

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:

Dr. OLIVER WINDRAM

01627437

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

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

As a non-destructive, and high-efficiency technology, spectroscopy has developed rapidly in crop management for precision agriculture. Many studies on spectroscopic detection of crop stress have made effective progress, including disease detection and water stress, etc. Various vegetation indices have also been established to indicate crops growth status. However, complex growing environment often cause crops to be affected by multiple stresses, which in turn affect crops yield and quality. Few studies have been conducted on this complex stress, especially the two physiological closely related abiotic stresses, heat and drought. Here we took broccoli as our research object, and collected the hyperspectral reflectance data of its leaves under heat and drought, as well as their combination by spectrophotometer. Then, we explored the spectral characteristics by various vegetation indices, and the results showed that a single vegetation index could hardly be significant in all treatments. Next, we train different machine learning models to predict these complex stresses, and the highest AUC can reach 0.9494 by logistic regression. Besides, we have developed a tool that can dynamically visualize the significant different wavelength between treatments. And additionally, we also try to combine spectroscopy and deep learning to build a system for predicting the shelf-life of broccoli heads, for now, it can track broccoli heads with ResNeXt to 97.2% accuracy and segment them with Unet to 99.2% accuracy, while the signal related to the shelf-life is still in progress.

Keywords: Broccoli; Machine Learning; Spectral feature; Heat stress; Drought stress

26

Word count: 2191.

27 Contents

28	1 INTRODUCTION	4
29	2 MATERIALS AND METHODS	8
30	2.1 Data Collection	8
31	2.2 Vegetation Indices Calculations	9
32	2.3 Machine Learning	10
33	2.3.1 Model Fitting	10
34	2.3.2 Feature Engineering	12
35	2.4 Computer Vision	13
36	2.5 Workflow	14
37	3 RESULTS	14
38	3.1 Data	14
39	3.2 Vegetation Indices	16
40	3.3 Dimensionality Reduction	17
41	3.4 Feature Selcetion	18
42	3.5 Models	19
43	3.6 Shelf Life Prediction	20
44	4 DISCUSSION	22
45	5 Supplementary Information	30
46	5.1 Experiment Apparatus	30

1 INTRODUCTION

The environment in which crops growth and reproduction is highly dynamic, crops are constantly exposed to various stresses, including abiotic stimuli such as humidity, light intensity and temperature, and biotic stimuli, like pathogens. Correspondingly, to cope with these stresses, plants have evolved a highly dynamic response mechanism, from gene expression regulation to changes within various secondary metabolites, which in turn alter their physiological characteristics, such as enzyme activity, stomatal aperture, photosynthetic rate, transpiration rate and so on. Rapid detection of stresses through these physiological changes of crops is particularly important for obtaining high yield and high quality products. Traditional detection methods usually require skilled growers who are capable to notice the subtle color changes or a slight droop or curl of plants leaves, which indicate the stress. However, it's generally subjective and time-consuming. In contrast, spectral detection and spectral imaging, as a non-destructive, accurate and efficient method, has developed rapidly both in research and practical application.

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

79 Previous studies on the application of spectral techniques in plant stress detec-
80 tion have made extensive progress, involving various crops and stress. In biotic
81 stresses, plant disease detection is the most studied. Early in 1982, the diffuse
82 reflectance spectra of potato tubers in the visible and near-infrared bands were
83 measured and analyzed in an attempt to detect the presence of disease before
84 its effects were visible (Muir et al., 1982). On top of that, spectral researches
85 also progress in the detection of various plant diseases, such as panicle blast,
86 brown planthopper, the bacterial leaf in rice (Kobayashi et al., 2001, Prasannaku-
87 mar et al., 2013, Yang, 2010) and yellow rust, powdery mildew in wheat (Bravo
88 et al., 2003, Cao et al., 2013). In abiotic stress, diverse researches focus on wa-
89 ter stress. Among them, many studies are based on canopy temperature based
90 Crop Water Stress Index (CWSI) measured from infrared thermometry (Alcha-
91 natis et al., 2010, Aladenola and Madramootoo, 2014, Bellvert et al., 2016). The
92 principle of CWSI theory is that transpiration cools the surface of leaves. When
93 soil moisture in the root zone decreases, stomatal conductance and transpiration
94 are weakened, and then leaf temperature increases. Further, what makes this the-
95 ory popular is its linear relationship between canopy temperature and insufficient
96 temperature and vapor pressure, as well as the development of empirical methods
97 for quantifying crop water stress (Idso et al., 1981). However, though the canopy
98 temperature is very useful for water stress detection, it still has some physiological
99 concerns. In some plants, the diurnal fluctuation in stomatal conductance make
100 the relationship unclear between canopy temperature and stress levels (Zarco-
101 Tejada et al., 2012). Moreover, leaf temperature does not directly explain other
102 physiological changes, such as photosynthesis pigments or non-stomatal reduc-
103 tion of photosynthesis under water stress (Zarco-Tejada et al., 2013). Therefore,
104 various alternative vegetation indices (VIs) based on the visible and red edge
105 spectral region are developed to capture water stress related signals (Berni et al.,
106 2009, Zarco-Tejada et al., 2013, Wang et al., 2015, Rossini et al., 2013, Panigada
107 et al., 2014, Dangwal et al., 2016).

108

109 Although spectroscopic studies of biotic and abiotic stresses can achieve signifi-
110 cant detection under different models, the stress they detect is often single, while
111 in practical production, crops tend to suffer from multiple complex stresses during

112 their growth, such as heat and drought and their combinations, especially in the
113 context of global warming. More interestingly, the molecular and physiological
114 mechanisms by which plants respond to heat and drought stress have been ex-
115 tensively studied and they show lots of connection. oth of them can differentially
116 affect the RNA stability, alter the enzyme activity and disrupt the steady-state of
117 metabolic flux, which in most of cases can cause a common response, oxidative
118 damage (Kollist et al., 2018, Suzuki et al., 2012, Mittler et al., 2012, McClung and
119 Davis, 2010). Besides, photosynthesis is an important physiological phenomenon
120 affected by drought and heat stress. Drought can lead to stomatal closure and
121 reduces CO₂ uptake which makes plants more susceptible to photodamage. And
122 also it can induce negative changes in photosynthetic pigments, either increase or
123 decrease chlorophyll content(Lawlor and Cornic, 2002, Anjum et al., 2011, Din
124 et al., 2011). Similarly, exposure to high temperature can also result in a reduc-
125 tion in chlorophyll biosynthesis, thereby disturbing the photosynthetic pigment
126 components (Camejo et al., 2006). Despite many physiological links between
127 plant heat stress and water stress, and it is of great meaningful to detect them in
128 the pratical production, little research have been conducted on the relationship
129 between these two stress spectral characteristics and their detection. So it's a
130 great interest and also useful to see how well we can detect the heat stress and
131 drought and their combination from the crops by data mining their spectral char-
132 acteristics changes.

133

134 Alongside this, methodologically, many previous studies on crops stresses detec-
135 tion used statistical discriminant model, specifically, variety of VIs, as a simple
136 and effective algorithm for quantitative and qualitative evaluation of vegetation
137 cover, growth dynamics, and stress levels. But due to different spectral combina-
138 tions, instrumentation, platforms, and resolutions used, it's hard to have a uni-
139 fied mathematical expression that defines all VIs, customized algorithms needed
140 to developed and tested against specific application requirements(Xue and Su,
141 2017). What'more, the prediction is often unsatisfactory and the generalization
142 ability somewhat insufficient. Compared with traditional statistical discriminant
143 model, machine learning methods, which developed rapidly in recent year, can
144 generally improve the speed and accuracy of prediction. The main difference be-

145 tween machine learning and statistics lies in their purpose. Statistical models are
146 more designed to infer the relationship between variables, while machine learn-
147 ing models are intended to make the most accurate prediction possible.

