



DEPARTMENT OF LIFE SCIENCES

---

# **Broccoli Drought and Heat Complex Stress Detection and Shelf Life Prediction Based on Spectrometry and Machine Learning**

---

*AUTHOR:*

XIAOSHENG LUO

*SUPERVISOR:*

*CID:*

01627437

Dr. OLIVER WINDRAM

August 27, 2019

A thesis submitted in partial fulfilment of the requirements for the degree

of Master of Research at Imperial College London

Formatted in the style of Methods in Ecology and Evolution

Submitted for the MRes in Computational Methods in Ecology and Evolution

## ***Declaration***

The images of the broccoli head on the conveyor belt used to construct and training the neural network was originally provided by Nathan E. Barlow (Imperial College London), and then the optimized image video was collected by the author and the supervisor, Dr. Oliver Windram (Imperial College London). Besides, Dr. Oliver Windram was mainly responsible for setting up the direction of this project. Acquisition of experimental data, data cleaning, data analysis, method modification, model training, parameter tuning and writing were exclusively performed by the author himself.

1

## 2 Abstract

3 As a non-destructive, and high-efficiency technology, spectroscopy has developed  
4 rapidly in crop management for precision agriculture. Many studies on spectro-  
5 scopic detection of crop stress have made effective progress, including disease  
6 detection and water stress, etc. Various vegetation indices have also been es-  
7 tablished to indicate crops growth status. However, complex growing conditions  
8 often cause crops to be affected by multiple stresses, which in turn affect crops  
9 yield and quality. Few studies have been conducted on this complex stress, espe-  
10 cially the two physiological closely related abiotic stresses, heat and drought.

11

12 Here we took broccoli as our research subject, and collected the hyperspectral  
13 reflectance data of its leaves under heat and drought, as well as their combination  
14 by spectrophotometer. Then, we explored the spectral characteristics by various  
15 vegetation indices, and the results showed that a single vegetation index could  
16 hardly be significant in all treatments. Next, we trained four machine learning  
17 models, including Logistic regression, Support Vector Machine, Random forest  
18 and XGBoost to predict these complex stresses, and the highest AUC can reach  
19 0.9494 by logistic regression.

20

21 Besides, we have developed a visualization method that can dynamically and  
22 more intuitively see the significant differences of spectral characteristics between  
23 different stress treatments. The strategy is that first, it sets a appropriate wave-  
24 length width to reduce the redundancy of hyperspectral information, then ran-  
25 domly searches and accumulates the significant statistical analysis results with  
26 the width, and finally outputs the important band peaks dynamically. It can be  
27 helpful for spectral image band selection.

28

29 And additionally, we also try to combine spectroscopy and deep learning to build a  
30 system for predicting the shelf-life of broccoli heads. For now, it can track broccoli  
31 heads with ResNeXt to 97.2% accuracy and segment them with Unet to 99.2%  
32 accuracy, while the signal related to the shelf-life is still in progress.

33

34   **Keywords:** Broccoli; Machine Learning; Spectral feature; Heat stress; Drought  
35   stress

36   Word count: 5786.

# <sup>37</sup> **Contents**

<sup>38</sup> <b>1 INTRODUCTION</b>	<b>5</b>
<sup>39</sup> <b>2 MATERIALS AND METHODS</b>	<b>8</b>
<sup>40</sup> <b>2.1 Workflow</b> . . . . .	<b>8</b>
<sup>41</sup> <b>2.2 Data Collection</b> . . . . .	<b>9</b>
<sup>42</sup> <b>2.3 Vegetation Indices Calculations</b> . . . . .	<b>10</b>
<sup>43</sup> <b>2.4 Machine Learning</b> . . . . .	<b>11</b>
<sup>44</sup> <b>2.4.1 Model Fitting</b> . . . . .	<b>12</b>
<sup>45</sup> <b>2.4.2 Feature Engineering</b> . . . . .	<b>14</b>
<sup>46</sup> <b>2.5 Computer Vision</b> . . . . .	<b>14</b>
<sup>47</sup> <b>3 RESULTS</b>	<b>15</b>
<sup>48</sup> <b>3.1 Data</b> . . . . .	<b>15</b>
<sup>49</sup> <b>3.2 Vegetation Indices</b> . . . . .	<b>17</b>
<sup>50</sup> <b>3.3 Dimensionality Reduction</b> . . . . .	<b>18</b>
<sup>51</sup> <b>3.4 Feature Selection</b> . . . . .	<b>20</b>
<sup>52</sup> <b>3.5 Models</b> . . . . .	<b>21</b>
<sup>53</sup> <b>3.6 Shelf Life Prediction</b> . . . . .	<b>23</b>
<sup>54</sup> <b>4 DISCUSSION</b>	<b>25</b>
<sup>55</sup> <b>5 Limitations and Future Research</b>	<b>28</b>
<sup>56</sup> <b>6 ACKNOWLEDGEMENTS</b>	<b>28</b>
<sup>57</sup> <b>7 CODE &amp; DATA ACCESSIBILITY</b>	<b>29</b>
<sup>58</sup> <b>8 Supplementary Information</b>	<b>36</b>

59 **1 INTRODUCTION**

60 The environment in which crops growth is highly dynamic, crops are constantly  
61 exposed to various stresses, including abiotic stimuli such as humidity, light in-  
62 tensity and temperature, and biotic stimuli, like pathogens. Correspondingly, to  
63 cope with these stresses, plants have evolved a highly dynamic response mecha-  
64 nism, from gene expression regulation to changes in various secondary metabo-  
65 lites, which in turn alter their physiological characteristics, such as enzyme activ-  
66 ity, stomatal aperture, photosynthetic rate and transpiration rate (Carter, 1993).  
67 Rapid detection of stresses through these physiological changes of crops is partic-  
68 ularly important for obtaining high yield and high quality products. Traditional  
69 detection methods usually require experts to be able to detect the subtle color  
70 changes or a slight droop or curl of plants leaves, which indicate the stress. How-  
71 ever, it's generally subjective and time-consuming. In contrast, spectral detection  
72 and spectral imaging, as a non-destructive, accurate and efficient method, has  
73 developed rapidly both in research and practical application (Xue and Su, 2017).

74

75 The spectral reflectance characteristics and mechanism of leaves have been well  
76 summarized (Knippling, 1970, Gates et al., 1965). When the leaves receive solar  
77 radiation, only part of the incident energy is reflected and the rest is transmitted  
78 or absorbed for photosynthesis. The typical reflectance spectrum of a leaf is that  
79 in the visible region (0.4 – 0.7 $\mu$ m), the reflectance of leaves is generally very low,  
80 especially in red (around 0.63 – 0.70 $\mu$ m) and blue (around 0.45 – 0.52 $\mu$ m), this  
81 absorption characteristic is mainly caused by plant pigments. Specifically, the ab-  
82 sorption of red primarily contributed from chlorophyll, and the absorption of blue  
83 is also involved in carotenes and xanthophylls (Gates et al., 1965, Rabideau et al.,  
84 1946). Furthermore, the high reflection in near-infrared (0.7–1.3 $\mu$ m) is caused by  
85 the internal cellular structure (Mestre, 1935, Willstätter and Mieg, 1907). Leaf  
86 cuticular wax is transparent and hardly reflects solar radiation. Radiation can  
87 be transmitted through the epidermis, then dispersed and multiple reflected and  
88 refracted in the mesophyll cells and air cavity, where different refracted index  
89 between air (1.0) and hydrated cell walls (1.4) account for these effect (Sinclair  
90 et al., 1968).

91

92 Previous studies on the application of spectral techniques in plant stress detec-  
93 tion have made extensive progress, involving various crops and stress. In biotic  
94 stresses, plant disease detection is the most studied. Early in 1982, the diffuse  
95 reflectance spectra of potato tubers in the visible and near-infrared bands were  
96 measured and analyzed in an attempt to detect the presence of disease before  
97 its effects were visible ([Muir et al., 1982](#)). On top of that, spectral research also  
98 progress in the detection of various plant diseases, such as panicle blast, brown  
99 planthopper, the bacterial leaf in rice ([Kobayashi et al., 2001](#), [Prasannakumar  
et al., 2013](#), [Yang, 2010](#)) and yellow rust, powdery mildew in wheat ([Bravo et al.,  
2003](#), [Cao et al., 2013](#)). In abiotic stress, diverse research focus on water stress.  
100 Among them, many studies are based on canopy temperature based Crop Wa-  
101 ter Stress Index (CWSI) measured from infrared thermometry ([Alchanatis et al.,  
2010](#), [Aladenola and Madramootoo, 2014](#), [Bellvert et al., 2016](#)). The principle of  
102 CWSI theory is that transpiration cools the surface of leaves. When soil moisture  
103 in the root zone decreases, stomatal conductance and transpiration are weak-  
104 ened, and then leaf temperature increases. What makes this theory popular is  
105 its linear relationship between canopy temperature and air temperature and va-  
106 por pressure, as well as the development of empirical methods for quantifying  
107 crop water stress ([Idso et al., 1981](#)). However, though the canopy temperature  
108 is very useful for water stress detection, it still has some physiological concerns.  
109 In some plants, the diurnal fluctuation in stomatal conductance make the re-  
110 lationship unclear between canopy temperature and stress levels ([Zarco-Tejada  
et al., 2012](#)). Moreover, leaf temperature does not directly explain other physi-  
111 ological changes, such as photosynthesis pigments or non-stomatal reduction of  
112 photosynthesis under water stress ([Zarco-Tejada et al., 2013](#)). Therefore, vari-  
113 ous alternative vegetation indices (VIs) based on the visible and red edge spectral  
114 region are developed to capture water stress related signals ([Berni et al., 2009](#),  
115 [Zarco-Tejada et al., 2013](#), [Wang et al., 2015](#), [Rossini et al., 2013](#), [Panigada et al.,  
2014](#), [Dangwal et al., 2016](#)).

121

122 Although spectroscopic studies of biotic and abiotic stresses can achieve signifi-  
123 cant detection under different models, the stress they detect is often single, while

in practical production, crops tend to suffer from multiple complex stresses during their growth, such as heat and drought and their combinations, especially in the context of global warming. More importantly, the molecular and physiological mechanisms by which plants respond to heat and drought stress have been extensively studied and they show lots of relevance. Both of them can differentially affect the RNA stability, alter the enzyme activity and disrupt the steady-state of metabolic flux, which in most of cases can cause a common response, oxidative damage (Kollist et al., 2018, Suzuki et al., 2012, Mittler et al., 2012, McClung and Davis, 2010). Moreover, photosynthesis is an important physiological phenomenon affected by drought and heat stress. Drought can lead to stomatal closure and reduces CO<sub>2</sub> uptake which makes plants more susceptible to photo damage, also it can induce negative changes in photosynthetic pigments, either increase or decrease chlorophyll content (Lawlor and Cornic, 2002, Anjum et al., 2011, Din et al., 2011). Similarly, exposure to high temperature can also result in a reduction in chlorophyll biosynthesis, thereby disturbing the photosynthetic pigment components (Camejo et al., 2006). Although many physiological links between plant heat stress and water stress, and it is of great meaningful to detect them in the practical production, little research have been conducted on the relationship between these two stress spectral characteristics and their detection. So it's a great interest and also useful to see how well we can detect the heat stress and drought and their combination from the crops by data mining their spectral characteristics changes.

