



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> 1 INTRODUCTION

<sup>39</sup> The environment in which crops growth is highly dynamic, crops are constantly  
<sup>40</sup> exposed to various stresses, including abiotic stimuli such as humidity, light in-  
<sup>41</sup> tensity and temperature, and biotic stimuli, like pathogens. Correspondingly, to  
<sup>42</sup> cope with these stresses, plants have evolved a highly dynamic response mecha-  
<sup>43</sup> nism, from gene expression regulation to changes in various secondary metabo-  
<sup>44</sup> lites, which in turn alter their physiological characteristics, such as enzyme ac-  
<sup>45</sup> tivity, stomatal aperture, photosynthetic rate and transpiration rate (?). Rapid  
<sup>46</sup> detection of stresses through these physiological changes of crops is particularly  
<sup>47</sup> important for obtaining high yield and high quality products. Traditional detec-  
<sup>48</sup> tion methods usually require experts to be able to detect the subtle color changes  
<sup>49</sup> or a slight droop or curl of plants leaves, which indicate the stress. However,  
<sup>50</sup> it's generally subjective and time-consuming. In contrast, spectral detection and  
<sup>51</sup> spectral imaging, as a non-destructive, accurate and efficient method, has devel-  
<sup>52</sup> oped rapidly both in research and practical application (?).

<sup>53</sup>

<sup>54</sup> The spectral reflectance characteristics and mechanism of leaves have been well  
<sup>55</sup> summarized (??). When the leaves receive solar radiation, only part of the in-  
<sup>56</sup> cident energy is reflected and the rest is transmitted or absorbed for photosyn-  
<sup>57</sup> thesis. The typical reflectance spectrum of a leaf is that in the visible region  
<sup>58</sup> ( $0.4 - 0.7\mu\text{m}$ ), the reflectance of leaves is generally very low, especially in red  
<sup>59</sup> (around  $0.63 - 0.70\mu\text{m}$ ) and blue (around  $0.45 - 0.52\mu\text{m}$ ), this absorption charac-  
<sup>60</sup> teristic is mainly caused by plant pigments. Specifically, the absorption of red  
<sup>61</sup> primarily contributed from chlorophyll, and the absorption of blue is also in-  
<sup>62</sup> volved in carotenes and xanthophylls (??). Furthermore, the high reflection in  
<sup>63</sup> near-infrared ( $0.7 - 1.3\mu\text{m}$ ) is caused by the internal cellular structure (??). Leaf  
<sup>64</sup> cuticular wax is transparent and hardly reflects solar radiation. Radiation can  
<sup>65</sup> be transmitted through the epidermis, then dispersed and multiple reflected and  
<sup>66</sup> refracted in the mesophyll cells and air cavity, where different refracted index  
<sup>67</sup> between air (1.0) and hydrated cell walls (1.4) account for these effect (?).

<sup>68</sup>

<sup>69</sup> Previous studies on the application of spectral techniques in plant stress detec-

tion have made extensive progress, involving various crops and stress. In biotic stresses, plant disease detection is the most studied. Early in 1982, the diffuse reflectance spectra of potato tubers in the visible and near-infrared bands were measured and analyzed in an attempt to detect the presence of disease before its effects were visible (?). On top of that, spectral research also progress in the detection of various plant diseases, such as panicle blast, brown planthopper, the bacterial leaf in rice (???) and yellow rust, powdery mildew in wheat (??). In abiotic stress, diverse research focus on water stress. Among them, many studies are based on canopy temperature based Crop Water Stress Index (CWSI) measured from infrared thermometry (???). The principle of CWSI theory is that transpiration cools the surface of leaves. When soil moisture in the root zone decreases, stomatal conductance and transpiration are weakened, and then leaf temperature increases. What makes this theory popular is its linear relationship between canopy temperature and air temperature and vapor pressure, as well as the development of empirical methods for quantifying crop water stress (?). However, though the canopy temperature is very useful for water stress detection, it still has some physiological concerns. In some plants, the diurnal fluctuation in stomatal conductance make the relationship unclear between canopy temperature and stress levels (?). Moreover, leaf temperature does not directly explain other physiological changes, such as photosynthesis pigments or non-stomatal reduction of photosynthesis under water stress (?). Therefore, various alternative vegetation indices (VIs) based on the visible and red edge spectral region are developed to capture water stress related signals (??????).

Although spectroscopic studies of biotic and abiotic stresses can achieve significant 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,

103 oxidative damage (????). Moreover, photosynthesis is an important physiological  
104 phenomenon affected by drought and heat stress. Drought can lead to stomatal  
105 closure and reduces CO<sub>2</sub> uptake which makes plants more susceptible to photo  
106 damage, also it can induce negative changes in photosynthetic pigments, either  
107 increase or decrease chlorophyll content (???). Similarly, exposure to high tem-  
108 perature can also result in a reduction in chlorophyll biosynthesis, thereby dis-  
109 turbing the photosynthetic pigment components (?). Although many physiologi-  
110 cal links between plant heat stress and water stress, and it is of great meaningful  
111 to detect them in the practical production, little research have been conducted  
112 on the relationship between these two stress spectral characteristics and their  
113 detection. So it's a great interest and also useful to see how well we can detect  
114 the heat stress and drought and their combination from the crops by data mining  
115 their spectral characteristics changes.

116

117 Alongside this, methodologically, many previous studies on crops stresses detec-  
118 tion used statistical discriminant model, especially various VIs. VIs are quite sim-  
119 ple and effective algorithms for quantitative and qualitative evaluation of veg-  
120 etation cover, growth dynamics, and stress levels. But due to different spectral  
121 combinations, instrumentation, platforms, and resolutions used, it's hard to have  
122 a unified mathematical expression that defines all VIs, customized algorithms  
123 needed to developed and tested against specific application requirements (?).  
124 In this way, the prediction is often unsatisfactory and the generalization ability  
125 somewhat insufficient. Compared with traditional statistical discriminant model,  
126 machine learning methods, which developed rapidly in recent years, can gen-  
127 erally improve the speed and accuracy of prediction. The main difference be-  
128 tween machine learning and statistics lies in their purpose. Statistical models are  
129 more designed to infer the relationship between variables, while machine learn-  
130 ing models are intended to make the most accurate prediction possible.