148

149 Overall, here we explore the ability to use hyperspectral techniques to detect and
150 differentiate more complex stresses of crops, that is, the combinations of drought
151 and heat, which are highly correlated with each other. Firstly, we calculated
152 some widely known VIs for statistical testing in anticipation of obtaining simple
153 and effective stress-related signals. Then, in order to approach the upper limit of
154 accuracy for detecting different stresses, we apply the machine learning strategy,
155 using linear classifiers, such as Logistic regression (LG), linear support vector ma-
156 chine (SVM) and tree models, like random forest (RF) and XGBoost for training
157 a robust classifier. On top of this, in order to increase the interpretability of the
158 model, great efforts had been made in feature engineering, including dimension
159 reduction, statistical filter method and model-based embedded method, the com-
160 parison and selection of features are carried out for model training, and finally
161 achieved considerable prediction accuracy. In particular, during these process, in
162 order to better visualize the statistical differences of hyperspectral features under
163 different stress comparisons dynamically, meanwhile be able to adjust the wave-
164 length width to reduce the redundancy of hyperspectral information and provide
165 a reference for specific band selection, a simple visualization search tool has been
166 developed.

167

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

2 MATERIALS AND METHODS

2.1 Data Collection

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

Table 1: Apparatus used in the experiment.

Item	Description
Camera Machine vision lens	Ximea 1.3 MP NIR Enhanced Camera MQ013RG-ON 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

During stress treatment, broccolis were randomly divided into four groups with 8 individuals in each group. They were treated under control, heat stress (27°C), drought stress (without watering) and the combination of heat and drought stress. Leaf reflectance spectroscopy data was collected by the Ocean Optics FLAME-S-XR1 spectrophotometer in a completely dark room, while the spectral image were collected by Ximea cameras with a specific bandpass filter (FEL0800, FBH650-40) (Table 1) under the illumination of the corresponding wavelength of the LED lamp, four days of data were collected until the leaves show a distinct dehydration drooping phenotype. And in the open-air greenhouse, spectral images are collected in a grow tent.

191

192 Broccoli heads for shelf-life prediction come from POLLYBELL FARMS LTD. They
 193 are divided into two groups, 18 in each, one of which is stored in cold storage
 194 for a some time, and the other is harvested freshly. They were placed naturally
 195 at room temperature and spectral image data were collected every day through
 196 the cameras with bandpass filter (FBH520-40, FBH650-40,FBH850-40) until they
 197 decay significantly.

198 2.2 Vegetation Indices Calculations

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

Table 2: Various vegetation indices

Name	Abbrev.	Equation	References
Ratio vegetation index	RVI	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	$(R_{531} - R_{570}) / (R_{531} + R_{570})$	(?)
Water index	WI	R_{900} / R_{970}	(?)
Structure independent pigment index	SIPI	$(R_{800} - R_{445}) / (R_{800} + R_{680})$	(?)

R_λ is the reflectance at wavelength λ ; 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.

2.3 Machine Learning

Machine learning generally includes several steps in practical operation, such as data collection and preprocessing, model selection, training, evaluation and repeatedly fine-tuning until a good prediction effect is achieved (Figure 2). In the data preprocessing stage, Z-score standardization is applied to simplify the calculation and the categorical data is one-hot encoded. In the strategy of training algorithms, firstly, logistic regression (LG), support vector machine (SVM), random forest (RF) and XGBoost algorithm were trained to fit the raw data and obtained the baseline score, and then the performance of the models were optimized by feature engineering and parameters adjustment. Most of the code used in this process is based on the API provided by sklearn(Pedregosa et al., 2011).

2.3.1 Model Fitting

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

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

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

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

SVM: the main idea of SVM is to find the decision boundary with the largest

classification interval between two different categories, which means that a hyperplane separating classes in the feature space is defined by the principle of maximum margin between the closest different data points, also known as support vectors. For simple linear separability problems, it can be described as an optimization problem by mathematical formulas as follows:

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

The minimized item represents the distance from the support vectors to the decision boundary with sign, known as geometry margin. By scaling w and b so that (x_j, y_j) as the point to get the minimum value, $y_j (w^T x_j + b) = 1$, so the other sample points are naturally greater than or equal to 1. Derived all the way and we can finally got a methemathical optimization problem:

$$\begin{aligned} \min_{w,b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } y_i (w^T x_i + b) \geq 1, \quad i = 1, 2, \dots, m \end{aligned}$$

Further, in order to allow the SVM to ignore some noise, a slack variable ($\xi_i \geq 0$) is introduced to allow some wrong classification, that is, allow some data points' functional margin less than 1, correspondingly, a penalty term is needed to add to the objective function to limit the slack variable, and here is the basic linear separable SVM:

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

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

Ensemble methods: Random Forests and XGBoost(Chen and Guestrin, 2016), both are based on decision tree model, they use the methods of bagging and boosting respectively, which can help to prevent high variance and high bias. Random forest mainly consists of two stochastic processes, random sampling of samples and features to construct many decision trees that are independent of

each other. The final prediction results are summarized by voting strategy. As for XGBoost, it's an algorithm developed from gradient boosted decision trees and designed for speed and performance. It's widely used in many competitions and achieved good grades.

In the process of training machine learning algorithms, when multi-classification is performed on logistic regression and support vector machines, "one to the rest" strategy is applied. All the models are validated by 6-fold cross-validation, and ROC_AUC is used as the metric for model evaluation.

2.3.2 Feature Engineering

Generally, data and features determine the upper bound of machine learning, whereas models and algorithms only approximate this upper bound. The purpose of feature engineering is to extract effective features and remove redundant features from the original data. It basically includes feature extraction, feature construction and feature selection [2](#). Separately, feature extraction mainly uses dimension reduction methods such as Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA). Feature construction is to construct new features based on previous expertise. Feature selection methods can be roughly divided into three types:

- Filter: scoring each feature according to divergence, correlation, etc., and then set a threshold for selection feature.
- Embedded: use some machine learning algorithms and models to train and get the coefficients of each feature, select features according to the coefficient, kind of similar to the filter method, but models are trained to determine the pros and cons of features. Specifically, multi-method ensemble selection ([Feilhauer et al., 2015](#)) was modified from the regression problem and later adapted to the classification problem.
- Wrapper: recursive elimination feature method, due to the high computational complexity and the long execution time of the algorithm, it is not adopted here.

In addition, in order to find a suitable bandpass filter for the camera, a search box

with specific bandwidth was used to repeatedly and randomly select features, and the importance of features is sorted by simple ANOVA and TukeyHSD significance test. Finally, the graph is plotted by accumulating significant ($p < 0.05$) bandwidth features.