Alongside this, methodologically, many previous studies on crops stresses detection used statistical discriminant model, especially various VIs. VIs are quite simple and effective algorithms for quantitative and qualitative evaluation of vegetation cover, growth dynamics, and stress levels. But due to different spectral combinations, instrumentation, platforms, and resolutions used, it's hard to have a unified mathematical expression that defines all VIs, customized algorithms needed to developed and tested against specific application requirements (Xue and Su, 2017). In this way, the prediction is often unsatisfactory and the generalization ability somewhat insufficient. Compared with traditional statistical discriminant model, machine learning methods, which developed rapidly in re-

157 cent years, can generally improve the speed and accuracy of prediction. The main  
158 difference between machine learning and statistics lies in their purpose. Statistical  
159 models are more designed to infer the relationship between variables, while  
160 machine learning models are intended to make the most accurate prediction pos-  
161 sible.

162

163 Overall, here we explore the ability to use hyperspectral techniques to detect  
164 and differentiate more complex stresses of crops, that is, the combinations of  
165 drought and heat, which are highly correlated with each other. Firstly, we cal-  
166 culated some widely known VIs for statistical testing in anticipation of obtaining  
167 simple and effective stress-related signals. Then, in order to approach the up-  
168 per limit of accuracy for detecting different stresses, we apply machine learning  
169 strategy, using linear classifiers, such as logistic regression (LG), linear support  
170 vector machine (SVM) and tree models, like random forests (RF) and XGBoost for  
171 training a robust classifier. On top of this, to increase the interpretability of the  
172 model, great efforts had been made in feature engineering. Specifically, we first  
173 extract features through dimensionality reduction of Principal Component Anal-  
174 ysis (PCA) and Linear Discriminant Analysis (LDA), and explore the importance  
175 of features under potential dimensions. Then statistical filter method and model-  
176 based embedded method are applied for feature selection. In particular, during  
177 these process, in order to better visualize the statistical differences of hyperspec-  
178 tral features under different stress comparisons dynamically, meanwhile be able  
179 to adjust the wavelength width to reduce the redundancy of hyperspectral infor-  
180 mation and provide a reference for specific band selection, a simple visualization  
181 search tool has been developed. Finally, the features obtained by these methods  
182 are combined to conduct model training, and the prediction effects of drought,  
183 heat stress and their combination stress were analyzed by model comparison.

184

185 Additionally, in order to optimize the customization of supermarket broccoli prod-  
186 ucts' shelf-life, a computer vision strategy based on deep learning and spectral  
187 detection technology is being developed. Now it can track and segment broc-  
188 coli heads on the conveyor belt through convolution neural network, but futher  
189 spectral signals related to shelf-life still need to be analyzed.

190 **2 MATERIALS AND METHODS**

191 **2.1 Workflow**

192 The project mainly includes two parts (Figure 1). The laboratory part is to grow  
193 broccolis under control conditions and then perform individual and combined  
194 stress treatments, collect spectral images and leaf reflectance spectrum data to  
195 explore the signals that can effectively distinguish among them and construct a  
196 robust machine learning classifier. The application part is to construct a detection  
197 system which can predict the shelf life of broccoli on the conveyor belt through  
198 computer vision methods and spectral images under specific bandwidth, which is  
199 selected based on the results obtained in the laboratory.

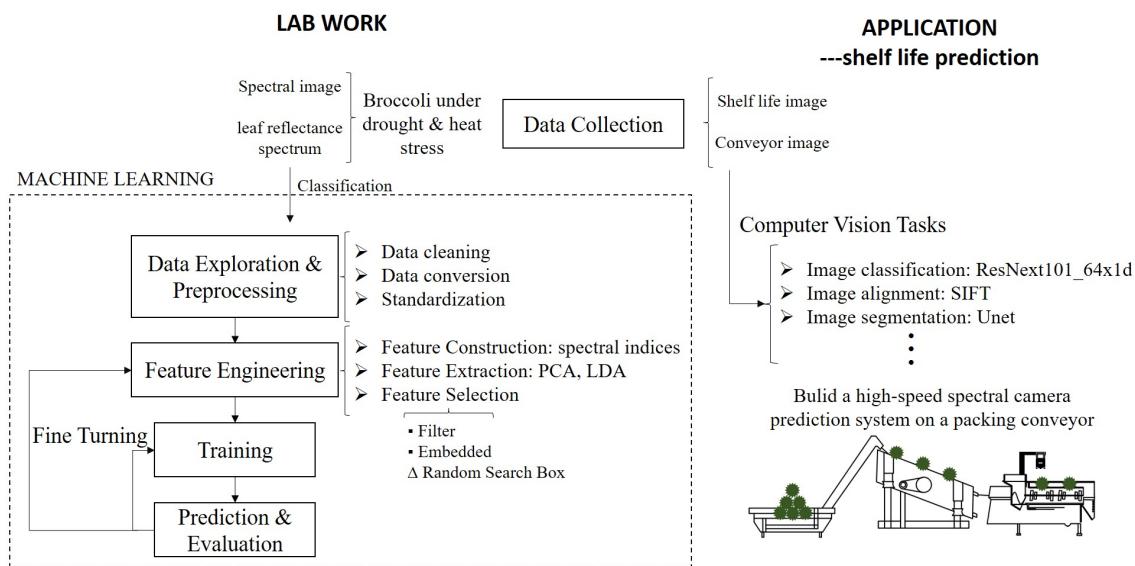


Figure 1: Project workflow

200 **2.2 Data Collection**

201 Broccolis were grown in Control Environment room and greenhouse. The normal  
202 growth temperature is controlled at 23 °C and the humidity is controlled at 60%,  
203 with long daylight (16 hours illumination, 8 hours darkness) treatment, and wa-  
204 ter is poured every 3 days from the tray to the soil. The whole growth cycle of  
205 broccoli takes about three months, during which it needs to be transferred to a  
206 suitable pot according to the size of broccoli.

207

208 During stress treatment, broccolis were randomly divided into four groups with

209 8 individuals in each group. They were treated under control, heat stress ( $27^{\circ}\text{C}$ ),  
 210 drought stress (without watering for three days) and the combination of heat and  
 211 drought stress. Leaf reflectance spectroscopy data was collected by the Ocean  
 212 Optics FLAME-S-XR1 spectrophotometer in a completely dark room, while the  
 213 spectral image were collected by Ximea cameras with a specific bandpass filters  
 214 (FEL0800, FBH650-40) (Table 1) under the illumination of the corresponding  
 215 wavelength of the LED lamp. Four days of data were collected until the leaves  
 216 show a distinct dehydration drooping phenotype. In the open-air greenhouse,  
 217 spectral images are collected with the same apparatus in a grow tent, in order to  
 218 eliminate the effect of unstable solar radiation.

**Table 1: Apparatus used in the experiment.**

Item	Description
Camera	Ximea 1.3 MP NIR Enhanced Camera MQ013RG-ON
Machine vision lens	MVL12M23-12 mm EFL, f/1.4, for 2/3" " C-Mount Format Cameras, with Lock.
Band pass filter	FEL0800: Ø25.0 mm Premium Longpass Filter, Cut-On Wavelength: 800 nm. FBH520: Ø25.0 mm Hard-Coated Bandpass Filters, Blocking Regions (OD >5): 200 - 485 nm, 556 - 1200 nm. FBH650: Ø25.0 mm Hard-Coated Bandpass Filters, Blocking Regions (OD >5), 200 - 611 nm, 690 - 1200 nm. FBH850: Ø25.0 mm Hard-Coated Bandpass Filters, Blocking Regions (OD >5), 200 - 805 nm, 896 - 1200 nm.
LED controller	Intelligent LED Solutions 12-Channel Light Controller
LED	12 Die LED array Full Spectrum 360-955nm

219 Broccoli heads for shelf-life prediction come from POLLYBELL FARMS LTD. They  
 220 are divided into two groups, 18 in each, one of which is stored in cold storage for  
 221 two weeks, and the other is harvested freshly. They were placed naturally at room  
 222 temperature and spectral image data were collected every day through the cam-  
 223 eras with bandpass filter (FBH520-40, FBH650-40,FBH850-40) and hyperspectral  
 224 data was collected by spectrophotometer as well until they decay significantly.

## 225 2.3 Vegetation Indices Calculations

226 The names, abbreviations, calculation formulas and citations of the various veg-  
 227 etation indices used in this study are as follows, mainly includes commonly used  
 228 remote sensing indices, chlorophyll-related indices, and indices related to water  
 229 stress indications.

**Table 2: Various vegetation indices**

Name	Abbrev.	Equation	References
Ratio vegetation index	RV	$R_n/R_r$	(Jordan, 1969) (Pearson and Miller, 1972)
Normalized difference vegetation index	NDVI	$(R_{800} - R_{680}) / (R_{800} + R_{680})$	(Rouse Jr et al., 1974) (Tucker, 1979)
Enhanced vegetation index	EVI	$2.5(R_n - R_r)(R_n + 6 \cdot R_r - 7.5 \cdot R_b + 1)$	(Huete et al., 2002)
Chlorophyll vegetation index	CVI	$R_n \cdot R_r / R_g^2$	(Vincini et al., 2008)
Chlorophyll index - green	CI-G	$R_n / R_g - 1$	(Gitelson et al., 2003)
Chlorophyll index - red edge	CI-RE	$R_n / R_{re} - 1$	(Gitelson et al., 2003)
Photochemical reflectance index	PRI	$(R_{531} - R_{570}) / (R_{531} + R_{570})$	(Gamon et al., 1992)
Water index	WI	$R_{900} / R_{970}$	(Zarco-Tejada et al., 2003)
Structure independent pigment index	SPI	$(R_{800} - R_{445}) / (R_{800} + R_{680})$	(Peñuelas and Inoue, 1999)

$R_\lambda$  is the reflectance at wavelength  $\lambda$ ;  $n, re, b, g$  and  $r$  represent NIR (760–900 nm), RE (700–730 nm), blue (450–520 nm), green (520–600 nm) and red (630–690 nm) respectively.

## 230 2.4 Machine Learning

231 Machine learning generally includes several steps in practical operation, such as  
 232 data collection and preprocessing, model selection, training, evaluation and re-  
 233 peatedly fine-tuning until a good prediction effect is achieved (Figure 1). In the  
 234 data preprocessing stage, Z-score standardization is applied to simplify the cal-  
 235 culation and the categorical data is one-hot encoded. In the strategy of training  
 236 algorithms, firstly, LG, SVM, RF and XGBoost algorithm were trained to fit the  
 237 raw data and obtained the baseline score, and then the performance of the mod-

238 els were optimized by feature engineering and parameters adjustment. Most of  
239 the code used in this process is based on the API provided by scikit-learn ([Pe-](#)  
240 [dregosa et al., 2011](#)).