131

132 Overall, here we explore the ability to use hyperspectral techniques to detect  
133 and differentiate more complex stresses of crops, that is, the combinations of  
134 drought and heat, which are highly correlated with each other. Firstly, we cal-  
135 culated some widely known VIs for statistical testing in anticipation of obtaining

simple and effective stress-related signals. Then, in order to approach the upper limit of accuracy for detecting different stresses, we apply machine learning strategy, using linear classifiers, such as logistic regression (LG), linear support vector machine (SVM) and tree models, like random forests (RF) and XGBoost for training a robust classifier. On top of this, to increase the interpretability of the model, great efforts had been made in feature engineering. Specifically, we first extract features through dimensionality reduction of Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA), and explore the importance of features under potential dimensions. Then statistical filter method and model-based embedded method are applied for feature selection. In particular, during these process, in order to better visualize the statistical differences of hyperspectral features under different stress comparisons dynamically, meanwhile be able to adjust the wavelength width to reduce the redundancy of hyperspectral information and provide a reference for specific band selection, a simple visualization search tool has been developed. Finally, the features obtained by these methods are combined to conduct model training, and the prediction effects of drought, heat stress and their combination stress were analyzed by model comparison.

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

## 2 MATERIALS AND METHODS

### 2.1 Workflow

The project mainly includes two parts (Figure ??). The laboratory part is to grow broccolis under control conditions and then perform individual and combined stress treatments, collect spectral images and leaf reflectance spectrum data to explore the signals that can effectively distinguish among them and construct a robust machine learning classifier. The application part is to construct a detection

166 system which can predict the shelf life of broccoli on the conveyor belt through  
 167 computer vision methods and spectral images under specific bandwidth, which is  
 168 selected based on the results obtained in the laboratory.

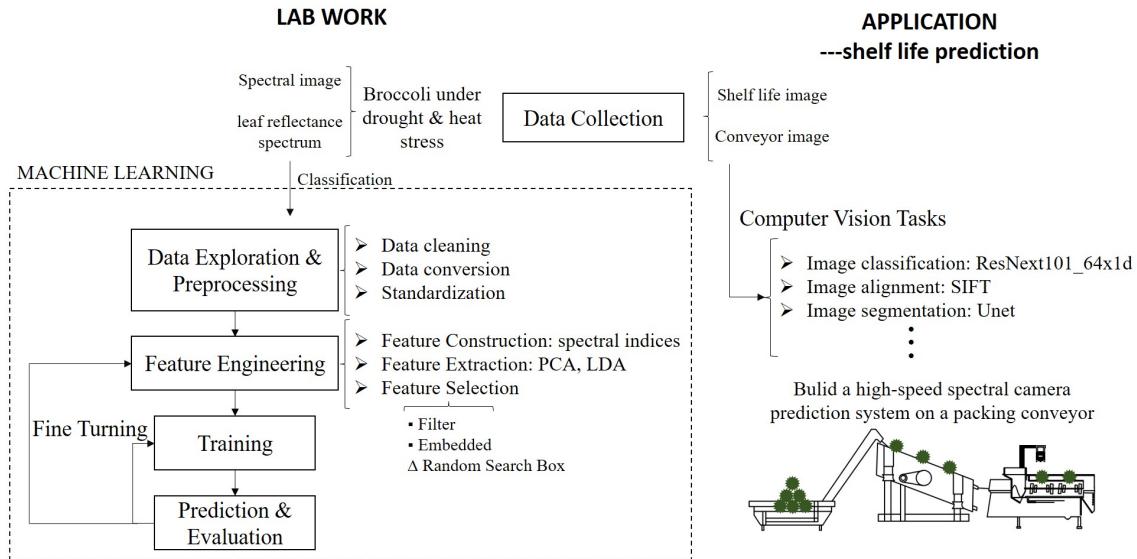


Figure 1: Project workflow

## 169 2.2 Data Collection

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

176

177 During stress treatment, broccolis were randomly divided into four groups with  
 178 8 individuals in each group. They were treated under control, heat stress (27°C),  
 179 drought stress (without watering for three days) and the combination of heat and  
 180 drought stress. Leaf reflectance spectroscopy data was collected by the Ocean  
 181 Optics FLAME-S-XR1 spectrophotometer in a completely dark room, while the  
 182 spectral image were collected by Ximea cameras with a specific bandpass filters  
 183 (FEL0800, FBH650-40) (Table ??) under the illumination of the corresponding  
 184 wavelength of the LED lamp. Four days of data were collected until the leaves  
 185 show a distinct dehydration drooping phenotype. In the open-air greenhouse,

<sup>186</sup> spectral images are collected with the same apparatus in a grow tent, in order to  
<sup>187</sup> 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 LED	Intelligent LED Solutions 12-Channel Light Controller 12 Die LED array Full Spectrum 360-955nm

<sup>188</sup> Broccoli heads for shelf-life prediction come from POLLYBELL FARMS LTD. They  
<sup>189</sup> are divided into two groups, 18 in each, one of which is stored in cold storage for  
<sup>190</sup> two weeks, and the other is harvested freshly. They were placed naturally at room  
<sup>191</sup> temperature and spectral image data were collected every day through the cam-  
<sup>192</sup> eras with bandpass filter (FBH520-40, FBH650-40,FBH850-40) and hyperspectral  
<sup>193</sup> data was collected by spectrophotometer as well until they decay significantly.

### <sup>194</sup> 2.3 Vegetation Indices Calculations

<sup>195</sup> The names, abbreviations, calculation formulas and citations of the various veg-  
<sup>196</sup> etation indices used in this study are as follows, mainly includes commonly used  
<sup>197</sup> remote sensing indices, chlorophyll-related indices, and indices related to water  
<sup>198</sup> stress indications.

**Table 2: Various vegetation indices**

Name	Abbrev.	Equation	References
Ratio vegetation index	NDVI	$R_n / R_r$	(?) (?)
Normalized difference vegetation index	NDVI	$(R_{800} - R_{680}) / (R_{800} + R_{680})$	(?) (?)
Enhanced vegetation index	EVI	$2.5(R_n - R_r)(R_n + 6 \cdot R_r - 7.5 \cdot R_b + 1)$	(?)
Chlorophyll vegetation index	CVI	$R_n \cdot R_r / R_g^2$	(?)
Chlorophyll index - green	CI-G	$R_n / R_g - 1$	(?)
Chlorophyll index - red edge	CI-RE	$R_n / R_{re} - 1$	(?)
Photochemical reflectance index	PRI	$(R_{531} - R_{570}) / (R_{531} + R_{570})$	(?)
Water index	WI	$R_{900} / R_{970}$	(?)
Structure independent pigment index	SICI	$(R_{800} - R_{445}) / (R_{800} + R_{680})$	(?)

$R_\lambda$  is the reflectance at wavelength  $\lambda$ ;  $n, r, e, b, g$  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.