2.4 Computer Vision

The project involves computer vision tasks such as image classification, segmentation, and image alignment. Specifically, image alignment was performed by classic Scale-invariant feature transform (SIFT)(Lowe et al., 1999), Image classification and segmentation are mainly accomplished by the transfer learning of convolution neural network (Figure 1). More specifically, Image classification

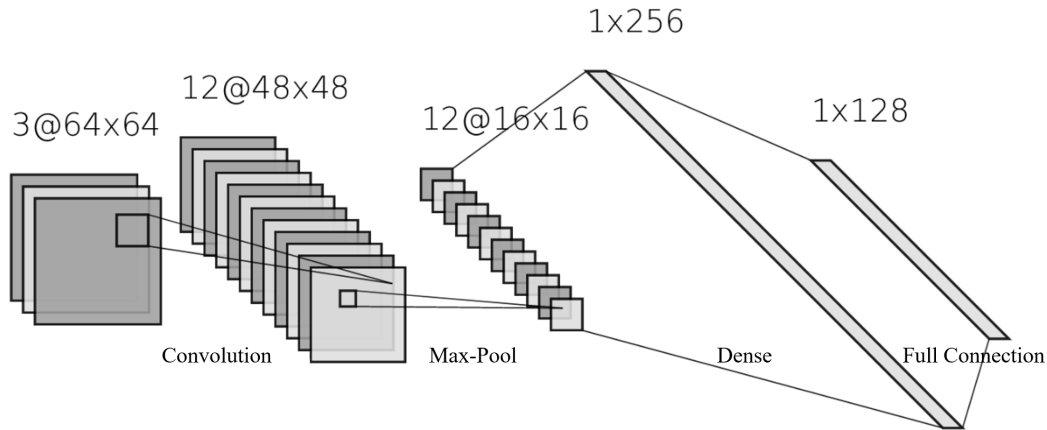


Figure 1: Structure of a simple convolution neural network

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 back is the height and width of pixels.

was implemented by ResNeXt (Xie et al., 2016), developed by UC San Diego and Facebook AI Research, while Image segmentation was implemented by Unet (Ronneberger et al., 2015). The training was conducted at P1000 on the HPC of Imperial College London, the optimizer for the neural network is Adam, the cost function is cross-entropy, and cyclical learning rates (Smith, 2017) was used, Most of the code is based on the API provided by Pytorch, fastai library (Howard et al.,

2018), and opencv.

2.5 Workflow

The project mainly includes two parts (Figure 2). 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 system which can predict the shelf life of broccoli on the conveyor belt through computer vision methods and spectral images under specific bandwidth, which is selected based on the results obtained in the laboratory.

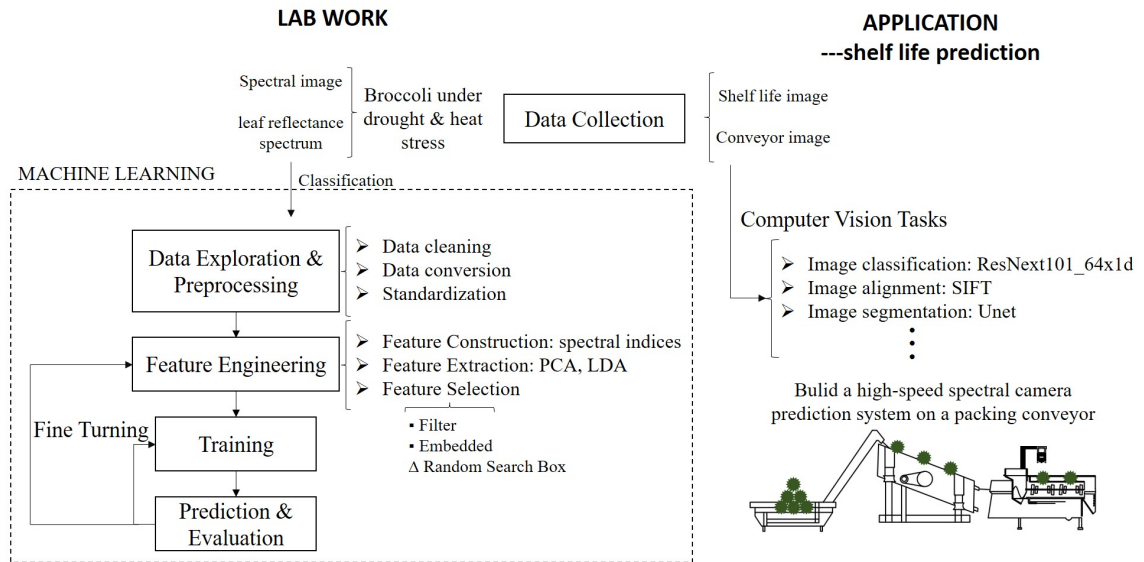


Figure 2: Project workflow

3 RESULTS

3.1 Data

We select the data of the day when the broccoli just appeared phenotype under the stress (Figure 3) to ensure the treatment effect. 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.

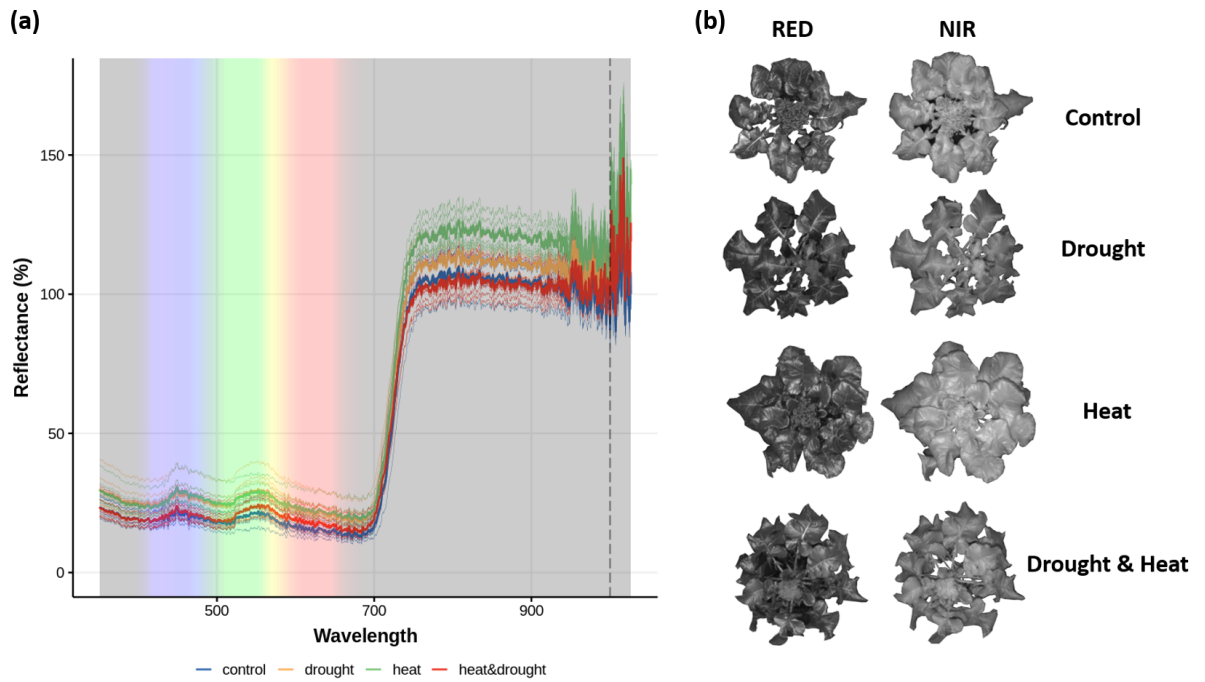


Figure 3: Spectral data of broccoli under heat and drought stress

(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 averaged by the hyperspectral scan of all the leaves of each sample, and the thick line is the average of all samples in different treatments. The right side of the dashed line was discarded in 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.