#### 241 2.4.1 Model Fitting

242 Logistic regression: the binomial logistic regression model is a classification model,  
243 which is represented by the conditional probability distribution  $P(Y|X)$ , in the  
244 form of parameterized logistic distribution. Here, the value of  $X$  is a real number,  
245 and the random variable  $Y$  takes a value of 0 or 1, then we estimate the model  
246 parameters by supervised learning. The binomial logistic regression model is the  
247 conditional probability distribution as follows:

$$P(Y = 1|x) = \frac{\exp(w \cdot x)}{1 + \exp(w \cdot x)}$$

$$P(Y = 0|x) = \frac{1}{1 + \exp(w \cdot x)}$$

248 Here,  $x$  is the input vector,  $w$  is the weight vector and  $Y$  is the output vec-  
249 tor,  $Y \in \{0, 1\}$ ,  $x = (x^{(1)}, x^{(2)}, \dots, x^{(n)}, 1)^T$ ,  $w = (w^{(1)}, w^{(2)}, \dots, w^{(n)}, b)^T$ . By comparing  
250 the probability of  $P(Y = 1|x)$  and  $P(Y = 0|x)$  can finally determine the category.  
251 The cost function of logistic regression can be derived by the method of maxi-  
252 mum likelihood estimation, which is known as the average of cross-entropy loss.  
253 Meanwhile, in order to prevent over-fitting, the L1 or L2 regularization terms was  
254 added during the optimization process.

256

257 SVM: the main idea of SVM is to find the decision boundary with the largest  
258 classification interval between two different categories, which means that a hy-  
259 perplane separating classes in the feature space is defined by the principle of  
260 maximum margin between the closest different data points, also known as sup-  
261 port vectors. For simple linear separability problems, it can be described as an  
262 optimization problem by mathematical formulas as follows:

$$\max_{w,b} \left[ \min_{x_i} \frac{y_i (w \cdot x_i + b)}{\|w\|} \right]$$

263 The minimized item represents the distance from the support vectors to the deci-  
264 sion boundary with sign, known as geometry margin. After derivation and trans-  
265 formation, and allow the SVM to ignore some noise, that is, allow some data  
266 points' functional margin less than 1, a slack variable ( $\xi_i \geq 0$ ) is introduced to al-  
267 low some wrong classification. correspondingly, a penalty term is needed to add  
268 to the objective function to limit the slack variable, and here is the basic linear  
269 separable SVM:

$$\begin{aligned} & \min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t. } & y_i (w^T x_i + b) \geq 1 - \xi_i \quad (i = 1, 2, \dots, m) \\ & \xi_i \geq 0 \quad (i = 1, 2, \dots, m) \end{aligned}$$

270 Finally, the problem can be resolved by Lagrange Duality and SMO algorithm  
271 (Platt, 1998). In addition, kernel mapping has also been tried to verify the model's  
272 performance under nonlinearly separable conditions.

273

274 Ensemble methods: Random Forests and XGBoost (Chen and Guestrin, 2016),  
275 both are based on decision tree model, they use the strategy of bagging and boost-  
276 ing respectively, which can help to prevent high variance and high bias. Random  
277 forest mainly consists of two stochastic processes, random sampling of samples  
278 and features to construct many decision trees that are independent of each other.  
279 The final prediction results are summarized by voting strategy. As for XGBoost,  
280 it's an algorithm developed from gradient boosted decision trees and designed  
281 for speed and performance. It's widely used in many competitions and achieved  
282 good grades.

283

284 Metric: ROC curve (receiver operating characteristic curve) is a graph showing  
285 the performance of a classification model at all classification thresholds. The curve  
286 plots the points of (True positive Rate, False Positive Rate) at different classifica-  
287 tion thresholds. Lowering the classification threshold classifies more items as pos-  
288 itive, thus increasing both False Positives and True Positives. Area under the ROC  
289 Curve (AUC) measures the entire two-dimensional area underneath the entire  
290 ROC curve. AUC is desirable for its scale-invariant and classification-threshold-  
291 invariant.

292

293 When multi-classification is performed on logistic regression and SVM, "one to the  
294 rest" strategy is applied. All the models are validated by 6-fold cross-validation,  
295 and ROC and AUC is used as the metric for model evaluation.

296 **2.4.2 Feature Engineering**

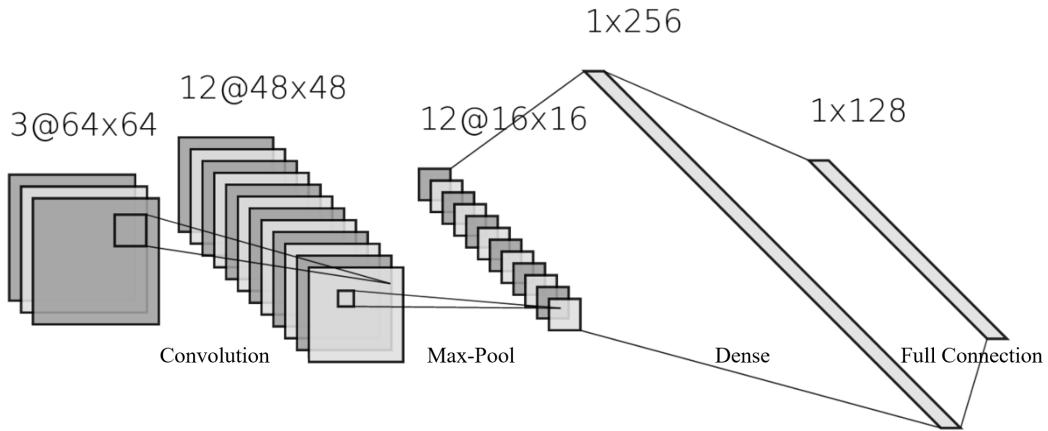
297 Generally, data and features determine the upper bound of machine learning,  
298 whereas models and algorithms only approximate this upper bound. The pur-  
299 pose of feature engineering is to extract effective features and remove redundant  
300 features, also it can make features machine readable and contextually relevant.  
301 It basically includes feature extraction, feature construction and feature selection  
302 (Figure 1). Separately, feature extraction mainly uses dimension reduction meth-  
303 ods such as PCA and LDA. Feature construction is to construct new features based  
304 on previous expertise. Feature selection methods can be roughly divided into  
305 three types:

- 306     • Filter: scoring each feature according to divergence, correlation, etc., and  
307       then set a threshold for selection feature.
- 308     • Embedded: use some machine learning algorithms and models to train and  
309       get the coefficients of each feature, select features according to the coef-  
310       ficient, kind of similar to the filter method, but models are trained to de-  
311       termine the pros and cons of features. Specifically, multi-method ensemble  
312       selection ([Feilhauer et al., 2015](#)) was modified from the regression problem  
313       and later adapted to the classification problem.
- 314     • Wrapper: recursive elimination feature method, due to the high computa-  
315       tional complexity and the long execution time of the algorithm, it is not  
316       adopted here.

317 In addition, in order to find a suitable bandpass filter for the camera, a search box  
318 with specific bandwidth was used to repeatedly and randomly select features, and  
319 the importance of features is sorted by simple ANOVA and TukeyHSD significance  
320 test. Finally, the graph is plotted by accumulating significant ( $p < 0.05$ ) bandwidth  
321 features.

## 322 2.5 Computer Vision

323 The project involves computer vision tasks such as image classification, segmen-  
324 tation, and image alignment. Specifically, image alignment was performed by  
325 classic Scale-invariant feature transform (SIFT) ([Lowe et al., 1999](#)), Image clas-  
326 sification and segmentation are mainly accomplished by the transfer learning of  
327 convolution neural network (Figure 2).



**Figure 2: Structure of a simple convolution neural network**

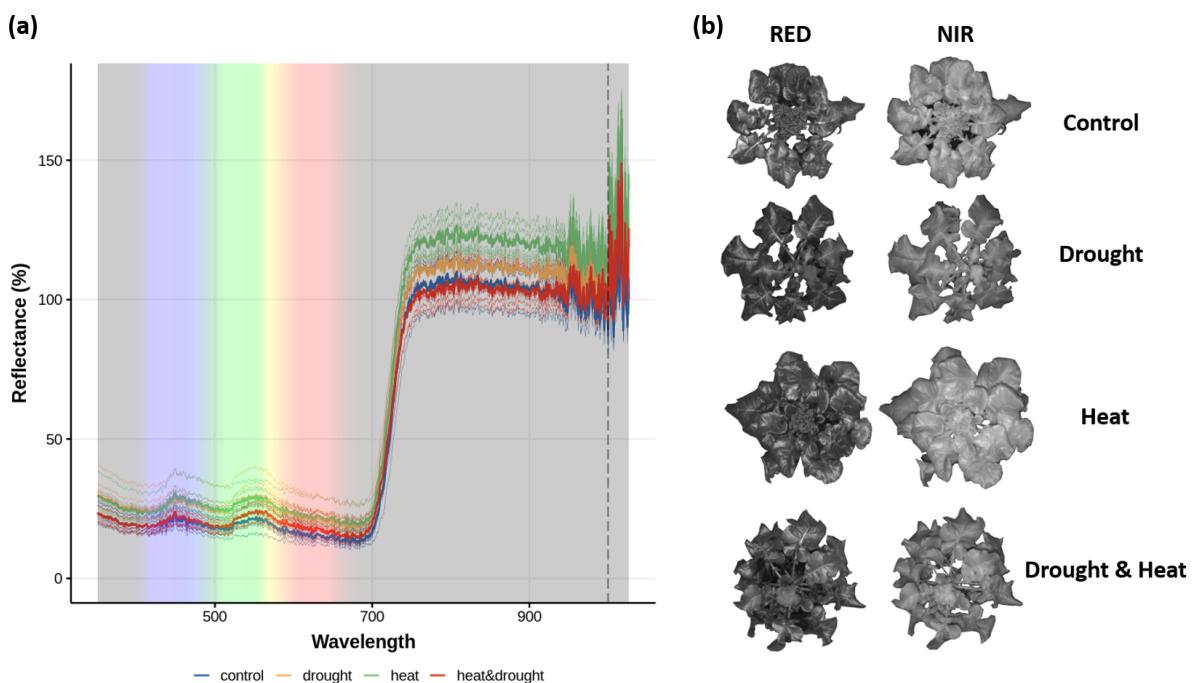
An image was taken as input (for example, a RGB image, normally three channels), then through the calculation with the multiple kernels' parameters and activation functions (usually ReLu) in convolution layer, and the downsampling process in pooling layer, can achieve the purpose of weight sharing and parameter reduction. Finally, the results are expanded and classified by the fully connected layer and the softmax function. In the figure, the number in front of @ is the number of channels, and the number followed by is the height and width of pixels.

