

Development of Inter-Leaves Weed and Plant Regions Identification Algorithm using Histogram of Oriented Gradient and K-Means Clustering

Dheeman Saha

Department of Electrical Engineering
and Computer Science, South Dakota
State University
Brookings, SD, 57007, USA
Dheeman.Saha@sdstate.edu

George Hamer

Department of Electrical Engineering
and Computer Science, South Dakota
State University
Brookings, SD, 57007, USA
George.Hamer@sdstate.edu

Ji Young Lee

Department of Electrical Engineering
and Computer Science, South Dakota
State University
Brookings, SD, 57007, USA
Jiyoung.Lee@sdstate.edu

ABSTRACT

This paper proposes a weed detection mechanism, where the carrot leaves are segmented from the weeds (mostly *Chamomile*). In the early stage, both weeds and carrot leaves are intermixed with each other and have similar color texture. This makes it difficult to identify without the help of the domain experts. Therefore, it is essential to remove the weed regions so that the carrot plants can grow without any interruptions. The process of identifying the weeds become more challenging when both plant and weed regions overlap (inter-leaves). The proposed method takes account of this problem and breaks down the identification mechanism into three major components: Image Segmentation, Feature Extraction, and Classification. In the Image Segmentation stage, K-Means clustering is applied to select the images that will be used for the identification purpose. Next, in the Feature Extraction stage structural information of the weed and leaves will be extracted from the lower unit images. Furthermore, to extract the information from the Region of Interest (ROI), Histogram of Oriented Gradient (HoG) is used to locate and label all the weed and carrot leaves regions. In the Classification stage, Support Vector Machine (SVM) analyzes all the information and labels the regions. This method of weed detection is effective as it automates the identification process and fewer herbicides will be used, which in-turn benefits the environment. The proposed method successfully classifies the plant regions at a success rate of 92% using an open dataset and outperformed some of the previous approaches.

CCS Concepts

•Computing methodologies → Feature selection; Support vector machines;

Keywords

Image Segmentation; K-Means Clustering; Automatic Greenness Identification; Histogram of Oriented Gradient (HoG); Support Vector Machine (SVM); Precision Agriculture (PA)

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the Owner/Author.

RACS '17, September 20–23, 2017, Krakow, Poland

© 2017 Copyright is held by the owner/author(s).

ACM ISBN 978-1-4503-5027-3/17/09.

<https://doi.org/10.1145/3129676.3129700>

1. INTRODUCTION

Most research in Precision Agriculture (PA) focuses on weed control management where crop production can be maximized with less use of herbicides. It has been estimated that spot control of weeds has the potential to reduce the amount of chemical applied by as much as 80% [1]. However, these weeds regions are manually identified by farmers, which is a tedious and time-consuming task. So, an automated system needs to be developed where weed control can be performed round the clock and crop production is not interrupted.

Moreover, in agricultural fields, the weed normally grows close-to-crop or between intra-row which need to be regulated to avoid substantial yield loss [2]. Therefore, appropriate classification mechanisms are required to extract information from the overlapping regions.

Furthermore, machine learning algorithms play essential roles to separate the weed leaves from the plants, like SVM [3, 4] is used to identify the weed portions that are situated between rows of crops in a cell-based approach. Individual cells in the grid are analyzed to determine whether to spray herbicide or not. This cell-based analysis is not suitable when the high precision treatment is required as it is computationally inefficient.

Although such approaches successfully identified the weed regions, none of them are suitable for weed identification in the overlapping regions, which are extremely common in real scenarios. If the weed regions around these areas can be located then more crops can be considered for the cultivation purpose. The proposed method is an extension of the previous work [5], where information from the inter-leaves regions are extracted for the classification purpose.

The primary focus in our method is to develop an automated system which identifies the weed regions with minimal help from the domain experts. To accomplish this goal, the following contributions are made in the paper:

- The K-Means clustering algorithm is applied to the input image to determine the image that consists of only plants and weeds.
- The Histogram of Oriented Gradient (HoG) is used a feature extractor, which helps to locate and label the weed and plant regions.
- Features are extracted of both plants and weeds with the help of Morphological Operations.