On the other hand, in the hyperspectral data of the leaves (Figure 3 (a)), it is difficult to get a distinct discriminant pattern from the perspective of data distribution, because the samples of different treatments are cross-covered. However, the overall trend of the broccoli leaves reflectance spectrum can still be clearly seen. There are small peaks at the blue and green-yellow junctions in the visible region, and it is well known that a small valley in the red band and strong reflection rate shifting in the near-infrared band. In particular, the average reflectance of the heat treatment appears to be significantly higher than other treatments.

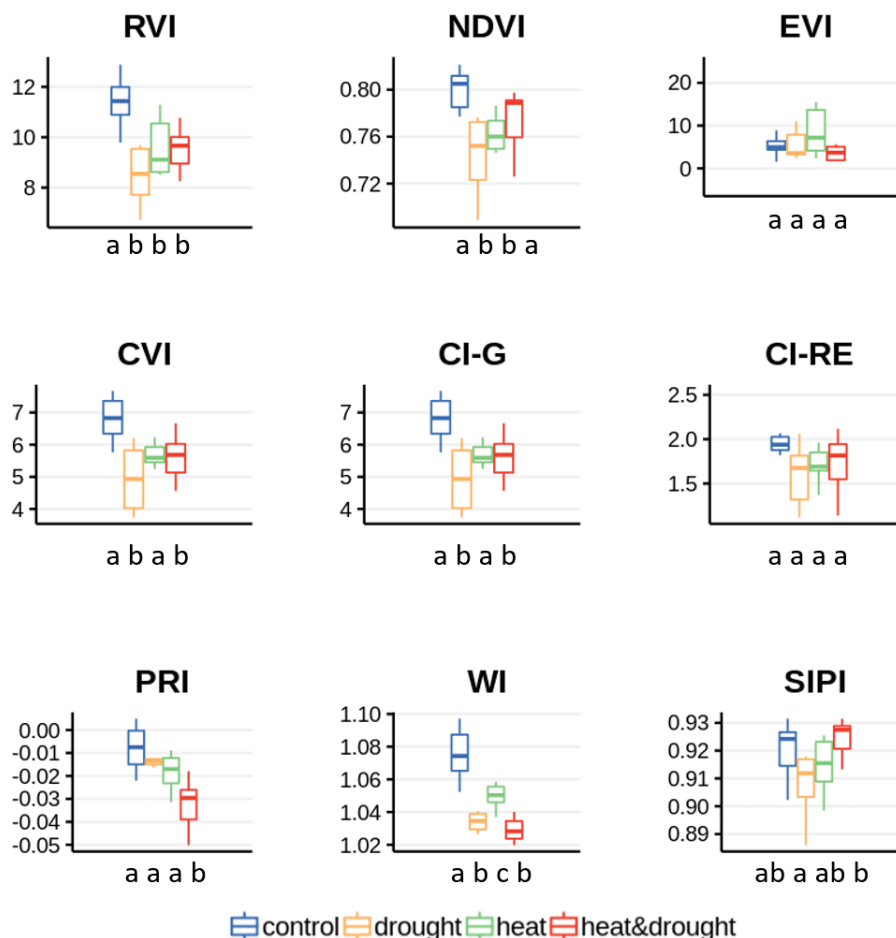


Figure 4: Various vegetation indices under different stresses

There are 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 ($p < 0.05$).

328 The spectral indices related to vegetation cover and chlorophyll content as well
 329 as water stress were calculated under different stress (Figure 4). The results
 330 show that significant differences between every two treatments can't be simply
 331 obtained from a single index. The results with significant differences also show
 332 different forms of band feature calculation methods, which is difficult to establish
 333 a unified pattern, suggesting that more complex models are needed.

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 variable, speed up computation and extract effective information hidden in the data.

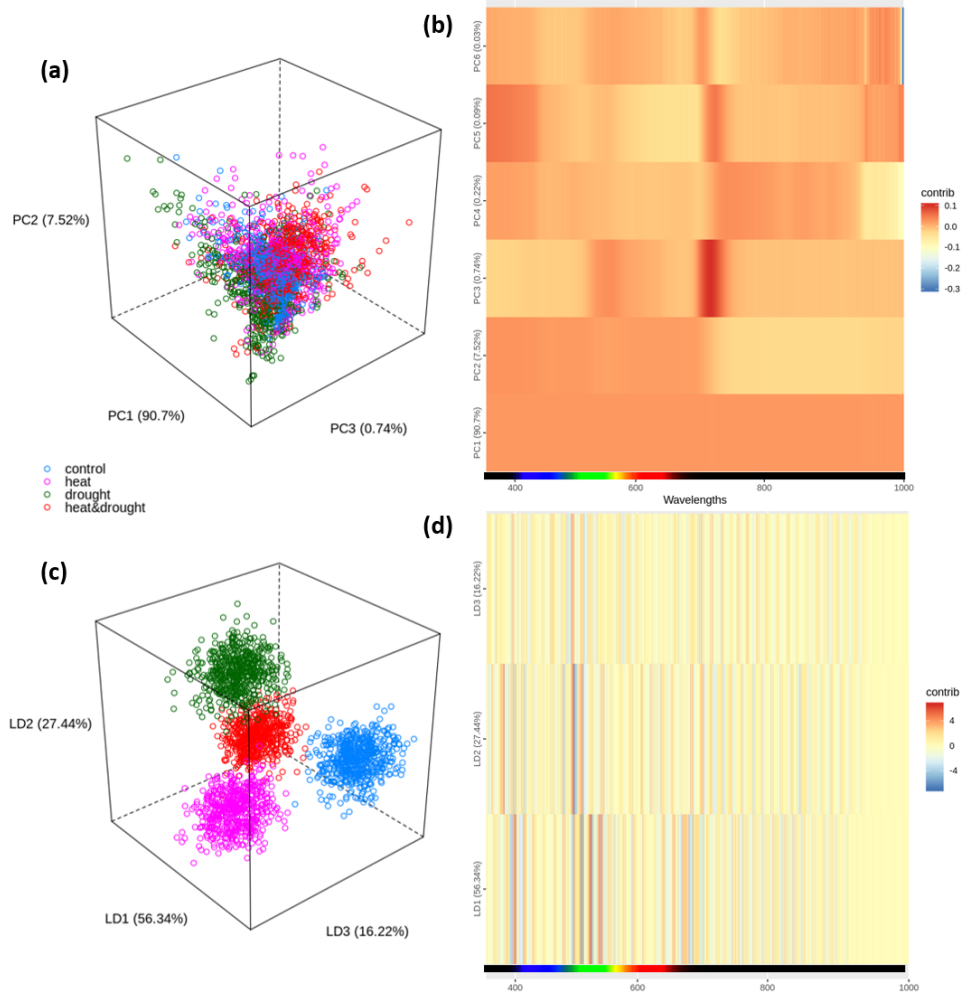


Figure 5: Feature dimensionality reduction by PCA and LDA

(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. Surprisingly, PC1 explains 90.7% of the variance, and the contribution of each wavelength seems to con-

344 tribute equally to it, possibly due to the systematic errors. In PC2, the visible light
 345 region contributes a larger variance, and there are two specific narrow bands in
 346 PC3 that are positively correlated with it. On the other hand, LDA, the supervised
 347 dimension reduction method, can completely separate the different treatments.
 348 In these components, the green (around $520nm$) and red edge (around $680nm$)
 349 wavelength may be important to separate each stress treatment in the principal
 350 components.