328 More specifically, Image classification was implemented by ResNeXt ([Xie et al.,](#)  
329 [2016](#)), developed by UC San Diego and Facebook AI Research, while Image seg-  
330 mentation was implemented by Unet ([Ronneberger et al., 2015](#)).The training  
331 was conducted on GPU P1000, the optimizer for the neural network is Adam  
332 ([Kingma and Ba, 2014](#)), the cost function is cross-entropy, and cyclical learning  
333 rates ([Smith, 2017](#)) was used, Most of the code is based on the API provided by  
334 Pytorch, fastai library ([Howard et al., 2018](#)), and opencv.

### 3 RESULTS

#### 3.1 Data

To study the spectral characteristics of drought stress, heat stress, and their combination, Leaf reflectance non-image hyperspectral data (Figure 3a) and spectral images of two channels (Figure 3b) were collected. We select the data of the day when the broccoli just appeared phenotype under the stress to ensure the treatment effect and maximize model performance. Under the control and heat conditions, the broccolis have no obvious phenotype, while under drought and combined stress treatment, the broccoli leaves are a little drooping due to dehydration, and the combined stress is slightly more obvious than the drought (Figure 3b).



**Figure 3: Spectral data of broccoli under heat and drought stress**

(a) is the hyperspectral data of broccoli leaves detected by spectrophotometer in dark environment, the vertical axis represent their relative reflectivity. The thin line is the average of 80 hyperspectral scans per sample, and the thick line is the average of all samples in different treatments. The right side of the dashed line was discarded is subsequent processing due to abnormal signal fluctuation. (b) is the spectral images taken by the camera with red bandpass filter ( $CWL = 650$  nm,  $FWHM = 40$  nm) and near infrared bandpass filter ( $> 800$  nm), under the illumination of the corresponding band of LED.

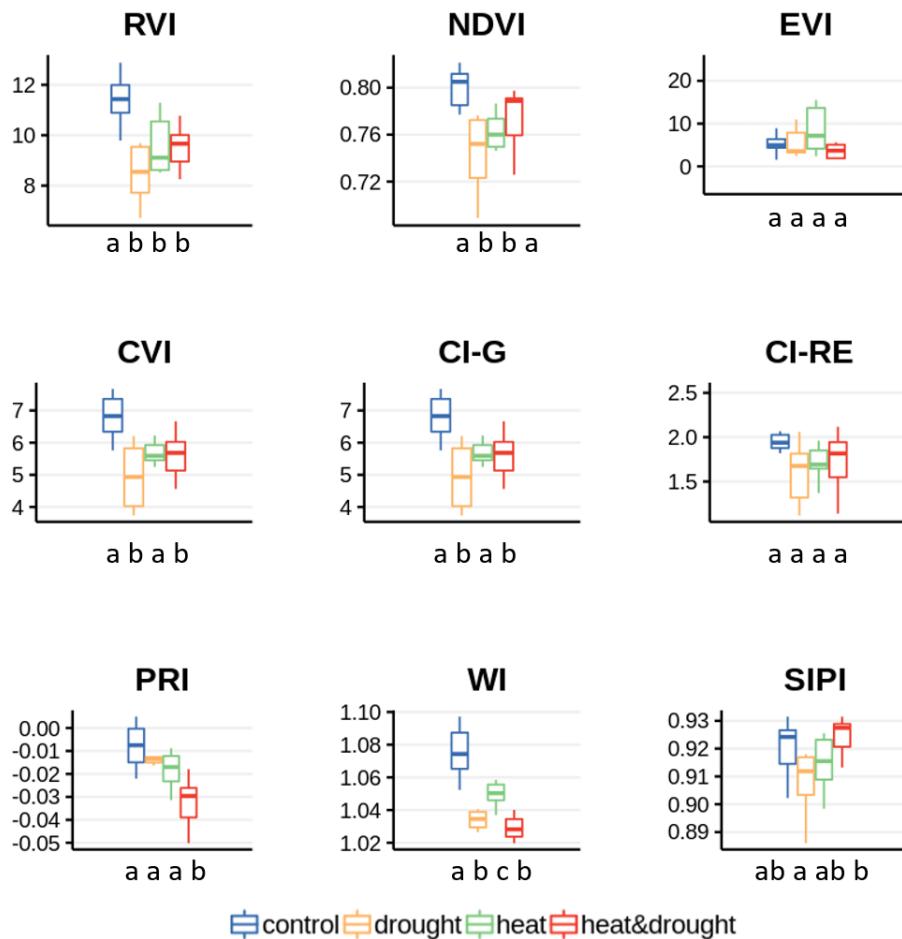
For the hyperspectral data of the leaves (Figure 3a), we scanned the leaves of

28 samples (7 samples per treatment), collected data every 10 scans, and finally got 80 data points per sample. The plots between treatments overlap to each other, so there is no clustering to get a simple discriminant pattern. However, the overall trend of the broccoli leaves reflectance spectrum can still be clearly seen. Small peaks at the blue and green-yellow junctions in the visible region, and well-known small valley in the red and strong reflection rate shifting in the near-infrared. These are mainly caused by the pigment and cell structure of broccoli (Gates et al., 1965). It's also worth pointing out that the average reflectance of the heat treatment appears to be generally higher than other treatments in the near-infrared.

## 3.2 Vegetation Indices

For the sake of looking for a simple and effective algorithm to distinguish four different treatments, several common spectral indices were calculated under different stress (Figure 4). Specifically, basic vegetation index RVI, is based on the phenomenon that leaves absorb relatively more red than infrared light. It's widely used for green biomass estimations (Jordan, 1969). NDVI is the most widely used VI to characterize canopy growth and vigor (Karnieli et al., 2010). EVI was introduced to correct soil and atmospheric effects on NDVI (Xue and Su, 2017). As for CVI, CI-G and CI-RE, they are to some extent related to chlorophyll content (Gitelson et al., 2003, Vincini et al., 2008). and for PRI, WI and SIFI have been proven to be effective of water status in some species (Ihuoma and Madramootoo, 2017, Katsoulas et al., 2016).

The results show that significant differences between every two treatments can't be simply obtained from a single index. The results with significant differences also show different forms of band feature calculation methods, which is difficult to establish a unified pattern, suggesting that more complex models are needed.

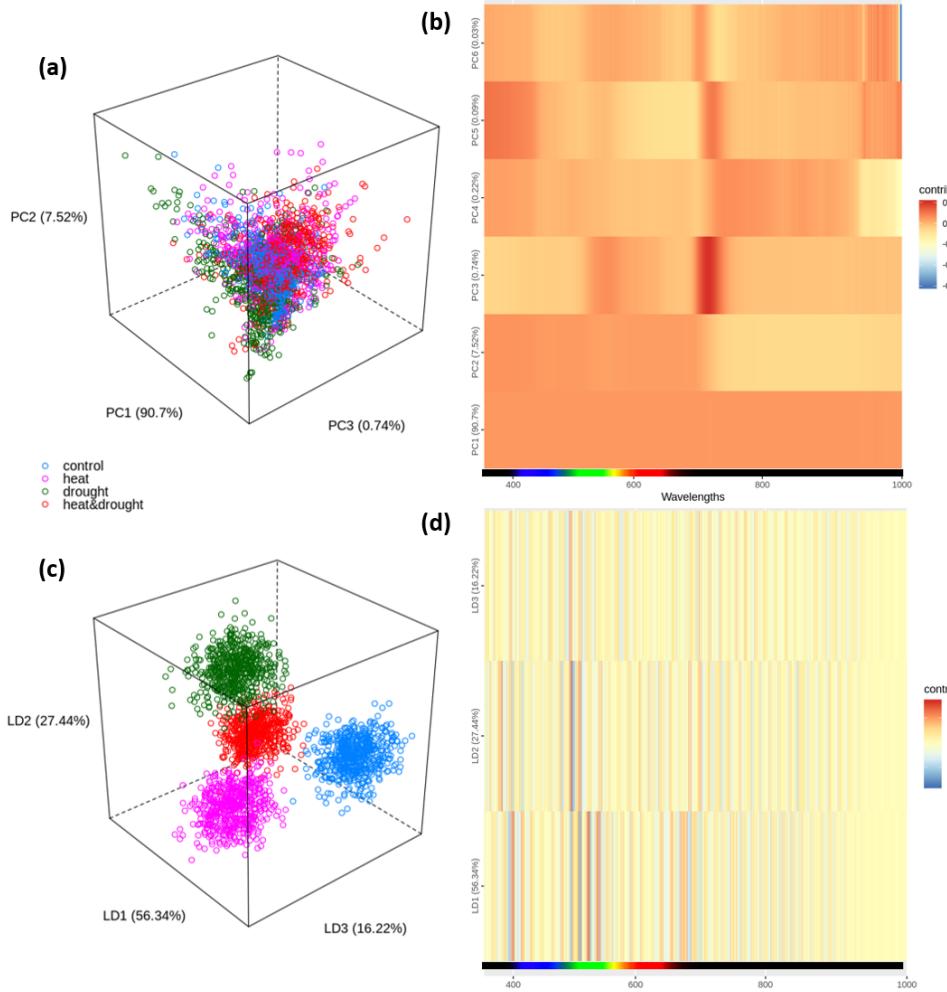


**Figure 4: Various vegetation indices under different stresses**

Seven samples per treatment, the equation for calculating vegetation index refers to Tabel 2. Normality test was implemented by shapiro test, and homogeneity of variance test was implemented by Barlett test before using ANOVA and Tukey's HSD for post-hoc analysis. Different letters under the x-axis represent significant differences ( $p < 0.05$ ).

### 3.3 Dimensionality Reduction

There may be multi-collinearity between hyperspectral features, that is variables may be correlated. Meanwhile, too many variables may hinder the pattern for model fitting, and it may also involve a lot of redundant information. Therefore, dimensionality reduction was used to reduce variables, speed up computation and extract effective information hidden in the data.



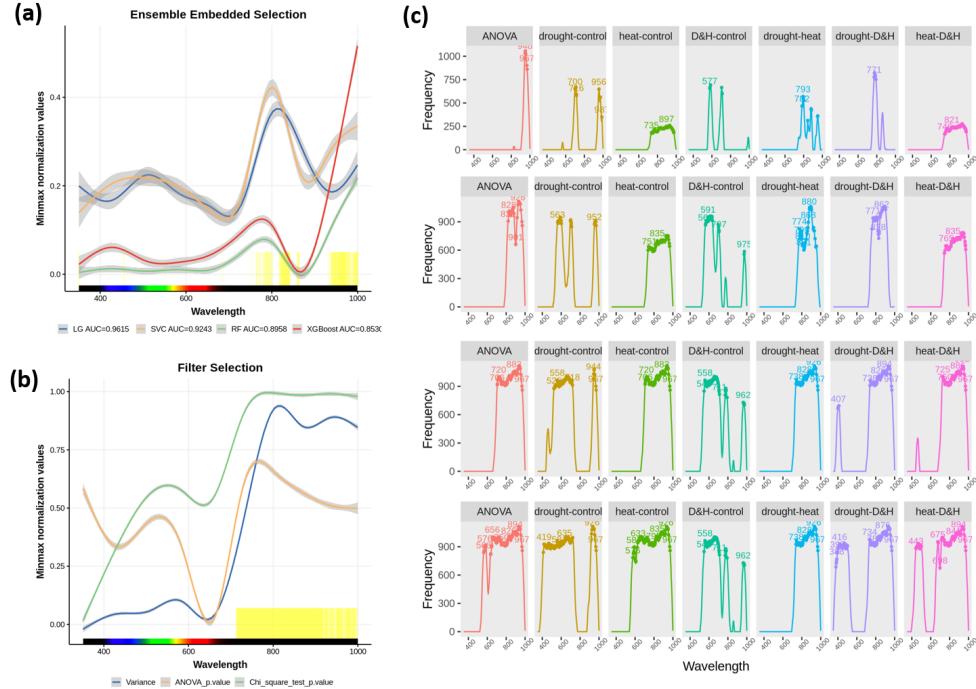
**Figure 5: Feature dimensionality reduction by PCA and LDA**

(a) and (b) are the distribution of the samples in the first three principal component dimensions, in which the values in the coordinate axis are the explanatory rates of the overall differences; Heat maps (c) and (d) indicate the correlation between each wavelength and each component.

