

Parameter Selection in Weed/Crop Discrimination

Evan McGinnis
Graduate Student, Biosystems
Engineering Dept
INFO 529
evanmc@email.arizona.edu

Abstract

This paper is an analysis of various parameters extracted from vegetation images. These parameters are used in KNN, Logistic Regression, Random Forest, and Support Vector Machine to discriminate between crop (cultivars of lettuce) and weeds.

Contents

1	Introduction	2
2	Image Segmentation	2
3	Feature Extraction	4
3.1	Length Width Ratio	4
3.2	Shape Index	5
3.3	Normalized Distance from Cropline	5
3.4	Hue	5
3.5	Saturation	5
3.6	YIQ	6
3.7	Compactness	6
3.8	Elongation	7
3.9	Eccentricity	7
3.10	Roundness	7
3.11	Convexity	7
3.12	Solidity	8
4	Feature Selection	8
4.1	Univariate	8
4.2	Low Variance	8
4.3	Feature Importance	8
4.4	Recursive Elimination	9
4.5	Principal Component Analysis	9
4.6	Feature Importance	9
5	Crop/Weed Discrimination	10
5.1	KNN	10
5.2	Logistic Regression	10
5.3	Support Vector Machine	10
5.4	Random Forest	10
5.5	Boosted Gradient	10
6	Conclusions	10
7	References	11

List of Figures

1	Before and after segmentation	4
2	Normalized Distance to Cropline (Source: author)	5
3	YIQ (Source: [5])	6
4	Examples of compactness Source: [Wirth2004-l]	6
5	Examples of elongation Source: [Wirth2004-l]	7
6	Illustration of major and minor axis Source: [Wirth2004-l]	7
7	Illustration of convexity Source: [Wirth2004-l]	8
8	Illustration of solidity Source: [Wirth2004-l]	8

1 Introduction

This paper is part of a larger effort to detect weeds in imagery gathered under field conditions and treat those weeds in real-time. More precisely, weed detection with an image must be complete with 100s of milliseconds to treat weeds at commercially viable speeds. While such concerns are beyond the scope of this paper, we will instead focus on selecting object parameters that result in the most accurate predictions of the classification of vegetation into one of two categories: weed or crop. It is of no concern what specific species of weed or crop cultivar – all that is important here is a binary classification of weed or not weed. The dataset examined in this paper was gathered in an agricultural research field in Yuma, Arizona by capturing images while walking beside a planting bed. The processing of this image set follows this procedure:

1. Segment images (remove pixels that are not vegetation)
2. Label vegetation as weed or crop
3. Extract various parameters from vegetation images
4. Evaluate significance of parameters to assigned class
5. Evaluate predictions with various models using selected parameters

2 Image Segmentation

The portions of the images that did not contain pixels with vegetation present were discarded. These images were segmented using various visible light indices [6]. As this process is not the primary subject of this paper, it will be given only superficial mention here. Various approaches to image segmentation are summarized in Table 1.

Table 1: Visible light indices ([1, 2])

Index	Formula	Comment
Triangular Greeness	$R_{green} - \alpha R_{red} - \beta R_{blue}$ $\alpha = \frac{2(\lambda_{blue} - \lambda_{green})}{(\lambda_{blue} - \lambda_{red})}$ $\beta = \frac{2(\lambda_{green} - \lambda_{red})}{(\lambda_{blue} - \lambda_{red})}$	Corrects for camera calibration using the peak sensitivity
Normalized Difference	$128 * \left(\left(\frac{(G-R)}{(G+R)} \right) + 1 \right)$	The NDI index produces a near-binary image.
Excess Green	$R = \frac{R}{R_{max}}$ $G = \frac{G}{G_{max}}$ $B = \frac{B}{B_{max}}$	ExG provided a clear contrast between plants and soil
Excess Red	$1.3R - G$	inspired by the fact that there are 4% blue, and 32% green, compared with 64% red cones in the retina of the human eye
Color Index of Vegetation Extraction	$0.441R - 0.811G + 0.385B + 18.78745$	This method was proposed to separate green plants from soil background in order to evaluate the crop growing status.
Excess Green - Excess Red	$ExG - ExR$	ExG used to extract the plant region and ExR used to eliminate the background noise (soil and residue) where green-red material (stems, branches, or petioles) may exist
Normalized Green-Red Difference	$\frac{(G-R)}{(G+R)}$	The method of NGRDI was used to overcome the differences in exposure settings selected by the digital camera when acquiring aerial photography of the field.
Vegetative Index	$\frac{G}{R^a B^{(1-a)}}, a = 0.667$	VEG has a significant advantage because it is robust to lighting change.
Com1	$ExG + CIVE + ExGR + VEG$	TODO
Modified Excess Green	$1.262G - 0.884R = 0.311B$	TODO
Combined Indices 2	$0.36ExG + 0.47CIVE + 0.17VEG$	Uses weighting factors to emphasize strengths of various approaches