## 199 2.4 Machine Learning

200 Machine learning generally includes several steps in practical operation, such as  
 201 data collection and preprocessing, model selection, training, evaluation and re-  
 202 peatedly fine-tuning until a good prediction effect is achieved (Figure ??). In the  
 203 data preprocessing stage, Z-score standardization is applied to simplify the cal-  
 204 culation and the categorical data is one-hot encoded. In the strategy of training  
 205 algorithms, firstly, LG, SVM, RF and XGBoost algorithm were trained to fit the  
 206 raw data and obtained the baseline score, and then the performance of the mod-  
 207 els were optimized by feature engineering and parameters adjustment. Most of  
 208 the code used in this process is based on the API provided by scikit-learn (?).

209 **2.4.1 Model Fitting**

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

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

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

216 Here,  $x$  is the input vector,  $w$  is the weight vector and  $Y$  is the output vec-  
217 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  
218 the probability of  $P(Y = 1|x)$  and  $P(Y = 0|x)$  can finally determine the category.  
219 The cost function of logistic regression can be derived by the method of maxi-  
220 mum likelihood estimation, which is known as the average of cross-entropy loss.  
221 Meanwhile, in order to prevent over-fitting, the L1 or L2 regularization terms was  
222 added during the optimization process.

224  
225 SVM: the main idea of SVM is to find the decision boundary with the largest  
226 classification interval between two different categories, which means that a hy-  
227 perplane separating classes in the feature space is defined by the principle of  
228 maximum margin between the closest different data points, also known as sup-  
229 port vectors. For simple linear separability problems, it can be described as an  
230 optimization problem by mathematical formulas as follows:

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

231 The minimized item represents the distance from the support vectors to the deci-  
232 sion boundary with sign, known as geometry margin. After derivation and trans-  
233 formation, and allow the SVM to ignore some noise, that is, allow some data  
234 points' functional margin less than 1, a slack variable ( $\xi_i \geq 0$ ) is introduced to al-  
235 low some wrong classification. correspondingly, a penalty term is needed to add

236 to the objective function to limit the slack variable, and here is the basic linear  
237 separable SVM:

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

238 Finally, the problem can be resolved by Lagrange Duality and SMO algorithm (?).  
239 In addition, kernel mapping has also been tried to verify the model's performance  
240 under nonlinearly separable conditions.

241

242 Ensemble methods: Random Forests and XGBoost (?), both are based on decision  
243 tree model, they use the strategy of bagging and boosting respectively, which can  
244 help to prevent high variance and high bias. Random forest mainly consists of two  
245 stochastic processes, random sampling of samples and features to construct many  
246 decision trees that are independent of each other. The final prediction results are  
247 summarized by voting strategy. As for XGBoost, it's an algorithm developed from  
248 gradient boosted decision trees and designed for speed and performance. It's  
249 widely used in many competitions and achieved good grades.

250

251 Metric: ROC curve (receiver operating characteristic curve) is a graph showing  
252 the performance of a classification model at all classification thresholds. The curve  
253 plots the points of (True positive Rate, False Positive Rate) at different classifica-  
254 tion thresholds. Lowering the classification threshold classifies more items as pos-  
255 itive, thus increasing both False Positives and True Positives. Area under the ROC  
256 Curve (AUC) measures the entire two-dimensional area underneath the entire  
257 ROC curve. AUC is desirable for its scale-invariant and classification-threshold-  
258 invariant.

259

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

263 **2.4.2 Feature Engineering**

264 Generally, data and features determine the upper bound of machine learning,  
265 whereas models and algorithms only approximate this upper bound. The pur-  
266 pose of feature engineering is to extract effective features and remove redundant  
267 features, also it can make features machine readable and contextually relevant.  
268 It basically includes feature extraction, feature construction and feature selec-  
269 tion (Figure ??). Separately, feature extraction mainly uses dimension reduction  
270 methods such as PCA and LDA. Feature construction is to construct new features  
271 based on previous expertise. Feature selection methods can be roughly divided  
272 into three types:

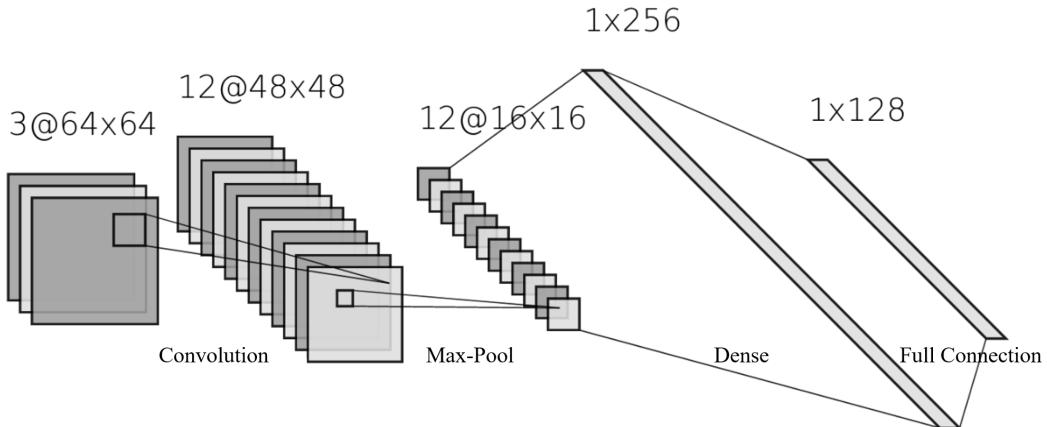
- 273 • Filter: scoring each feature according to divergence, correlation, etc., and  
274 then set a threshold for selection feature.
- 275 • Embedded: use some machine learning algorithms and models to train and  
276 get the coefficients of each feature, select features according to the coef-  
277 ficient, kind of similar to the filter method, but models are trained to de-  
278 termine the pros and cons of features. Specifically, multi-method ensemble  
279 selection (?) was modified from the regression problem and later adapted  
280 to the classification problem.
- 281 • Wrapper: recursive elimination feature method, due to the high computa-  
282 tional complexity and the long execution time of the algorithm, it is not  
283 adopted here.

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

289 **2.5 Computer Vision**

290 The project involves computer vision tasks such as image classification, segmen-  
291 tation, and image alignment. Specifically, image alignment was performed by  
292 classic Scale-invariant feature transform (SIFT) (?), Image classification and seg-

293 mention are mainly accomplished by the transfer learning of convolution neural  
294 network (Figure ??).



