (75%).

### 124 3. METHODS

125 The methodology followed is illustrated in Figure 1. First, multispectral and RGB images are acquired by UAS over the  
126 soybean breeding field, together with measurement from height fixed bars, spectral responses from reflectance targets and  
127 GPS (Global Positioning System) data from GCPs on field. After that, a photogrammetric pipeline was carried out, obtaining  
128 orthomosaics coming from MSI and point clouds from RGB data. Features from each plot are extracted to perform a RF and  
129 XGBoost model, training the learning process and validating it with destructive grain yield measurements, with the main goal  
130 being to predict the plots grain yield based on imagery data.

131  Figure 1. Workflow.

#### 132 3.1. UAS imagery

133 Proper flight planning is crucial to guarantee the imagery acquisition reaches the theoretical parameters, produces high quality  
134 images, achieves optimization of existing resources as well as minimizes the capture time.

135 Once the study area is defined, Sensefly software determines the flight strips, the camera orientation and the image acquisition  
136 regarding the restricted forward and side overlap and guaranteeing the scale for the required GSD (Ground Sample Distance),  
137 2.54 cm (1 inch) for RGB and 5.08 cm (2 inches) for MSI, based on the onboard sensor. Due to the proportion of spatial  
138 resolution of both flights, their combination in a single product is easier and there is no need for additional resampling  
139 operations. The parameters that define image capture are determined during flight execution depending on light conditions,  
140 wind and flight speed.

#### 141 3.2. Photogrammetric pipeline

142 Firstly, a topographic survey was performed that allows for the absolute georeferencing and scaling of the model. For this  
143 purpose, accuracy targets such as GCP were placed along the study area so as to be detectable in the acquired images. Once  
144 the aerial imagery had been captured, a standard photogrammetric pipeline was performed by image-based modelling  
145 techniques. Each dataset was handled by a framework based on camera calibration [27], image orientation and dense point  
146 cloud extraction [28]. The Pix4Dmapper software package (Pix4D SA, Lausanne, Switzerland) was employed for image  
147 processing, producing orthomosaics and 3D point clouds. In addition, the GCPs' measurements were employed in retrieving

148 the camera's interior parameters and correcting for any systematic error or block deformation. At this point, it is worth  
149 mentioned that the parameter's extraction from multispectral imagery is done through orthomosaic (i) while from RGB,  
150 geometric parameters are extracted based on 3D point clouds (ii).

151 (i) Images gathered by the Parrot Sequoia Multispectral sensor generated datasets for each flight that included Green, Red,  
152 Red Edge and NIR information. This sensor is a radiometric self-calibrating system. It incorporates an integrated irradiance  
153 sensor (Sunshine sensor) that allows irradiance values to be synchronized with the onboard GPS, IMU and magnetometer.  
154 Moreover, the relative influence of the atmosphere is minimal because the atmospheric column spanned by the radiation is  
155 unimportant and can be neglected in the calculations [29]. To radiometrically check this calibration, at the same time to the  
156 aerial data acquisition, a radiometric campaign on field was carried out over reflectance targets. Finally, the orthomosaics for  
157 each band are accurately geo-referenced to EPSG 32616, WGS84 CRS and the bands are merged, considering the parallax,  
158 using the Geospatial Data Abstraction Library (GDAL).

159 (ii) For the RGB data, geometric variables based on the generated point cloud, with a spatial resolution >100 points/m<sup>2</sup> and  
160 mesh calculations allows plant height estimations [4] and canopy volume, characterizing crop geometry with a high detail and  
161 accuracy (3.4.2 section).

### 162 **3.3. Point cloud processing**

163 Generated point clouds per each RGB flight are used to extract the soybean height and canopy volume, critical for biomass  
164 estimation [30]. In order to compute these absolute values, the reference dataset was used as explained below. These point  
165 clouds possibly enclose outliers owing to the massive and automated nature of the photogrammetric processing. To filter  
166 isolated clusters, a statistical analysis on each point's neighbourhood is performed by assuming a Gaussian distribution of  
167 neighbours' distances [31]. Afterward, to guarantee fully registered point clouds, the Iterative Closest Point algorithm [32] is  
168 used, getting an assumable mean error among ground points from the obtained point clouds. Afterwards, point clouds were  
169 filtered by a common bounding box, with the aim to derive physiological crop dynamics. A deviation point cloud of height  
170 variations between the reference dataset (where the plants do not emerge yet) and the studied datasets was computed.  
171 Consequently, an accurate cloud-to-cloud distance was derived, giving a local approximation model to the reference cloud by  
172 a quadric surface. These point cloud-based plant heights were calibrated by a comparison to 5 fixed bars randomly placed in  
173 the study area by measuring the height with a ruler to obtain field surveyed ground truth at the same time as the flights were  
174 performed.

175 The next step was the triangulation of these point clouds-based plant height. The meshing algorithm chosen was 3D Delaunay  
176 triangulation [33]. These meshes have to be refined to remove the errors generated during the automated process, through the  
177 approximation of Attene [34].

### 178 **3.4. Plot feature extraction**

179 We extracted different features per plot grouped in radiometric (through the multispectral orthomosaic) and geometric (based  
 180 on the point cloud by RGB data) parameters.

### 181 3.4.1. Radiometric features

182 Individual plot boundaries need to be extracted and defined separately from images with an assigned plot ID that defines their  
 183 genomic type by a field-map based plot extraction. First, we created a SPH file from the field map using QGIS open source  
 184 software. The script starts from the top right and builds the first polygon using the defined plot size and skips the gap between  
 185 plots and generates the next one until it gets to the last plot on the bottom left. One advantage is that it can be generalized to  
 186 other crop types as long as the field map is provided and the plots are planted in regular distance and have a consistent size  
 187 within a trial.

188 Once the individual plots are extracted, the ‘Triple S’ pipeline (Statistical computing of Segmented Soybean multispectral  
 189 imagery) was run. ‘Triple S’ [8] is an open source pipeline coded in Python that uses the GDAL library and Open Source  
 190 Computer Vision Library [35] running over Anaconda Prompt. From each plot, it generates the following information ordered  
 191 in a spreadsheet by the name of the plot file as follows: first, the image is classified in ground and soybean by  $k$ -means  
 192 clustering [36] using the near infrared band, which provides a bigger difference in the spectral response between end-members.  
 193 Once the image is filtered, the statistical parameters of the pixel-values of soybean end-member are calculated according to  
 194 Gaussian and robust models. Since, the possible presence of systematizms, and/or outliers, will hinder the fulfilment of the  
 195 hypothesis of a Gaussian distribution, statistics like the mean and the standard deviation will not provide a suitable analysis  
 196 [37]. For this reason, the following robust estimators are adopted in the present study: the median  $m$ , the normalized median  
 197 absolute deviation (NMAD) (Equation 1), the square root of the biweight midvariance (BwMv) (Equation 2), and the  
 198 interpercentile ranges (IPR):

$$199 \quad \text{NMAD} = 1.4826 \cdot \text{MAD} \quad (1)$$

$$200 \quad \text{BwMv} = \frac{n \sum_{i=1}^n a_i (x_i - m)^2 (1 - U_i^2)^4}{(\sum_{i=1}^n a_i (1 - U_i^2) (1 - 5U_i^2))^2} \quad (2)$$

$$201 \quad a_i = \begin{cases} 1, & \text{if } |U_i| < 1 \\ 0, & \text{if } |U_i| \geq 1 \end{cases} \quad (3)$$

$$202 \quad U = \frac{x_i - m}{9\text{MAD}} \quad (4)$$

203 being the median absolute deviation (MAD) (Equation 5), i.e. the median ( $m$ ) of the absolute deviations from  
 204 the data’s median ( $m_x$ ):

$$205 \quad \text{MAD} = m(|x_i - m_x|) \quad (5)$$



206 Please note that, for asymmetric distribution, will not be possible to provide a plus-minus range, therefore an absolute  
 207 interpercentile range at different confidence intervals will be provided (50 % also known as interquartile range, 90 % and 99  
 208 %), and additionally some percentile values such as 2.5 %, 25 %, 75 % and 97.5 %.

209 In the second step, canopy cover area (m<sup>2</sup>) was obtained by reading the coordinates in the metadata and relating it to the  
 210 number of soybean end-member pixels. The next step consists of acquiring the number of rows through an edge map that  
 211 determines if the row is completed. Canny algorithm [38] was used to obtain the edge map from the NIR band, in this case.  
 212 Finally, Principal Component Analysis (PCA) [39] computes the length of each row. The row length is the number of soybean  
 213 pixels along the first eigenvector of the covariance matrix [40]. Next, with median reflectance values, a bunch of VI  
 214 (Vegetation Index) are calculated as Table 3 indicates:

215 Table 3. VI used as inputs from the model

VI	Equation	Proposed by
NDVI	$(\text{NIR}-\text{R})/(\text{NIR}+\text{R})$	[41]
SAVI	$(1+\text{L})*(\text{NIR}-\text{R})/(\text{NIR}+\text{R}+\text{L})$	[42]
MSAVI	$(2*\text{NIR}+1-((2*\text{NIR}+1)^2-8*(\text{NIR}-\text{R})*(\text{NIR}-\text{R}))^{0.5})/2$	[43]
GESAVI	$(\text{NIR}-\text{a})*(\text{R}-\text{b})/(\text{R}+\text{z})$	[44]
CI <sub>re</sub>	$(\text{NIR}/\text{RE})-1$	[45]
CI <sub>g</sub>	$(\text{NIR}-\text{G})-1$	[45]
VARI	$(\text{G}-\text{R})/(\text{G}+\text{R})$	[44]
RVI	$(\text{NIR}/\text{R})$	[47]
DVI	$(\text{NIR}-\text{R})$	[48]
RDVI	$(\text{NIR}-\text{R})/(\text{NIR}+\text{R})^{0.5}$	[49]
TVI	$0.5*(120*(\text{NIR}-\text{G})-200*(\text{R}-\text{G}))$	[50]

216

### 217 3.4.2. Geometric features

218 In order to extract the point cloud from each plot, the commonly used file-based solution Rapidlasso LAStools [51] was used;  
 219 specifically, the tool named ‘lasclip’ using the SHP file already generated based on the field map.