These indices are used to create a mask that is then applied to the original source image to permit vegetation to show while masking details that are not relevant (ground pixels, stones, and other items that may appear in field conditions) The intent here is to remove all pixels that are not relevant to the task of distinguishing between crop and weed.¹

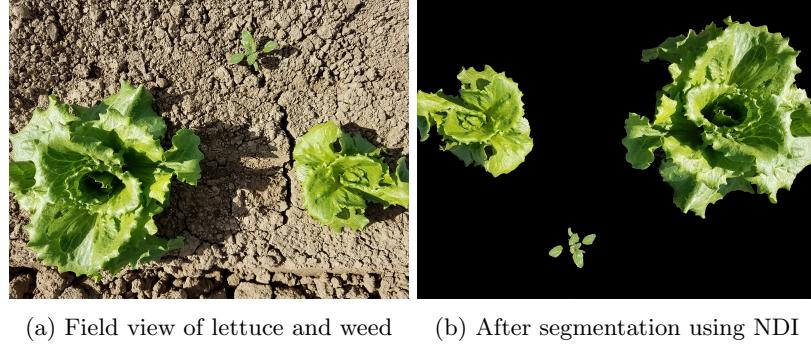


Figure 1: Before and after segmentation

The segmented image has discarded ground pixels while retaining most of the pixels that will be used, but a close examination reveals that pixels in the stems of the weed are also eliminated, as they are less green than the rest of the plant. While they are not eliminated, pixels in the area of the deep shadows of the vegetation may affect attempts to classify objects based on color attributes. It is also envisioned that in the field images will be acquired under controlled, not ambient lighting conditions. For the purposes of this paper, images will use the NDI segmentation approach.

3 Feature Extraction

The library used for feature extraction is the OpenCV toolkit, so the discussion of various features may inadvertently slip into using OpenCV terms. Some basic shape descriptors are used below, and the most fundamental ones that apply here are the bounding box, the rectangular box that completely encloses the object, and the convex hull (and convex area), the smallest set of straight lines that completely contains an object. The concept of a *centroid* is also important to this discussion. A centroid is the center of mass of an object, and a concept that will be used in this analysis.

The segmented images produced are then processed by a first identifying the separate objects (often called blobs) within the image and then computing various aspects of each of those objects. Weeds or crop may tend to exhibit these features to an extent that they can be used to distinguish between the two. Weeds, for instance, may be more elongated than crop, or may tend to have saturation differences that may not be readily apparent to the casual observer.

3.1 Length Width Ratio

The ratio of width to length is not – as the name might imply – a simple ratio, but is expressed as:

$$S = \begin{bmatrix} Var(X) & Cov(XY) \\ Cov(XY) & Var(Y) \end{bmatrix}, \lambda = \frac{eig_1(S)}{eig_2(S)}$$

¹Observers will note that the segmentation photo is rotated. This bears further investigation, but will have no impact on the features examined here.

Where $eig_1(S)$ and $eig_2(S)$ are the maximum eigenvalues of the matrix S , with λ representing the ratio. [3]

3.2 Shape Index

The shape index of an object is a metric expressing the relationship between an object's perimeter and its area.

$$\alpha = \frac{e}{4\sqrt{A}}$$

3.3 Normalized Distance from Cropline

The cropline in a planting is, simply, the line along the bed where crop can be expected. Under field conditions the cropline will appear in the same spot in each photo. The image set here, however, was manually acquired by walking along the crop row and capturing images. Unfortunately, this means that the crop line location will differ from one picture to another. For this image set, the cropline is defined as the line that intersects the centroid of the objects with the largest area. Crops will most often have a distance from the cropline very close to zero. Weeds, on the other hand, may have a distance close to zero if they appear within the line of crop, but often appear far from the crop line. Figure 2 illustrates the concept of a

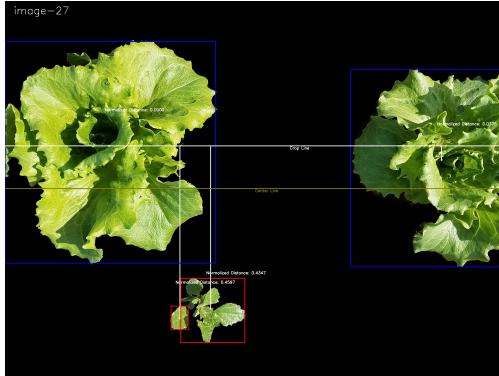


Figure 2: Normalized Distance to Cropline (Source: author)

cropline and the normalized distance of vegetation from it. In this image we see two growths of lettuce that are very close to the cropline (distances here are in pixels, but the units are not significant. This could be expressed in millimeters) at 0 and 0.375 and a weed lying 0.4347 units from the cropline. The line marked *Center Line* is for reference purposes and can be ignored for now. There are two additional items that are worth noting about this image: the dots connecting the plant to the cropline are the *centroids* mentioned earlier, and the colored bounding boxes signify the class of the object, something we will return to in a later section.