351 3.4 Feature Selction

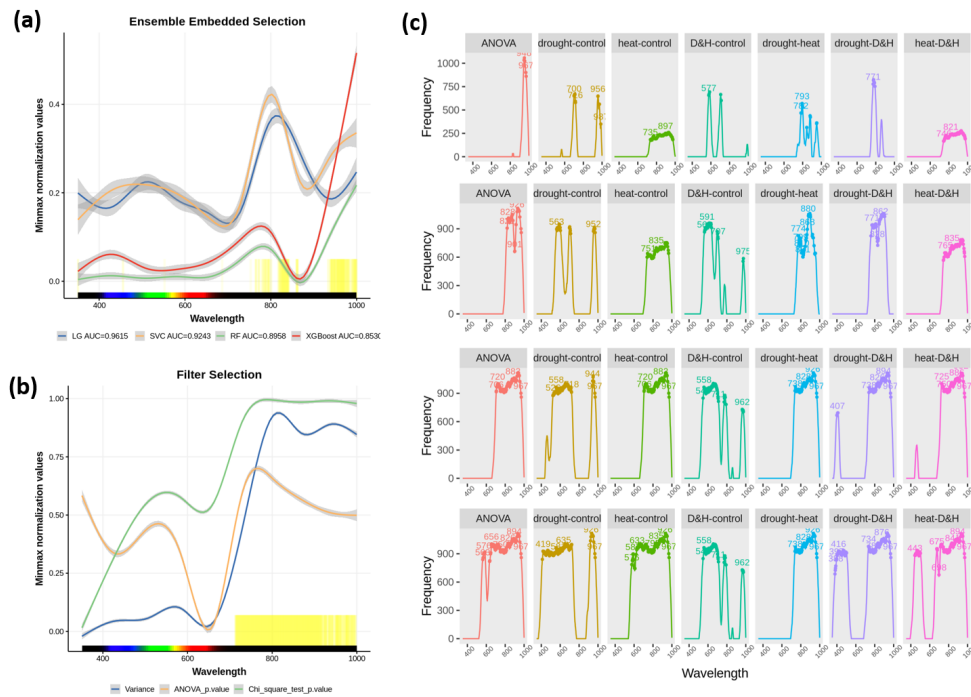


Figure 6: Feature Selection

In (a) and (b), the gray shadows show the confidence intervals, and the yellow vertical lines indicate the important features upon the setting threshold. (c) 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.

352 Through further data mining, we try to use different methods to select the wave-
 353 lengths that can effectively distinguish different treatments. Filter and Embedded
 354 method are commonly used feature selection methods in machine learning. The
 355 results of the Filter show that the variance, chi-square test and the ANOVA of each
 356 wavelength are coincident in the NIR region. And through the modified ensem-
 357 ble method, the correlation coefficient or feature importance obtained by fitting

the LG, SVM, RF and XGBoost are used to establishing the threshold according to their respective goodness of fitting, results show that two NIR fragments and other fragments in some bands may be important.

However, although the results obtained by the above methods are in good agreement with our existing knowledge, the detail of features importance in the pairwise relationship of different treatment remain unclear, also for embedded method, the generalization ability could be constrained by the model. Moreover, in practical application, bandwidth is usually limited to a certain value, continuous band selection is more useful. So here, we create a random search box sorting method. Through simple statistical analysis, it can dynamically display the feature importance under specific bandwidth, and detail the feature importance between different treatments which can provide a feature engineering basis for more sophisticated and complex models.

3.5 Models

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

Results auc shows that in the fitting of the original data, the linear classifier LG and linear SVM fit well, while tree-based model performance is general, mainly because the number of features is too large. which makes the tree model easier to overfit. After removing the noise from PCA, the performance of the tree model is improved obviously, but the effect of LDA is not ideal. Furthermore, through the result of feature selection in the previous step, the models were trained at 409 ~ 429nm, 558 ~ 577nm and 700 ~ 896nm bands, and the performance of the tree model is further improved.

Table 3: ROC_AUC of 4 machine learning models under different data

	Logit	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

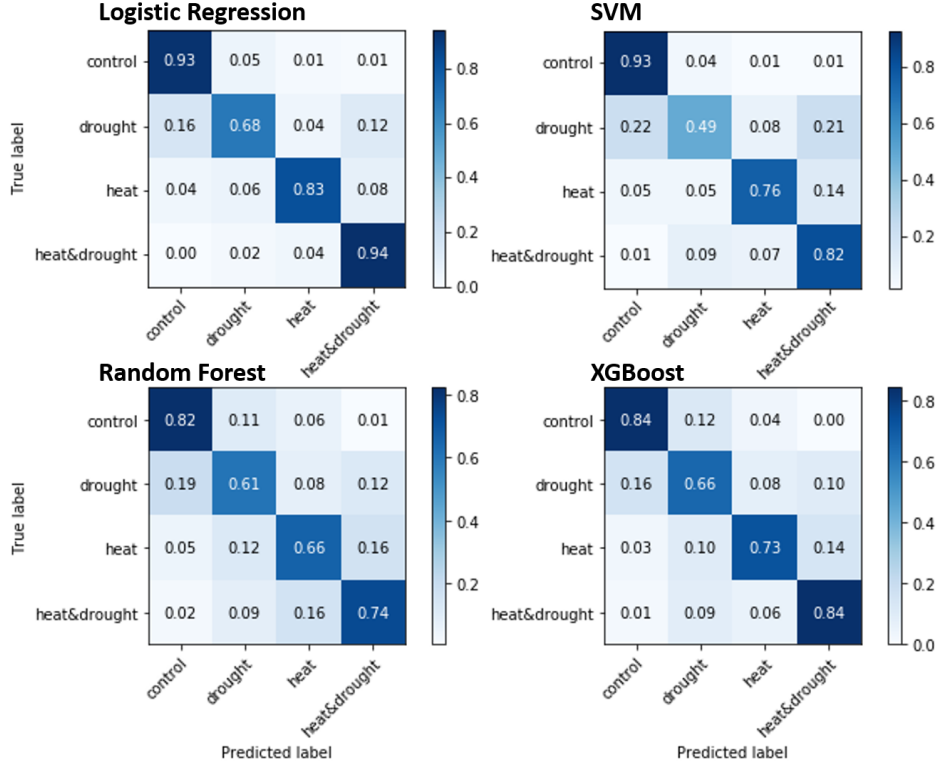


Figure 7: The confusion matrix

Next, we observe their confusion matrix when they have their best AUC. The prediction performance of the four models is similar for different treatments. Among them, the control and combined stress treatments have the best predictive effect, followed by the heat stress, while the drought stress is the worst, and its mispredictions tend to appear more in control and combined stress.

3.6 Shelf Life Prediction

To predict the shelf life on the broccoli packaging conveyor, we collected thousands of broccoli spectral images (Figure 8,(a)) with high-speed band-specific filter cameras and specific band LED lights, and then we marked 5708 broccoli

395 images to tell whether the broccoli head was in the middle of the lens. And
 396 through data augmentation, in the case of limited computing resource, a simple
 397 classifier can be constructed by transfer learning of ResNext101_64 convolution
 398 neural network, which can easily achieve an accuracy of 97.2%. By setting a lower
 399 threshold of predicting probability, then we can select the highest probability of
 400 broccoli head images between the two valleys to track broccolis. After that, for the
 401 sake of removing the influence of the background, we labeled 492 masks for broc-
 402 coli images segmentation, the Unet was trained to reach 99.2% accuracy. Next,
 403 through the SIFT operator, the images of different channels can be aligned for
 404 subsequent analysis.