The results of unsupervised dimensionality reduction PCA show that, when the variables are mapped to the linear-independent direction of maximum variance, clustering between different treatment is not effective (Figure 5a). Surprisingly, PC1 explains 90.7% of the variance, and the contribution of each wavelength seems to contribute equally to it, possibly due to the systematic errors. In PC2, the visible light region contributes a larger variance, and there are two specific narrow bands in PC3 that are positively correlated with it (Figure 5b). On the other hand, LDA, the supervised dimension reduction method, can completely separate the different treatments (Figure 5c). In these components, the green (around 520nm) and red edge (around 680nm) wavelength may be important to separate each stress treatment in the principal components (Figure 5d).

### 391 3.4 Feature Selection

392 In order to remove irrelevant features, reduce over-fitting in the process of model  
 393 training, and find explanatory wavelengths that can effectively distinguish different  
 394 treatments. We explored several feature selection methods as followed.



**Figure 6: Feature Selection**

For visualization purposes, all coefficients are min-max normalized. In (a) and (b), the gray shadows show the confidence intervals, and the yellow vertical lines indicate the important features upon the setting threshold (Filter: the coincidence of the top 50% importance features of each filter. Embedded: AUC weighted average of coefficients of each model). (c) From top to bottom shows the significant difference of wavelengths between different treatments dynamically, the earlier the peak appears, the more significant it is, the number on the graph represents where the peak is.

395 Filter and Embedded methods are commonly used feature selection methods in  
 396 machine learning. The results of the Filter method (Figure 6b) show that the  
 397 variance, chi-square test and the ANOVA of each wavelength are coincident upon  
 398 threshold in the near-infrared region. As for the modified ensemble method ([Feilhauer et al., 2015](#)), the correlation coefficient or feature importance obtained by  
 399 fitting the LG, SVM, RF and XGBoost are weighted with their AUC scores, then  
 400 set their average plus standard deviation as threshold. And the features upon the  
 401 threshold are mainly include two near-infrared fragments and other narrow frag-  
 402 ments in visible light (Figure 6a).

404

405 However, although we have acquired some features through these two methods,  
406 it's still not convincing enough for us to explain the relationship between features  
407 and the treatments. For the embedded method, the generalization ability could  
408 be constrained by the algorithm itself. As for Filter method, it mainly focuses  
409 on the correlation between individual features and treatments. The advantage  
410 of this method is that it is efficient in computing and robust to over-fitting prob-  
411 lems, but it tends to choose redundant features because they do not consider the  
412 correlation between features. Moreover, for both methods, the detail of features  
413 importance in the pairwise relationship of different treatments remain unclear.  
414 And more significant, if these hyperspectral signals are applied to practice, con-  
415 tinuous bandwidth signals make more sense for spectral cameras.

416

417 So here, we developed a tool to better visualize and select significant continu-  
418 ous wavelengths. It works like this. First, we set a bandwidth (here, 40 nm),  
419 then randomly and iteratively search for the wavelength of this bandwidth in the  
420 hyperspectral signals, and take the average of the signal for variance analysis.  
421 Record the significant band ( $p < 0.05$ ), and finally count the number of times the  
422 significant band appears and dynamically display based on the magnitude of the  
423 significance (Figure 6c). Here, based on the results, we roughly selected 409-429  
424 nm, 558-577 nm, 700-896 nm and 926-967 nm wavelengths for latter model  
425 training.

### 426 3.5 Models

427 Four kinds of machine learning model were applied to fit the different treatments'  
428 hyperspectral data, each treatment has 7 samples, 80 scans per sample and each  
429 scan is the average of 10 scans. A total of  $4 \times 7 \times 80$  data was used for model fitting,  
430 with 6-fold cross-validation and ROC-AUC as an evaluation metric.

431

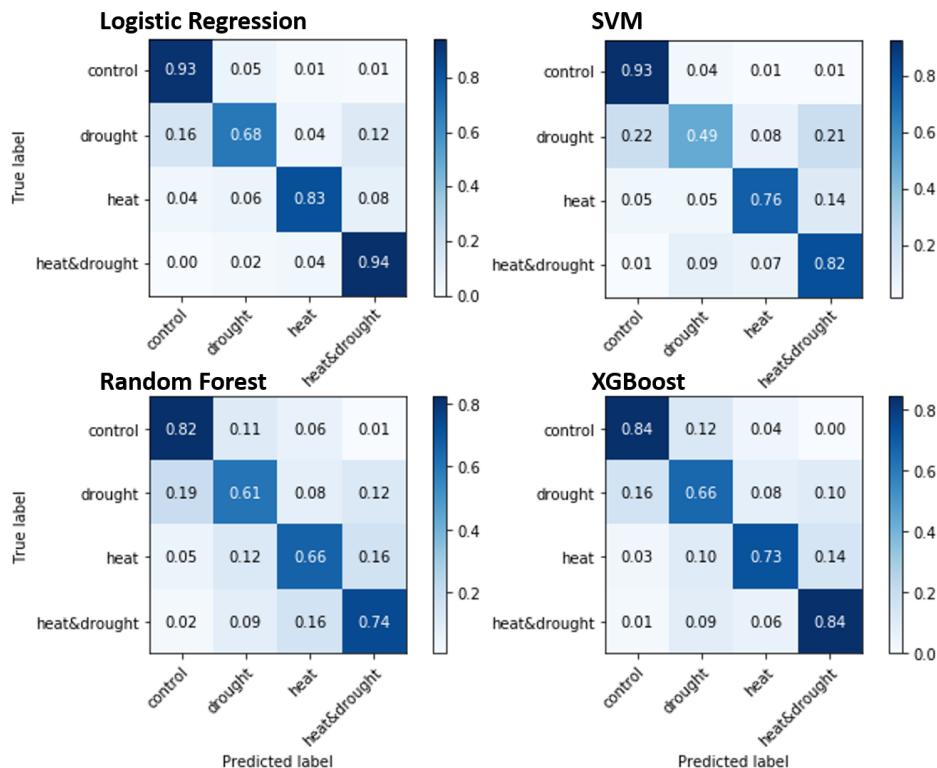
432 Results (Table 3) shows that in the fitting of the original data, the linear classifier  
433 LG and linear SVM fit well, while tree-based model performance is relatively poor,  
434 mainly because the number of features is too large. Too much redundant infor-  
435 mation makes the tree model easier to overfit. After removing some noise from

436 PCA, the performance of the tree models are improved, but the effect of LDA is  
 437 less pronounced. Furthermore, through the result of feature selection in the pre-  
 438 vious step, the models were trained with 409-429 nm, 558-577 nm, 700-896 nm  
 439 and 926-967 nm wavebands, and the performance of the tree model is further  
 440 improved.

**Table 3: AUC of 4 machine learning models under different features**

	LG	SVM	RF	XGBoost
raw	0.9494	0.9122	0.7794	0.8235
LDA	0.6826	0.6817	0.6710	0.6655
PCA	0.8509	0.8507	0.8285	0.9005
Selection	0.9108	0.8758	0.7438	0.8026
Selection+PCA	0.9155	0.8754	0.8767	0.9096

Selection stands for 409-429 nm, 558-577 nm, 700-896 nm and 926-967 nm wavelengths.



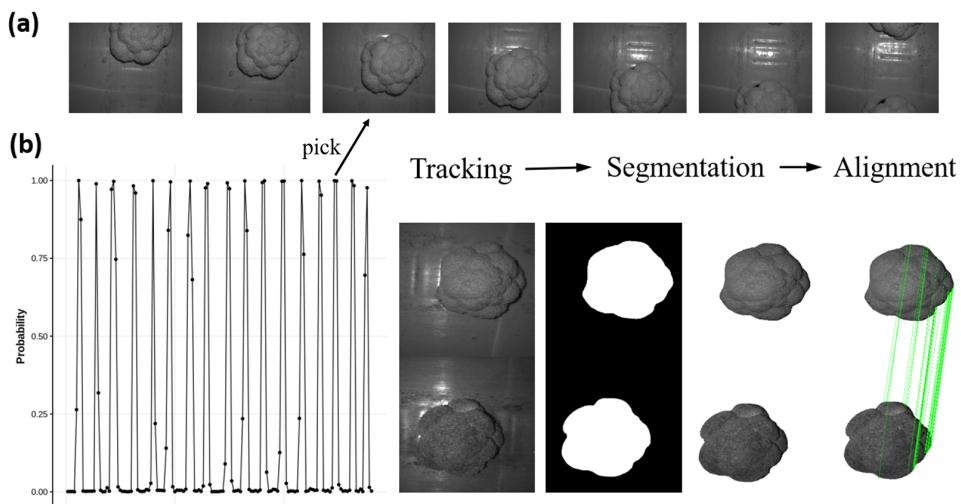
**Figure 7: The confusion matrix**

The prediction effect when the four models have the best AUC. True label refers to the original category of the data, and the predicted label is the category predicted by the model. The number in the box represents the accuracy, calculated by the number of predicted samples divided by the total number of samples for each category

441 Next, we observe their confusion matrix (Figure 7) when they have their best

442 AUC. The prediction performance of the four models is similar for different treat-  
 443 ments. Among them, the control and combined stress treatments have the best  
 444 predictive effect, followed by the heat stress, while surprisingly, the drought stress,  
 445 it has visually obvious symptom (droopy leaves), is the worst, and its mispredic-  
 446 tions tend to appear more in control and combined stress.