3.4 Hue

The hue of an object is define TODO: Insert an actual definition. In this case, the image is converted to the Hue Saturation and Intensity colorspace and the mean value for the hue is taken.

3.5 Saturation

In this case, the image is converted to the Hue Saturation and Intensity colorspace and the mean value for the saturation is taken.

3.6 YIQ

The YIQ model of color us used by the NTSC color TV system. Y represents the luma information, I and Q the chrominance information. The processing employed here is to convert the image to the YIQ color space and take the mean value for the I, or in-phase component. The I component is the feature of interest here,

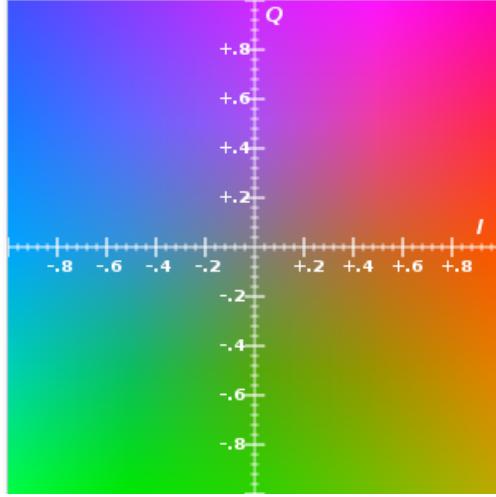


Figure 3: YIQ (Source: [5])

and conversion of RGB to YIQ is achieved with this transformation:

$$\begin{bmatrix} Y \\ I \\ Q \end{bmatrix} \approx \begin{bmatrix} 0.299 & 0.587 & 0.114 \\ 0.5959 & -0.2746 & -0.3213 \\ 0.2115 & -0.5227 & 0.3112 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix}$$

3.7 Compactness

The compactness of an object is defined as the ratio of its area to the area of a circle with the same perimeter as the original object, and is given by this equation [6]:

$$\frac{4\pi \text{ area}}{\text{perimeter}^2}$$

The most compact object is a circle, whose value is computed as 1. Objects with irregular boundaries will have values larger than 1.



Figure 4: Examples of compactness Source: [Wirth2004-1]

3.8 Elongation

Elongation is the ratio of the length to the width of an object's bounding box [6]:

$$\frac{width_{bounding}}{length_{bounding}}$$

This produces a metric between 0 (more elongated) and 1 (roughly circular or square).

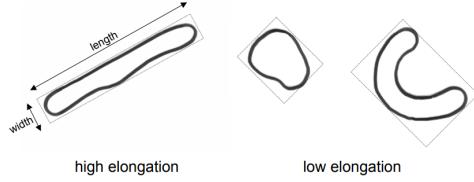


Figure 5: Examples of elongation Source: [Wirth2004-1]

3.9 Eccentricity

The eccentricity (or ellipticity) of an object is the ratio of the length of the minor axis to the length of the major axis.

$$\frac{length_{minor-axis}}{length_{major-axis}}$$

The major axis of an object is expressed as the (x,y) endpoints of the longest line that can be drawn through an object. The minor axis is the longest line that can be drawn through an object while remaining perpendicular to the major axis. [6]

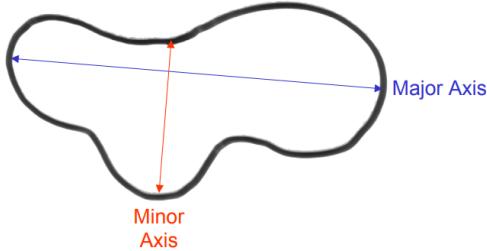


Figure 6: Illustration of major and minor axis Source: [Wirth2004-1]

3.10 Roundness

The roundness of an object is an expression varying between 1 (perfectly circular) and 0 (departure from circularity).

$$\frac{4\pi \text{ area}}{(\text{convex perimeter})^2}$$

3.11 Convexity

The convexity of an object is the amount an object differs from a convex object, expressed as the ratio of an object's convex perimeter to the perimeter [6]:

$$\frac{\text{convex perimeter}}{\text{perimeter}}$$

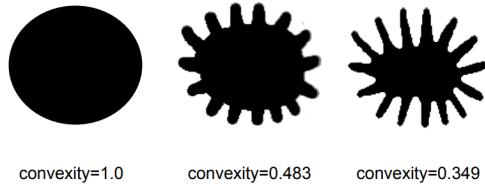


Figure 7: Illustration of convexity Source: [Wirth2004-1]

3.12 Solidity

The solidity of an object varies between 1 (completely solid) and 0, an indication that the object has irregular boundaries.