- The classification of weeds and plants are done with Support Vector Machine (SVM) and can classify the plant leaves at a success rate of 92%.

2. PROPOSED METHOD

This section discusses each of the components that are used for the detection purpose. The breakdown of the three major components is shown in Figure 1.

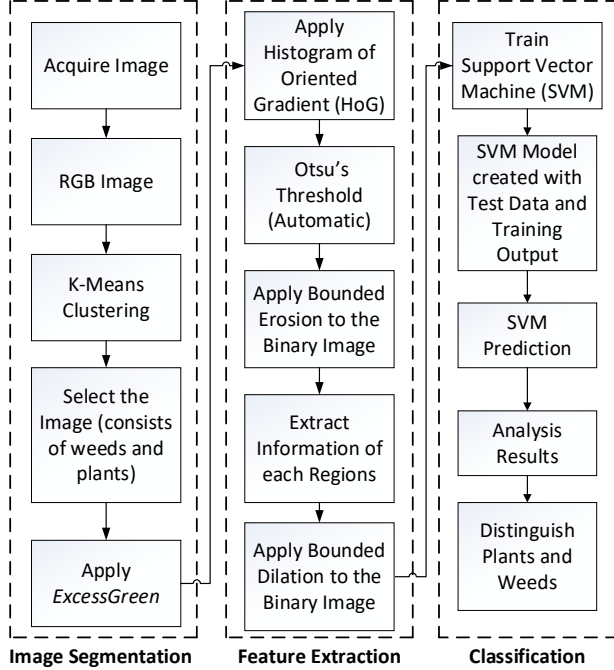


Figure 1. The flowchart of the proposed method for classifying and decision-making process.

The proposed method consists of three subsections: Image Segmentation, Feature Extraction, and Classification. All the subsections play a critical role in identifying the plant leaves from weeds. In the following sub sections, the method will be discussed in detail.

2.1 Image Segmentation

Initially, in the Image Segmentation stage, the images are loaded into the system from the database. The main task of this section is to pick the image that consists of plants and weeds where other unnecessary items are eliminated like soil, pesticides, algae, etc. To achieve this task, we did not make use of the traditional approach mentioned in Shi *et al.* [6], where individual bands of colors are examined to identify the greenness of the regions.

Instead, we used K-Means clustering algorithm [7], where each of the cluster regions is analyzed and then the desired image is stored in the new dataset. In this algorithm, the number of observations n is broken down into k clusters. The clusters develop around the mean value which changes until all the pixel values are considered.

One thing to consider is that the clustering of the K-Means algorithm is random. Therefore, randomness needs to be eliminated so that the images which consist of only plants and weeds are selected. This step can be accomplished with the help of the algorithm stated in Figure 2.

```

if (Total_Red_1 < Total_Red_2 && Total_Red_1 < Total_Red_3)
    if (Total_Blue_1 < Total_Blue_2 && Total_Blue_1 < Total_Blue_3)
        if (Total_Green_1 < Total_Green_2 && Total_Green_1 < Total_Green_3)
            input_image = segmented_image(1);
        end
    end
end
elseif (Total_Red_2 < Total_Red_1 && Total_Red_2 < Total_Red_3)
    if (Total_Blue_2 < Total_Blue_1 && Total_Blue_2 < Total_Blue_3)
        if (Total_Green_2 < Total_Green_1 && Total_Green_2 < Total_Green_3)
            input_image = segmented_image(2);
        end
    end
end
elseif (Total_Red_3 < Total_Red_1 && Total_Red_3 < Total_Red_2)
    if (Total_Blue_3 < Total_Blue_1 && Total_Blue_3 < Total_Blue_2)
        if (Total_Green_3 < Total_Green_1 && Total_Green_3 < Total_Green_2)
            input_image = segmented_image(3);
        end
    end
end
end
end
  
```

Figure 2. The algorithm that is used to eliminate the K-Means clustering randomness.