220 Next, geometric features were extracted from the point cloud-based plant height and mesh from each plot; specifically,  
 221 maximum and mean height and the standard deviation as a quantification of the height variability from the point cloud. From  
 222 the mesh obtained as a triangulation of the point cloud, the canopy volume of each plot was calculated.

### 223 3.5. ML models: RF and XGBoost

224 Once the plot features were extracted, the yield prediction model was performed. Specifically, machine learning algorithms  
 225 develop an accurate prediction model from the training dataset. The analysis of optical sensor data often contains noise, this  
 226 issue can be compensated for by adding an appropriate quantity of characteristic training data [51]. From all ML methods,  
 227 assembly algorithms integrate a high number of individually weak but complementary predictors, to create a robust estimator.  
 228 This amalgamation could be done as either bagging or as boosting. Furthermore, tree learning algorithms do not involve linear

229 interactions between features (perfect for this type of data). For this study, RF as bagging and XGBoost algorithm as boosting  
 230 were chosen. A brief description of these both algorithms follows.

231 RF is one of the most known algorithms belonging to model aggregation ideas, introduced by [52]. The basics of RF theory  
 232 cover the convergence theorem and generalization error bound. More specifically, it is an ensemble machine learning method  
 233 [53] based on constructing a multitude of decision trees at training time, sampled independently and with the same distribution.  
 234 At each node, a given number of input variables are randomly chosen and the best split is calculated within this subset. No  
 235 pruning step is performed so all the trees of the forest are maximal trees. Another advantage of RF is that it is useful not only  
 236 in regression and classification problems, but also in the selection of variables. The out-of-bag (OOB) sample is the dataset  
 237 not used to generate the actual tree. It is used to estimate the prediction error as well as to assess variable importance in order  
 238 to perform the variable selection.

239 XGBoost, on the other hand, is a scalable nonlinear machine learning algorithm for tree boosting developed by [54]. This  
 240 method implies a computationally effective improvement of gradient boosting decision tree implementation where a new weak  
 241 learner is built to be maximally correlated with the negative gradient of the loss function related to the whole assembly for  
 242 each iteration [55]. Specifically, XGBoost speeds up the boosted tree construction operating in parallel and suggests a new  
 243 distributed algorithm for tree searching. The importance of each feature to the training model is considered when the boosted  
 244 trees are constructed to intelligently obtain the appropriated feature scores. Another characteristic is that XGBoost additionally  
 245 offer the possibility of penalizing the complexity of the trees.

246 To sum up, ML approaches aim to find a relationship between an input  $X = \{x_1, x_2, \dots, x_N\}$  and an output  $Y$  in the training  
 247 dataset and apply it to a testing dataset to assess the quality of the model. Thus, for both ML processes, *scikit-learn* [56] Python  
 248 libraries were implemented. The study area is randomly divided into a training and a testing zone, with a range of 15% using  
 249 split function imported from *sklearn.metrics* library. The random state is fixed to always obtain the same result. In the case of  
 250 RF, the maximum depth of a tree was set to 5 (default 6) to decrease the complexity of the model. The number of boosted trees  
 251 was set to 1000, commonly less than a thousand. For the XGBoost model, the learning rate was intentionally set to 0.06,  
 252 slighter than the default value (0.3), to head up to a more precise generalization [57]. The number of boosted trees was also  
 253 set to 1000 and the subsample to 0.8 to reduce the risk of over-fitting, making the training dataset more robust to the noise  
 254 generating randomness. The accuracy (ACC) is calculated as follows (Equation 6):

$$255 \quad ACC = 100 - \frac{(100 * \sum_{i=1}^{n_{test}} (\frac{x_{pred,test}^i - x_{act,test}^i}{x_{pred,test}^i}))}{n_{test}} \quad (6)$$

256 where  $x_{pred,test}^i$  is the predicted GY of the  $i^{th}$  plot from the testing dataset,  $x_{act,test}^i$  is the measured GY of the  $i^{th}$  plot  
 257 from the testing dataset used as the actual value and  $n_{test}$  is the total number of testing samples within the study area.

## 258 4. EXPERIMENTAL RESULTS

### 259 4.1. MSI results by 2D image processing

260 Images gathered by the Parrot Sequoia Multispectral sensor generate datasets for each flight that included Green (G), Red (R),  
 261 Red Edge (RE) and Near InfraRed (NIR) information. The weather conditions when the flights were done was clear and free  
 262 of clouds (during noon time). Data was separately processed per band by a photogrammetric pipeline to obtain the orthomosaic  
 263 required for GIS integration, considering the parallax. At the same time to the aerial data acquisition, a radiometric campaign  
 264 on field was carried out to radiometrically check the calibration of the sensor. Thus, calibration targets were placed in the study  
 265 area and measured by the spectroradiometer, obtaining a mean difference in reflectance between the measured target in field  
 266 and in the orthomosaic to less than 3.02% per band. In addition, to accurately reflect the breeding field planting configuration,  
 267 a script was developed to overlay defined plot sizes with known spacing and eliminate border effects by changing the plot size.  
 268 This automated plot extraction allows us to analyse each plot consisting, in total, of 900 individual plots with variable size.  
 269 Figure 2 illustrates the color composite of the multispectral orthomosaic (NIR+R+G) (a) and the automatic plot extraction over  
 270 a randomly selected area (b). Figure 2.c shows how Triple S was used for July 23<sup>rd</sup>, 2018 (DAP 61) to compute canopy cover,  
 271 row number and length for one random plot. As a brief analysis, we can see how the outliers influence the values, making  
 272 differences between mean and median value. The standard deviation represents the spatial variability in reflectance with no  
 273 correlation found along time per band once the outliers are removed. The threshold is the value obtained using *K-means* ( $k=2$   
 274 in this case: vegetation and ground) to mask the soybean member using NIR band (band 4).

275 Figure 2. Colour composite MSI mapping over the study area on July 23<sup>rd</sup> 2018 (DAP 61) (DL (500047.3,  
 276 4480849.5); UR (500364.0, 4480968.0); EPSG 32616) (a), a detail of field-map based plot extraction (b) and Triple S  
 277 software run over a random plot for July 23<sup>rd</sup>, 2018 (DAP 61) (c).

278 The statistics of variables from MSI analysis by plot are presented in Table 4 for the different study dates, July 23<sup>rd</sup>, 2018 and  
 279 August 1<sup>st</sup>, 2018, respectively: CC (canopy cover) and soybean reflectance by band. The length of row parameter was rejected  
 280 because of the lack of variation enough within the plot, also being influenced by the plot cut and the filter applied (*k-means*  
 281 clustering).

282 Table 4. Statistics of canopy cover and soybean reflectance by band of soybean class per plot from MSI analysis at DAP 61  
 283 and 70: mean, standard deviation (Std), median, normalized median absolute deviation (NMAD), square root of the biweight  
 284 midvariance (BwMv), percentiles at 2.5 % (P2.5%), 25 % (Q25%), 75 % (Q75%) and 97.5 % (P97.5%), interquartile range  
 285 (IQR) and interpercentile range at 90 % (IPR90%) and 99 % (IPR99%) confidence interval.

	<i>Parameter (%)</i>	<i>Mean</i>	<i>Std</i>	<i>Median</i>	<i>NMAD</i>	<i>BwMv</i>	<i>P2.5%</i>	<i>Q25%</i>	<i>Q75%</i>	<i>P97.5%</i>	<i>IQR</i>	<i>IPR90%</i>	<i>IPR99%</i>
	Canopy Cover	79.54	20.29	85.45	6.66	7.60	3.37	80.09	89.47	98.32	9.37	70.54	99.85
7/23/2018 (DAP 61)	Green	6.21	0.72	6.19	0.79	0.73	4.93	5.65	6.71	7.60	1.06	2.34	3.31
	Red	2.53	0.27	2.51	0.23	0.26	2.02	2.37	2.68	3.08	0.31	0.88	1.51
	Red Edge	31.84	2.61	32.03	2.26	2.44	26.29	30.42	33.45	36.91	3.02	8.25	16.49
	Near InfraRed	55.15	7.02	55.42	5.80	6.83	39.45	51.70	59.48	69.23	7.78	24.46	35.83
	Canopy Cover	86.90	5.48	87.77	3.62	3.79	69.42	85.22	90.06	94.35	4.83	19.35	31.97

8/01/2018 (DAP 70)	Green	5.90	0.39	5.85	0.39	0.39	5.25	5.61	6.15	6.76	0.54	1.25	1.96
	Red	2.62	0.18	2.60	0.18	0.18	2.33	2.48	2.73	3.03	0.24	0.57	1.00
	Red Edge	32.22	1.34	32.22	1.28	1.35	29.63	31.31	33.06	34.96	1.75	4.52	6.89
	Near InfraRed	55.59	2.27	55.58	2.19	2.22	50.99	54.12	57.08	60.08	2.96	7.59	13.14

286

287 It can be seen with the canopy cover parameter, the breach of the normality hypothesis causes the statistical dispersion to be  
288 overestimated, compared to robust values (NMAD, BwMv, percentile (P) and IPR).