**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.

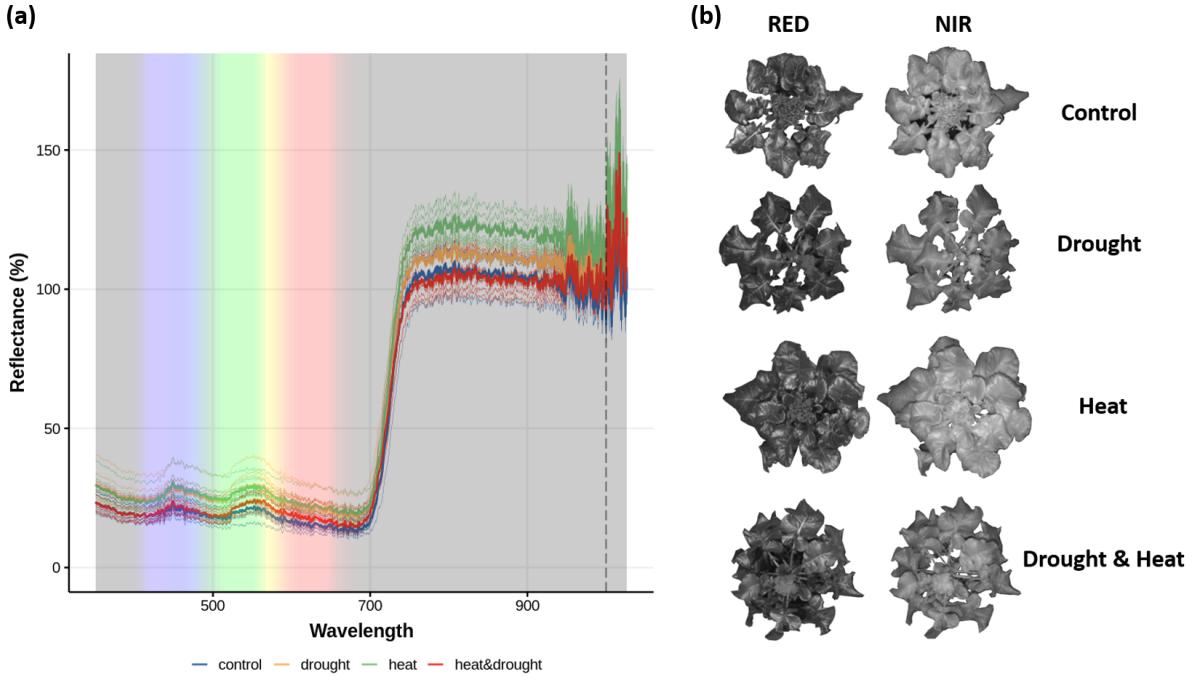
295 More specifically, Image classification was implemented by ResNeXt (?), devel-  
296 oped by UC San Diego and Facebook AI Research, while Image segmentation was  
297 implemented by Unet (?). The training was conducted on GPU P1000, the opti-  
298 mizer for the neural network is Adam (?), the cost function is cross-entropy, and  
299 cyclical learning rates (?) was used, Most of the code is based on the API provided  
300 by Pytorch, fastai library (?), and opencv.

### 301 **3 RESULTS**

#### 302 **3.1 Data**

303 To study the spectral characteristics of drought stress, heat stress, and their com-  
304 bination, Leaf reflectance non-image hyperspectral data (Figure ??a) and spec-  
305 tral images of two channels (Figure ??b) were collected. We select the data of  
306 the day when the broccoli just appeared phenotype under the stress to ensure  
307 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 ??b).



**Figure 3: hyperspectral non-imaging data and spectral iamge data of brocoli 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 ( $\text{CWL} = 650 \text{ nm}$ ,  $\text{FWHM} = 40 \text{ nm}$ ) and near infrared bandpass filter ( $> 800 \text{ nm}$ ), under the illumination of the corresponding band of LED.

For the hyperspectral data of the leaves (Figure ??a), 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 (?). It's also worth pointing out that the average reflectance of the heat treatment

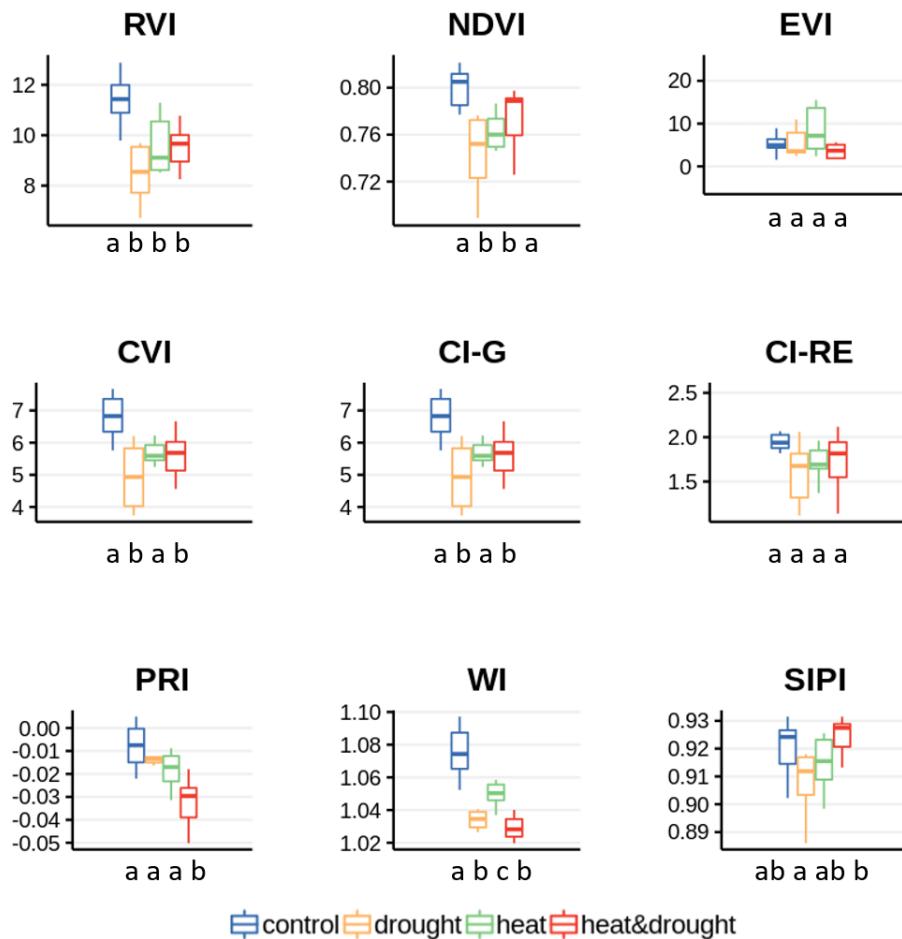
<sup>321</sup> appears to be generally higher than other treatments in the near-infrared.