The algorithm stated in Figure 2 evaluates the band of colors in each cluster. It has been pre-analyzed that the desired cluster consists of least number of pixel values in each RGB image. After performing the clustering algorithm, the outcome of the three clusters is shown in Figure 3.

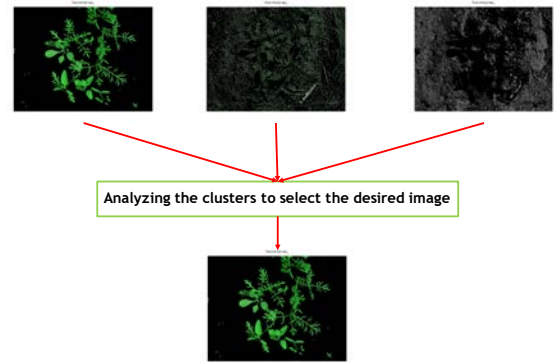


Figure 3. The algorithm mentioned in Figure 2 is used to select the desired image.

The final image is stored in the database and will be considered as an input image for weed identification process. Subsequently, *ExcessGreen* in Equation 1 is executed where more weight is assigned to the green pixel and other color bands are eliminated.

$$ExcessGreen = 2 * (g) - (r) - (b) \quad (1)$$

Afterward, Otsu's Thresholding [8] which is an automatic thresholding is applied to convert the RGB image into a binary image.



Figure 4(a). Otsu's Thresholding

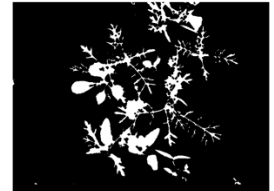


Figure 4(b). Applying Bounded Erosion



Figure 4(c). Applying Bounded Dilation

The next step is to identify the overlapping regions from the binary image. As stated in our previous approach [5] the Morphological Operations are used to separate the weeds from the plants. The Morphological Operations are a class of techniques that analyze the shape information in monochrome images [13]. The features of both weed and plant regions are extracted using *Bounded Erosion*, where the erosion operation is carried out in a loop process. This continues until all the regions are separated and are possible to label the regions of both plants and weed as shown in Figure 4(b). The dilation operation is just the reverse operation of the erosion. *Bounded Dilation* is carried out until all the regions remerged to that of the initial binary image as shown in Figure 4(c). The labeling of the regions is carried out with the help of Histogram of Oriented Gradient (HoG), which will be discussed in the next section.

2.2 Feature Extraction

The set of features that are considered are listed in the following subsections.

2.2.1 Features related to Areas and Perimeter

The features are extracted after the *Bounded Erosion* operation as shown in Figure 4(b). Table 1 shows the list of features that are considered and they are *Area*, *Perimeter* and *Convex Area*. Both the erosion and dilation operations use the Structuring Elements (SE) which define the neighboring structure of an object. As shown in Figure 4(a), plant leaves are bit rounder and convex shaped where else weed leaves are much thinner and straighter. Therefore, the SE of 'disk' and 'line' are selected as the primary structure of the regions in the image. Features within interrogated images corresponding to these selected shapes are extracted and the values are stored for further evaluation. When all the feature values are extracted the *Bounded Dilation*, operation is carried out to restore the image back to the initial image.

Table 1. List of Features considered for the experiment

Feature Number	Feature	Description
$feature_1$	Area	Area of the pixels covered by leaves and weeds
$feature_2$	Perimeter	Length of the pixels covered by leaves and weeds
$feature_3$	Convex Area	Area of the Convex Hull for leaves and weeds

2.2.2 Histogram of Oriented Gradient (HoG)

The Histogram of Oriented Gradient (HoG) [9] is a feature descriptor which is used to identify and label objects from an image. The objects are identified based on the intensity gradient or edge direction. For our purpose, the HoG feature descriptor is used to allocate the shape information of the regions and label the regions of weeds and plants. The HoG descriptor makes use of the connected regions called cells which map the direction of each of

the gradient vector [9]. The cell size set for our experiment is a 4x4 window. This is because if we increase the cell size than a lot of detail of the image will not be considered and if the smaller cell size is considered then the computational costs will increase. Figure 5 shows the outcome of the HoG descriptor where the gradient direction of both weed and plants are mapped and are considered as a set of features along with the other extracted features for the Support Vector Machine (SVM).