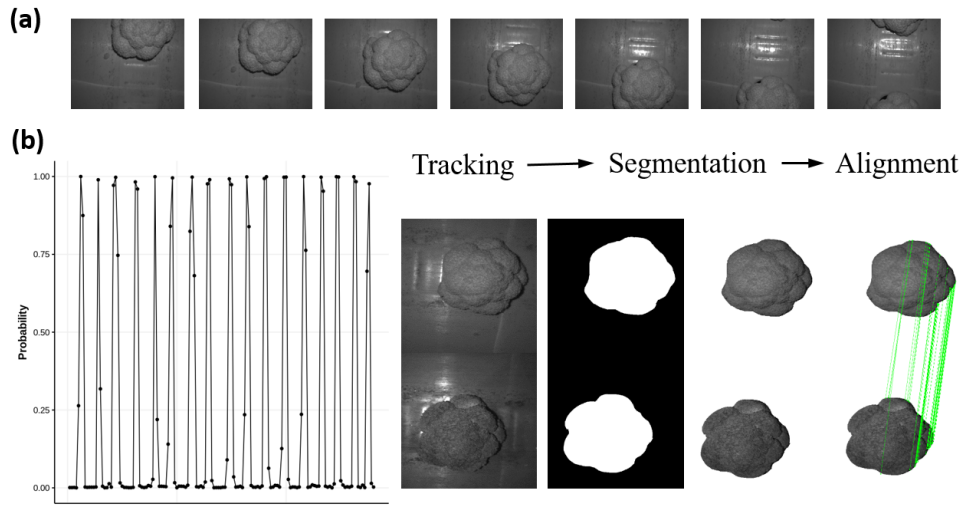


Figure 8: Processing of broccoli head images on conveyor belt

(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. If it did, then segment by Unet and matched by SIFT.

405 Next, we placed the new and stored broccoli heads at room temperature, let
 406 them decay naturally to capture the shelf-life related signals. With time elapsing,
 407 the reflectivity of the three selected channels can somewhat reflect the rotting
 408 changes. However, our focus is more on the first two days before the broccoli
 409 heads showed obvious decay phenotype. Among the three selected channels, the
 410 newly harvested broccoli heads could not be effectively distinguished from the
 411 stored one. Some vegetation indices were also tried, but more research is needed.

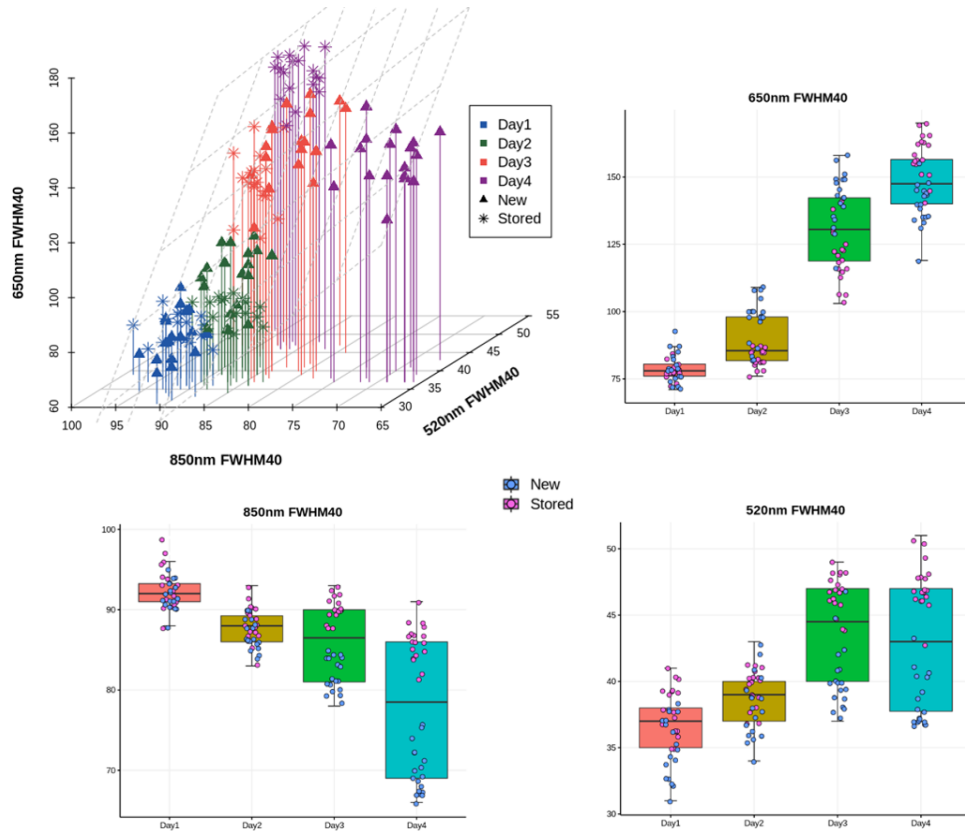


Figure 9: Spectral reflection signal of three channels when the broccoli heads naturally rot

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.

4 DISCUSSION

From the results of various VIs calculated in different treatments, it seems not surprising that not getting a unified index that can effectively distinguish these four treatments, because of the complex relationship among them. However, detailed analysis can still give us some inspiration. Firstly, indices based on red and near-infrared such as RVI, NDVI and EVI, can somewhat reflect the physiological characteristics of plants, however, they are more used in the field of remote sensing for various kinds of local, regional, and global scale models, including general circulation and biogeochemical models (Peterson et al., 1988, Huete et al., 2002). While in CVI, CI-G, CI-RE these indices of leaf chlorophyll content (Gitelson et al., 2003, Vincini et al., 2008), they seem to have a consistent result, except CI-RE, which

423 doesn't consider the green channel. Results show that in the water-deficient envi-
424 ronment (drought, heatdrought), the chlorophyll content of broccoli leaves was
425 significantly affected. This water stress phenomenon has been extensively sup-
426 ported, and it's also well understood that lack of water hinders nutrient transport
427 in plants and thus affects photosynthetic pigments synthesis. As for the water
428 index (WI), it reflects water absorption in the mesophyll andd had been shown
429 to have a good indication of water content in many crops(Wang et al., 2015, ?,
430 ?). The result here is also very satisfactory, it can be seen that there were sig-
431 nificant differences between the two treatments except drought and combined
432 stress. This indicates that in heat stress, the water content of the plant leaves is
433 also affected, even they are well watered, and this effect is not sufficient to su-
434 perimpose the significance in the combined stress relative to the drought stress.
435 The functional basis of the PRI is related on its sensitivity to rapid changes in
436 carotenoids through the de-epoxidation of the xanthophyll pigments(?). It can
437 serve as an indirect means for water stress detection due to the effects of water
438 stress on the efficiency of photosynthesis. Researchers have demonstrated the
439 sensitivity of PRI to short-term crop water stress detection(??Zarco-Tejada et al.,
440 2013), and to the long-term change of carotenoid/chlorophyll ratio(?). Here we
441 found that PRI is sensitive to the combination of drought and heat stress in broc-
442 coli, which may imply the cumulative effect of stress on photosynthetic efficiency
443 and pigments. And Of course, more studies is needed to support these inference.