#### 289 4.2. RGB results by 3D modelling

290 RGB data generates 3D point clouds. The point cloud from June 7<sup>th</sup> (DAP 15) was used as a terrain reference. It contains  
291 1,613,588 points while the one from July 23<sup>rd</sup> (DAP 61) has 5.74% more points for the same study area, 1,711,892. The one  
292 from August 1<sup>st</sup> (DAP 70) has 1,699,878 points. Please note that the variation of the spatial resolution of the computed point  
293 clouds for DAP 61 and 70 is due to the texture changes, which affects (among other factors) the densification operation. The  
294 three flights reach the same GSD. The next step was the registration of the point cloud from DAP 61 and DAP 70 against the  
295 one from DAP 15 using the ICP algorithm [58] on terrain points. Firstly, the coarse registration was done by manually picking  
296 similar GCP. Secondly, the ICP algorithm finds that affine transformation matrix that minimizes the distances between closet  
297 points from terrain points of the two point clouds considered. Once the alignment was done, the height value was checked  
298 against 5 height fixed bars randomly placed over the study area, reaching a difference of less than 2.46 cm for the study date  
299 of July 23<sup>rd</sup> (DAP 61) and 2.21 cm for the study date of August 1<sup>st</sup> (DAP 70). On the other hand, the deviation point cloud  
300 from July 23<sup>rd</sup> (DAP 61) reaches the following statistical parameters: a minimum height of 0 m, a maximum of 1.244 m, a  
301 mean of 0.578 m and a standard deviation of 0.614 (Figure 3.a); while the one from August 1<sup>st</sup> (DAP 70) has a minimum  
302 height of 0 m, a maximum of 1.476 m, a mean of 0.798 m and a standard deviation of 0.803.

303 Figure 3.b analyses two particular plots from July 23<sup>rd</sup>, 2018 where the visual differences in quantifying the canopy volume  
304 could be appreciated. Calibrated point clouds are converted into meshes by applying a 3D Delaunay triangulation and refined:  
305 filling of holes through algorithms of planar triangulation, repairing of meshing gaps by threshold algorithms and removal of  
306 topological and geometric noise by anti-aliased Laplacians filters. The grid was chosen as 45 cm as a trade-off between spatial  
307 resolution that affects the accuracy and computational cost. Finally, these meshes give us the value of the canopy volume per  
308 plot.

309 Figure 3. Deviation point cloud over Soybean from July 23<sup>rd</sup> (DAP 61) using June 7<sup>th</sup> (DAP 15) as reference in *meters* (a)  
310 and canopy volume calculation of two random plots from the same date (b).

311 The statistics of variables from RGB analysis by plot are presented in Table 5 for the different study dates, July 23<sup>rd</sup>, 2018 and  
312 August 1<sup>st</sup>, 2018, respectively: CV (canopy volume), H max (maximum height) and variation of these parameters within the  
313 plot (OCV and OHmax). From these results, we can affirm that these variations (OCV and OHmax) can be assumed as equal.

314 Table 5. Statistics of CV and H max by band per plot from RGB analysis at DAP 61 and 70: mean, standard deviation (Std),  
 315 median, normalized median absolute deviation (NMAD), square root of the biweight midvariance (BwMv), percentiles at 2.5  
 316 % (P2.5%), 25 % (Q25%), 75 % (Q75%) and 97.5 % (P97.5%), interquartile range (IQR) and interpercentile range at 90 %  
 317 (IPR90%) and 99 % (IPR99%) confidence interval.

	<i>Parameter</i>	<i>Mean</i>	<i>Std.</i>	<i>Median</i>	<i>NMAD</i>	<i>BwMv</i>	<i>P 2.5%</i>	<i>Q 25%</i>	<i>Q 75%</i>	<i>P 97.5%</i>	<i>IQR</i>	<i>IPR 90%</i>	<i>IPR 99%</i>
7/23/2018 (DAP 61)	CV (dm <sup>3</sup> )	1282.7 8	218.37	1253.85	199.04	215.52	917.37	1135.57	1413.23	1754.78	277.66	729.35	1119.80
	Hmax (cm)	92.33	16.64	87.58	10.89	12.36	73.02	81.70	97.16	139.61	15.46	56.37	84.95
	OCV (% dm <sup>3</sup> )	19.65	4.98	19.18	4.96	4.77	11.78	15.91	22.57	31.18	6.66	16.03	27.06
	OHmax (% cm)	19.14	4.84	18.66	4.63	4.61	11.44	15.54	21.73	30.78	6.20	15.66	25.57
8/01/2018 (DAP 70)	CV (dm <sup>3</sup> )	1496.7 9	242.17	1487.83	219.58	233.32	1065.43	1327.34	1630.02	2079.86	302.69	819.91	1388.60
	Hmax (cm)	107.57	16.95	104.02	11.76	13.62	84.53	96.76	113.45	154.90	16.70	55.02	88.09
	OCV (% dm <sup>3</sup> )	21.28	5.95	20.92	6.18	5.95	11.65	16.74	25.03	34.52	8.29	19.13	29.36
	OHmax (% cm)	20.84	5.74	20.49	5.95	5.68	11.55	16.48	24.54	33.87	8.06	18.92	30.03

318

319 In this case, the Gaussian values of the central tendency and dispersion of the parameters do not differ markedly as in the  
 320 previous case (Table 6). However, the normality condition is not met in any of the previous 18 cases, with the results of the  
 321 Robust Jarque-Bera test [59] for a significance level of 5%.

322

### 323 4.3. ML model results

324 In this study, we developed tree learning models via RF and XGBoost for soybean yield prediction by UAS-based imagery.  
 325 To sum up, we used 840 plots with a rate of 15% to check the model: 714 trained plots and 126 tested plots. The features used  
 326 are 60 between both dates, 12 coming from the RGB analysis (canopy volume, maximum height and their standard deviations  
 327 within each plot from DAP 61, DAP 70 and from the point cloud that represents the increment from DAP 61 to DAP 70) and  
 328 48 from the MSI coming from DAP 61 and 70, containing canopy cover value, 24 parameters from each band (mean, median,  
 329 standard deviation) and 22 VI (GESAVI, NDVI, SAVI, MSAVI, Clre, Clg, VARI, RVI, DVI, RDVI and TVI). As a result,  
 330 we achieve an accuracy over **90.72%** by RF and **91.36%** by XGBoost computed as Equation 6 indicates.

331 The features which represent more than 71% of the importance in each model are shown in the Figure 4.a by RF and 4.b by  
 332 XGBoost. Analysing this importance parameter, we can see that the Clg index for the DAP 70 is the most related feature while  
 333 TVI and DVI are negligible regarding Grain Yield in both models. Clg represents the canopy chlorophyll content using G and  
 334 NIR band.

335 Figure 4. Features importance for more than 71% by RF (a) and XGBoost (b).

336 To quantify how the sensors contribute to the accuracy of the fusion models, both models were run using only RGB features,  
 337 increasing the MAE (Mean Absolute Error) in 36.99% by RF and 31.72% by XGBoost. When only MSI features are used, the  
 338 MAE increases in 8.97% by RF and 14.74 by XGBoost; clearly showing how multispectral features are more related to yield  
 339 than geometric measurements based on RGB data.

340 To analyse when the images should be captured, we run the models only with features provided by DAP 61, the MAE increases  
 341 in 10.49% by RF and 12.74% by XGBoost. When the models are run with features from DAP 70, the MAE increases in 3.16%  
 342 by RF and 5.95% by XGBoost. These results affirm that the images from DAP 70 better predict the yield than the images  
 343 captured on DAP 61.

344

## 345 5. VALIDATION RESULTS AND DISCUSSION

346 In this section, an accurate analysis of the predicted values from the ML models is carried out. Figure 5 show the absolute  
 347 errors for the actual GY sorted from smallest to largest per plot along the training dataset (Figure 5.a) and testing dataset  
 348 (Figure 5.b). In both process, XGBoost and RF, the error is larger when the actual GY values are more extreme are. As  
 349 expected, RF works better in fixing the training dataset than the testing, compared with XGBoost. However, we can assume  
 350 that both ML approaches achieve the same total accuracy generating the regression model.

351 Machine learning models are able to accurately fit the training data. As a disadvantage, they are susceptible to overfitting when  
 352 small or large datasets with an insufficient level of variation [60]. For this reason, the validation errors along time were  
 353 compared against the trained errors verifying that the validation errors do not increment while the trained errors decrease.