### 447 3.6 Shelf Life Prediction

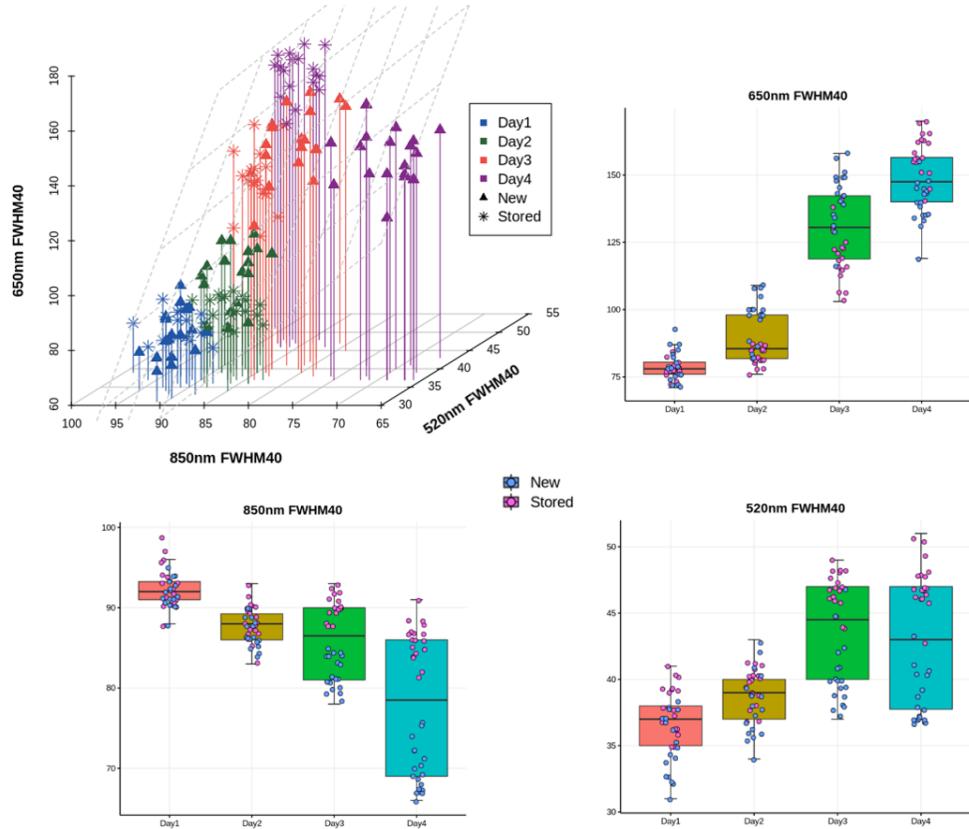


**Figure 8: Processing of broccoli head images on conveyor belt**

(a) from left to right is a broccoli time series images on a conveyor belt in the near-infrared channel; (b) is the ResNeXt prediction to track whether broccoli is in the middle of the lens, the closer the prediction probability is to 1, the more likely the broccoli head is in the middle of the lens. Select the image with the largest predicted probability between the two valleys, then segment by Unet and matched by SIFT.

448 To predict the shelf life on the broccoli packaging conveyor, we collected thou-  
 449 sands of broccoli spectral images (Figure 8a) with high-speed band-specific filter  
 450 cameras and specific waveband LED lights, and then we labeled 5708 broccoli im-  
 451 ages to tell whether the broccoli head was in the middle of the lens, so that we can  
 452 capture the spectral signal of the whole head. And through data augmentation,  
 453 with limited computing resources, a simple classifier can be constructed by trans-  
 454 fer learning of ResNext101\_64 (Xie et al., 2016) convolution neural network,  
 455 which can easily achieve an accuracy of 97.2%. By setting a lower threshold of  
 456 predicting probability, then we can select the highest probability of broccoli head  
 457 images between the two valleys to track broccolis (Figure 8b). After that, for

458 the sake of removing the influence of the background, we labeled 492 masks for  
 459 broccoli images segmentation, the Unet ([Ronneberger et al., 2015](#)) was trained  
 460 to reach 99.2% accuracy. Next, through the powerful SIFT operator ([Lowe et al.,](#)  
 461 [1999](#)), the images of different channels can be aligned for subsequent shelf-life  
 462 related signals analysis.



**Figure 9: Spectral reflection signal of three channels when the broccoli heads naturally rot**

Spectral images were collected by a camera with 3 waveband filters (green, red and near-infrared) under stable LED illumination. By segmenting the background and using the median to represent the signal of the entire broccoli head. The spectral reflectance of the broccoli head changes significantly with its natural decay. Among the early spectral signals which we are more concerned, 520 nm and 850 nm bands seem to be more effective to distinguish the new and stored heads. However, in general, their non-linear variations in different signals are elusive, and further modeling is needed.

463 Next, we placed the new and stored broccoli heads at room temperature, let them  
 464 decay naturally to capture the shelf-life related signals. With time elapsing, the re-  
 465 flectivity of the three selected channels can somewhat reflect the rotting changes.  
 466 However, our focus is more on the first two days before the broccoli heads showed  
 467 obvious decay phenotype. Among the three selected channels (green, red and

468 near-infrared), the newly harvested broccoli heads could not be effectively dis-  
469 tinguished from the stored one. Further statistical modeling or deep learning  
470 strategies need to be explored

## 471 4 DISCUSSION

472 From the results of various VIs calculated in different treatments (Figure 4), it  
473 seems not surprising that not getting a unified index that can effectively distin-  
474 guish these four treatments, because of the complex relationship among them.  
475 However, detailed analysis can still give us some inspiration. Firstly, indices based  
476 on red and near-infrared such as RVI, NDVI and EVI, can somewhat reflect the  
477 physiological characteristics of plants, however, they are more used in the field  
478 of remote sensing for various kinds of local, regional, and global scale models,  
479 including general circulation and biogeochemical models (Peterson et al., 1988,  
480 Huete et al., 2002). While in CVI, CI-G, CI-RE these indices of leaf chlorophyll  
481 content (Gitelson et al., 2003, Vincini et al., 2008), they seem to have a consis-  
482 tent result, except CI-RE, which doesn't consider the green channel. Results show  
483 that in the water-deficient environment (drought, heat&drought), the chlorophyll  
484 content of broccoli leaves was significantly affected (CVI, CI-G). This water stress  
485 phenomenon has been extensively supported, and it's also well understood that  
486 lack of water hinders nutrient transport in plants and thus affects photosynthetic  
487 pigments synthesis. As for the water index (WI), it reflects water absorption in  
488 the mesophyll andd had been shown to have a good indication of water content in  
489 many crops (Wang et al., 2015, Dawson et al., 1999, Peñuelas et al., 1997). The  
490 result here is also very satisfactory, it can be seen that there were significant dif-  
491 ferences between the two treatments except drought and combined stress. This  
492 indicates that in heat stress, the water content of the plant leaves is also affected,  
493 even though they are well watered, and this effect is not sufficient to superimpose  
494 the significance in the combined stress relative to the drought stress. The func-  
495 tional basis of the PRI is related on its sensitivity to rapid changes in carotenoids  
496 through the de-epoxidation of the xanthophyll pigments (Magney et al., 2016). It  
497 can serve as an indirect means for water stress detection due to the effects of wa-  
498 ter stress on the efficiency of photosynthesis. Researchers have demonstrated the

499 sensitivity of PRI to short-term crop water stress detection ([Gamon et al., 1997](#),  
500 [Suárez et al., 2010](#), [Zarco-Tejada et al., 2013](#)), and to the long-term change of  
501 carotenoid / chlorophyll ratio ([Mänd et al., 2010](#)). Here, we found that PRI is  
502 sensitive to the combination of drought and heat stress in broccoli, which may  
503 imply the cumulative effect of stress on photosynthetic efficiency and pigments.  
504 And of course, more studies is needed to support these inference.

505

506 In the process of machine learning model training for hyperspectral data, the  
507 results of feature engineering show that the important features are mainly con-  
508 centrated in the green and near infrared ([Figure 5,6](#)). In detail, because the  
509 hyperspectral bandwidth is small, there will inevitably be a lot of redundant in-  
510 formation. Data dimensionality reduction can effectively extract important infor-  
511 mation, shorten model training time and reduce over-fitting. Here, dimen-  
512 sional-  
513 ity reduction by PCA can effectively de-correlate and remove the linear relation-  
514 ship between dimensions, but it does not consider the classification information.  
515 Therefore, after dimensionality reduction, the loss of information will be mini-  
516 mized, but classification may become more difficult ([Figure 5a](#)). The data points  
517 in the graph are not clustering, and from the contribution of each component to  
518 the principal component ([Figure 5b](#)), it's truly hard to get useful information.  
519 Another commonly used dimension reduction method is LDA, which seeks to dis-  
520 tinguish data points as easily as possible after dimension reduction. After dimen-  
521 sionality reduction, the sample data has the largest inter-class distance and the  
522 smallest intra-class variance in the new dimension space, and the data has the best  
523 separability in the low dimension space ([Figure 5c](#)). It can almost reach 100%  
524 classification, so that the contribution of each wavelength may better explain the  
525 important features for classification, here are the green and red edges([Figure 5d](#)).

526

526 Then, as for feature selection, several methods of experimentation had a good  
527 consistency, that is, infrared waveband information might be relatively important  
528 for classification ([Figure 6](#)). It may be suggested that changes in mesophyll cell  
529 structure, such as membrane structure, are more likely to affect spectrum reflec-  
530 tion in leaves under heat and drought stress, while changes in pigments are less  
531 important to distinguish between them. In particular, through the dynamic visu-

532 alization of the random feature search box, we can more clearly see the impor-  
533 tance of features to the relationship between them (Figure 6c). The water stress  
534 and control group showed the most significant difference around 700-900 nm,  
535 which was basically consistent with the WI. Wavelength around 700-800 nm may  
536 be important for distinguishing between heat stress and control. As for combined  
537 stress, it seems similar to the water stress. And between the combined stress and  
538 the individual stress, there is a difference around 420 nm.

539

540 Finally, the results of the four machine learning models show that linear clas-  
541 sifier performs well when the data dimension is relatively large, while the tree  
542 model does not (Table 3). This is understandable because regularization is used  
543 in training linear classifiers, which can effectively deal with multiple collinearity  
544 problems and reduce the weight of redundant information. The tree model can  
545 also be improved after dimensionality reduction. According to the statistical anal-  
546 ysis results, we empirically select 409-429, 558-577, 700-896 and 926-967 nm  
547 waveband information for training, so we can see that the performance of the XG-  
548 Boost model has been improved effectively, its AUC can reach 0.9096. By showing  
549 the confusion matrix (Figure 7), it is not surprising that the control group and  
550 the combined stress group can achieve the best distinction. But surprisingly, the  
551 heat stress group can be more effectively distinguished from other stresses, as  
552 opposed to the drought group which may have more phenotype. The erroneous  
553 distinction of drought group mostly appears in the difference to control group,  
554 which may be explained to some extent that some of the leaves do not reach the  
555 threshold at which the drought can be detected.

556

557 As for the prediction of broccoli shelf life, due to the increasingly mature computer  
558 vision technology based on deep learning, and relatively stable environment and  
559 large data generated in production. It is easy to obtain high accuracy through a  
560 large number of data labeling and transfer learning. What is important is that for  
561 the capture of broccoli shelf life related signals, the more challenging is the signal  
562 difference in the early fresh period. Although we can get a signal that changes  
563 significantly with the decay of broccoli, how to predict its shelf life in the early  
564 stage remains to be further studied.