### <sup>322</sup> 3.2 Vegetation Indices

<sup>323</sup> For the sake of looking for a simple and effective algorithm to distinguish four  
<sup>324</sup> different treatments, several common spectral indices were calculated under dif-  
<sup>325</sup> ferent stress (Figure ??). Specifically, basic vegetation index RVI, is based on the  
<sup>326</sup> phenomenon that leaves absorb relatively more red than infrared light. It's widely  
<sup>327</sup> used for green biomass estimations (?). NDVI is the most widely used VI to char-  
<sup>328</sup> acterize canopy growth and vigor (?). EVI was introduced to correct soil and  
<sup>329</sup> atmospheric effects on NDVI (?). As for CVI, CI-G and CI-RE, they are to some  
<sup>330</sup> extent related to chlorophyll content (??). and for PRI, WI and SIFI have been  
<sup>331</sup> proven to be effective of water status in some species (??).

<sup>332</sup>

<sup>333</sup> The results show that significant differences between every two treatments can't  
<sup>334</sup> be simply obtained from a single index. The results with significant differences  
<sup>335</sup> also show different forms of band feature calculation methods, which is difficult  
<sup>336</sup> to establish a unified pattern, suggesting that more complex models are needed.

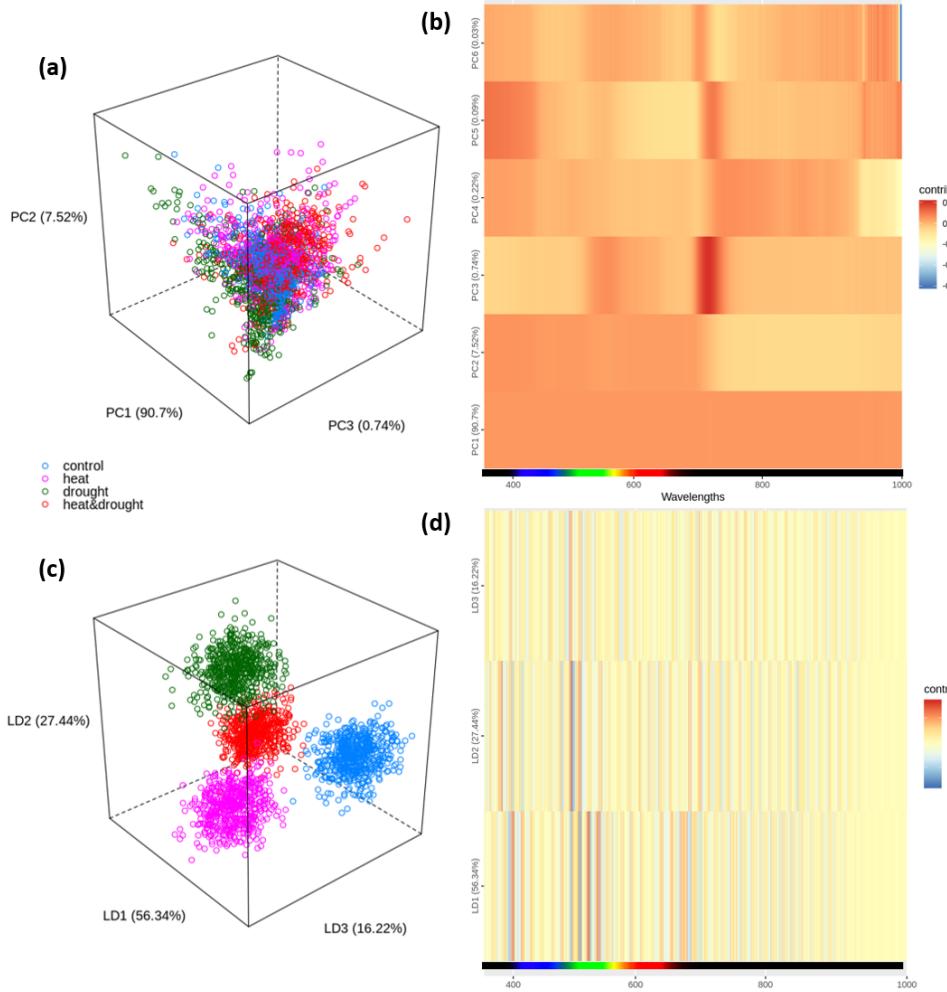


**Figure 4: Various vegetation indices under different stresses**

Seven samples per treatment, the equation for calculating vegetation index refers to Tabel ???. 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.



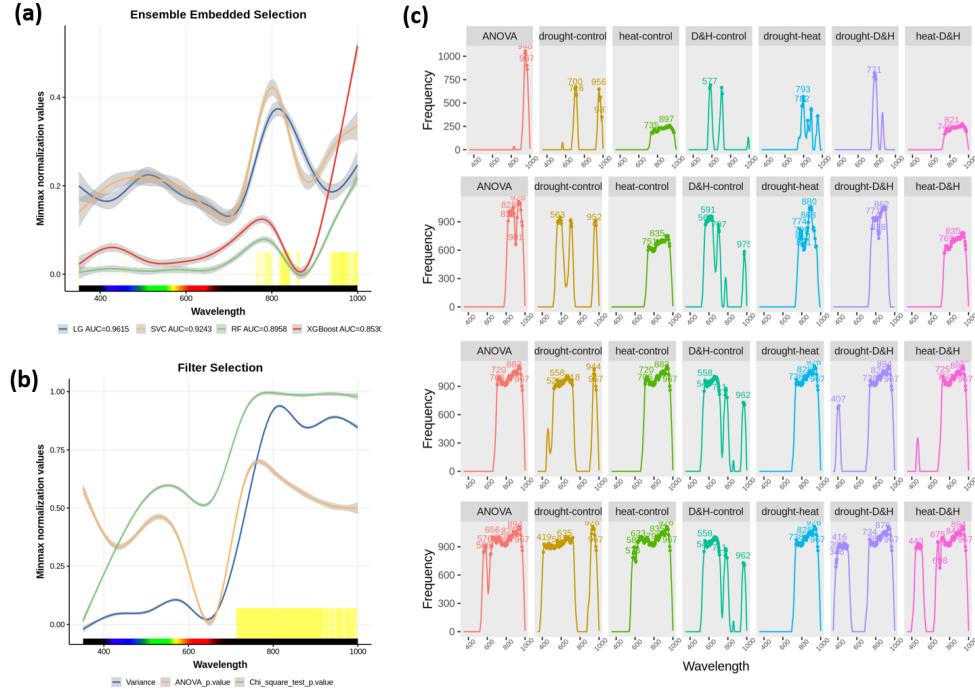
**Figure 5: Feature dimensionality reduction by PCA and LDA**

(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 ??a). 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 ??b). On the other hand, LDA, the supervised dimension reduction method, can completely separate the different treatments (Figure ??c). 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 ??d).

### 354 3.4 Feature Selection

355 In order to remove irrelevant features, reduce over-fitting in the process of model  
 356 training, and find explanatory wavelengths that can effectively distinguish different  
 357 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.