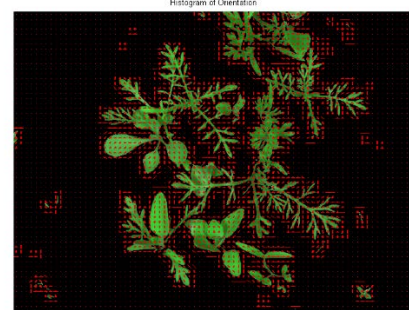


Figure 5. Histogram of Oriented Gradient (HoG) locates the weed and plants locations.

2.3 Classification

The Support Vector Machine (SVM) is a quintessential mechanism for classification. In general, the classifier is constructed in the N-dimensional hyper-plane that optimally separates the two classes [10]. The SVM consists of two phases: *training* and *testing*. In the *training* phase, the dataset is $\{x_i, y_i\}$, where $i = 1, 2, \dots, l$ and $x_i \in R^d$ where R is the Vector Space and d is the dimension of input training data. All the extracted features are the Vector Space where $y_i \in \{-1, +1\}$ comprise the output labels -1 and +1, indicating the identified class. The label -1 indicates as weed class whereas +1 indicates as plant class. The SVM becomes significantly more accurate in classification as the margin along the separating hyper-plane increases [10]. The equation of the plane that maximizes the margin is given in Equation 2 where w is the weight vector, b is the intercept term and x is the feature dataset.

$$f(x) = \text{sign}(w \cdot x + b) \quad (2)$$

The performance of SVM depends on the type of kernel functions that will be used for the classification purpose [14]. Kernel functions like Polynomial, Radial Basis Kernel and Sigmoid Kernel functions provides higher classification accuracy but come with a computational cost. Deng *et al.* [11] used Linear Kernel function and stated that linear classifier performance is cost effective than other kernel functions. Thus, these kernel functions evaluations suggested us to move ahead with Linear Kernel.

3. RESULTS AND ANALYSIS

3.1 Input Images Type

The test dataset consisted of 60 images from organic carrot fields collected from Haug *et al.* [12]. The specification of each of the images is similar to the ones mentioned in paper [5].

The proposed method is simulated with MATLAB 2014 simulator in 64-bit Windows 10 Operating System with 16 GB RAM and Intel Core i7 CPU at 2.93 GHz.

3.2 Evaluation of the Results

Since the dataset is limited the images are broken down into 40 training images and 20 test images. The training dataset images

consist of wide variety of weed to plant ratios. This variation will help to determine the efficiency of the SVM Model.

The Equation 3 *Plant Leaves Identification* is used to determine the detection accuracy of the plant leaves and the follow information are generated.

True Positive (TP): The number of locations that are correctly identified as leaves.

True Negative (TN): The number of locations that are correctly identified as weed.

False Positive (FP): The number of locations that are not correctly identified as leaf, and instead, incorrectly as weed.

False Negative (FN): The number of locations that are not correctly identified as weed, and instead, incorrectly as plant.

$$\text{Plant Leaves Identification} = \frac{TP}{TP+FP} * 100\% \quad (3)$$

The overall success rate of the plant leaves is about 92.01%. As stated in Figure 6 the proposed system showed improvement when compared with previous weed detection methods in the papers [5, 12], where the average accuracy mentioned are around 85% and 90%. Higher accuracy of the carrot leaves has been achieved as the feature descriptor, HoG locates and labels the regions in both the overlapping and non-overlapping portions. This inclusion of the labeled features along with other features increased the accuracy of the object identification and minimized any false identification.