354 To quantitatively assess the models' performance, different errors were computed. Table 6 shows the values of error metrics  
 355 from both models in (kg/ha) evaluated for the training and the testing dataset. A 95% confidence level was applied to these  
 356 estimated errors. As a reference value, the mean GY measured per plot is 3783.409 kg/ha for all the dataset; 3777.45 kg/ha for  
 357 the training dataset and 3817.16 kg/ha for the testing dataset. The Mean Bias Error (MBE), the Absolute Mean Bias Error  
 358 (AMBE), the Root Mean Square Error (RMSE), the Relative Error (RE) and the Absolute Error (AE) were computed as follows  
 359 (Equations 7-11):

$$360 \quad MBE = \frac{\sum_{i=1}^n (x_{pred}^i - x_{act}^i)}{n} \quad (7)$$

$$361 \quad AMBE = \frac{\sum_{i=1}^n |(x_{pred}^i - x_{act}^i)|}{n} \quad (8)$$

$$362 \quad RMSE = \sqrt{\frac{\sum_{i=1}^n (x_{pred}^i - x_{act}^i)^2}{n}} \quad (9)$$

$$363 \quad RE = 100 * \frac{\sum_{i=1}^n (x_{pred}^i - x_{act}^i) / x_{act}^i}{n} \quad (10)$$

$$364 \quad AE = 100 * \frac{\sum_{i=1}^n |(x_{pred}^i - x_{act}^i) / x_{act}^i|}{n} \quad (11)$$

365 where  $x_{pred}^i$  is the predicted GY of the  $i^{th}$  plot,  $x_{act}^i$  is the measured GY of the  $i^{th}$  plot used as the actual value and  $n$  is  
366 the total number of samples within the study area. The NMAD was defined in section 3.4.1 (see Equation 1)  
367 In addition, the Nash and Sutcliffe index,  $\eta$  is also computed (Equation 12); used in modelling to characterize the error related  
368 to the spatial heterogeneity:

$$369 \quad \eta = 1 - \frac{\sum_{i=1}^n (x_{pred}^i - x_{act}^i)^2}{\sum_{i=1}^n (x_{pred}^i - \bar{x}_{act})^2} \quad (12)$$

370 where  $\bar{x}_{act}$  is the actual average GY.

371 Some of these evaluation metrics have been extensively used to analysis the power of regression models [61].

372 Table 6. Error metrics of both models in (kg/ha) at 95% confidence interval evaluated in training and testing dataset: MBE  
373 (Mean Bias Error), AMBE (Absolute Mean Bias Error), RMSE (Root Mean Square Error), NMAD (normalized median  
374 absolute deviation), RE (Relative Error), AE (Absolute Error) and  $\eta$  (the Nash and Sutcliffe index).

<b>Dataset</b>	<b>Model</b>	<b>MBE</b>	<b>AMBE</b>	<b>RMSE</b>	<b>NMAD</b>	<b>RE</b>	<b>AE</b>	<b><math>\eta</math></b>
<b>Training</b>	<i>RF</i>	13.61	140.25	181.19	167.48	1.14 %	4.03 %	0.80
	<i>XGBoost</i>	30.39	240.45	303.99	292.12	1.98 %	6.87 %	0.21
<b>Testing</b>	<i>RF</i>	-4.17	325.33	410.24	384.62	1.37 %	9.06 %	-2.46
	<i>XGBoost</i>	-7.15	306.76	394.66	353.04	1.18 %	8.55 %	-1.52

375  
376 Smaller values of MBE, AMBE, RMSE, NMAD, RE and AE and larger values of  $\eta$  ( $-\infty < \eta \leq 1$ ) indicate better precision and  
377 accuracy of the prediction model. With these results, we can affirm that XGBoost performs better than RF for this type of  
378 data, probably dealing better with overfitting.

379 Figure 5 shows the scatter plots of the measured vs. predicted GY values from the training (Figure 5.c) and testing dataset  
380 (Figure 5.d) in both models, RF and XGBoost. In both cases is fit a linear function according to a bisquare weighting. For  
381 the computation the outliers are discarded according to the studentized residuals at for a significance level of 0.05 for a two  
382 tails distribution. The coefficients, the regression values ( $R^2$ ) and the highest studentized residual are shown in Table 7. The  
383  $i$ -th studentized residual ( $sr_i$ ) is computed as the division of the residual ( $r_i$ ) of the  $i$ -th observation by the exact residual  
384 standard deviation [62] (Equation 13):

$$385 \quad sr_i = \frac{r_i}{\sqrt{MS_{Res} (1-h_{ii})}} \quad (13)$$

386 being  $MS_{Res}$  the mean squared error of the regression fit calculated by removing the  $i$ -th observation, and  $h_{ii}$   
387 is leverage value for the  $i$ -th observation ( $i$ -th element of the diagonal of the hat matrix).

388 As shown by [61], studentized residual is generally recommended instead normalized residual for least squares fit, since any  
389 point with a large residual and a large  $h_{ii}$  is potentially highly influential. If the absolute value of a studentized residuals is  
390 greater than a critical threshold, then the observation is marked as outlier. The critical threshold is defined from a  $t$ -

391 distribution with  $n-p-1$  degrees of freedom; being  $n$  de number of observations and  $p$  the number of fit coefficients. A total  
 392 of 47 and 32 outliers were detected for the training RF and XGBoost models respectively; and 2 and 1 for the testing RF and  
 393 XGBoost models respectively.

394 Figure 5. Prediction errors and actual Grain Yield (GY) sorted smallest to largest per plot along the training dataset (a) and  
 395 testing dataset (b) (please note that the errors and GY are plotted in the primary and secondary axis, respectively); scatter  
 396 plots of the measured against the predicted grain yield ( $\text{kg}\cdot\text{ha}^{-1}$ ) by RF and XGBoost from training (c) and testing dataset  
 397 (d). In both cases is drawn the line corresponding to the robust linear fit at 95 % of confidence.

398 Table 7. Robust linear fit coefficient,  $R^2$  value, highest studentized residuals mad RMSE & NMAD values of the fitting.

Dataset	Model	a	b	$R^2$	Max Studentized Residual	RMSE	NMAD
<i>Training</i>	<i>RF</i>	1.372	-1420.5	0.9728	1.95	94.06	102.40
	<i>XGBoost</i>	1.429	-1638.8	0.7787	1.95	262.03	263.05
<i>Testing</i>	<i>RF</i>	1.433	-1614.7	0.3828	1.88	399.87	360.08
	<i>XGBoost</i>	1.290	-1069.1	0.4183	1.89	387.11	370.03

399 A brief checkup about how different genotypes affect our GY prediction is introduced in Figure 6, where the AE (Absolute  
 400 Error) from the testing dataset is grouped by families (families within 4-6 predicted values in the testing dataset), potentially  
 401 being PI404188A the best family predicted.

402 Figure 6. Errors from the testing dataset grouped by family.

## 403 6. CONCLUSIONS

404 This paper demonstrates the great potential of UAS to predict soybean yield from multi-sensor data fusion as a rapid, accurate  
 405 and cost-effective tool for automated high throughput phenotyping. Specifically, this study evaluates the power of high spatial  
 406 resolution optical data, combined with regression models based on machine learning approaches (RF and XBOOST) to  
 407 effectively obtain high correlations with yield in breeding trials.

408 Although data fusion is able to increase the accuracy in phenotype prediction, future researches should address the efficiency  
 409 of different sensor combinations. The sensor cost and the accuracy improvement should be assessed for each study.  
 410 Additionally, this workflow can be successfully used for other HTPPs (High Throughput Phenotyping Platforms) and other  
 411 crops planted in breeding nurseries. Even so, more comprehensive studies are necessary, including studies on different crop  
 412 species at different phenotypic stages. Furthermore, UAS approaches for precision farming are in constant evolution and  
 413 represents an extremely dynamic sector. In this context, this research is our contribution as a methodology for yield prediction  
 414 in soybean from UAS-based multi-sensor data fusion by machine learning approaches.

## 415 DECLARATIONS

### 416 Availability of data and materials



417 The datasets used and analysed during the current study are available from the corresponding author on reasonable request.

#### 418 **Competing interests**

419 The authors declare that they have no competing interests.

#### 420 **Funding**

421 The research is partly funded by the project ‘Development of Analytical Tools for Drone-based Canopy Phenotyping in Crop  
422 Breeding’ from the American Institute of Food and Agriculture.

#### 423 **Author’s contributions**

424 M.H. conceived the experiments; M.H. performed the experiments; M.H. and P.R. analysed the data; M.H. was a major  
425 contributor in writing the manuscript; P.R. wrote and edited the manuscript; KR supervised the research and acquired funding;  
426 all authors approved the final manuscript.

#### 427 **Acknowledgements**

428 Authors would like to thank Miguel Lopez, Fabiana Freitas, David Schlueter, Ryan Ferguson, Aaron Shwarts, Smit Stuart  
429 and Keith Cherkauer for their collaboration during the experimental phase of this research.