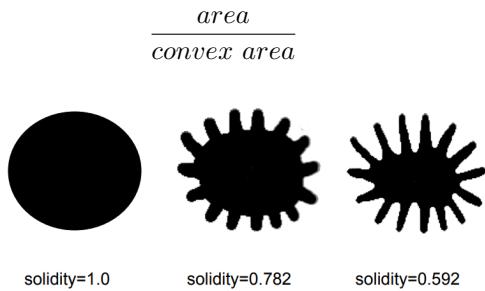


Figure 8: Illustration of solidity Source: [Wirth2004-1]

4 Feature Selection

The features described in the previous section were generated for a set of segmented images, resulting in each object being described by these attributes. Several technique were then used to explore the relationship between these attributes and the labeled class. Specifically, only the most important variables will be selected in the predictions, and those parameters with a weak association will be dropped.

4.1 Univariate

In the univariate scheme, the features with the strongest association with the class are selected. In this scheme, scikit-learn supports a suite of statistical tests, and we will use the ANOVA F-value method for feature evaluation.

4.2 Low Variance

4.3 Feature Importance

The Feature Importance scheme uses Random Forest and Extra Trees approaches to rank features for importance.

Applying this to the dataset created results in these scores, where more important features are assigned higher scores. Table 2 shows the results of the scores assigned to the various features. In this table, we can see that the *YIQ*, the *Normalized Distance*, and the *Compactness* are the most important features. Interestingly, color features play prominently here, far exceeding the scores assigned to structural features such as the *Solidity* of an object. Structural features such as *Eccentricity* have scores low enough that it is likely the case that they could be omitted from the model without a negative impact on accuracy.

Feature	Importance
Length-Width Ratio	0.02998853
Shape Index	0.05357191
Distance	0.05081641
Normalized Distance	0.12969820
Hue	0.00542753
Saturation	0.08911404
YIQ Mean	0.38947503
Compactness	0.09122537
Eccentricity	0.00761004
Roundness	0.04467927
Convexity	0.05499904
Solidity	0.05339462

Table 2: Feature Importance Scores

4.4 Recursive Elimination

The Recursive Feature Elimination scheme works backwards from a full model (all attributes) to remove attributes, build a model, and evaluate the results to characterize a feature's contribution to the prediction.

4.5 Principal Component Analysis

Principal Component Analysis (PCA) involves re-projecting the data (a change in basis) as a technique to reduce the dimensionality of the data [4]

4.6 Feature Importance

Selection	Feature	Length-Width Ratio	Shape Index	Distance	Normalized Distance	Hue	Saturation	YIQ Mean	Compactness	Eccentricity	Roundness	Convexity	Solidity
Univariate	23.6	120.5	70.1	138.6	1.1	134.9	517.8	115.0	0.3	38.0	1.0	34.1	
Variance	2.114	0.008	9150.389	0.027	117.093	1350.474	>0.001	0.0376	0.174	0.174	0.312	0.014	
Recursive	3	2	12	1	9	10	11	4	7	5	8	6	
Importance	0.039	0.070	0.027	0.260	0.033	0.108	0.252	0.087	0.015	0.067	0.011	0.032	

Table 3: Feature Selection using various approaches.

5 Crop/Weed Discrimination

5.1 KNN

5.2 Logistic Regression

5.3 Support Vector Machine

5.4 Random Forest

5.5 Boosted Gradient

6 Conclusions

Table 4: Learning Results

Method	Train	Test
KNN	0	0
Logistic Regression	0	0
SVM	0	0
Random Forest	0	0

7 References

- [1] E Hamuda, M Glavin, and E Jones. “A survey of image processing techniques for plant extraction and segmentation in the field”. In: *Comput. Electron. Agric.* 125 (2016), pp. 184–199. ISSN: 0168-1699. DOI: 10.1016/j.compag.2016.04.024.
- [2] E Raymond Hunt et al. “A visible band index for remote sensing leaf chlorophyll content at the canopy scale”. In: *Int. J. Appl. Earth Obs. Geoinf.* 21 (Apr. 2013), pp. 103–112. ISSN: 0303-2434. DOI: 10.1016/j.jag.2012.07.020.
- [3] F Lin et al. “Detection of corn and weed species by the combination of spectral, shape and textural features”. In: *Sustainability (Switzerland)* 9.8 (2017). ISSN: 2071-1050. DOI: 10.3390/su9081335.
- [4] Andreas C Müller and Sarah Guido. *Introduction to Machine Learning with Python*. 1st ed. O'Reilly Media, Nov. 2016. ISBN: 9781449369415.
- [5] various. *YIQ*. <https://en.wikipedia.org/wiki/YIQ>.
- [6] Michael A Wirth. *Shape Analysis and Measurement*. 2004.