## 565    5   Limitations and Future Research

566    In the laboratory work, because planting broccolis requires a lot of time and space  
567    investment, limited by this, we have not been able to obtain large-scale image data  
568    and multibody repeated hyperspectral data. For the training of machine learning  
569    models and deep neural networks, they require a large amount of data, so more  
570    training data still needs to be acquired to train robust models.

571

572    The discussion of the specific molecular and physiological mechanisms for the  
573    results, most of them are extended through similar studies, and there could be  
574    differences between species and platforms. Further, for the pigment or physio-  
575    logical changes under complex stress of plants still need to be studied. To link the  
576    significant different spectral characteristics of complex stress with the changes of  
577    plant metabolites could be the direction of exploration.

578

579    Besides, for the relatively poor drought stress prediction effects. Error predic-  
580    tion is more likely to occur in control and combination stress (Figure 7). The  
581    results could be due to the different degree of water stress on leaves. Collecting  
582    quantitative stress level data as dependent variable, turning classification prob-  
583    lem into a regression problem could have better model performance.

584

585    While for predicting the shelf life of broccolis, it is feasible to make predictions  
586    from statistical signal differences, but as we see, predicting the shelf life of healthy  
587    broccoli heads with subtle difference requires more advanced technology and  
588    more data, and deep learning strategy is still the direction we need to develop.

## 589    6 ACKNOWLEDGEMENTS

590    I'm very greatful to my supervisor Dr. Oliver Windram for his patient guidance  
591    on this project and meticulous feedback on the writing. I'm also greatful to Chris  
592    Adam for his guidance on how to use the spectrophotometer. Besides, many  
593    thanks for Sarah Blanford from Sainsbury's and James Brown from POLLYBELL  
594    FARMS LTD for their support on this project.

## 595 7 CODE & DATA ACCESSIBILITY

596 All the code used for this project can be obtained from:

597 <https://github.com/Luoxsh6/CMEECourseWork>,

598 and the data from:

599 <https://imperialcollegelondon.box.com/s/k4ckuprv1if7kxqzb4lsiqhuxp22zwee>

## 600 References

601 Aladenola, O. and Madramootoo, C. (2014), ‘Response of greenhouse-grown bell  
602 pepper (*capsicum annuum* l.) to variable irrigation’, *Canadian journal of plant  
603 science* **94**(2), 303–310.

604 Alchanatis, V., Cohen, Y., Cohen, S., Moller, M., Sprinstin, M., Meron, M., Tsipris,  
605 J., Saranga, Y. and Sela, E. (2010), ‘Evaluation of different approaches for esti-  
606 mating and mapping crop water status in cotton with thermal imaging’, *Preci-  
607 sion Agriculture* **11**(1), 27–41.

608 Anjum, S., Wang, L., Farooq, M., Hussain, M., Xue, L. and Zou, C. (2011), ‘Brassi-  
609 nolide application improves the drought tolerance in maize through modulation  
610 of enzymatic antioxidants and leaf gas exchange’, *Journal of Agronomy and Crop  
611 Science* **197**(3), 177–185.

612 Bellvert, J., Marsal, J., Girona, J., Gonzalez-Dugo, V., Fereres, E., Ustin, S. and  
613 Zarco-Tejada, P. (2016), ‘Airborne thermal imagery to detect the seasonal evolu-  
614 tion of crop water status in peach, nectarine and saturn peach orchards’, *Remote  
615 Sensing* **8**(1), 39.

616 Berni, J. A., Zarco-Tejada, P. J., Suárez, L. and Fereres, E. (2009), ‘Thermal and  
617 narrowband multispectral remote sensing for vegetation monitoring from an  
618 unmanned aerial vehicle’, *IEEE Transactions on geoscience and Remote Sensing*  
619 **47**(3), 722–738.

620 Bravo, C., Moshou, D., West, J., McCartney, A. and Ramon, H. (2003), ‘Early dis-  
621 ease detection in wheat fields using spectral reflectance’, *Biosystems Engineering*  
622 **84**(2), 137–145.

- 623 Camejo, D., Jiménez, A., Alarcón, J. J., Torres, W., Gómez, J. M. and Sevilla,  
624 F. (2006), ‘Changes in photosynthetic parameters and antioxidant activities  
625 following heat-shock treatment in tomato plants’, *Functional Plant Biology*  
626 **33**(2), 177–187.
- 627 Cao, X., Luo, Y., Zhou, Y., Duan, X. and Cheng, D. (2013), ‘Detection of powdery  
628 mildew in two winter wheat cultivars using canopy hyperspectral reflectance’,  
629 *Crop Protection* **45**, 124–131.
- 630 Carter, G. A. (1993), ‘Responses of leaf spectral reflectance to plant stress’, *Amer-*  
631 *ican Journal of Botany* **80**(3), 239–243.
- 632 Chen, T. and Guestrin, C. (2016), Xgboost: A scalable tree boosting system, in  
633 ‘Proceedings of the 22nd acm sigkdd international conference on knowledge  
634 discovery and data mining’, ACM, pp. 785–794.
- 635 Dangwal, N., Patel, N., Kumari, M. and Saha, S. (2016), ‘Monitoring of water  
636 stress in wheat using multispectral indices derived from landsat-tm’, *Geocarto  
637 International* **31**(6), 682–693.
- 638 Dawson, T., Curran, P., North, P. and Plummer, S. (1999), ‘The propagation of  
639 foliar biochemical absorption features in forest canopy reflectance: A theoretical  
640 analysis’, *Remote Sensing of Environment* **67**(2), 147–159.
- 641 Din, J., Khan, S., Ali, I., Gurmani, A. et al. (2011), ‘Physiological and agronomic  
642 response of canola varieties to drought stress’, *J Anim Plant Sci* **21**(1), 78–82.
- 643 Feilhauer, H., Asner, G. P. and Martin, R. E. (2015), ‘Multi-method ensemble se-  
644 lection of spectral bands related to leaf biochemistry’, *Remote Sensing of Envi-*  
645 *ronment* **164**, 57–65.
- 646 Gamon, J., Penuelas, J. and Field, C. (1992), ‘A narrow-waveband spectral in-  
647 dex that tracks diurnal changes in photosynthetic efficiency’, *Remote Sensing of  
648 environment* **41**(1), 35–44.
- 649 Gamon, J., Serrano, L. and Surfus, J. (1997), ‘The photochemical reflectance  
650 index: an optical indicator of photosynthetic radiation use efficiency across  
651 species, functional types, and nutrient levels’, *Oecologia* **112**(4), 492–501.

- 652 Gates, D. M., Keegan, H. J., Schleter, J. C. and Weidner, V. R. (1965), ‘Spectral  
653 properties of plants’, *Applied optics* **4**(1), 11–20.
- 654 Gitelson, A. A., Gritz, Y. and Merzlyak, M. N. (2003), ‘Relationships between leaf  
655 chlorophyll content and spectral reflectance and algorithms for non-destructive  
656 chlorophyll assessment in higher plant leaves’, *Journal of plant physiology*  
657 **160**(3), 271–282.
- 658 Howard, J. et al. (2018), ‘fastai’, <https://github.com/fastai/fastai>.
- 659 Huete, A., Didan, K., Miura, T., Rodriguez, E. P., Gao, X. and Ferreira, L. G.  
660 (2002), ‘Overview of the radiometric and biophysical performance of the modis  
661 vegetation indices’, *Remote sensing of environment* **83**(1-2), 195–213.
- 662 Idso, S., Jackson, R., Pinter Jr, P., Reginato, R. and Hatfield, J. (1981), ‘Normalizing  
663 the stress-degree-day parameter for environmental variability’, *Agricultural  
664 meteorology* **24**, 45–55.
- 665 Ihuoma, S. O. and Madramootoo, C. A. (2017), ‘Recent advances in crop water  
666 stress detection’, *Computers and Electronics in Agriculture* **141**, 267–275.
- 667 Jordan, C. F. (1969), ‘Derivation of leaf-area index from quality of light on the  
668 forest floor’, *Ecology* **50**(4), 663–666.
- 669 Karnieli, A., Agam, N., Pinker, R. T., Anderson, M., Imhoff, M. L., Gutman, G. G.,  
670 Panov, N. and Goldberg, A. (2010), ‘Use of ndvi and land surface temperature  
671 for drought assessment: Merits and limitations’, *Journal of climate* **23**(3), 618–  
672 633.
- 673 Katsoulas, N., Elvanidi, A., Ferentinos, K. P., Kacira, M., Bartzanas, T. and Kittas,  
674 C. (2016), ‘Crop reflectance monitoring as a tool for water stress detection in  
675 greenhouses: A review’, *biosystems engineering* **151**, 374–398.
- 676 Kingma, D. P. and Ba, J. (2014), ‘Adam: A method for stochastic optimization’,  
677 *arXiv preprint arXiv:1412.6980* .
- 678 Knipling, E. B. (1970), ‘Physical and physiological basis for the reflectance of visi-  
679 ble and near-infrared radiation from vegetation’, *Remote sensing of environment*  
680 **1**(3), 155–159.

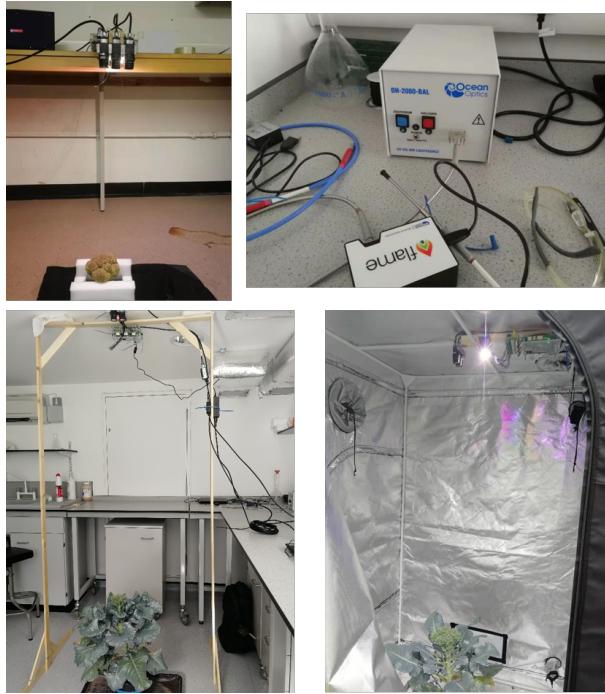
- 681 Kobayashi, T., Kanda, E., Kitada, K., Ishiguro, K. and Torigoe, Y. (2001), 'De-  
682 tection of rice panicle blast with multispectral radiometer and the potential of  
683 using airborne multispectral scanners', *Phytopathology* **91**(3), 316–323.
- 684 Kollist, H., Zandalinas, S. I., Sengupta, S., Nuhkat, M., Kangasjärvi, J. and Mittler,  
685 R. (2018), 'Rapid responses to abiotic stress: priming the landscape for the  
686 signal transduction network', *Trends in plant science* .
- 687 Lawlor, D. W. and Cornic, G. (2002), 'Photosynthetic carbon assimilation and  
688 associated metabolism in relation to water deficits in higher plants', *Plant, cell  
& environment* **25**(2), 275–294.
- 690 Lowe, D. G. et al. (1999), Object recognition from local scale-invariant features.,  
691 in 'iccv', Vol. 99, pp. 1150–1157.
- 692 Magney, T. S., Vierling, L. A., Eitel, J. U., Huggins, D. R. and Garrity, S. R. (2016),  
693 'Response of high frequency photochemical reflectance index (pri) measure-  
694 ments to environmental conditions in wheat', *Remote sensing of Environment*  
695 **173**, 84–97.
- 696 Mänd, P., Hallik, L., Peñuelas, J., Nilson, T., Duce, P., Emmett, B. A., Beier, C., Es-  
697 tiarte, M., Garadnai, J., Kalapos, T. et al. (2010), 'Responses of the reflectance  
698 indices pri and ndvi to experimental warming and drought in european shrub-  
699 lands along a north–south climatic gradient', *Remote Sensing of Environment*  
700 **114**(3), 626–636.
- 701 McClung, C. R. and Davis, S. J. (2010), 'Ambient thermometers in plants: from  
702 physiological outputs towards mechanisms of thermal sensing', *Current Biology*  
703 **20**(24), R1086–R1092.
- 704 Mestre, H. (1935), 'The absorption of radiation by leaves and algae', *Cold Spring  
705 Harbor Symposia on Quantitative Biology* pp. 191–209.
- 706 Mittler, R., Finka, A. and Goloubinoff, P. (2012), 'How do plants feel the heat?',  
707 *Trends in biochemical sciences* **37**(3), 118–125.
- 708 Muir, A., Porteous, R. and Wastie, R. (1982), 'Experiments in the detection of  
709 incipient diseases in potato tubers by optical methods', *Journal of Agricultural  
710 Engineering Research* **27**(2), 131–138.