444

445 In the process of machine learning model training for hyperspectral data, the
446 results of feature engineering show that the important features are mainly con-
447 centrated in the green and near infrared. In detail, because the hyperspectral
448 bandwidth is small,there will inevitably be a lot of redundant information. Data
449 dimensionality reduction can effectively extract important information, shorten
450 model training time and reduce over-fitting. Here, dimensionality reduction by
451 PCA can effectively de-correlate and remove the linear relationship between di-
452 mensions, but it does not consider the classification information. Therefore, after
453 dimensionality reduction, the loss of information will be minimized, but classifi-
454 cation may become more difficult. The data points in the graph are not scattered,
455 and from the contribution of each component to the principal component, it is

456 true. It's hard to get useful information. Here, dimensionality reduction by PCA
457 can effectively de-correlate and remove the linear relationship between dimen-
458 sions, but it does not consider the classification information. Therefore, after
459 dimensionality reduction, the loss of information can reduce to the lowest level,
460 but it may not be helpful in classification ???. The sample points in the graph are
461 not centralized, and from the contribution of each wavelength to the principal
462 components, it is really difficult to obtain useful information. Another commonly
463 used dimension reduction method is LDA, which seeks to distinguish data points
464 as easily as possible after dimension reduction. After dimensionality reduction,
465 the sample data has the largest inter-class distance and the smallest intra-class
466 variance in the new dimension space, and the data has the best separability in
467 the low dimension space ???. It can almost reach 100% classification, so that the
468 contribution of each wavelength may better explain the important features for
469 classification, here are the green and red edges.

470

471 Then, as for feature selection, several methods of experimentation have a good
472 consistency, that is, infrared band information might be relatively important for
473 classification. It may be suggested that changes in mesophyll cell structure, such
474 as membrane structure, are more likely to affect spectrum reflection in leaves
475 under heat and drought stress, while changes in pigments are less important to
476 distinguish between them. In particular, through the dynamic visualization of the
477 random feature search box, we can more clearly see the importance of features
478 to the relationship between them. The water stress and control group showed the
479 most significant difference around 700 ~ 900nm, which was basically consistent
480 with the WI. Wavelength around 700 ~ 800nm may be important for distinguish-
481 ing between heat stress and control. As for combined stress, it seems similar
482 to the water stress. And between the combined stress and the individual stress,
483 there is a difference around 420nm.

484

485 Finally, the results of the four machine learning models show that linear clas-
486 sifier performs well when the data dimension is relatively large, while the tree
487 model does not. This is understandable because regularization is used in training
488 linear classifiers, which can effectively deal with multiple collinearity problems

489 and reduce the weight of redundant information. The tree model can also be im-
490 proved after dimensionality reduction. According to the statistical analysis results,
491 we empirically select 409-429, 557-558, 700-896 nm band information for train-
492 ing, so we can see that the performance of the XGBoost model has been improved
493 effectively, its AUC can reach 0.9096. By showing the confusion matrix, it is not
494 surprising that the control group and the combined stress group can achieve the
495 best distinction. But surprisingly, the heat stress group can be more effectively
496 distinguished from other stresses, as opposed to the drought group which may
497 have more phenotype. The erroneous distinction of drought group mostly ap-
498 pears in the difference to control group, which may be explained to some extent
499 that some of the leaves do not reach the threshold at which the drought can be
500 detected.

501

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

510

511

512 References

- 513 Aladenola, O. and Madramootoo, C. (2014), 'Response of greenhouse-grown bell
514 pepper (*capsicum annuum* l.) to variable irrigation', *Canadian journal of plant*
515 *science* **94**(2), 303–310.
- 516 Alchanatis, V., Cohen, Y., Cohen, S., Moller, M., Sprinstin, M., Meron, M., Tsipris,
517 J., Saranga, Y. and Sela, E. (2010), 'Evaluation of different approaches for esti-
518 mating and mapping crop water status in cotton with thermal imaging', *Preci-*
519 *sion Agriculture* **11**(1), 27–41.

- 520 Anjum, S., Wang, L., Farooq, M., Hussain, M., Xue, L. and Zou, C. (2011), 'Brassi-
521 nolide application improves the drought tolerance in maize through modulation
522 of enzymatic antioxidants and leaf gas exchange', *Journal of Agronomy and Crop*
523 *Science* **197**(3), 177–185.
- 524 Bellvert, J., Marsal, J., Girona, J., Gonzalez-Dugo, V., Fereres, E., Ustin, S. and
525 Zarco-Tejada, P. (2016), 'Airborne thermal imagery to detect the seasonal evolu-
526 tion of crop water status in peach, nectarine and saturn peach orchards', *Remote*
527 *Sensing* **8**(1), 39.
- 528 Berni, J. A., Zarco-Tejada, P. J., Suárez, L. and Fereres, E. (2009), 'Thermal and
529 narrowband multispectral remote sensing for vegetation monitoring from an
530 unmanned aerial vehicle', *IEEE Transactions on geoscience and Remote Sensing*
531 **47**(3), 722–738.
- 532 Bravo, C., Moshou, D., West, J., McCartney, A. and Ramon, H. (2003), 'Early dis-
533 ease detection in wheat fields using spectral reflectance', *Biosystems Engineering*
534 **84**(2), 137–145.
- 535 Camejo, D., Jiménez, A., Alarcón, J. J., Torres, W., Gómez, J. M. and Sevilla,
536 F. (2006), 'Changes in photosynthetic parameters and antioxidant activities
537 following heat-shock treatment in tomato plants', *Functional Plant Biology*
538 **33**(2), 177–187.
- 539 Cao, X., Luo, Y., Zhou, Y., Duan, X. and Cheng, D. (2013), 'Detection of powdery
540 mildew in two winter wheat cultivars using canopy hyperspectral reflectance',
541 *Crop Protection* **45**, 124–131.
- 542 Chen, T. and Guestrin, C. (2016), Xgboost: A scalable tree boosting system, in
543 'Proceedings of the 22nd acm sigkdd international conference on knowledge
544 discovery and data mining', ACM, pp. 785–794.
- 545 Dangwal, N., Patel, N., Kumari, M. and Saha, S. (2016), 'Monitoring of water
546 stress in wheat using multispectral indices derived from landsat-tm', *Geocarto*
547 *International* **31**(6), 682–693.
- 548 Din, J., Khan, S., Ali, I., Gurmani, A. et al. (2011), 'Physiological and agronomic
549 response of canola varieties to drought stress', *J Anim Plant Sci* **21**(1), 78–82.

550 Feilhauer, H., Asner, G. P. and Martin, R. E. (2015), 'Multi-method ensemble se-
 551 lection of spectral bands related to leaf biochemistry', *Remote Sensing of Envi-*
 552 *ronment* **164**, 57–65.

553 Gates, D. M., Keegan, H. J., Schleter, J. C. and Weidner, V. R. (1965), 'Spectral
 554 properties of plants', *Applied optics* **4**(1), 11–20.

555 Gitelson, A. A., Gritz, Y. and Merzlyak, M. N. (2003), 'Relationships between leaf
 556 chlorophyll content and spectral reflectance and algorithms for non-destructive
 557 chlorophyll assessment in higher plant leaves', *Journal of plant physiology*
 558 **160**(3), 271–282.

559 Howard, J. et al. (2018), 'fastai', <https://github.com/fastai/fastai>.

560 Huete, A., Didan, K., Miura, T., Rodriguez, E. P., Gao, X. and Ferreira, L. G.
 561 (2002), 'Overview of the radiometric and biophysical performance of the modis
 562 vegetation indices', *Remote sensing of environment* **83**(1-2), 195–213.

563 Idso, S., Jackson, R., Pinter Jr, P., Reginato, R. and Hatfield, J. (1981), 'Normaliz-
 564 ing the stress-degree-day parameter for environmental variability', *Agricultural*
 565 *meteorology* **24**, 45–55.

566 Jordan, C. F. (1969), 'Derivation of leaf-area index from quality of light on the
 567 forest floor', *Ecology* **50**(4), 663–666.