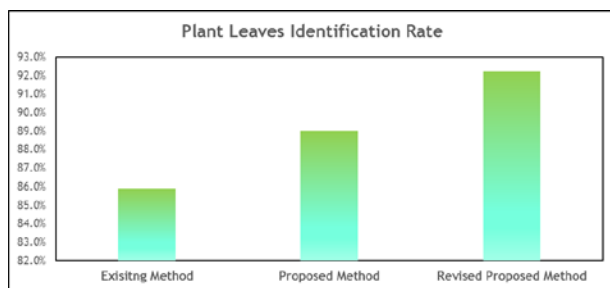


Figure 6. Plant Leaves Identification Rate.

4. CONCLUSIONS

The improved detection mechanism consists of two additional approaches to that of the previous approach stated in the paper [5]. Initially, in the pre-processing step, the K-Means clustering algorithm selects the image that consists of plant and weed regions only and avoided unnecessary tasks like handling noises from the images. This new dataset is considered for further evaluation purpose, where Morphological Operations are carried out to separate the plant leaves from weeds. Furthermore, HoG is used to locate and label the regions of weeds and plants even in the overlapping regions. Then the extracted features from the images are used to train the SVM Model to classify the weeds from the plants. The system can identify plant regions with a success rate of 92.01%. This detection method is significant as weeds normally grow close-to-crop or between intra-row which need to be regulated to avoid substantial yield loss [2].

In future research, larger dataset will be evaluated to bolster the credibility of the proposed method. Additionally, extensive weeds and plants segmentation research will be carried out to improve the detection accuracies.

5. REFERENCES

- [1] Timmermann, C., Gerhards, R., and Kühbauch, W. 2003. The economic impact of site-specific weed control. *Precision Agriculture*. 4, 3, 249-260.
- [2] Slaughter, D., Giles, D., and Downey, D. 2008. Autonomous robotic weed control systems: A review. *Computers and electronics in agriculture*. 61, 1, 63-78.
- [3] Tellaeche, A., Pajares, G., Burgos-Artizzu, X. P., and Ribeiro, A. 2011. A computer vision approach for weeds identification through Support Vector Machines. *Applied Soft Computing*. 11, 1, 908-915.
- [4] Shi, L., Duan, Q., Ma, X., and Weng, M. 2011. The research of support vector machine in agricultural data classification. In *Computer and Computing Technologies in Agriculture V*. Springer Berlin Heidelberg, 265-269.
- [5] Saha, D., Hanson, A., & Shin, S. Y. 2016, October. Development of Enhanced Weed Detection System with Adaptive Thresholding and Support Vector Machine. In *Proceedings of the International Conference on Research in Adaptive and Convergent Systems*. 85-88.
- [6] Shi, L., Duan, Q., Ma, X., and Weng, M. 2011. The research of support vector machine in agricultural data classification. In *Computer and Computing Technologies in Agriculture V*. Springer Berlin Heidelberg, 265-269.
- [7] Hartigan, J. A., and Wong, M. A. 1979. Algorithm AS 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*. 28, 1, 100-108.
- [8] Otsu, N. 1975. A threshold selection method from gray-level histograms. *Automatica*. 1, 23-27, 285-296.
- [9] Dalal, N., and Triggs, B. 2005. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition, CVPR 2005. IEEE Computer Society Conference*. 1, 886-893.
- [10] T. Joachims. 1999. Svm: Support vector machine. SVM-Light Support Vector Machine. 19, 4, University of Dortmund. DOI = <http://svmlight.joachims.org/>.
- [11] Deng, W., Huang, Y., Zhao, C., and Wang, X. 2014. Discrimination of Crop and Weeds on Visible and Visible/Near-Infrared Spectrums Using Support Vector Machine, Artificial Neural Network and Decision Tree. *Sensors & Transducers*. 26, 26.
- [12] Haug, S., and Ostermann, J. 2014. A crop/weed field image dataset for the evaluation of computer vision based precision agriculture tasks. *Computer Vision-ECCV 2014 Workshops*. (September 2014), 105-116.
- [13] Comer, M. L., and Delp, E. J. 1999. Morphological operations for color image processing. *Electronic Image*, 8, 3, 279-39.
- [14] Hsu, C. W., Chang, C. C., and Lin, C. J. 2003. A practical guide to support vector classification.