358 Filter and Embedded methods are commonly used feature selection methods in  
 359 machine learning. The results of the Filter method (Figure ??b) show that the  
 360 variance, chi-square test and the ANOVA of each wavelength are coincident upon  
 361 threshold in the near-infrared region. As for the modified ensemble method (?),  
 362 the correlation coefficient or feature importance obtained by fitting the LG, SVM,  
 363 RF and XGBoost are weighted with their AUC scores, then set their average plus  
 364 standard deviation as threshold. And the features upon the threshold are mainly  
 365 include two near-infrared fragments and other narrow fragments in visible light  
 366 (Figure ??a).

367

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

379

380 So here, we developed a tool to better visualize and select significant continu-  
381 ous wavelengths. It works like this. First, we set a bandwidth (here, 40 nm),  
382 then randomly and iteratively search for the wavelength of this bandwidth in the  
383 hyperspectral signals, and take the average of the signal for variance analysis.  
384 Record the significant band ( $p < 0.05$ ), and finally count the number of times the  
385 significant band appears and dynamically display based on the magnitude of the  
386 significance (Figure ??c). Here, based on the results, we roughly selected 409-  
387 429 nm, 558-577 nm, 700-896 nm and 926-967 nm wavelengths for latter model  
388 training.

### 389 3.5 Models

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

394

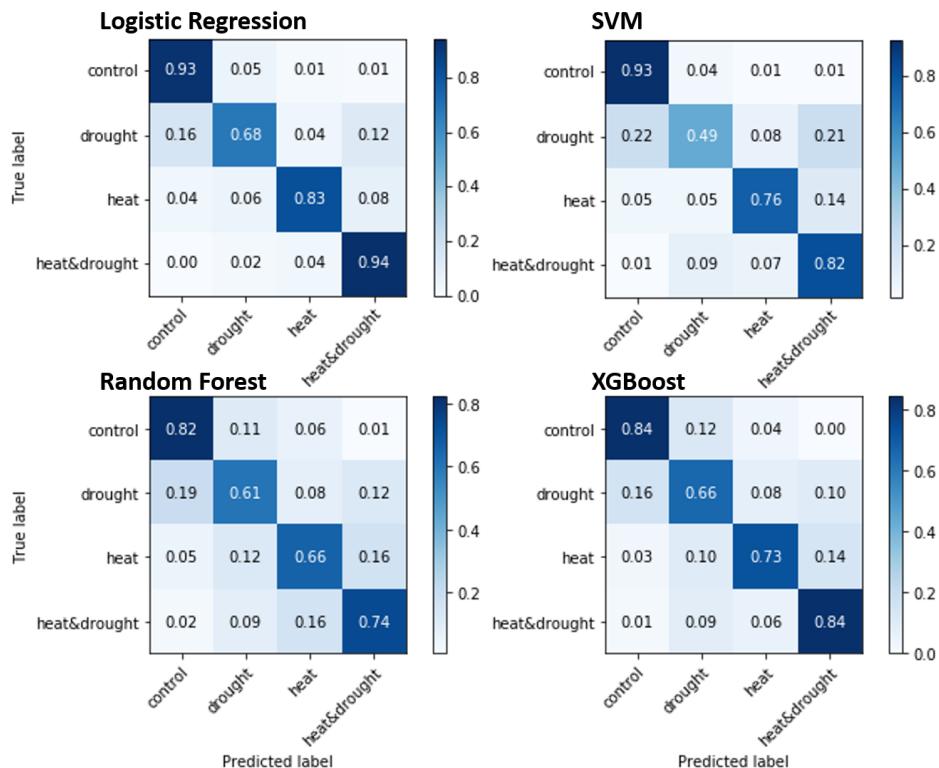
395 Results (Table ??) shows that in the fitting of the original data, the linear classi-  
396 fier LG and linear SVM fit well, while tree-based model performance is relatively  
397 poor, mainly because the number of features is too large. Too much redundant  
398 information makes the tree model easier to overfit. After removing some noise

399 from PCA, the performance of the tree models are improved, but the effect of  
 400 LDA is less pronounced. Furthermore, through the result of feature selection in  
 401 the previous step, the models were trained with 409-429 nm, 558-577 nm, 700-  
 402 896 nm and 926-967 nm wavebands, and the performance of the tree model is  
 403 further 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.