- 711 Panigada, C., Rossini, M., Meroni, M., Cilia, C., Busetto, L., Amaducci, S.,  
712 Boschetti, M., Cogliati, S., Picchi, V., Pinto, F. et al. (2014), 'Fluorescence, pri  
713 and canopy temperature for water stress detection in cereal crops', *International*  
714 *Journal of Applied Earth Observation and Geoinformation* **30**, 167–178.
- 715 Pearson, R. L. and Miller, L. D. (1972), Remote mapping of standing crop biomass  
716 for estimation of the productivity of the shortgrass prairie, in 'Remote sensing  
717 of environment, VIII', p. 1355.
- 718 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,  
719 Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A.,  
720 Cournapeau, D., Brucher, M., Perrot, M. and Duchesnay, E. (2011), 'Scikit-learn:  
721 Machine learning in Python', *Journal of Machine Learning Research* **12**, 2825–  
722 2830.
- 723 Peñuelas, J. and Inoue, Y. (1999), 'Reflectance indices indicative of changes  
724 in water and pigment contents of peanut and wheat leaves', *Photosynthetica*  
725 **36**(3), 355–360.
- 726 Peñuelas, J., Pinol, J., Ogaya, R. and Filella, I. (1997), 'Estimation of plant wa-  
727 ter concentration by the reflectance water index wi (r900/r970)', *International*  
728 *Journal of Remote Sensing* **18**(13), 2869–2875.
- 729 Peterson, D. L., Aber, J. D., Matson, P. A., Card, D. H., Swanberg, N., Wessman, C.  
730 and Spanner, M. (1988), 'Remote sensing of forest canopy and leaf biochemical  
731 contents', *Remote Sensing of Environment* **24**(1), 85–108.
- 732 Platt, J. C. (1998), Sequential minimal optimization: A fast algorithm for training  
733 support vector machines, Technical report, ADVANCES IN KERNEL METHODS  
734 - SUPPORT VECTOR LEARNING.
- 735 Prasannakumar, N., Chander, S., Sahoo, R. and Gupta, V. (2013), 'Assessment of  
736 brown planthopper,(nilaparvata lugens)[stål], damage in rice using hyperspec-  
737 tral remote sensing', *International journal of pest management* **59**(3), 180–188.
- 738 Rabideau, G. S., French, C. S. and Holt, A. S. (1946), The absorption and reflec-  
739 tion spectra of leaves, chloroplast suspensions, and chloroplast fragments as  
740 measured in an ulbricht sphere.

- 741 Ronneberger, O., Fischer, P. and Brox, T. (2015), U-net: Convolutional networks  
742 for biomedical image segmentation, in ‘International Conference on Medical  
743 image computing and computer-assisted intervention’, Springer, pp. 234–241.
- 744 Rossini, M., Fava, F., Cogliati, S., Meroni, M., Marchesi, A., Panigada, C., Gia-  
745 rdino, C., Busetto, L., Migliavacca, M., Amaducci, S. et al. (2013), ‘Assessing  
746 canopy pri from airborne imagery to map water stress in maize’, *ISPRS Journal*  
747 of Photogrammetry and Remote Sensing **86**, 168–177.
- 748 Rouse Jr, J., Haas, R., Schell, J. and Deering, D. (1974), ‘Monitoring vegetation  
749 systems in the great plains with erts’.
- 750 Sinclair, T. R. et al. (1968), ‘Pathway of solar radiation through leaves’.
- 751 Smith, L. N. (2017), Cyclical learning rates for training neural networks, in ‘2017  
752 IEEE Winter Conference on Applications of Computer Vision (WACV)’, IEEE,  
753 pp. 464–472.
- 754 Suárez, L., Zarco-Tejada, P., Berni, J., González-Dugo, V. and Fereres, E. (2010),  
755 Orchard water stress detection using high-resolution imagery, in ‘XXVIII Interna-  
756 tional Horticultural Congress on Science and Horticulture for People  
757 (IHC2010): International Symposium on 922’, pp. 35–39.
- 758 Suzuki, N., Koussevitzky, S., Mittler, R. and Miller, G. (2012), ‘Ros and redox  
759 signalling in the response of plants to abiotic stress’, *Plant, Cell & Environment*  
760 **35**(2), 259–270.
- 761 Tucker, C. J. (1979), ‘Red and photographic infrared linear combinations for mon-  
762 itoring vegetation’, *Remote sensing of Environment* **8**(2), 127–150.
- 763 Vincini, M., Frazzi, E. and D’Alessio, P. (2008), ‘A broad-band leaf chlorophyll  
764 vegetation index at the canopy scale’, *Precision Agriculture* **9**(5), 303–319.
- 765 Wang, X., Zhao, C., Guo, N., Li, Y., Jian, S. and Yu, K. (2015), ‘Determining the  
766 canopy water stress for spring wheat using canopy hyperspectral reflectance  
767 data in loess plateau semiarid regions’, *Spectroscopy Letters* **48**(7), 492–498.

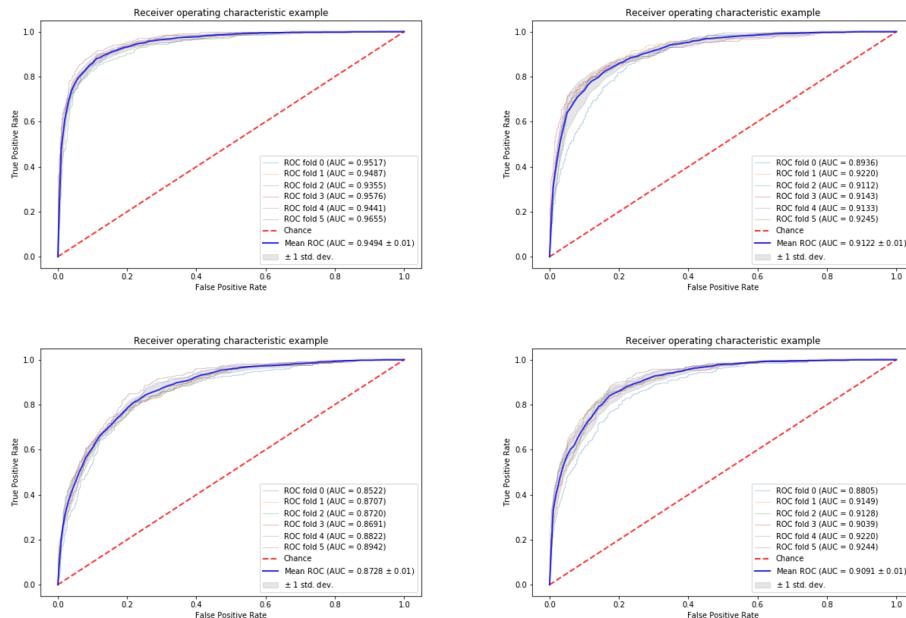
- 768 Willstätter, R. and Mieg, W. (1907), ‘Untersuchungen über chlorophyll; iv. ue-  
769 ber die gelben begleiter des chlorophylls’, *Justus Liebigs Annalen der Chemie*  
770 **355**(1), 1–28.
- 771 Xie, S., Girshick, R., Dollár, P., Tu, Z. and He, K. (2016), ‘Aggregated residual  
772 transformations for deep neural networks’, *arXiv preprint arXiv:1611.05431* .
- 773 Xue, J. and Su, B. (2017), ‘Significant remote sensing vegetation indices: A review  
774 of developments and applications’, *Journal of Sensors* **2017**.
- 775 Yang, C.-M. (2010), ‘Assessment of the severity of bacterial leaf blight in rice using  
776 canopy hyperspectral reflectance’, *Precision Agriculture* **11**(1), 61–81.
- 777 Zarco-Tejada, P. J., González-Dugo, V. and Berni, J. A. (2012), ‘Fluorescence,  
778 temperature and narrow-band indices acquired from a uav platform for wa-  
779 ter stress detection using a micro-hyperspectral imager and a thermal camera’,  
780 *Remote sensing of environment* **117**, 322–337.
- 781 Zarco-Tejada, P. J., González-Dugo, V., Williams, L., Suárez, L., Berni, J. A.,  
782 Goldhamer, D. and Fereres, E. (2013), ‘A pri-based water stress index combin-  
783 ing structural and chlorophyll effects: Assessment using diurnal narrow-band  
784 airborne imagery and the ccsi thermal index’, *Remote sensing of environment*  
785 **138**, 38–50.
- 786 Zarco-Tejada, P. J., Rueda, C. and Ustin, S. L. (2003), ‘Water content estimation in  
787 vegetation with modis reflectance data and model inversion methods’, *Remote  
788 Sensing of Environment* **85**(1), 109–124.

## 8 Supplementary Information



**Figure 10: Experiment Apparatus**

The tools used for images collection are basically consistent, spectral cameras with filters, spectrophotometer, LED lights and controllers. The figure shows the experiment apparatus for broccoli shelf life, broccoli in the control-environment room and in the grow tent.



**Figure 11: ROC for the best performance model**  
The ROC curve corresponding to the confusion matrix in the context



**Figure 12: Broccoli status for shelf-life prediction**

The RGB images of the naturally decaying broccoli, the images above is the broccoli stored for 2 weeks. The images below show the newly harvested broccoli, from left to right represent change of corresponding days.