#### 430 **Ethics approval and consent to participate**

431 Not applicable.

#### 432 **Consent for publication**

433 Not applicable.

## 434 **REFERENCES**

- 435 1. Furbank RT, Tester M. Phenomics-technologies to relieve the phenotyping bottleneck. *Trends Plant Sci.* 2011;  
436 16(12):635-644.
- 437 2. Araus JL, Kefauver SC, Zaman-Allah M, Olsen MS, Cairns JE. Translating high throughput phenotyping into  
438 genetic gain. *Trends Plant Sci.* 2018; 23(5):451-466.
- 439 3. Thenkabail PS, Lyon JG, Huete A. *Hyperspectral remote sensing of vegetation*. Boca Raton: CRC Press. 2011;  
440 1943–61.
- 441 4. Malambo L, Popescu SC, Murray SC, Putman E, Pugh NA, Horne DW, Vidrine M. Multitemporal field-based plant  
442 height estimation using 3D point clouds generated from small unmanned aerial systems high-resolution imagery.  
443 *International journal of applied earth observation and geoinformation.* 2018; 64, 31-42.
- 444 5. Roth L, Streit B. Predicting cover crop biomass by lightweight UAS-based RGB and NIR photography: an applied  
445 photogrammetric approach. *Precision Agriculture.* 2018; 19(1), 93-114.

- 446 6. Hassan MA, Yang M, Rasheed A, Jin X, Xia X, Xiao Y. Time-series multispectral indices from unmanned aerial  
447 vehicle imagery reveal senescence rate in bread wheat. *Remote Sens.* 2018; 10(6):809.
- 448 7. Whalley JL, Shanmuganathan S. Applications of image processing in viticulture: A review. 20<sup>th</sup> International  
449 Congress on Modelling and Simulation. 2013.
- 450 8. Herrero-Huerta M, Govindarajan S, Cherkauer K, Rainey K. Triple S: a new tool for soybean high throughput  
451 phenotyping from UAS-based multispectral imagery. *SPIE Defense + Commercial Sensing.* 2019; 1007-20.
- 452 9. Paulus S. Accessing the plant architecture in 3D for plant phenotyping-recent approaches and requirements. In  
453 *Precision agriculture'19.* 2019; 315-321. Wageningen Academic Publishers.
- 454 10. Khan SH, Hayat M, Bennamoun M, Sohel FA, Togneri R. Cost-sensitive learning of deep feature representations  
455 from imbalanced data. *IEEE transactions on neural networks and learning systems.* 2018; 29(8), 3573-3587.
- 456
- 457 11. Berni JAJ, Zarco-Tejada PJ, Suárez L, González-Dugo V, Fereres E. Remote sensing of vegetation from UAV  
458 platforms using lightweight multispectral and thermal imaging sensors. *Int. Arch. Photogramm. Remote Sens.*  
459 *Spatial Inform. Sci.* 2019; 38(6).
- 460
- 461 12. Bendig J, Bolten A, Bennertz S, Broscheit J, Eichfuss S, Bareth G. Estimating biomass of barley using crop surface  
462 models (CSMs) derived from UAV-based RGB imaging. *Remote Sens.* 2014; 6(11):10395.
- 463 13. Khan Z, Chopin J, Cai J, Eichi VR, Haefele S, Miklavcic S. Quantitative estimation of wheat phenotyping traits  
464 using ground and aerial imagery. *Remote Sens.* 2018; 10(6):950.
- 465 14. Iqbal F, Lucieer A, Barry K, Wells R. Poppy crop height and capsule volume estimation from a single UAS flight.  
466 *Remote Sens.* 2017; 9(7):647.
- 467 15. Hu P, Chapman SC, Wang X, Potgieter A, Duan T, Jordan D. Estimation of plant height using a high throughput  
468 phenotyping platform based on unmanned aerial vehicle and self-calibration: Example for sorghum breeding. *Eur.*  
469 *J. Agron.* 2018; 95:24-32.
- 470 16. Herrero-Huerta M, Felipe-García B, Belmar-Lizarán S, Hernández-López D, Rodríguez-Gonzálvez P, González-  
471 Aguilera D. Dense Canopy Height Model from a low-cost photogrammetric platform and LiDAR data. *Trees.* 2016;  
472 30(4), 1287-1301.
- 473 17. Wallace L, Lucieer A, Watson C, Turner D. Development of a UAV-LiDAR system with application to forest  
474 inventory. *Remote Sensing.* 2012; 4(6), 1519-1543.

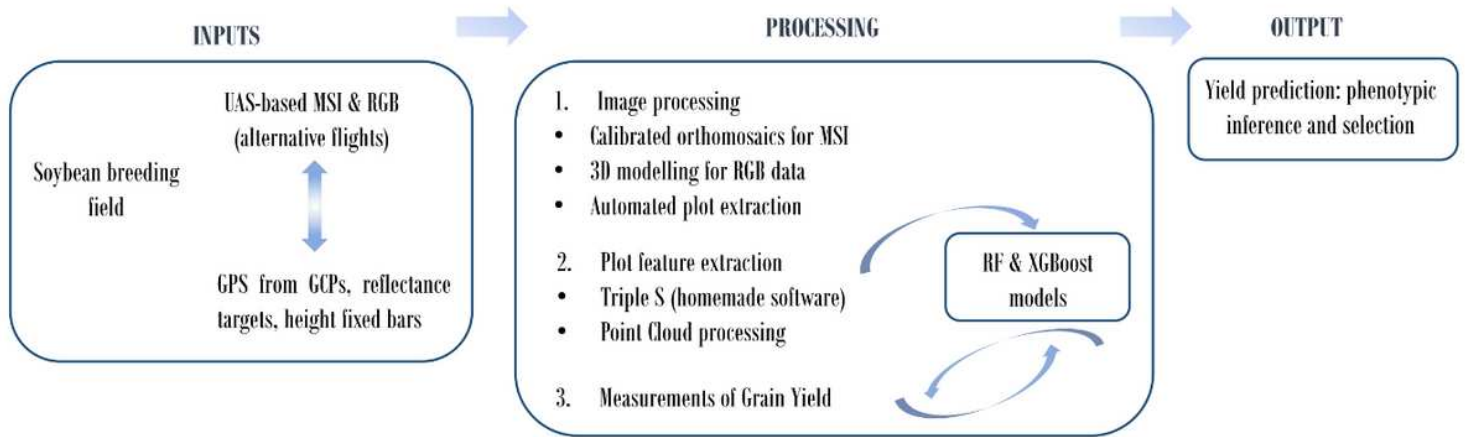
- 475 18. Namin ST, Esmacilzadeh M, Najafi M, Brown TB, Borevitz JO. Deep phenotyping: deep learning for temporal  
476 phenotype/genotype classification. *Plant methods*. 2018; 14(1), 66.
- 477 19. Awty-Carroll D, Clifton-Brown J, Robson P. Using k-NN to analyse images of diverse germination phenotypes and  
478 detect single seed germination in *Miscanthus sinensis*. *Plant methods*. 2018; 14(1), 5.
- 479 20. Ubbens J, Cieslak M, Prusinkiewicz P, Stavness I. The use of plant models in deep learning: an application to leaf  
480 counting in rosette plants. *Plant methods*. 2018; 14(1), 6.
- 481 21. Buxton H. Learning and understanding dynamics scene activity: a review. *Image and Vision Computing*. 2003;  
482 21(1), 125–136.
- 483 22. Maimaitijiang M, Ghulam A, Sidike P, Hartling S, Maimaitiyiming M, Peterson K, Burken J. Unmanned Aerial  
484 System (UAS)-based phenotyping of soybean using multi-sensor data fusion and extreme learning machine. *ISPRS  
485 Journal of Photogrammetry and Remote Sensing*. 2017; 134, 43-58.
- 486 23. Moreno-Torres JG, Raeder T, Alaiz-Rodríguez R, Chawla NV, Herrera F. A unifying view on dataset shift in  
487 classification. *Pattern Recogn*. 2012; 45(1):521–30.
- 488 24. Turner D, Lucieer A, Malenovsky Z, King DH, Robinson SA. Spatial co-registration of ultra-high resolution visible,  
489 multispectral and thermal images acquired with a micro-UAV over Antarctic moss beds. *Remote Sensing*. 2014; 6,  
490 4003–4024.
- 491 25. Lopez MA, Xavier A, Rainey KM. Phenotypic variation and genetic architecture for photosynthesis and water use  
492 efficiency in Soybean (*Glycine max* L. Merr). *Frontiers in plant science* 10 (2019): 680.
- 493 26. Takasu T. RTKLIB: Open Source Program Package for RTK-GPS. FOSS4G, Tokyo, Japan. 2009.
- 494 27. Remondino F, Fraser C. Digital camera calibration methods: Considerations and comparisons. *Int. Arch.  
495 Photogramm. Remote Sens. Spat. Inf. Sci*. 2006; 36, 266–272.
- 496 28. Herrero-Huerta M, González-Aguilera D, Rodríguez-Gonzalvez P, Hernández-López D. Vineyard yield estimation  
497 by automatic 3D bunch modelling in field conditions. *Computers and electronics in agriculture*. 2015; 110, 17-26.
- 498 29. Herrero-Huerta M, Hernández-López D, Rodríguez-Gonzalvez P, González-Aguilera D, González-Piqueras J.  
499 Vicarious radiometric calibration of a multispectral sensor from an aerial trike applied to precision agriculture.  
500 *Computers and Electronics in Agriculture*; 2014; 108, 28-38.