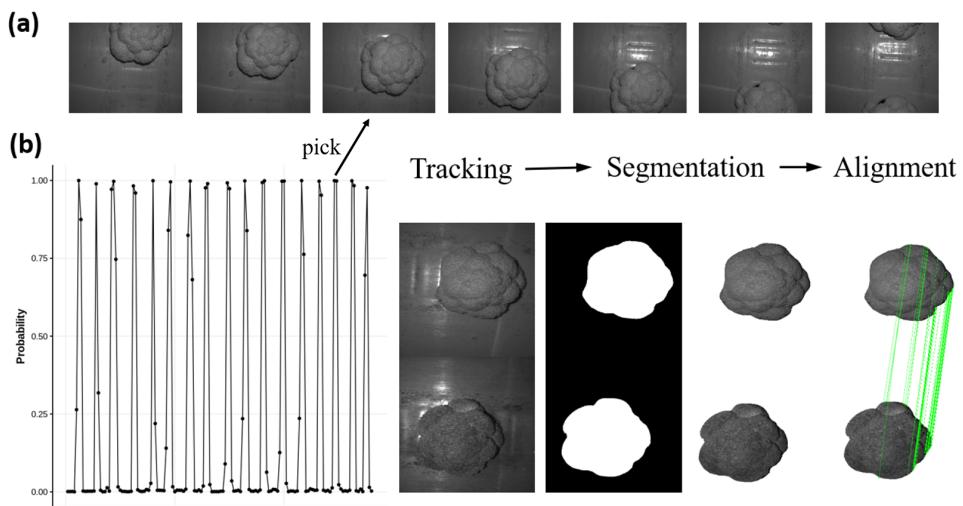
**Figure 7: The confusion matrix**

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

404 Next, we observe their confusion matrix (Figure ??) when they have their best

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

410 **3.6 Shelf Life Prediction**

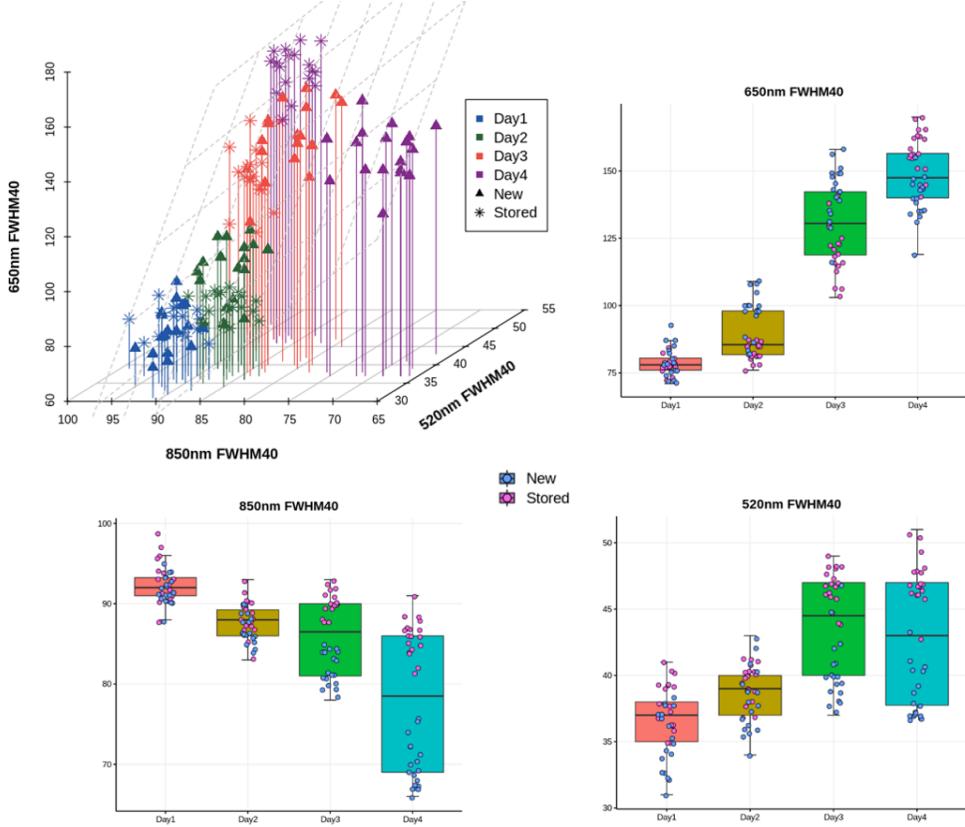


**Figure 8: Processing of broccoli head images on conveyor belt**

(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.

411 To predict the shelf life on the broccoli packaging conveyor, we collected thou-  
 412 sands of broccoli spectral images (Figure ??a) with high-speed band-specific filter  
 413 cameras and specific waveband LED lights, and then we labeled 5708 broccoli im-  
 414 ages to tell whether the broccoli head was in the middle of the lens, so that we  
 415 can capture the spectral signal of the whole head. And through data augmenta-  
 416 tion, with limited computing resources, a simple classifier can be constructed by  
 417 transfer learning of ResNext101\_64 (?) convolution neural network, which can  
 418 easily achieve an accuracy of 97.2%. By setting a lower threshold of predicting  
 419 probability, then we can select the highest probability of broccoli head images  
 420 between the two valleys to track broccolis (Figure ??b). After that, for the sake

421 of removing the influence of the background, we labeled 492 masks for broccoli  
 422 images segmentation, the Unet (?) was trained to reach 99.2% accuracy. Next,  
 423 through the powerful SIFT operator (?), the images of different channels can be  
 424 aligned for subsequent shelf-life 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.

425 Next, we placed the new and stored broccoli heads at room temperature, let them  
 426 decay naturally to capture the shelf-life related signals. With time elapsing, the re-  
 427 flectivity of the three selected channels can somewhat reflect the rotting changes.  
 428 However, our focus is more on the first two days before the broccoli heads showed  
 429 obvious decay phenotype. Among the three selected channels (green, red and  
 430 near-infrared), the newly harvested broccoli heads could not be effectively dis-

<sup>431</sup> distinguished from the stored one. Further statistical modeling or deep learning  
<sup>432</sup> strategies need to be explored

## <sup>433</sup> 4 DISCUSSION

<sup>434</sup> From the results of various VIs calculated in different treatments (Figure ??),  
<sup>435</sup> it seems not surprising that not getting a unified index that can effectively dis-  
<sup>436</sup> tinguish these four treatments, because of the complex relationship among them.  
<sup>437</sup> However, detailed analysis can still give us some inspiration. Firstly, indices based  
<sup>438</sup> on red and near-infrared such as RVI, NDVI and EVI, can somewhat reflect the  
<sup>439</sup> physiological characteristics of plants, however, they are more used in the field of  
<sup>440</sup> remote sensing for various kinds of local, regional, and global scale models, in-  
<sup>441</sup> cluding general circulation and biogeochemical models (??). While in CVI, CI-G,  
<sup>442</sup> CI-RE these indices of leaf chlorophyll content (??), they seem to have a consis-  
<sup>443</sup> tent result, except CI-RE, which doesn't consider the green channel. Results show  
<sup>444</sup> that in the water-deficient environment (drought, heat&drought), the chlorophyll  
<sup>445</sup> content of broccoli leaves was significantly affected (CVI, CI-G). This water stress  
<sup>446</sup> phenomenon has been extensively supported, and it's also well understood that  
<sup>447</sup> lack of water hinders nutrient transport in plants and thus affects photosynthetic  
<sup>448</sup> pigments synthesis. As for the water index (WI), it reflects water absorption in  
<sup>449</sup> the mesophyll andd had been shown to have a good indication of water content  
<sup>450</sup> in many crops (???). The result here is also very satisfactory, it can be seen that  
<sup>451</sup> there were significant differences between the two treatments except drought and  
<sup>452</sup> combined stress. This indicates that in heat stress, the water content of the plant  
<sup>453</sup> leaves is also affected, even though they are well watered, and this effect is not  
<sup>454</sup> sufficient to superimpose the significance in the combined stress relative to the  
<sup>455</sup> drought stress. The functional basis of the PRI is related on its sensitivity to rapid  
<sup>456</sup> changes in carotenoids through the de-epoxidation of the xanthophyll pigments  
<sup>457</sup> (?). It can serve as an indirect means for water stress detection due to the effects of  
<sup>458</sup> water stress on the efficiency of photosynthesis. Researchers have demonstrated  
<sup>459</sup> the sensitivity of PRI to short-term crop water stress detection (???), and to the  
<sup>460</sup> long-term change of carotenoid / chlorophyll ratio (?). Here, we found that PRI  
<sup>461</sup> is sensitive to the combination of drought and heat stress in broccoli, which may

462 imply the cumulative effect of stress on photosynthetic efficiency and pigments.  
463 And of course, more studies is needed to support these inference.