568 Knipling, E. B. (1970), 'Physical and physiological basis for the reflectance of visi-
 569 ble and near-infrared radiation from vegetation', *Remote sensing of environment*
 570 **1**(3), 155–159.

571 Kobayashi, T., Kanda, E., Kitada, K., Ishiguro, K. and Torigoe, Y. (2001), 'De-
 572 tection of rice panicle blast with multispectral radiometer and the potential of
 573 using airborne multispectral scanners', *Phytopathology* **91**(3), 316–323.

574 Kollist, H., Zandalinas, S. I., Sengupta, S., Nuhkat, M., Kangasjärvi, J. and Mittler,
 575 R. (2018), 'Rapid responses to abiotic stress: priming the landscape for the
 576 signal transduction network', *Trends in plant science* .

- 577 Lawlor, D. W. and Cornic, G. (2002), 'Photosynthetic carbon assimilation and
578 associated metabolism in relation to water deficits in higher plants', *Plant, cell
579 & environment* **25**(2), 275–294.
- 580 Lowe, D. G. et al. (1999), Object recognition from local scale-invariant features.,
581 in 'iccv', Vol. 99, pp. 1150–1157.
- 582 McClung, C. R. and Davis, S. J. (2010), 'Ambient thermometers in plants: from
583 physiological outputs towards mechanisms of thermal sensing', *Current Biology*
584 **20**(24), R1086–R1092.
- 585 Mestre, H. (1935), 'The absorption of radiation by leaves and algae', *Cold Spring
586 Harbor Symposia on Quantitative Biology* pp. 191–209.
- 587 Mittler, R., Finka, A. and Goloubinoff, P. (2012), 'How do plants feel the heat?',
588 *Trends in biochemical sciences* **37**(3), 118–125.
- 589 Muir, A., Porteous, R. and Wastie, R. (1982), 'Experiments in the detection of
590 incipient diseases in potato tubers by optical methods', *Journal of Agricultural
591 Engineering Research* **27**(2), 131–138.
- 592 Panigada, C., Rossini, M., Meroni, M., Cilia, C., Busetto, L., Amaducci, S.,
593 Boschetti, M., Cogliati, S., Picchi, V., Pinto, F. et al. (2014), 'Fluorescence, pri
594 and canopy temperature for water stress detection in cereal crops', *International
595 Journal of Applied Earth Observation and Geoinformation* **30**, 167–178.
- 596 Pearson, R. L. and Miller, L. D. (1972), Remote mapping of standing crop biomass
597 for estimation of the productivity of the shortgrass prairie, in 'Remote sensing
598 of environment, VIII', p. 1355.
- 599 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
600 Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A.,
601 Cournapeau, D., Brucher, M., Perrot, M. and Duchesnay, E. (2011), 'Scikit-learn:
602 Machine learning in Python', *Journal of Machine Learning Research* **12**, 2825–
603 2830.
- 604 Peterson, D. L., Aber, J. D., Matson, P. A., Card, D. H., Swanberg, N., Wessman, C.
605 and Spanner, M. (1988), 'Remote sensing of forest canopy and leaf biochemical
606 contents', *Remote Sensing of Environment* **24**(1), 85–108.

- 607 Platt, J. C. (1998), Sequential minimal optimization: A fast algorithm for training
608 support vector machines, Technical report, ADVANCES IN KERNEL METHODS
609 - SUPPORT VECTOR LEARNING.
- 610 Prasannakumar, N., Chander, S., Sahoo, R. and Gupta, V. (2013), ‘Assessment of
611 brown planthopper, (*nilaparvata lugens*) [stål], damage in rice using hyperspec-
612 tral remote sensing’, *International journal of pest management* **59**(3), 180–188.
- 613 Rabideau, G. S., French, C. S. and Holt, A. S. (1946), The absorption and reflec-
614 tion spectra of leaves, chloroplast suspensions, and chloroplast fragments as
615 measured in an ulbricht sphere.
- 616 Ronneberger, O., Fischer, P. and Brox, T. (2015), U-net: Convolutional networks
617 for biomedical image segmentation, in ‘International Conference on Medical
618 image computing and computer-assisted intervention’, Springer, pp. 234–241.
- 619 Rossini, M., Fava, F., Cogliati, S., Meroni, M., Marchesi, A., Panigada, C., Gia-
620 rdino, C., Busetto, L., Migliavacca, M., Amaducci, S. et al. (2013), ‘Assessing
621 canopy pri from airborne imagery to map water stress in maize’, *ISPRS Journal*
622 *of Photogrammetry and Remote Sensing* **86**, 168–177.
- 623 Rouse Jr, J., Haas, R., Schell, J. and Deering, D. (1974), ‘Monitoring vegetation
624 systems in the great plains with erts’.
- 625 Sinclair, T. R. et al. (1968), ‘Pathway of solar radiation through leaves’.
- 626 Smith, L. N. (2017), Cyclical learning rates for training neural networks, in ‘2017
627 IEEE Winter Conference on Applications of Computer Vision (WACV)’, IEEE,
628 pp. 464–472.
- 629 Suzuki, N., Koussevitzky, S., Mittler, R. and Miller, G. (2012), ‘Ros and redox
630 signalling in the response of plants to abiotic stress’, *Plant, Cell & Environment*
631 **35**(2), 259–270.
- 632 Tucker, C. J. (1979), ‘Red and photographic infrared linear combinations for mon-
633 itoring vegetation’, *Remote sensing of Environment* **8**(2), 127–150.
- 634 Vincini, M., Frazzi, E. and D’Alessio, P. (2008), ‘A broad-band leaf chlorophyll
635 vegetation index at the canopy scale’, *Precision Agriculture* **9**(5), 303–319.

- 636 Wang, X., Zhao, C., Guo, N., Li, Y., Jian, S. and Yu, K. (2015), ‘Determining the
637 canopy water stress for spring wheat using canopy hyperspectral reflectance
638 data in loess plateau semiarid regions’, *Spectroscopy Letters* **48**(7), 492–498.
- 639 Willstätter, R. and Mieg, W. (1907), ‘Untersuchungen über chlorophyll; iv. ue-
640 ber die gelben begleiter des chlorophylls’, *Justus Liebigs Annalen der Chemie*
641 **355**(1), 1–28.
- 642 Xie, S., Girshick, R., Dollár, P., Tu, Z. and He, K. (2016), ‘Aggregated residual
643 transformations for deep neural networks’, *arXiv preprint arXiv:1611.05431* .
- 644 Xue, J. and Su, B. (2017), ‘Significant remote sensing vegetation indices: A review
645 of developments and applications’, *Journal of Sensors* **2017**.
- 646 Yang, C.-M. (2010), ‘Assessment of the severity of bacterial leaf blight in rice using
647 canopy hyperspectral reflectance’, *Precision Agriculture* **11**(1), 61–81.
- 648 Zarco-Tejada, P. J., González-Dugo, V. and Berni, J. A. (2012), ‘Fluorescence,
649 temperature and narrow-band indices acquired from a uav platform for wa-
650 ter stress detection using a micro-hyperspectral imager and a thermal camera’,
651 *Remote sensing of environment* **117**, 322–337.
- 652 Zarco-Tejada, P. J., González-Dugo, V., Williams, L., Suárez, L., Berni, J. A.,
653 Goldhamer, D. and Fereres, E. (2013), ‘A pri-based water stress index combin-
654 ing structural and chlorophyll effects: Assessment using diurnal narrow-band
655 airborne imagery and the cwsí thermal index’, *Remote sensing of environment*
656 **138**, 38–50.

657 **5 Supplementary Information**

658 **5.1 Experiment Apparatus**