- 501 30. Tilly N, Aasen H, Bareth G. Fusion of plant height and vegetation indices for the estimation of barley biomass.  
502 Remote Sensing. 2015; 7(9), 11449-11480.
- 503 31. Herrero-Huerta M, Lindenbergh R, Rodríguez-González P. Automatic tree parameter extraction by a Mobile  
504 LiDAR System in an urban context. PloS one. 2018; 13(4), e0196004.
- 505 32. Besl PJ, McKay ND. A method for registration of 3-D shapes. Trans. Pattern Anal. Mach. Intell. 1992; 14, 239–256.
- 506 33. Golias NA, Dutton, RW. Delaunay triangulation and 3D adaptive mesh generation. Finite elements Anal. Des. 1997;  
507 25 (3), 331–341.
- 508 34. Attene M. A lightweight approach to repairing digitized polygon meshes. Vis Comput. 2010; 26(11):1393–1406.
- 509 35. Open Source Computer Vision Library, <http://sourceforge.net/projects/opencvlibrary/> (accessed February 2019).
- 510 36. Hartigan JA, Wong MA. Algorithm AS 136: A k-means clustering algorithm. Applied statistics, 1979; 100-108.
- 511 37. Nocerino E, Menna F, Remondino F, Toschi I, Rodríguez-González P. Investigation of indoor and outdoor  
512 performance of two portable mobile mapping systems. In Videometrics, Range Imaging, and Applications XIV.  
513 International Society for Optics and Photonics. 2017; doi: <https://doi.org/10.1117/12.2270761>
- 514 38. Canny JA. Computational Approach To Edge Detection. IEEE Trans. Pattern Analysis and Machine Intelligence.  
515 1986; 8(6):679–698.
- 516 39. Jolliffe IT. Graphical representation of data using principal components. Principal Component Analysis. Ed. New  
517 York, NY: Springer. 2002; 78–110.
- 518 40. Weinmann M, Jutzi B, Mallet C. Semantic 3D scene interpretation: a framework combining optimal neighborhood  
519 size selection with relevant features. Ann. Photogramm. Remote Sens. Spat. Inf. Sci. 2014; doi: 10.5194/isprsannals-  
520 II-3-181-2014.
- 521 41. Rouse JWJ, Haas RH, Schell JA, et al. Monitoring vegetation systems in the Great Plains with ERTS. Nasa Spec  
522 Publ. 1974; 351:309.
- 523 42. Huete AR. A soil-adjusted vegetation index (SAVI). Remote Sens Environ. 1988;25(3): 295–309.
- 524 43. Qi J, Chehbouni A, Huete A, Kerr Y, Sorooshian S. A modified soil adjusted vegetation index. Remote Sens.  
525 Environ. 1994; 48 (2), 119–126.

- 526 44. Gilabert MA, González-Piqueras J, García-Haro FJ, Meliá J. A generalized soil-adjusted vegetation index. *Remote*  
527 *Sens. Environ.* 2002; 82, 303–310.
- 528 45. Gitelson AA, Viña A, Ciganda V, et al. Remote estimation of canopy chlorophyll content in crops. *Geophys Res*  
529 *Lett.* 2005; 32(8):93–114.
- 530 46. Gitelson AA, Kaufman YJ, Stark R, et al. Novel algorithms for remote estimation of vegetation fraction. *Remote*  
531 *Sens Environ.* 2002; 80(1):76–87.
- 532 47. Jordan CF. Derivation of leaf-area index from quality of light on the forest floor. *Ecology.* 1969; 50(4):663–6.
- 533 48. Richardson AJ, Wiegand CL. Distinguishing vegetation from soil background information. *Photogramm Eng*  
534 *Remote Sens.* 1977; 43(12):1541–52.
- 535 49. Rougean, JL, Breon, FM. Estimating PAR absorbed by vegetation from bidirectional reflectance measurements.  
536 *Remote Sens. Environ.* 1995; 51, 375– 384.
- 537 50. Broge NH, Leblanc E. Comparing prediction power and stability of broadband and hyperspectral vegetation indices  
538 for estimation of green leaf area index and canopy chlorophyll density. *Remote Sens Environ.* 2001; 76(2):156–72.
- 539 51. Rapidlasso GmbH, 2019. <http://rapidlasso.com> (accessed 26.07.19).
- 540 52. Breiman L. Statistical modeling: the two cultures. *Statistical Science.* 2001; 16(3), 199–231.
- 541 53. Dietterich T. Ensemble methods in machine learning. *Lecture Notes Comput. Sci.* 2000; 1857, 1–15.
- 542 54. Chen T, Guestrin C. Xgboost: a scalable tree boosting system. In: *Proceedings of the 22Nd ACM SIGKDD*  
543 *international conference on knowledge discovery and data mining.* ACM. 2016; pp 785–794.
- 544 55. Natekin A, Knoll A. Gradient boosting machines, a tutorial. *Front neurorobotics.* 2013; 7:21.
- 545 56. Raschka S. *Python machine learning.* Packt Publishing Ltd, Birmingham. 2015.
- 546 57. Friedman JH. Stochastic gradient boosting. *Comput Stat Data Anal.* 2002; 38(4):367–378. Besl PJ, McKay ND.  
547 *Method for registration of 3-D shapes.* In *Robotics-DL tentative.* International Society for Optics and Photonics.  
548 1992; pp. 586-606.
- 549 58. Besl PJ, McKay ND. *Method for registration of 3-D shapes.* In *Robotics-DL tentative.* International Society for  
550 *Optics and Photonics.* 1992; pp. 586-606.

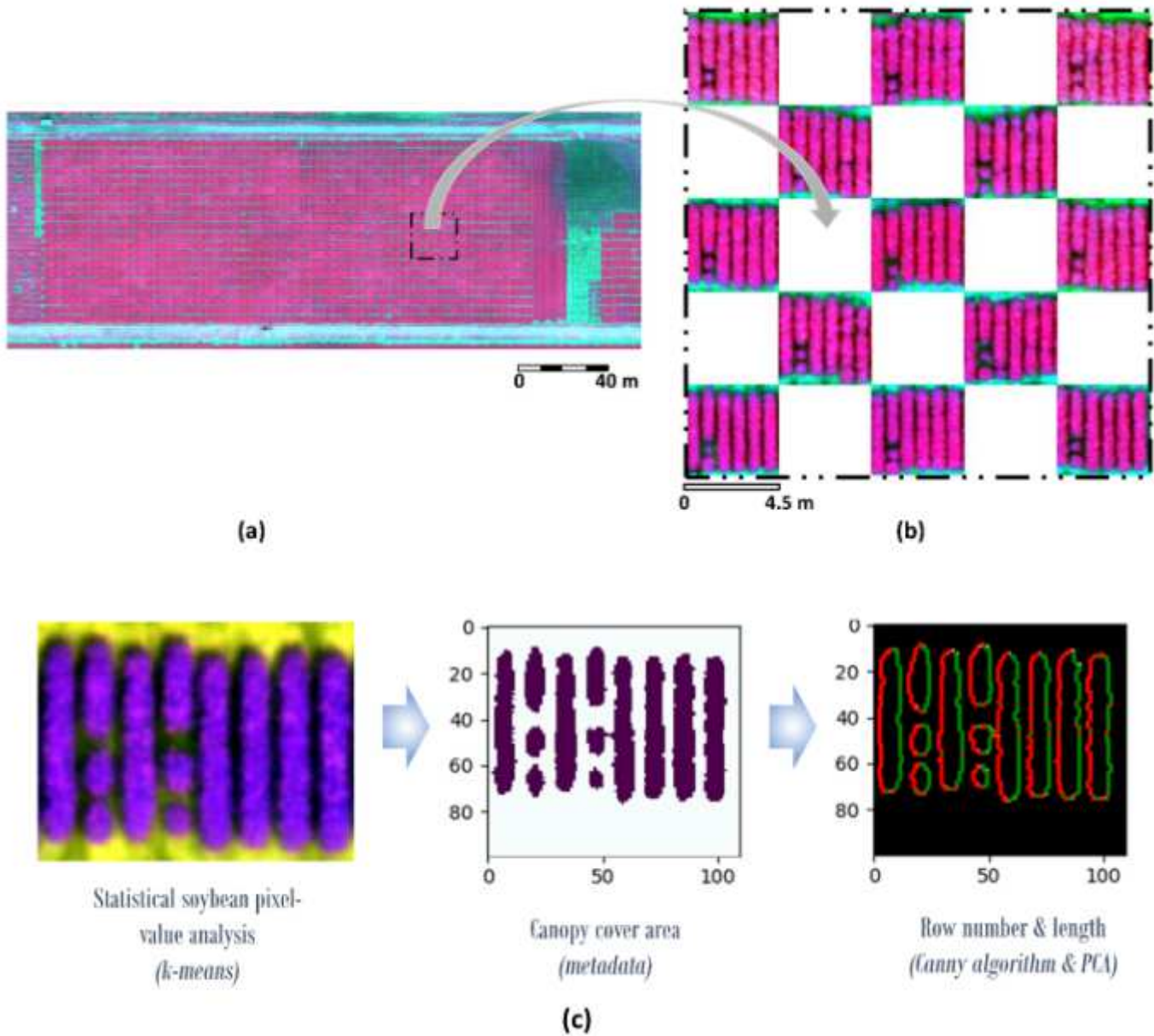
- 551 59. Gel YR, Gastwirth JL. A robust modification of the Jarque–Bera test of normality. *Economics Letters*. 2018; doi:  
552 <https://doi.org/10.1016/j.econlet.2007.05.022>.
- 553 60. Cawley GC, Talbot NL. On over-fitting in model selection and subsequent selection bias in performance evaluation.  
554 *J Mach Learn Res*. 2010; 11:2079–2107.
- 555 61. Elarab M, Ticlavilca AM, Torres-Rua AF, Maslova I, McKee M. Estimating chlorophyll with thermal and roadband  
556 multispectral high resolution imagery from an unmanned aerial system using relevance vector machines for precision  
557 agriculture. *Int. J. Appl. Earth Obs*. 2015; 43, 32–42.
- 558 62. Montgomery DC, Peck EA, Vining GG. *Introduction to linear regression analysis*, fifth ed., John Wiley & Sons,  
559 Inc., Hoboken, New Jersey, 2012.