464

465 In the process of machine learning model training for hyperspectral data, the  
466 results of feature engineering show that the important features are mainly con-  
467 centrated in the green and near infrared (Figure ??,??). In detail, because the  
468 hyperspectral bandwidth is small, there will inevitably be a lot of redundant in-  
469 formation. Data dimensionality reduction can effectively extract important in-  
470 formation, shorten model training time and reduce over-fitting. Here, dimen-  
471 sionality reduction by PCA can effectively de-correlate and remove the linear re-  
472 lationship between dimensions, but it does not consider the classification infor-  
473 mation. Therefore, after dimensionality reduction, the loss of information will  
474 be minimized, but classification may become more difficult (Figure ??a). The  
475 data points in the graph are not clustering, and from the contribution of each  
476 component to the principal component (Figure ??b), it's truly hard to get useful  
477 information. Another commonly used dimension reduction method is LDA, which  
478 seeks to distinguish data points as easily as possible after dimension reduction.  
479 After dimensionality reduction, the sample data has the largest inter-class dis-  
480 tance and the smallest intra-class variance in the new dimension space, and the  
481 data has the best separability in the low dimension space (Figure ??c). It can al-  
482 most reach 100% classification, so that the contribution of each wavelength may  
483 better explain the important features for classification, here are the green and  
484 red edges (Figure ??d).

485

486 Then, as for feature selection, several methods of experimentation had a good  
487 consistency, that is, infrared waveband information might be relatively impor-  
488 tant for classification (Figure ??). It may be suggested that changes in mesophyll  
489 cell structure, such as membrane structure, are more likely to affect spectrum re-  
490 flection in leaves under heat and drought stress, while changes in pigments are  
491 less important to distinguish between them. In particular, through the dynamic  
492 visualization of the random feature search box, we can more clearly see the im-  
493 portance of features to the relationship between them (Figure ??c). The water  
494 stress and control group showed the most significant difference around 700-900

495 nm, which was basically consistent with the WI. Wavelength around 700-800 nm  
496 may be important for distinguishing between heat stress and control. As for com-  
497 bined stress, it seems similar to the water stress. And between the combined  
498 stress and the individual stress, there is a difference around 420 nm.

499

500 Finally, the results of the four machine learning models show that linear clas-  
501 sifier performs well when the data dimension is relatively large, while the tree  
502 model does not (Table ??). This is understandable because regularization is used  
503 in training linear classifiers, which can effectively deal with multiple collinearity  
504 problems and reduce the weight of redundant information. The tree model can  
505 also be improved after dimensionality reduction. According to the statistical anal-  
506 ysis results, we empirically select 409-429, 558-577, 700-896 and 926-967 nm  
507 waveband information for training, so we can see that the performance of the XG-  
508 Boost model has been improved effectively, its AUC can reach 0.9096. By showing  
509 the confusion matrix (Figure ??), it is not surprising that the control group and  
510 the combined stress group can achieve the best distinction. But surprisingly, the  
511 heat stress group can be more effectively distinguished from other stresses, as  
512 opposed to the drought group which may have more phenotype. The erroneous  
513 distinction of drought group mostly appears in the difference to control group,  
514 which may be explained to some extent that some of the leaves do not reach the  
515 threshold at which the drought can be detected.

516

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

## 5 Limitations and Future Research

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

531

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

538

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

544

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

## 549 6 ACKNOWLEDGEMENTS

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

## 555 7 CODE & DATA ACCESSIBILITY

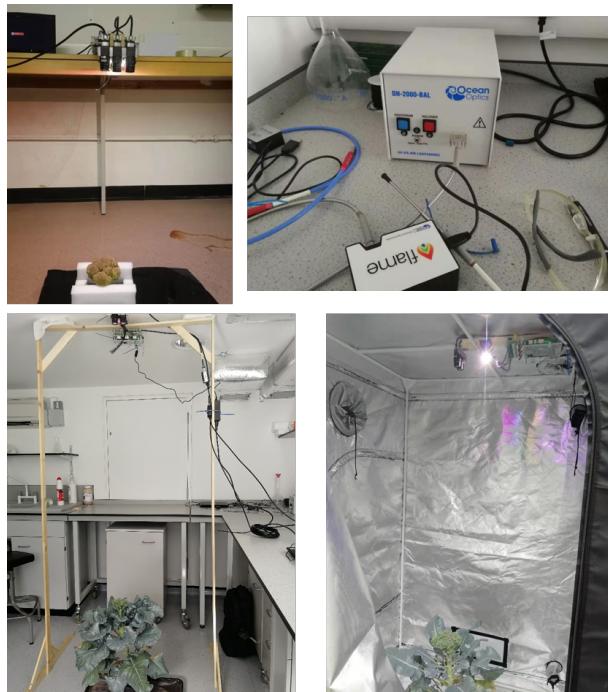
556 All the code used for this project can be obtained from:

557 <https://github.com/Luoсх6/CMEECourseWork>,

558 and the data from:

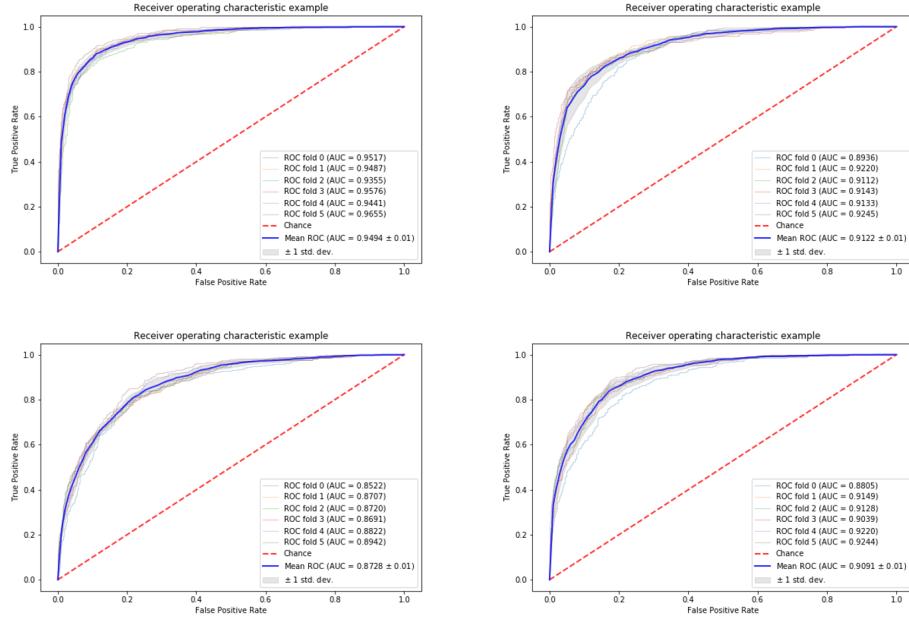
559 <https://imperialcollegelondon.box.com/s/k4ckuprv1if7kxqzb4lsiqhuxp22zwee>

## 560 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.