# Figures



**Figure 1**

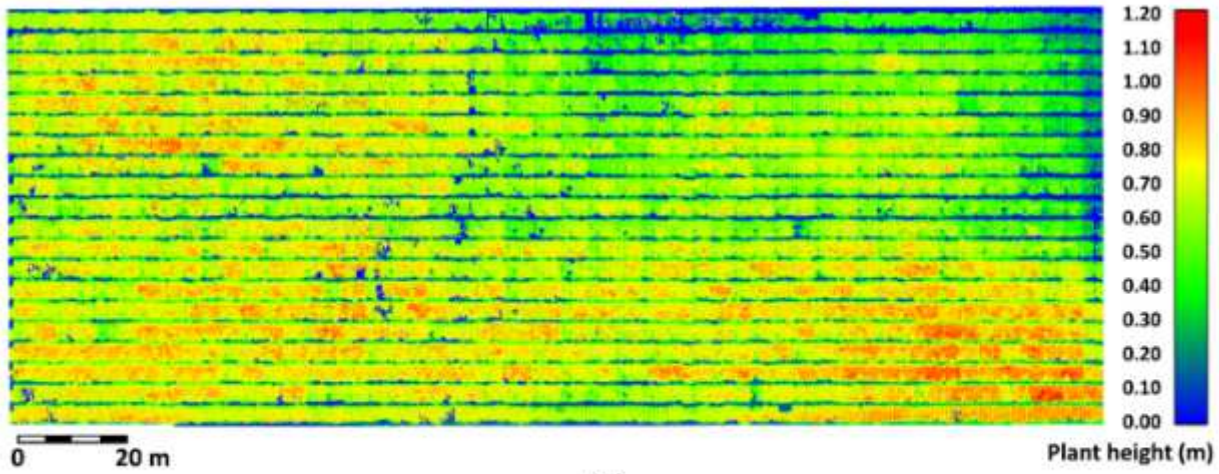
Workflow.



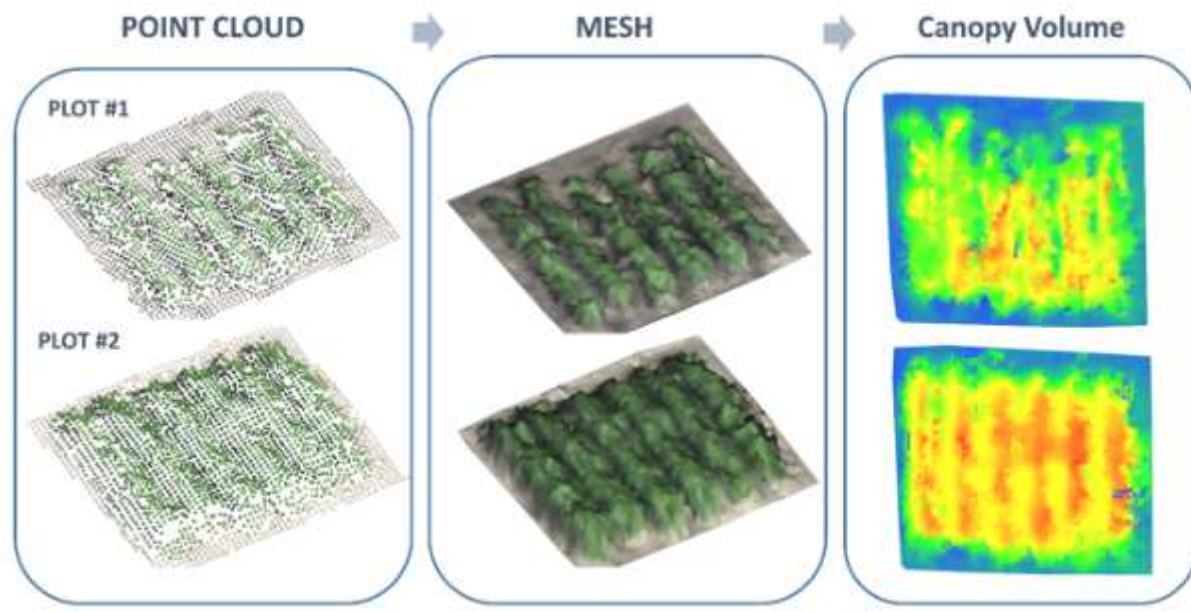
**Figure 2**

Colour composite MSI mapping over the study area on July 23rd 2018 (DAP 61) (DL (500047.3, 4480849.5); UR (500364.0, 4480968.0); EPSG 32616) (a), a detail of field-map based plot extraction (b) and Triple S software run over a random plot for July 23rd, 2018 (DAP 61) (c).





(a)



(b)

**Figure 3**

Deviation point cloud over Soybean from July 23rd (DAP 61) using June 7th (DAP 15) as reference in meters (a) and canopy volume calculation of two random plots from the same date (b).

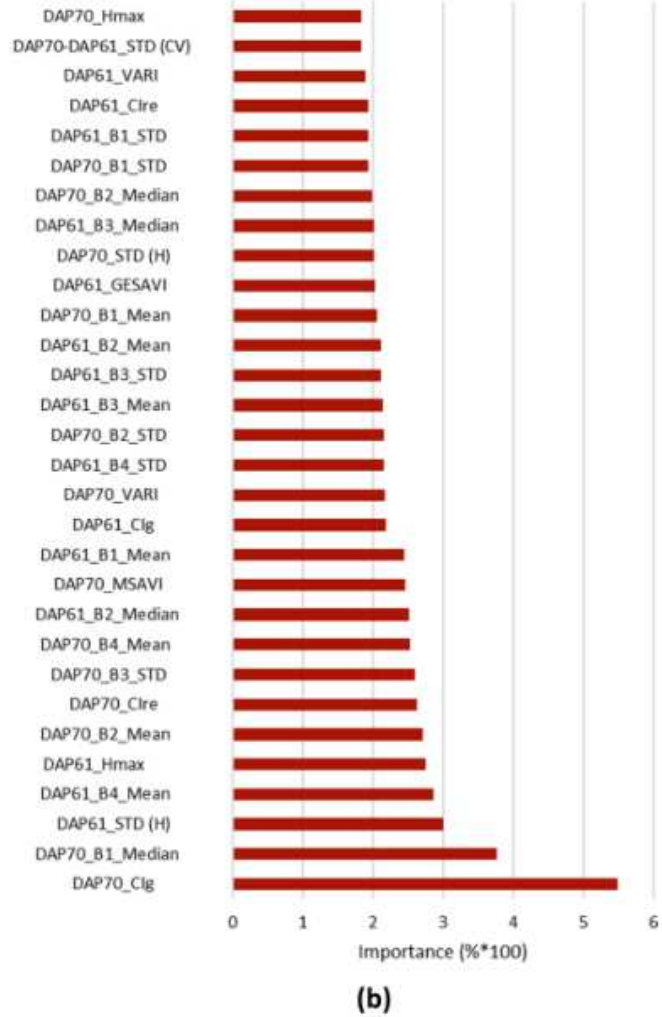
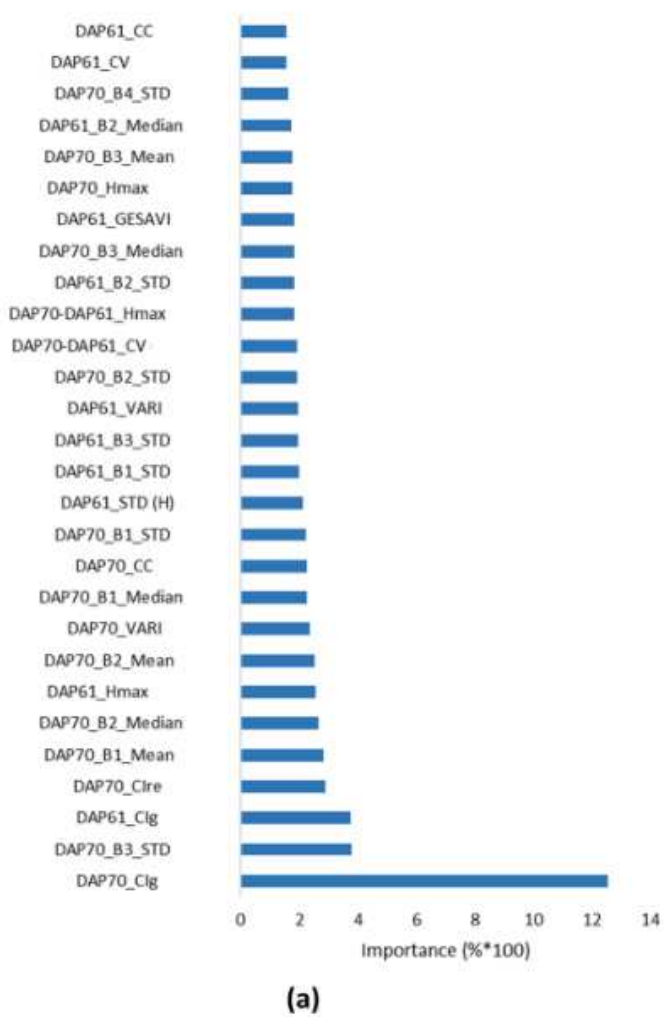
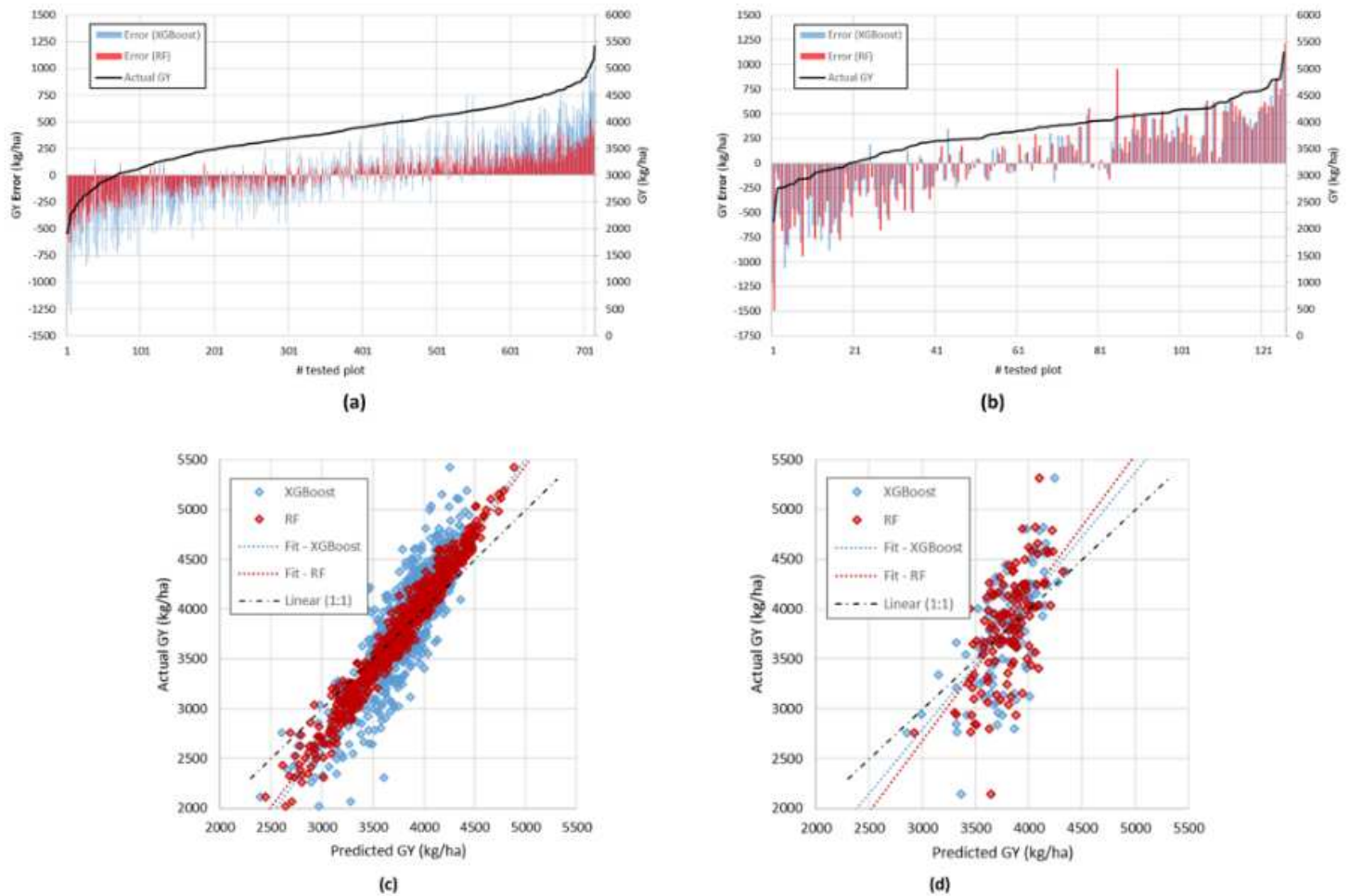


Figure 4

Features importance for more than 71% by RF (a) and XGBoost (b).



**Figure 5**

Prediction errors and actual Grain Yield (GY) sorted smallest to largest per plot along the training dataset (a) and testing dataset (b) (please note that the errors and GY are plotted in the primary and secondary axis, respectively); scatter plots of the measured against the predicted grain yield ( $\text{kg}\cdot\text{ha}^{-1}$ ) by RF and XGBoost from training (c) and testing dataset (d). In both cases is drawn the line corresponding to the robust linear fit at 95 % of confidence.

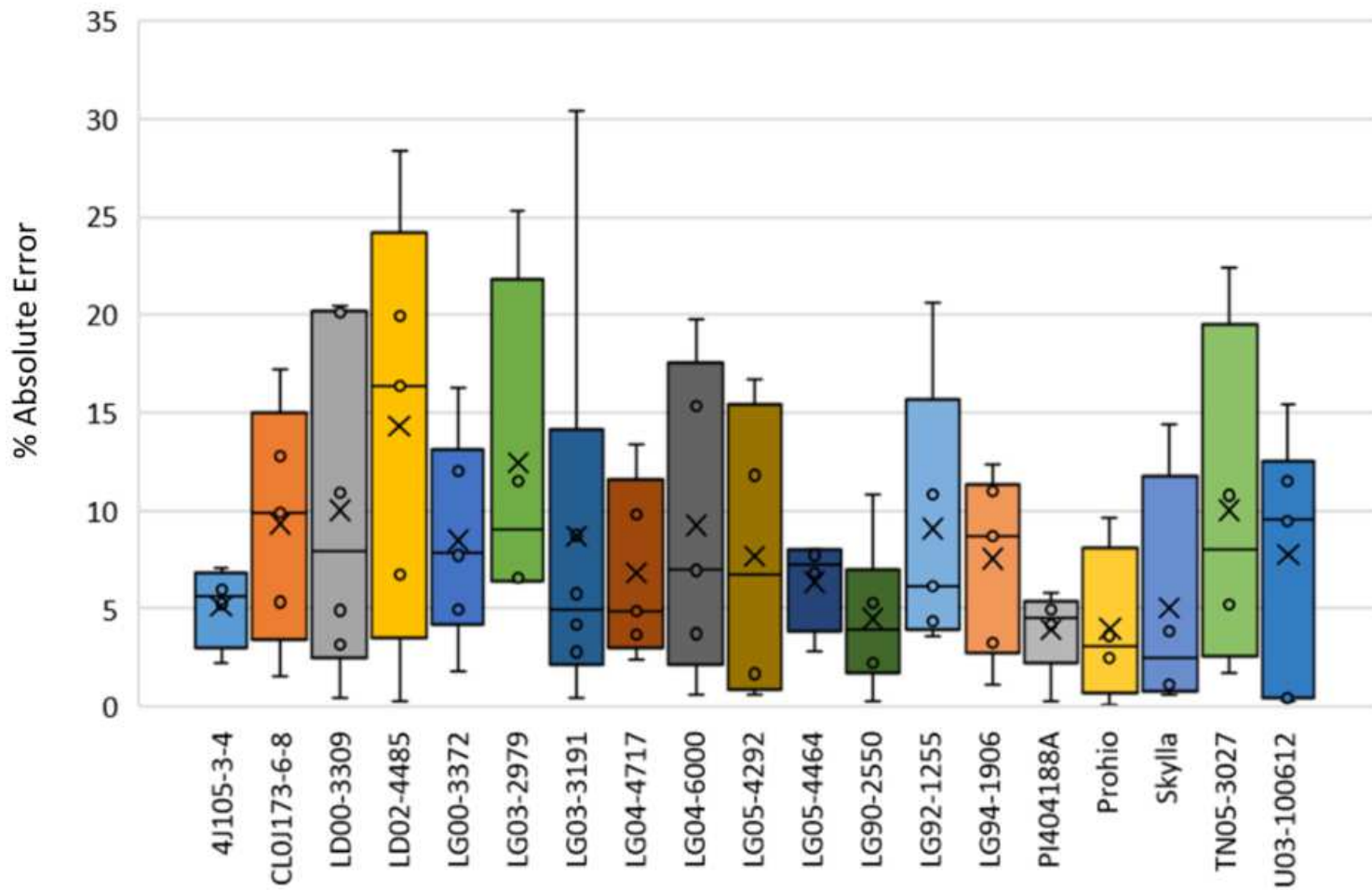


Figure 6

Errors from the testing dataset grouped by family